

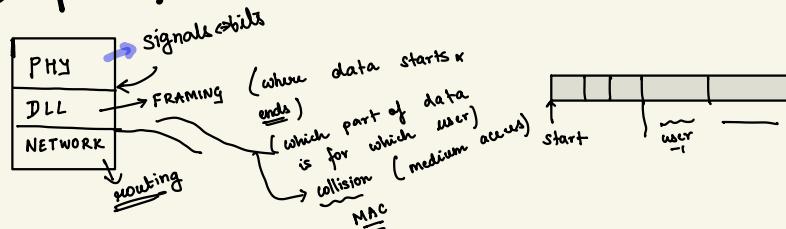


Class Notes

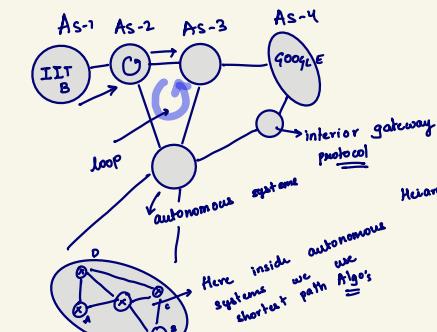
29/07/24 - Lec 1

REF. TEXT: PETERSON & DAVIE
COMP NETS, SYSTEM :

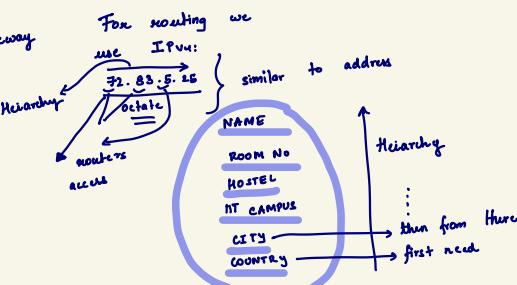
30/07/24 - Lec 2

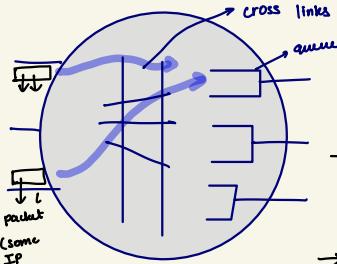


SIMPLEX: ONE WAY COMMUNICATION (e.g. radio, TV)
FULL-DUPLEX - (communicate in both directions simultaneously) → (a can send to b & b can " " a)
HALF-DUPLEX (com... " " " but not ") (only one signal a → b or b → a at one time)
↳ WIFI (read why)



outside these autonomous systems → BGP
tells other AS (autonomous system) which AS it can reach





Incoming packets mapped to same OLP Link:-

so we need queue !

So we can use FIFO

→ But if some client pays more
we give priority to his message

like prg :-

→ If queue is filled and extra packets come :

These packets are dropped (data drop)

(CONGESTION CONTROL) :-

↓
So ARPANET DID they will
only handle this issue at
end nodes not intermediate nodes
and only for L4 (messages)

MESSAGES :

L2: FRAME (1 unit of DATA)

L3: SYMBOLS (Like Napolean).

L3: Packet

L4: (TCP) → segment

(UDP) → data grame

End to end delay

of pkt

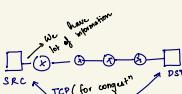
(sec to dst)

$$\leq \left(\sum_{i=1}^n Q_{max} + \sum_{i=1}^n \frac{d_i}{c_i} \right) t_{transmission}$$

BUFFERS

$$Q_{max} (bit/sec) \rightarrow \text{time for last packet to leave}$$

ith msmt



This detects congestion (not drops)

→ slows down. (So hate of incoming packets fallen down)

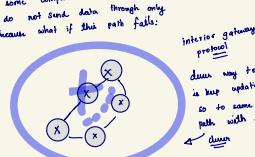
So we need to limit end to end delay

because delays could cause problems

For some companies:

They do not send data through only

path, because what if this path fails?



down way to handle congestion is keep updating weights of edges so to some destination it takes path with less congestion :-
↓ slower



L4: Transport layer

TCP: Transmission Control Protocol → Congestion control
→ Reliability
↓ file transfer

This checks for reliability

UDP → Does nothing

LAYER 5: Application layer

File (HTTP)

Email (SMTP)

VODP
(UDP over TCP)
(we don't have time for doing many packets read too often)

TEXT
MESSAGES

P2 P...

Layer 4:

TCP .. UDP ..

IF (natural notion)

L2: ATM, 4G, ethernet

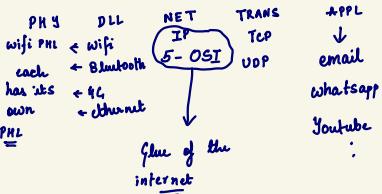
L1: Wavelength, optic fiber, ...

DESIGN PROTOCOLS IN MODULES

each subfunction handled by some protocol
we discussed above OSI-5 → APPL.
MUL.
NET.
TRANS.
PHYS.

DISADVANTAGES OF (DESIGNING/MODULARISATION)
 (1) EASE of development (we do not have to solve all problems)
 (2) Easier to debug (because we can work on one layer)
 (3) Many applications & technologies (at physical layer)
 working together :-
 (4) If we want make change at one layer
 we do not affect other layer.
 (5) ease of modification - one change 1 layer to address a problem

PROTOCOL LAYERING



Because of IP we can have many time
PHL & DLL.

(TCP) → RFC (req for comments)

This is how initially TCP was called. So while making a TCP we have to satisfy old RFC and then we can add our own protocols!!

For WiFi the protocol is IEEE 802.11[802.11ac...].

standards as standards increase rate of data increase

② DISAPP

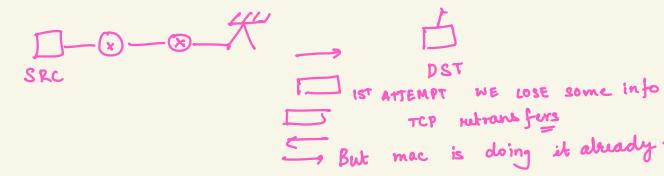
Redundancy of TASKS

① TCP : Handles retransmissions

(Reliable) ↑ (Because we TCP make

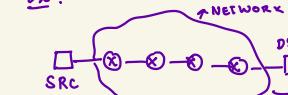
sure = no pkts drop)

But on the side mac retransmission is also going on



③ SUBOPTIMALITY:-

Ex: VOIP (lowest time delay)



We do not have choice of choosing our path. We get what network gives! No guarantee on QoS (Quality of Service: amount of pkt drops, Latency (Delay))

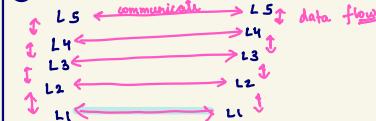
we could have some other network by which both could have less latency. But that do not know!
 Known as BEST-EFFORT.

ADVATAGES OF LAYERING:-

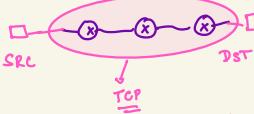
- ① DEVELOPMENT & DEBUGGING (EASE)
- ② EASY TO MODIFY ONE LAYER WITHOUT BREAKING THE SYSTEM. (we can work with one layer and not worry about the others).
- ③ Can have different choices at the different layers :- (compatibility)

DISADVANTAGES

① DEV I ② DEV 2



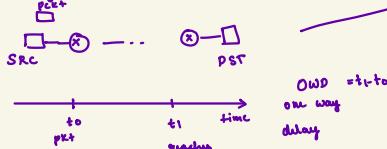
Higher layer do not communicate with low levels:-



Let us say a pk is dropped. It assumes that some routers in between must have had a queue full and this could not be reason, it just makes intelligent guess as these routers cannot communicate much with TCP. This drop could be due to some interference but TCP makes a wrong guess!

Latency Metrics :-

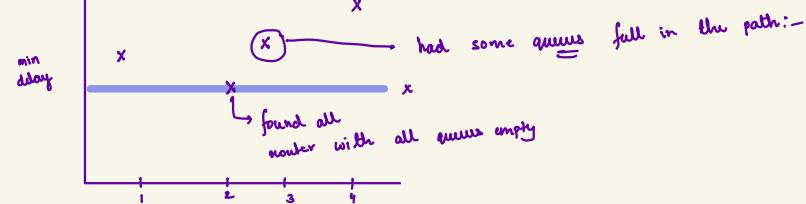
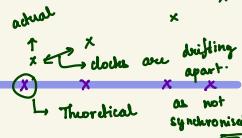
① One-Way Delay:



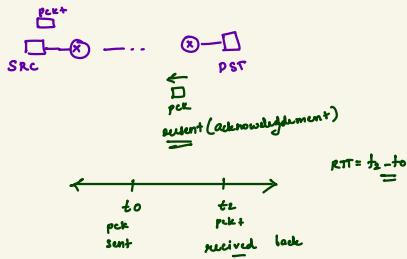
first issue :-

The clocks are not synchronised as they are not atomic clocks

OWD ?



② Round Trip Time



In VOIP we want
~few 100 ms at most

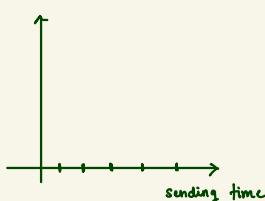
→ JITTER : Variability in OWD latencies :-

Pkts 1, 2, ...

OWD d₁, d₂, ...

e_k = |d_{k+1} - d_k| - jitter

$$\text{avg jitter} = \frac{1}{n-1} \sum_{k=1}^n e_k$$



I want the spacing between { Due to delays
packets to be same! } problem ho jati hai

One way communication

main problem hai

Jitter ke karne quality

ki maa bura hajali hai

video quality is bad because

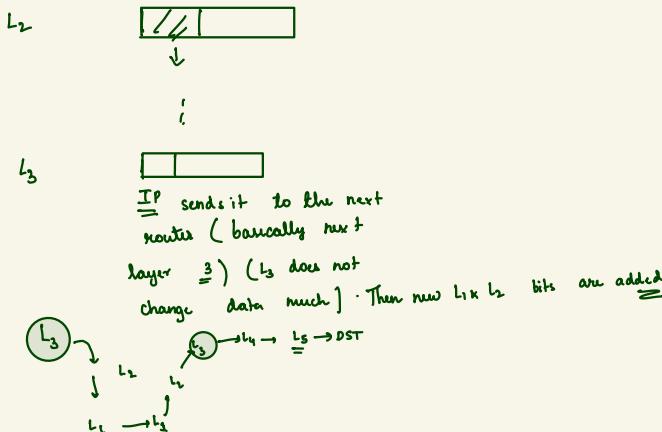
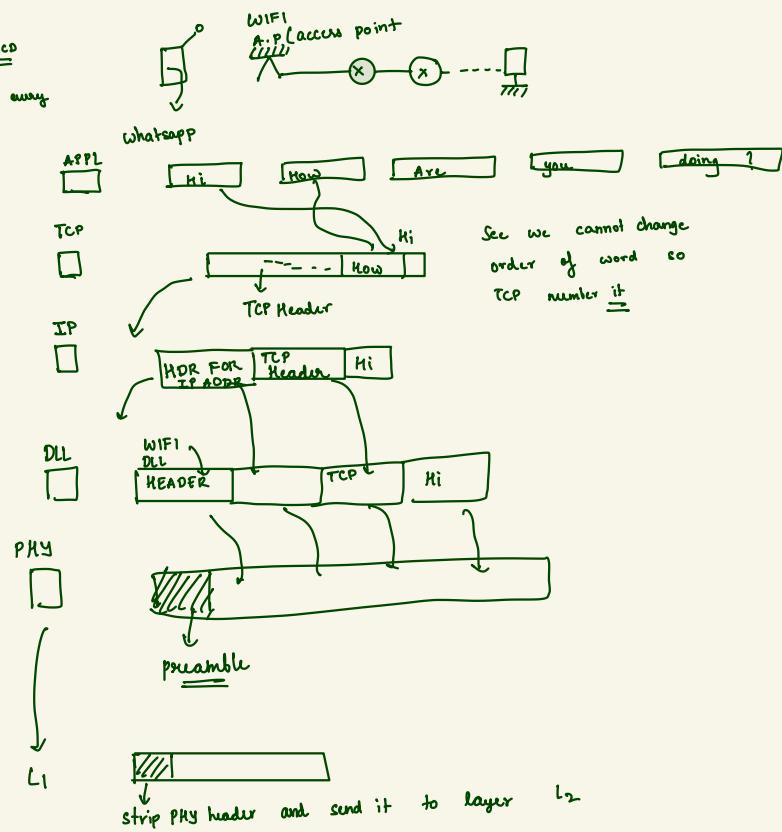
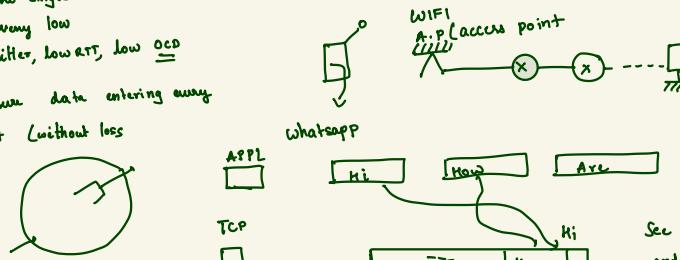
some information should come after it

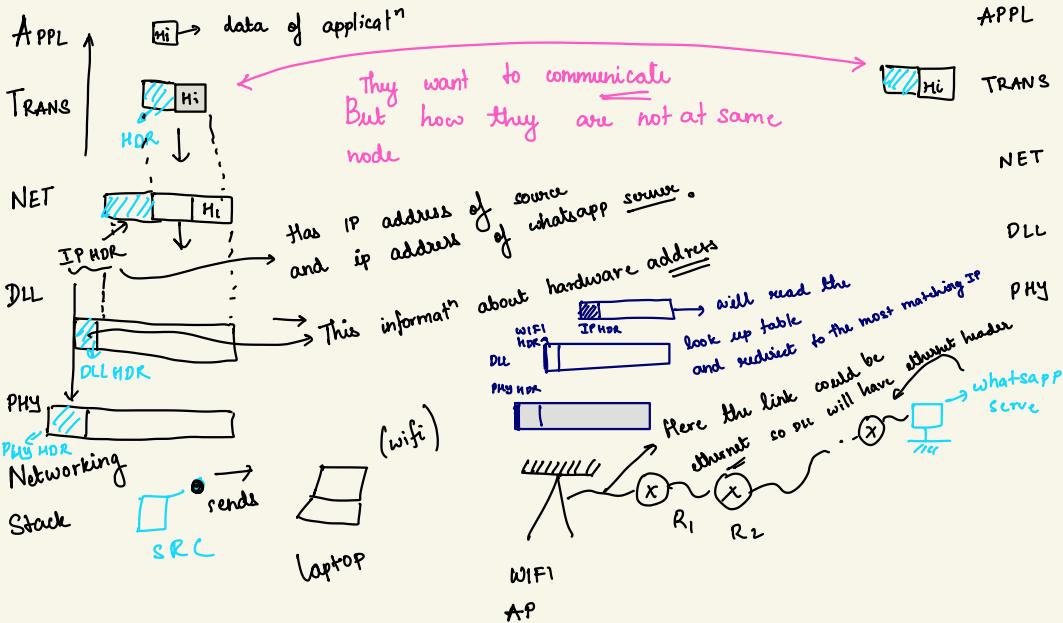
we less spacing now has more

Example in tennis (Remember Shalapur)

Telephone's Network

↳ Voice → This was the single application. We had very low delays, no loss, low jitter, low RTT, low latency.
 The wire telephone → We had to make sure data entering away counter should exit (without loss).





* Both laptop and wifi are in hearing range of this channel.
 So PHY layer of both take this message
 But DLL of laptop rejects this message
 due to Hardware address in DLL not meant for Laptop.

The DLL header depends on technology used

- Here we have WiFi Header
- Remember this encapsulation and deencapsulation happening till it reaches final IP, then completely deencapsulates.

- So if we want to send secret info, the wifi server should only have encryption key and you should not just "Mi" send encrypted version.

QoS: Quality of Service

Pat loss (Drop) rate :-

Latency (JWD, RTT, Jitter)

Throughput



Throughput: Data rate from a sender to a receiver per unit time }

→ We care sometimes because sometimes we need to push a large data



We do not care for QoS as in start we can have a delay, but then things work smoothly. But we care about jitter because we want to avoid buffering.

VOIP: Throughput?

↓
Claude Shannon

↳ "Information Theory"

Entropy: minimum number of bits required to encode the data

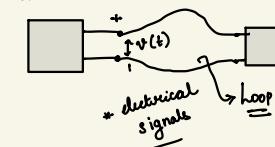
For voice we do not need much high throughput and if pkt drops then TCP backs off so we lose time

RTT ~ few 100 ms

Loss?

Physical layer

Wired Media



$$\text{loop} \rightarrow \text{magnetic field} \quad \star \frac{d\Phi}{dt} = \text{induced voltage}$$

→ We want to minimize the area to reduce the induced voltage

So we use twisted pairs:- → area minimizes



There are different types of twisted pairs:- Supports

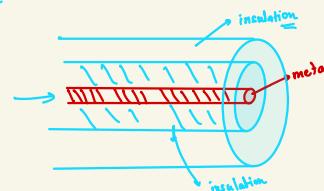
Cat 3: 10 Mbps upto 100m

Cat 5: 100 Mbps "

Cat 6: 1000 Mbps "

Coaxial Cable:-

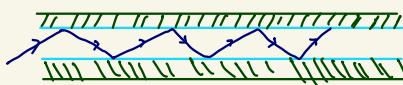
We can connect many PC's into it:-



Thin Net Coax: 100 Mbps, 200m
→ 0.2 inch in dia

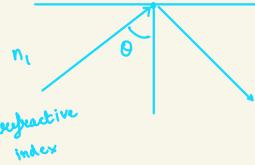
Thick Net Coax: 100 Mbps, 500m
→ 0.4 inch in dia

Optic fibre



n_2

Medium 2



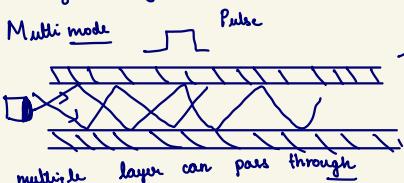
Medium 1

$$\theta_c = \sin^{-1} \left(\frac{n_2}{n_1} \right)$$

Single mode



If I turn this laser on.
Only a single layer passes through:-



multiple layer can pass through

SINGLE MODE AND MULTI MODE

in single mode single ray can get through
multi rays can't get through
This allows multiple rays to get in.

OPTIC fibre have high distance transmission ability

Pin

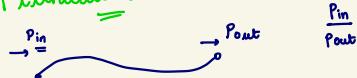
Pout

Single mode (only a single pulse comes same as previous)

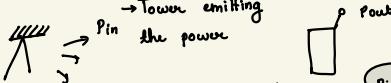
pulse of different layer get added

figuring out which is diff
So we can send much
through single mode

Absorption



We want to see the ratio :-



Very little is reaching to phone so $\frac{Pin}{Pout}$ very high
∴ we look at $10 \log_{10} \frac{Pin}{Pout}$ in dB

So a small diff in this high change in ratios

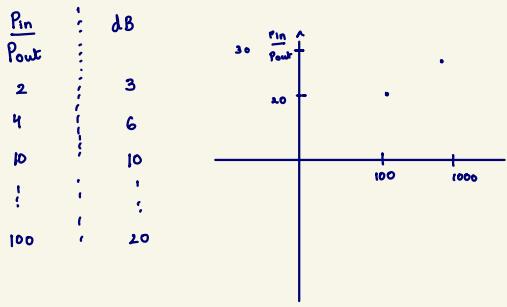
Ex $Pout = \frac{Pin}{2}$ } A 2 decrease in Pow
 $Att = 10 \log_{10} \frac{Pout}{Pin} \sim 3 \text{ dB}$ } is 3 decrease in Attenuation

Ex $Pout = \frac{Pin}{10}$ }

$Att = 10 \text{ dB}$

★

$$\begin{aligned}
 & \frac{Pout}{Pin} = 4 \quad 3 \text{ dB} \\
 & \frac{Pout}{Pin} = 2 \quad 6 \text{ dB} \\
 & 10 \log_{10} \frac{Pout}{Pin} = 3 \quad -① \\
 & 10 \log_{10} \frac{Pout}{Pin} = 3 - ② \\
 & ① + ② \\
 & 10 \log_{10} \frac{Pout}{Pin} = 10 \log \frac{Pout}{Pin} = 1
 \end{aligned}$$



Absolute Power in decibel Scale
1mW as a reference

Power P (watts) in dBm

$$10 \log_{10} \frac{P}{1 \text{ mW}} \rightarrow 10 \log_{10} \frac{P}{10^{-3}}$$

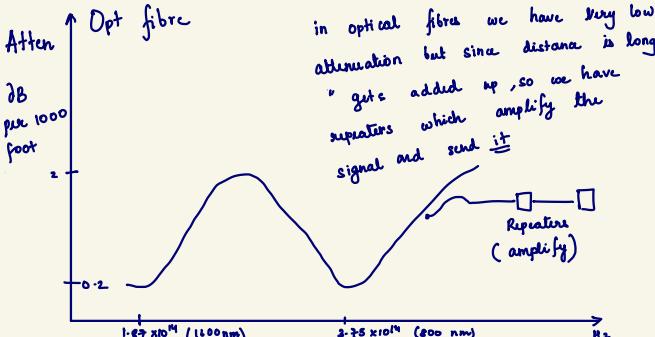
Ex $P = 1 \text{ mW} = 10^{-3} \text{ W}$
 $10 \log_{10} \frac{10^{-3}}{10^{-3}} = 0$

$\Rightarrow P = 2 \text{ mW}$
 $10 \log_{10} \frac{2 \times 10^{-3}}{10^{-3}} = 3 \text{ dB}$

So Government sets max transmission limit.

Also to make sure health hazards and other transmission

Rx Pow
 $P = -64 \text{ dBm}$ } absolute power
 Rx Power } small power we can still get wifi signals.
 Noise Power } noise + noise
 power Also noise can be present in circuit as not at 0 K (kelvin)

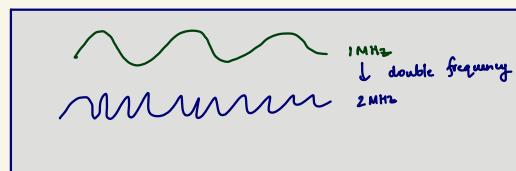
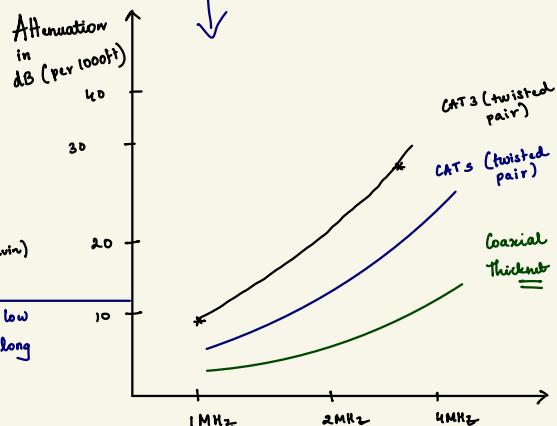
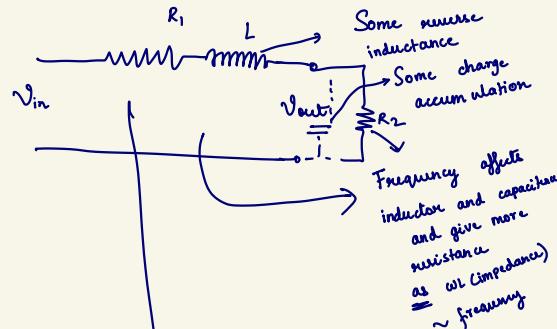


Attenuation can also be thought in terms of Amplitude²

Power $\propto (\text{Amp})^2$

Attenuation $L = 10 \log_{10} \left(\frac{A_{\text{in}}}{A_{\text{out}}} \right)^2$

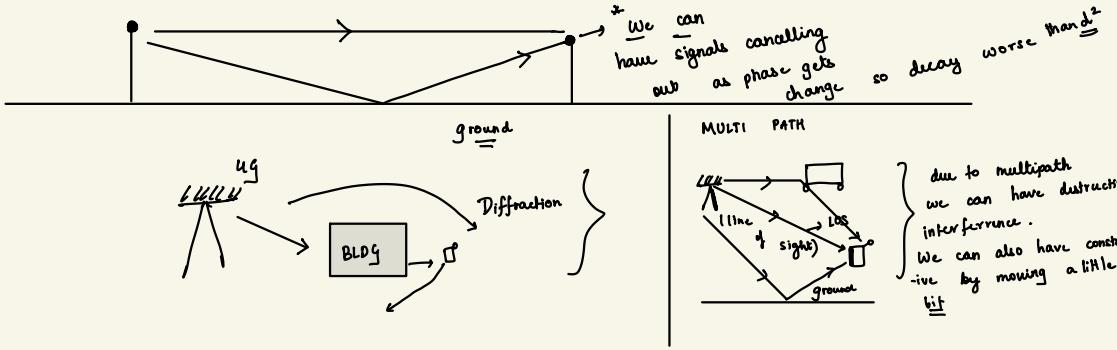
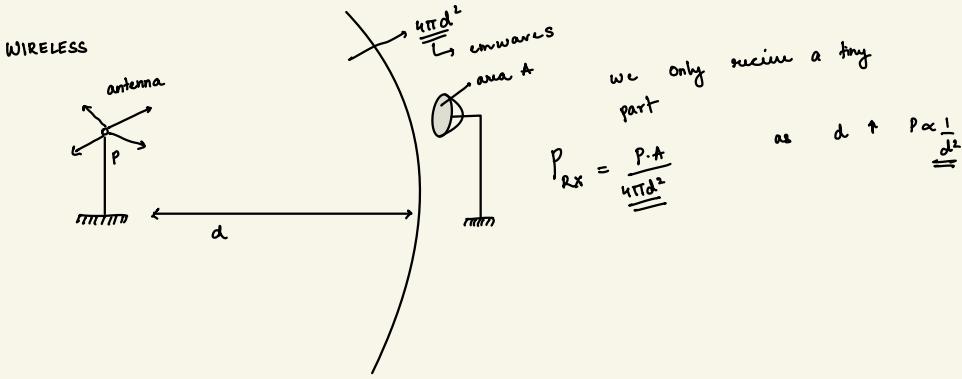
$= 20 \log_{10} \frac{A_{\text{in}}}{A_{\text{out}}} \equiv \text{dB}$



The more frequency ranges I have
I can send more signal

Entropy

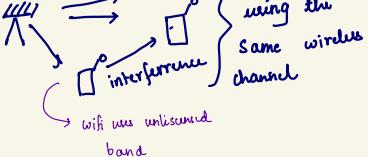
Shannon Capacity } max how many
↓ bits we
Max data rate can push through
for communication over
a channel



WIRELESS CHANNELS

in → out vs ---

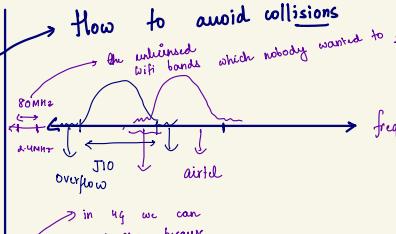
In wireless signals the attenuation is much higher than wired



two people using the same wireless channel

signals from/for them interfere due to collision.

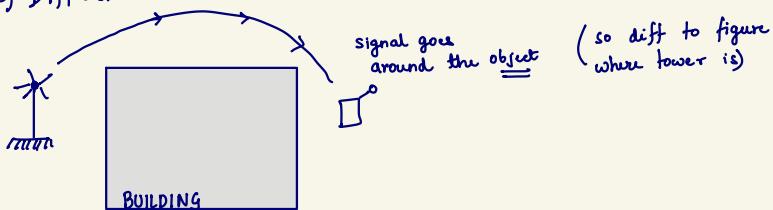
This happens in wifi



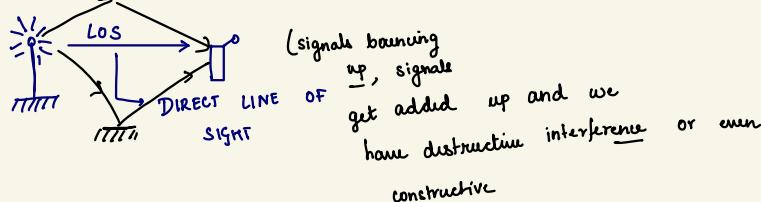
in 4g we can avoid this because the tower tells the phone when to communicate

in these bands we can transmit without paying government and cost is low
But since unlicensed many people use it and there could be a lot of interference

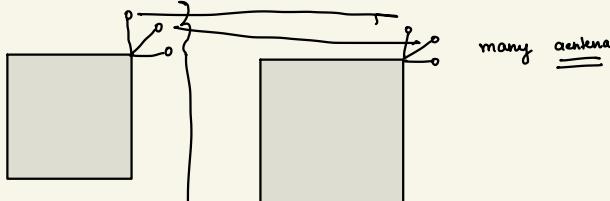
3) Diffraction:-



4) MULTIPATH

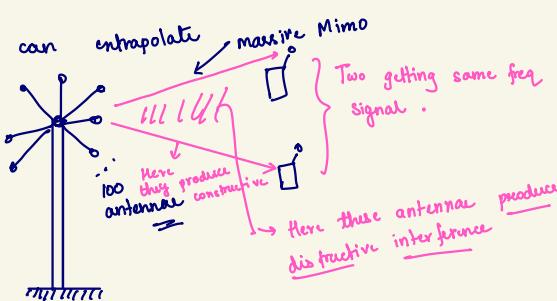


5) MIMO



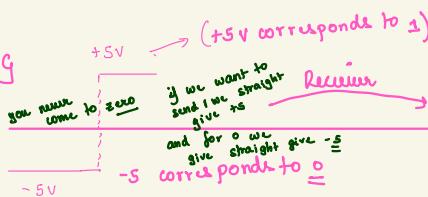
They increase probability of constructive interference because by having 3 we can have atleast one of them having constructive interference

So we can extrapolate massive Mimo



SIGNALING

WIRED
Non-Return to zero

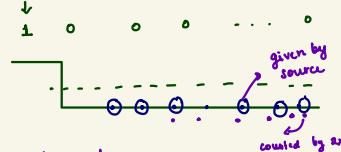


green is noise making it difficult to analyse signals

→ in return to zero we send signal zero by sending or then sending $-5V$.

ISSUES WITH NRZ

i) Bits



He has count how zeros he has received.
The signals are coming \propto Hz the receiver has a fast clock with ZN .
He will count more bits because there for a long time. So he counts more bits as he does not have synchronised clocks.

In Return to zero we know when zero comes because of a $0V$ gap.

No need to synchronised clock.

How does a receiver know when he starts getting message.

PREAMBLE

First he gets set bits which signify he is going to get message

e.g. 1010101010 ...

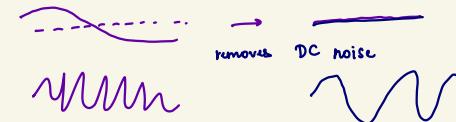
(2) Baseline Wander



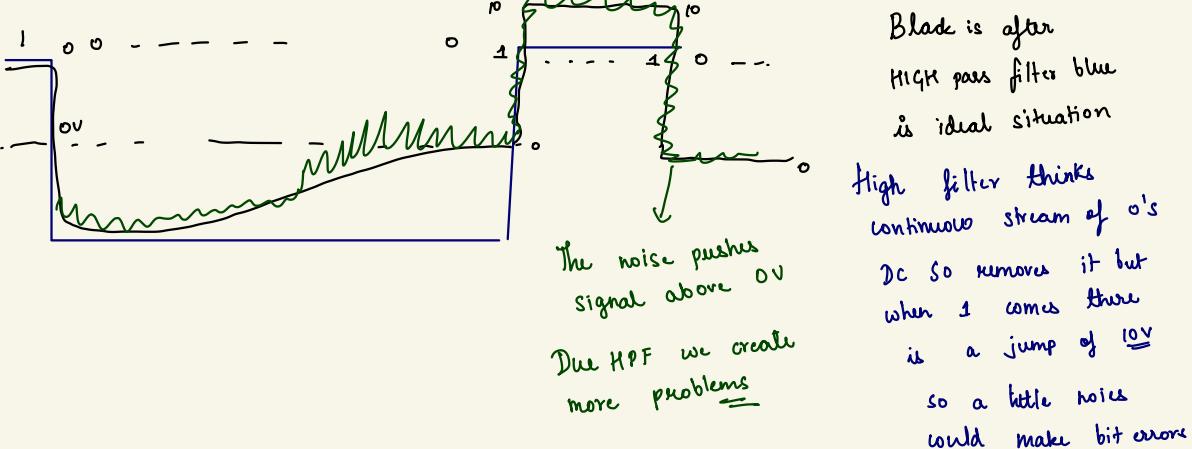
(A) Amplifier

So now user/receiver gets noisy signals so he can read this is 1 which is wrong

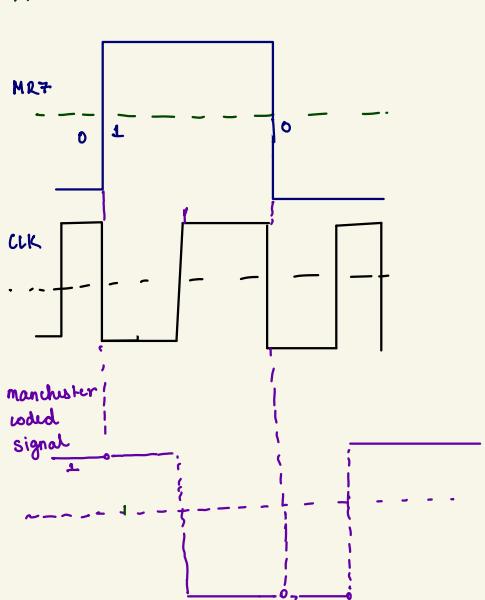
To avoid this we use a high pass filter → remove DC (offset) and low freq. signals



in phones we have band pass filter which removes low freq. and keeps frequency from high domain depending on "SIM CARD".



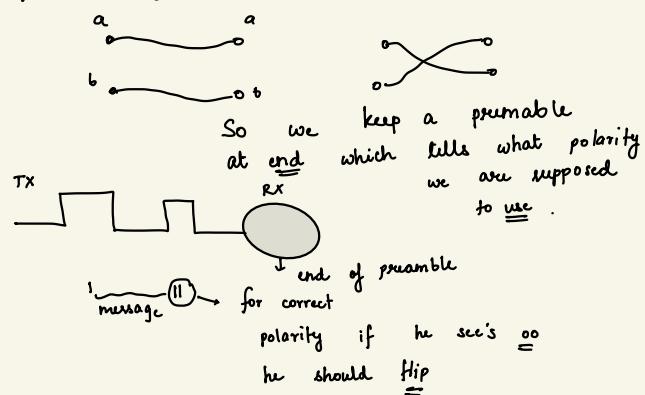
MANCHESTER CODING (used in Ethernet)

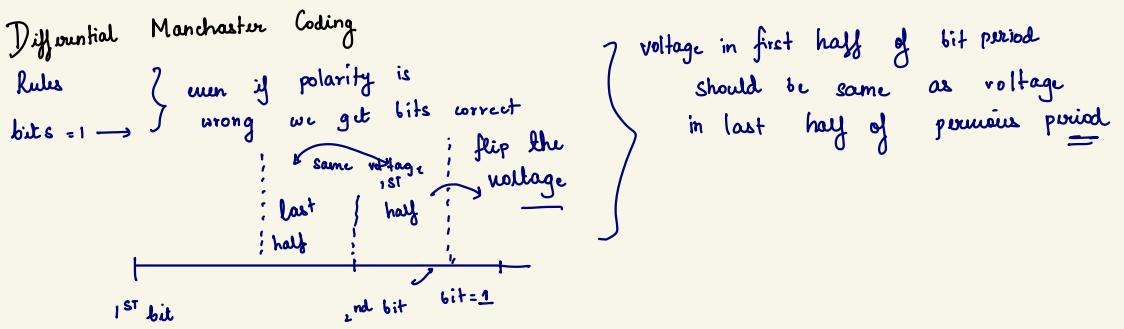


XOR	
0 0	0
1 0	1
0 1	1
1 1	0

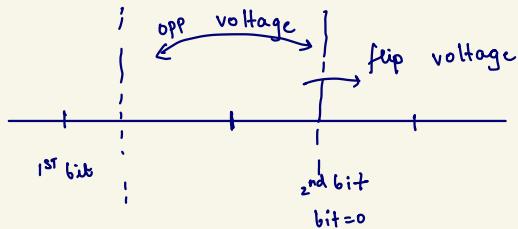
- (a) Signal transition at every bit period
- (b) Avg signal per bit is $\frac{1}{2}$

If polarity switches then message also flips

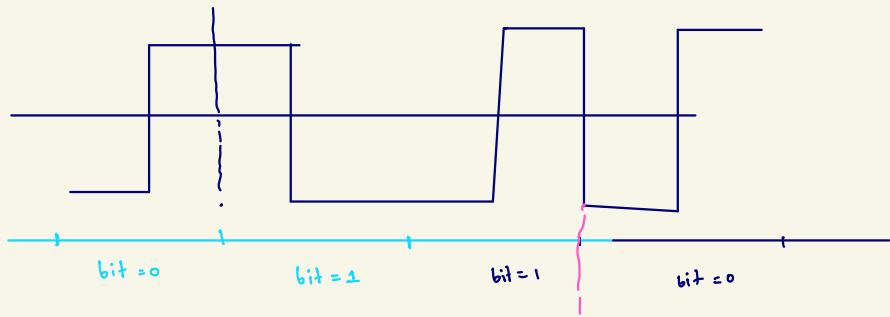




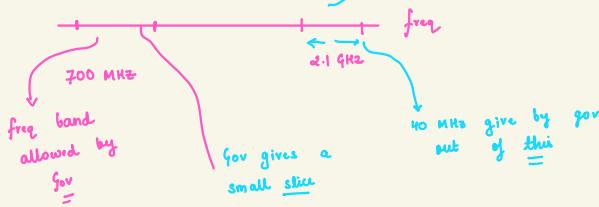
if $\text{bit} = 0 \rightarrow$ voltage in 1st half of bit is opp of that in last half of prev bit period



eg)



Modulation :-
in wireless channel

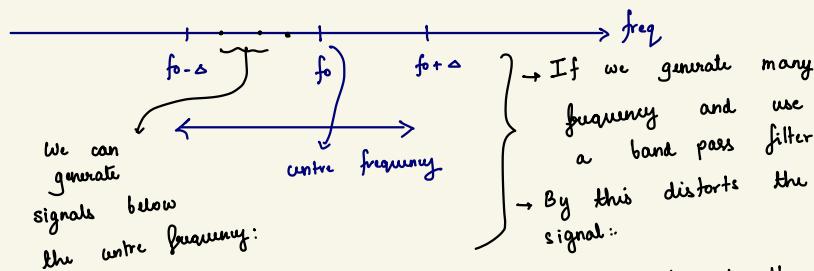


→ If we say Jio uses this band, we have to take fourier transform of the signal and it should lie in this band.

So there are international standards for signals.
So we can design signals which can help all people.

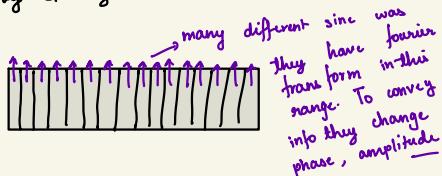
How to confine signal to a band

Default signal $A \cos(2\pi f t + \phi)$
amplitude frequency phase



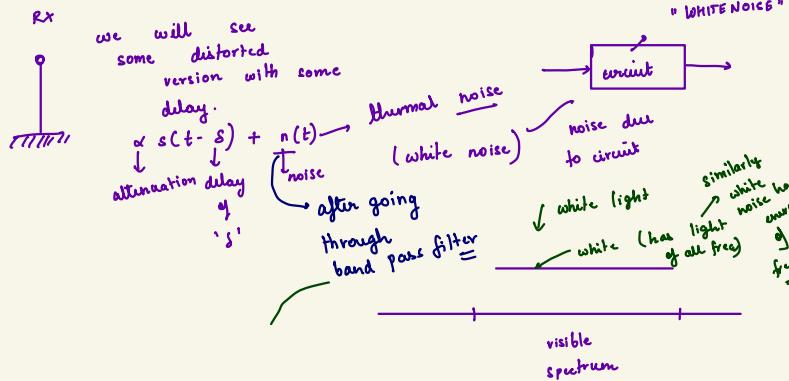
- If we generate many frequency and use a band pass filter
- By this distorts the signal.

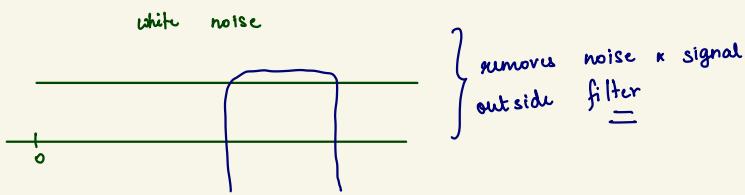
→ What sir suggested we can generate many signals in the same band by splitting a larger band into smaller bands



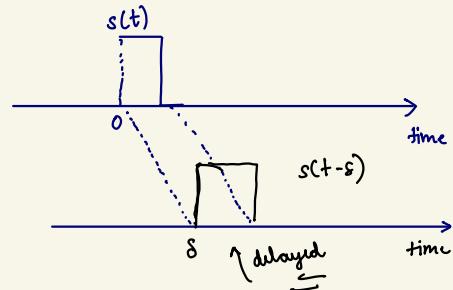
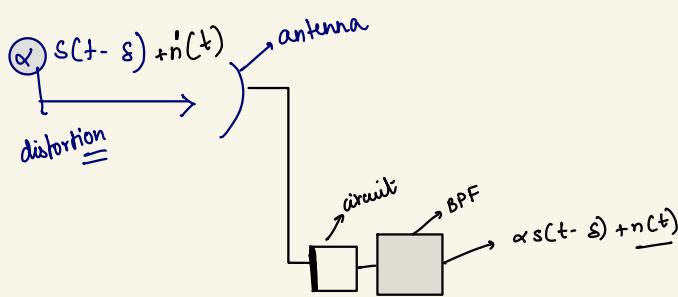
→ We want to directly be able to generate frequency of signal within this range

TX (transmitter)
 $\rightarrow A \cos(2\pi f t + \phi)$
 $s(t) \uparrow$



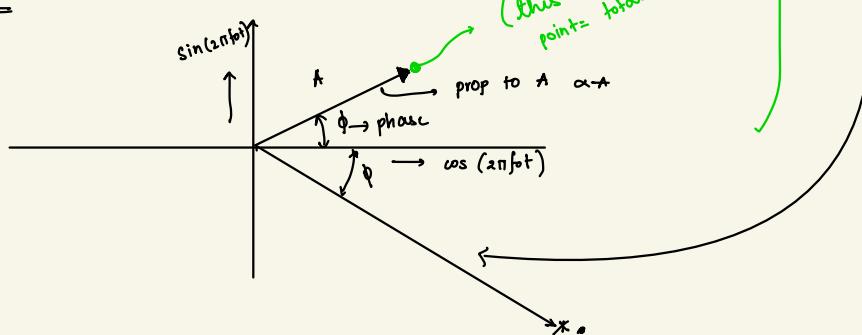


Band pass
to filter signal
of interest basically
bandwidth signal



$$A \cos(2\pi f_0 t - \phi) \rightarrow \text{Vector Space (?)}$$

We are converting sine wave signals into vector space
form



$$2A \cos(2\pi f_0 t + \phi)$$

In this vector space

$$a(t) = \dots$$

$$b(t) = \dots$$

$$\langle a(t), b(t) \rangle = \int_0^T a(t) b(t) dt$$

3 dim space

$$a = a_x \hat{e}_x + a_y \hat{e}_y + a_z \hat{e}_z$$

$$b = b_x \hat{e}_x + b_y \hat{e}_y + b_z \hat{e}_z$$

$$\frac{a \cdot b}{\|a\| \|b\|} = \cos \theta$$

$$\text{unit vectors}$$

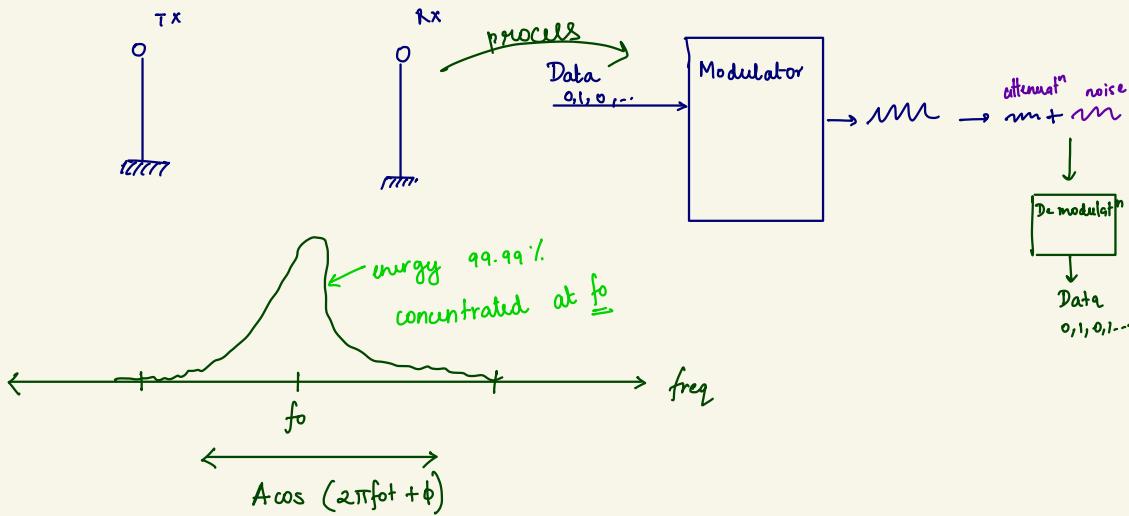
$$\sqrt{\frac{2}{T}} \cos(2\pi f_0 t)$$

$$T = \underline{11f_0}$$

$$\sqrt{\frac{2}{T}} \sin(2\pi f_0 t)$$

Q) inner prod of these = 0

$$\begin{aligned} \langle s_1(t), s_2(t) \rangle &= \frac{2}{T} \int_0^T \sin 2\pi f_0 t \cos 2\pi f_0 t dt \\ &= \frac{2}{T} \times \frac{1}{2} \int_0^T \sin 4\pi f_0 t dt \\ &= \frac{1}{T} \left[-\cos 4\pi f_0 t \right]_0^T = 0 \end{aligned}$$



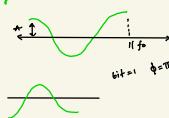
* This is not good as this cannot convey any information

We need change signal to send message

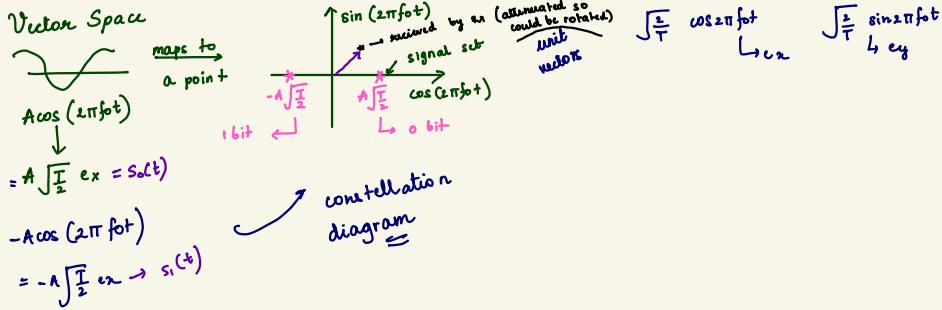
modulate the carrier signal
we can change the amplitude, phase, frequency

$$A \cos(2\pi f_0 t + \phi)$$

Ex) $\phi=0$; bit 0



These signals get attenuated } so we have to analysis to understand what is the



$$x(t) = \alpha s_i(t-s) + n(t)$$

will be rotated

delay rotates the way?

let us assume he corrects it:-

What will noise do?
Noise : white, Gaussian, Noise

Additive (AWGN)

$$\alpha s_i(t-s) + n(t)$$

my signal

noise added thus additive

energy of noise over range

We can model it as gaussian noise

$r(t) \rightarrow \text{received}$

$s_i(t)$ how to figure out which signal

$$x(t) = \alpha s_i(t) + n(t)$$

$$= s_{\text{ex}}(t) + n(t)$$

$$\langle a(t), b(t) \rangle = \int_0^T a(t) b^*(t) dt$$

$$s_i(t) = -A \cos 2\pi f_0 t$$

$$s_{\text{ex}} = \langle s_i(t), e_x \rangle = \int_0^T (-A \cos 2\pi f_0 t) (\sqrt{\frac{1}{2}} \cos 2\pi f_0 t) dt$$

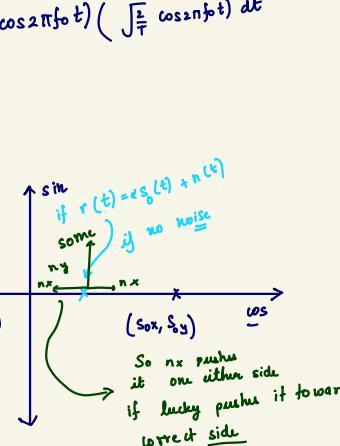
$$= -A \sqrt{\frac{1}{2}}$$

$$s_{i,y} = \langle s_i(t), e_y \rangle = 0$$

$$r(t) = r_{\text{ex}}, r_y ?$$

$$r_{\text{ex}}(t) = \langle r(t), e_x \rangle = \sin(2\pi f_0 t)$$

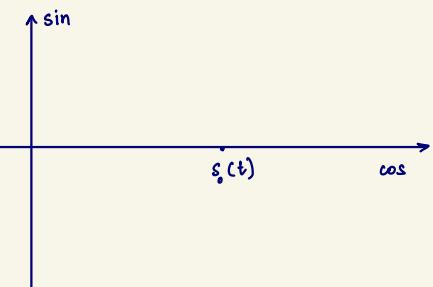
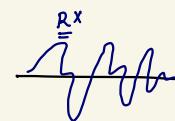
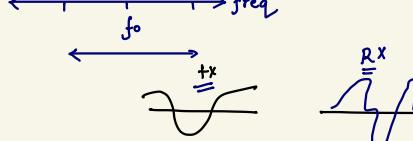
$$r_y(t) = \langle r(t), e_y \rangle = \sin(2\pi f_0 t) (s_{\text{ex}}, s_{\text{ey}})$$



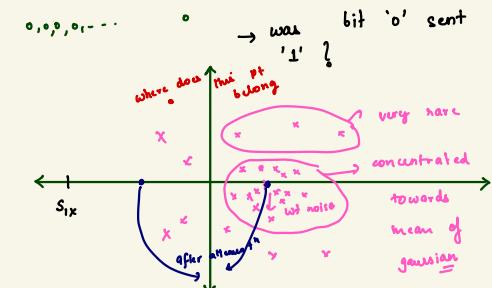
expected value

flat line

noise power / energy

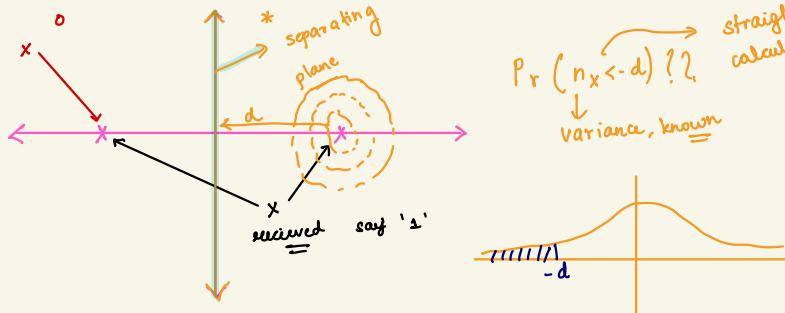


$n_x = \langle n(t), e_x \rangle$ each is identical
 $n_y = \langle n(t), e_y \rangle$ independent Gaussian RV

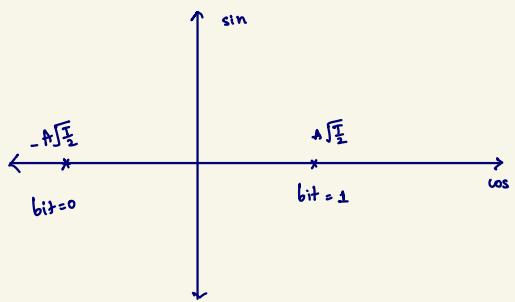


to check see closer to which pt of mean after attenuation.

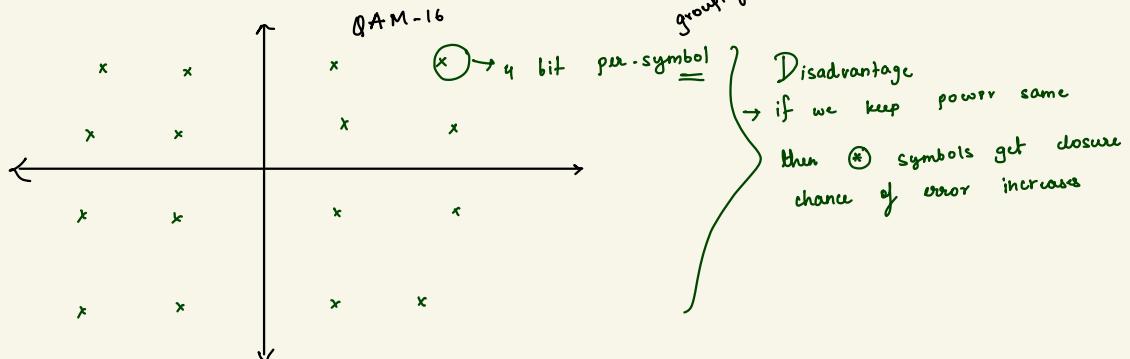
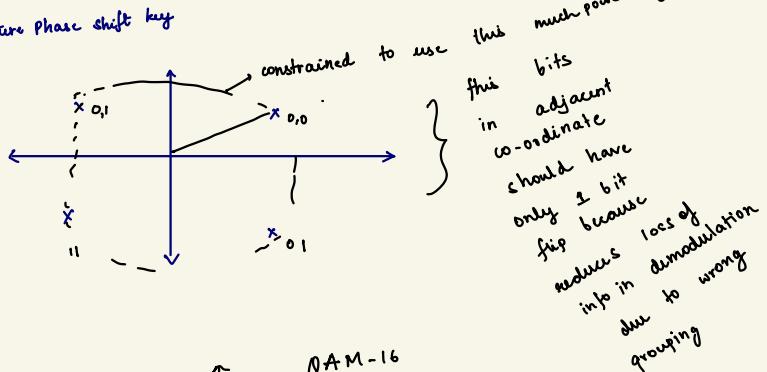
Suppose we know constellation after attenuation
and without noise?



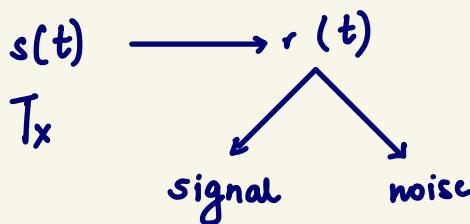
Binary phase Shift Key (shift)



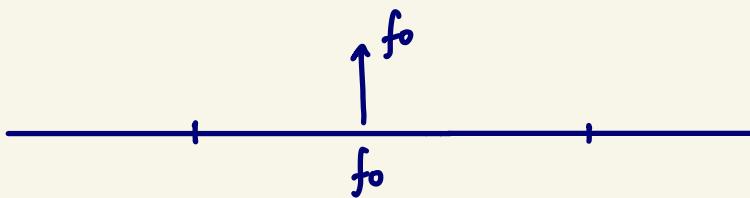
Quadrature phase shift key



PHY
↓
DLL

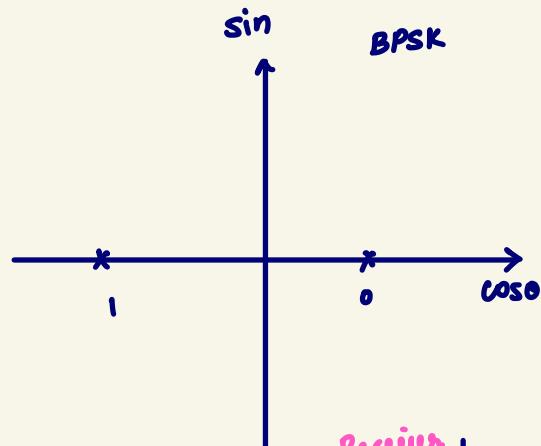


we want keep signal in a particular frequency band



* to convey information we need to modulate or change freq / amplitude

$$A \cos(2\pi f_0 t + \theta)$$



signal → channel → RX

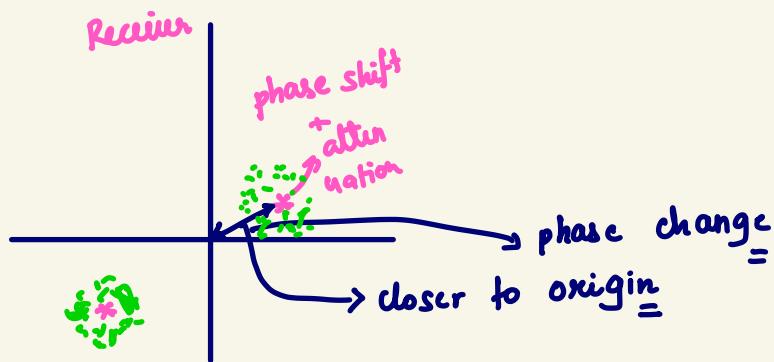
T_x

signal given by
source in blue

Receiver

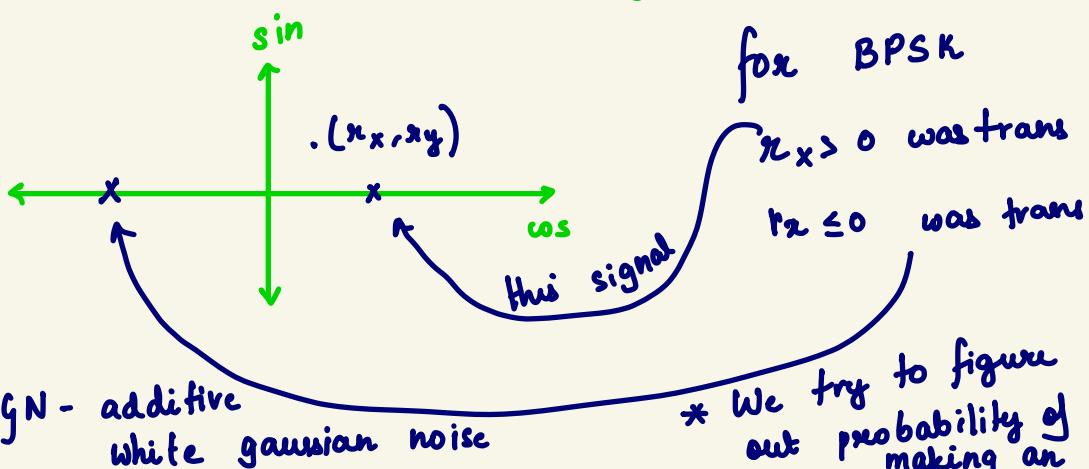
phase shift
+ attenuation

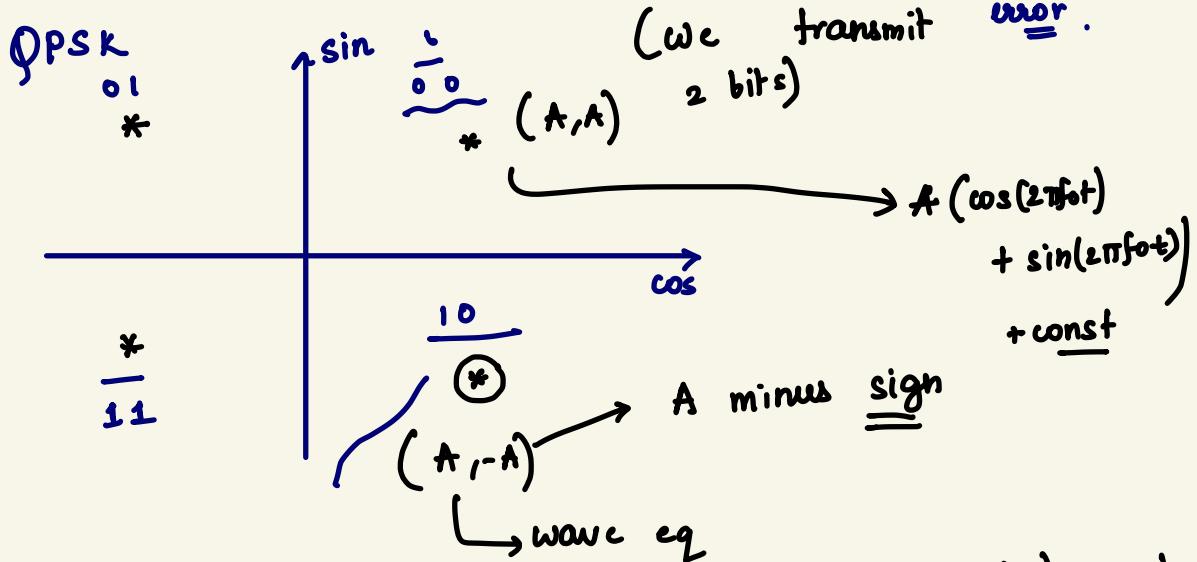




- Our constellation gets rotated. Also we will see signals around these points
- So receiver first rotates constellation back to correct angle with some error
- We do not know the signal yet (1 or 0)
So we take inner product of sent signal with e_x & e_y .

$$r_x = \langle r(t), e_x(t) \rangle \quad r_y = \langle r(t), e_y(t) \rangle$$



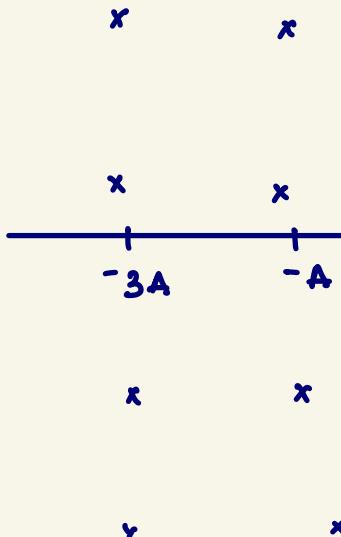


We have given convention and this type of bit assignment reduces error.

* if $(A, -A)$ was 11 instead of 10 then if due to error the signal of (A, A) was calculated to be $(A, -A)$ then we would 11 both bits are wrong. If 10 then we got one bit correct.

QAM-16 - (sending 16 bits)

We have 16 constellatⁿ pts, each pt can have 4 bits.

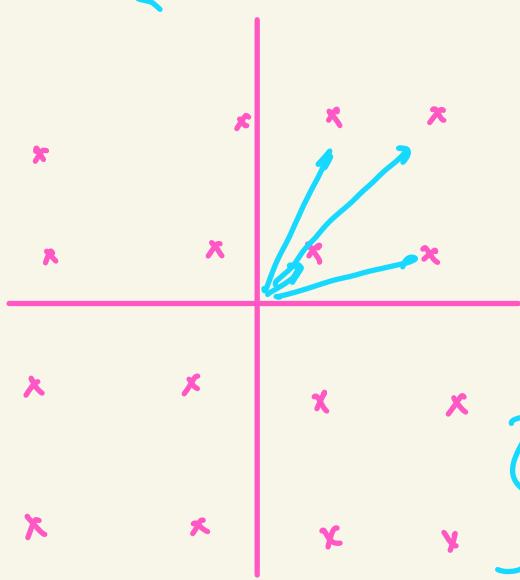
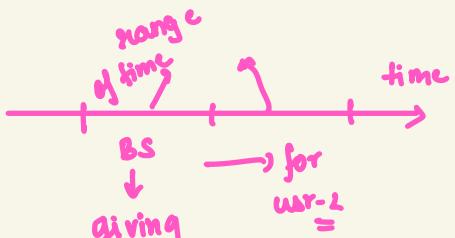
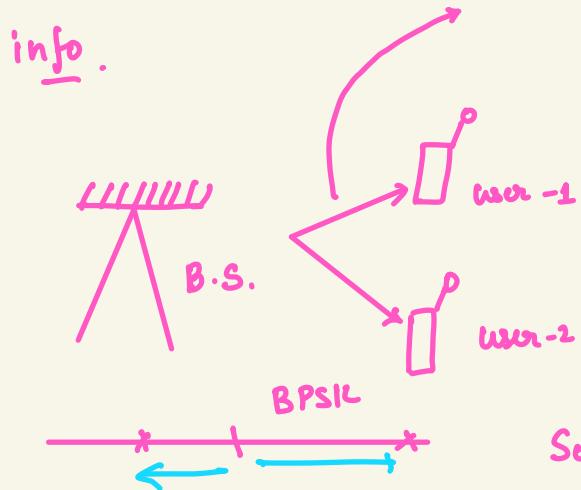


If we do not estimate channel correctly.
we get a rotated + shifted one then
we can map regions like we did in BPSK

But since not aligned the rotation (phase)
we can get lot of bit errors. So our phone has to keep doing this.

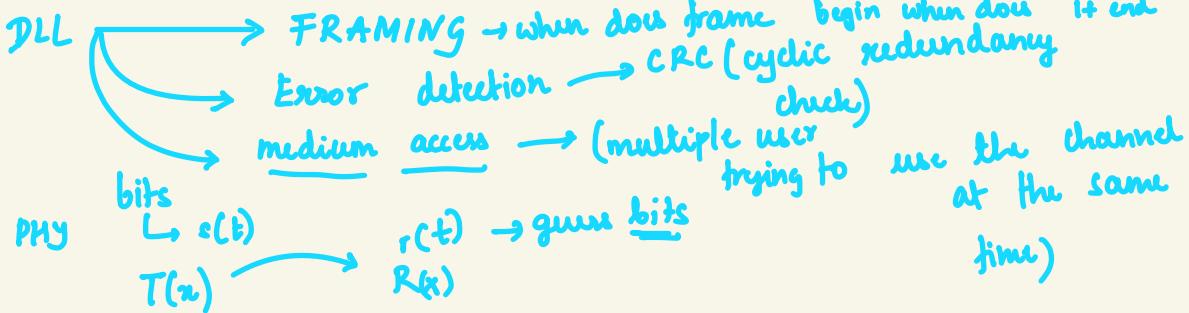
assuming this transported
gaussian centred at origin.
We could have errors and how to choose which constellatⁿ pt.
→ Choose the one with the least distance from centre

So our phone has to have so info.



So B.S. has to decide what type signal to user-1
So we need to know what level of attenuation,
phase errors I can get.

So we will be give some probability }
then we should select given attenuatn what errors we can make:- }
The energy for each type is different for signal in QPSK for a quadrant. In BPSK they have the same energy. But we can lot of bit errors in QPSK.

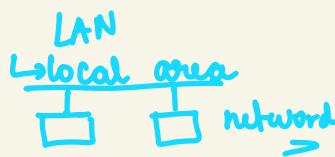
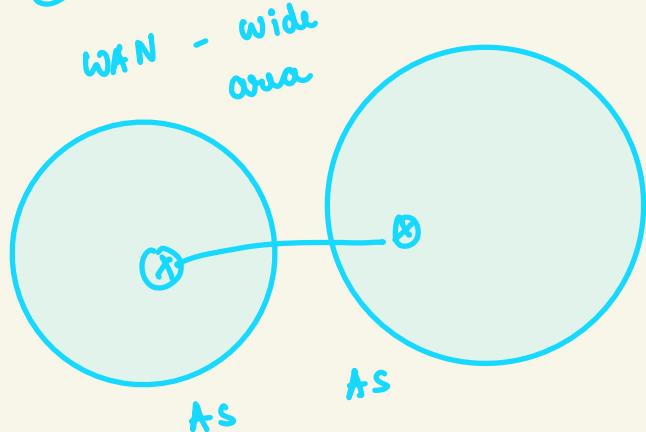


data 1011 01011 011 010 → which sub array of bits is IP address.

* How do we choose which frame has what data.

HDLC → High-level data link control

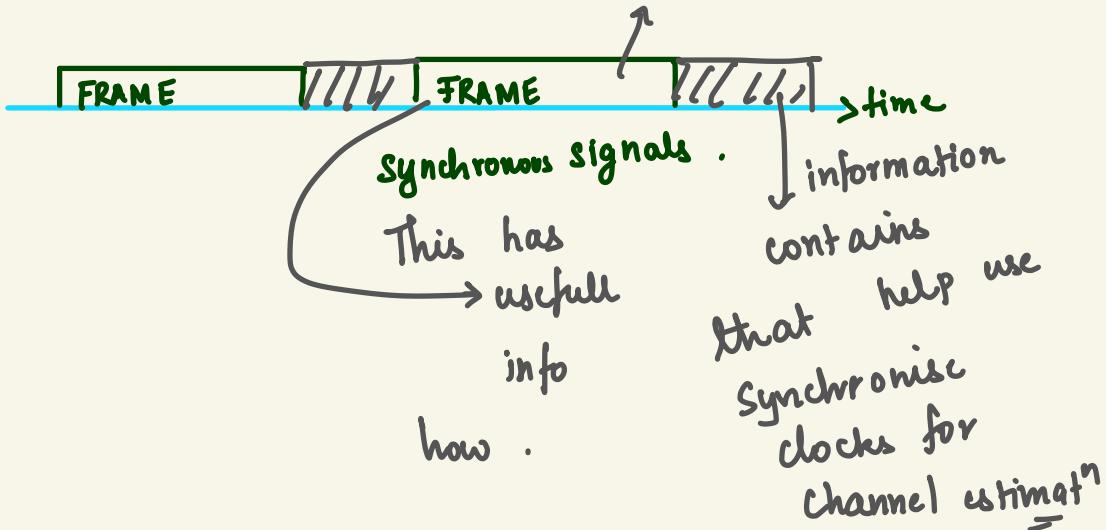
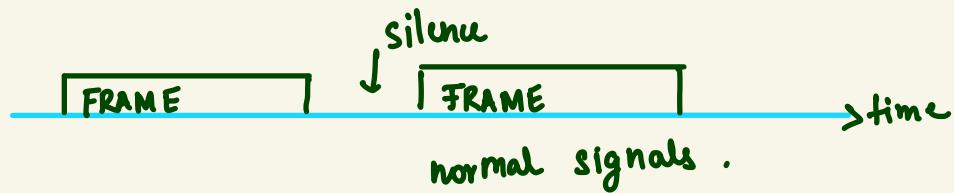
Ref: Peterson & Davie



Synchronous Mode (HDLC)

always transmitting *

} There is always some signal on line

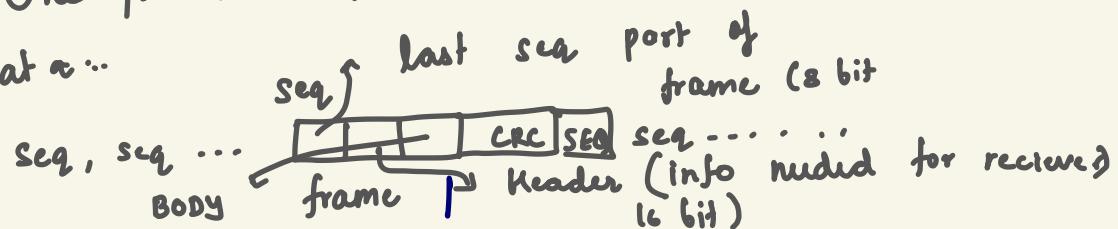


* How should the user when important info (FRAME) is there and when no useful info is being sent.

So HDLC has a DEFAULT SEQ : 0111 1110 ...

So we will have this sequence to synchronise

One problem if this sequence is in the data...

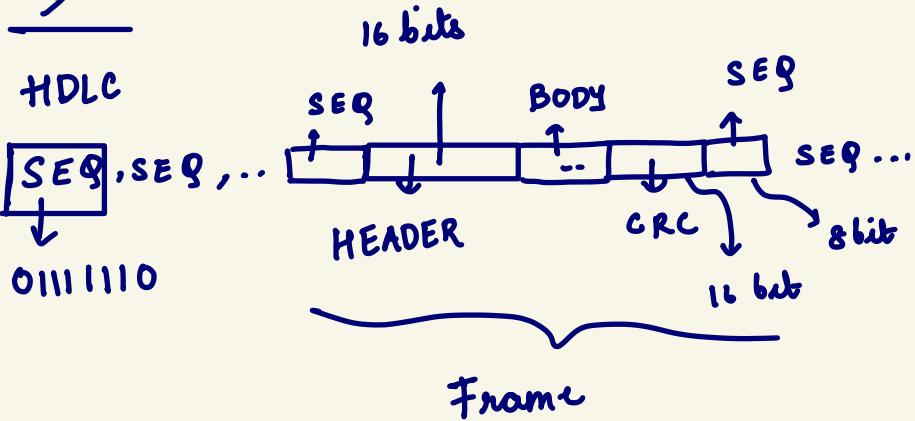


↓

this particular sequence could come in header to CRC then what to do

If our message/frame $\underline{0111110}$ so put some bit in our message to distort this.

DLL



what is seq come in body . Then what to do?
 → One way put length in header → does not hold
 chicken & egg?

→ So we can see seq has 6 consecutive ones. So where ever 5 ones in message insert a redundant zero. Whenever 5 ones put zero after it. Breaks the sequence:-

AT Transmitter

actual seq

00111110101010111010101011111101001111110

Ex √ This is header, body, crc

whenever 5 consec once add a zero Rx stops reading

001111010101010111100101011110 1110100111110

green are added zero.

Now receiver has these bits. They know there is bit stuffing. Whenever he see some one zero ignores '0'.

If he sees 6 ones then it belongs to seq.

→ If he sees 6 ones continuous and one one. Then there is an error and we discard the frame.

cyclic redundancy check

Header body $\xrightarrow{k \text{ bits}}$ CRC $\xrightarrow{\text{extra}}$

$n \text{ bits}$ \downarrow \nearrow extra bits

to check if there
is a bit error
in header and body

function of data of header & body

i) want to be able to easily detect a large type
of bit errors. (single bit, burst of consec. errors etc)

2) creation of CRC & verification of CRC should be
computable & efficient:-

3) for given K let shld work for any n .
→ shld be computable for any n :-

DATA WORDS

(n-bit)

sent by TX

01101 (correct)

01100 (flipped)

01100 (incorrect)

same wave

receiver gets
this one

Here 2^n
possibility

→ n-bit string

CODED WORDS
(n+k)-bits CRC is there

(we shift)
are mapping

Here 2^{n+k}
possibility

some possible messages
may not be mapped here.
so these messages are never
possible.

HAMMING DISTANCE

given 2 code words : a_1, a_2, \dots, a_m
 b_1, b_2, \dots, b_m

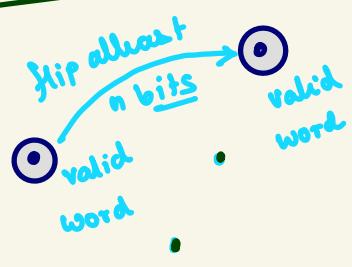
The number of bit in which they differ is
called the hamming distance :-

MIN HAMMING DISTANCE : Minimum of H.D between all pairs of code words of code :

We want design code words with maximum H.D

: error detection : If min H.D is N (of a coding scheme) then we can always detect atleast $n-1$ bit errors

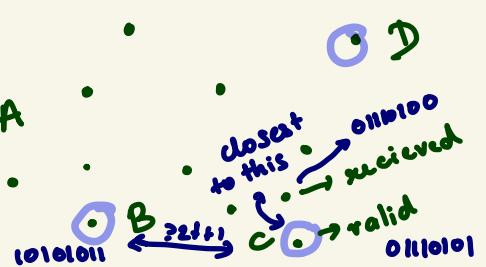
Because there are 2^n possibilities to flip n bits or more because min possible distance is $\geq n$



→ If we flip less than n bits then we go from valid to invalid

- If we want to error then it should fall on invalid code words. for that less than n bits

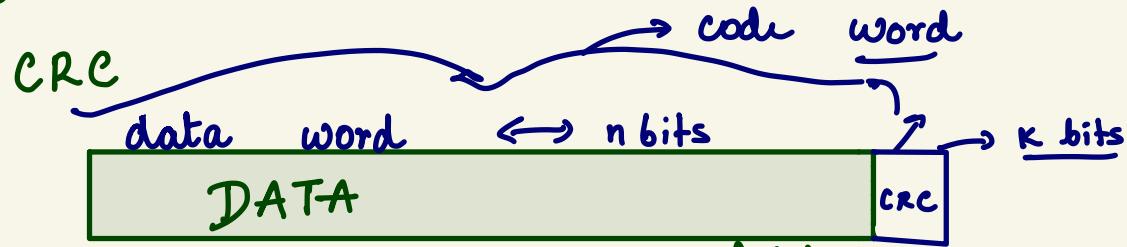
↳ error - correction



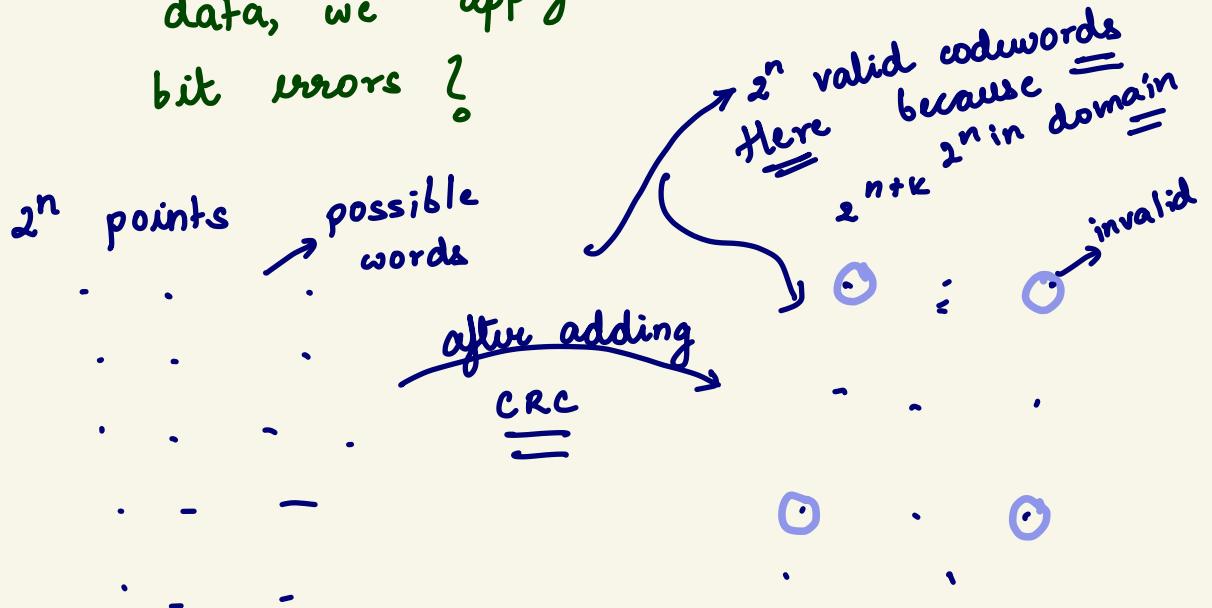
How one to choose correct can be caught at max $n-1$ errors

So choose this because we know we have 1 bit errors.
 Suppose min H.D $(2t+1)$ and # bit errors is ' t ' or
 less then by mapping received code to the max + valid
 code word correct all error

Galois Fields (To generate CRC)



we have added all fields to
 data, we apply crc to detect
 bit errors?



Minimum Hamming distance should be large to detect as many errors as possible.
 No matter how large data word code word should have a CRC.

Galois Fields
 BIT MUL 0,1 and

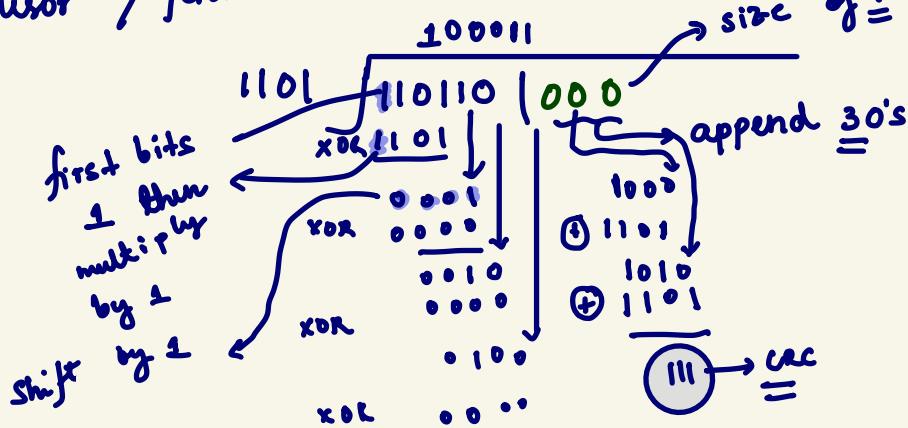
X	0	1
0	0	0
1	0	1

$+$	0	1
0	0	1
1	1	0

How generate CRC?

Ex) Data: 110110 $K=3$

Divisor / Generator: 1101 ($K+1$ bits size)

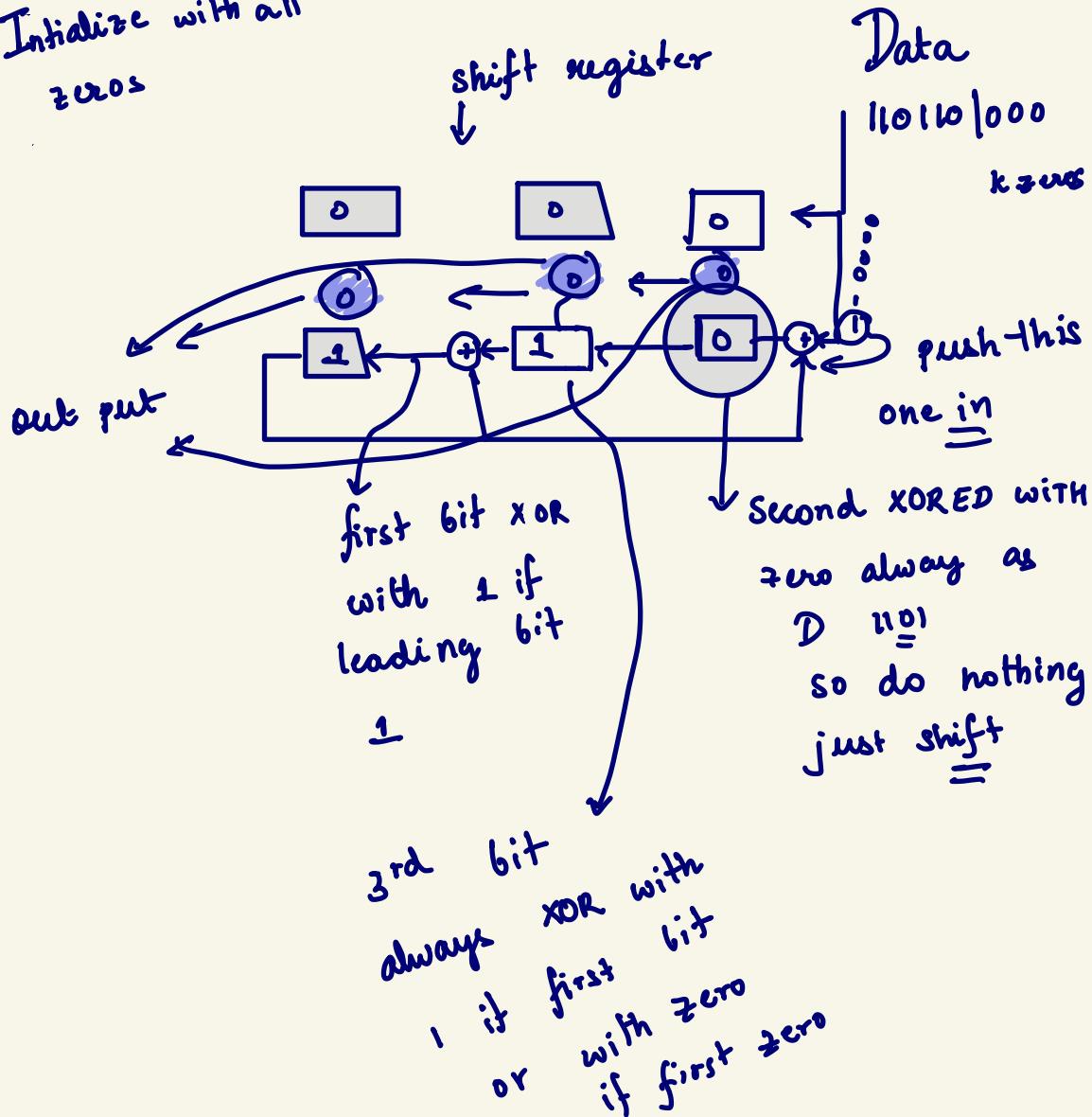


So code word

110110111

How to build a circuit

Initialize with all
zeros



We have code word at receiver
We get codeword as

$\frac{\text{DATA}}{=}$ | CRC

At receiver we can apply CRC process on DATA'
and see whether $\text{CRC}' = \underline{\text{CRC}''}$

Or we can feed DATA' + CRC' and see if
remainder $\equiv 0$ implements
 ↑ division
 by $C(x)$

same
circuit

OR

How to prove this method works
polynomial representation of bit strings

$$\begin{array}{rcl} 1101 & \rightarrow 1 \cdot x^3 + 1 \cdot x^2 + 0 \cdot x^1 + \underline{1 \cdot 1} & = x^3 + x^2 + 1 \\ x_3 & x^2 & x \\ & \downarrow & \downarrow \\ & 1101 & \end{array} \quad \xrightarrow{\substack{\text{divisor poly} \\ \underline{\underline{C(x)}}}} \quad$$

$$\begin{array}{rcl} (C(x)) (1+x) & \rightarrow & 1101 \\ & \xrightarrow{\substack{\text{OR} \\ \underline{\underline{C(x)}}}} & \end{array} \quad \begin{array}{rcl} (1+x) (x^3 + x^2 + 1) & \rightarrow & 101 \\ (x^4 + x^3 + x + x^3 + x^2 + 1) & \xrightarrow{\substack{\text{OR} \\ \underline{\underline{(x^4 + x^3 + x^2 + x + 1)}}}} & \end{array}$$

$$\begin{array}{rcl} \begin{array}{l} 1101 \\ \oplus 1101 \\ \hline 1011 \end{array} & \rightarrow & 11101 \\ & & \end{array}$$

$$= x^4 + x^2 + x + 1$$

Errors? error bits

$$110110111 \rightarrow P(x) \rightarrow \text{codeword}$$

$$+ \underbrace{000001001}_{\text{same same}} \rightarrow E(x) \rightarrow \text{error}$$

$$\text{we receive } 0 \mid \frac{P(x) + E(x)}{C(x)} = 0 ?$$

if zero we
are not detect
error

if non-zero
we caught the
error.

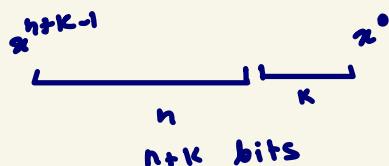
Types of errors

1) SINGLE BIT ERROR

$$E(x) = x^i \rightarrow \text{error function}$$

~~~

$\rightarrow$  for some  
 $0 \leq i \leq (n+k-1)$



$$\frac{P(x)}{C(x)} + \frac{E(x)}{C(x)} ?$$

If  $C(x) = x^k + \underbrace{\dots}_{\text{anything}} + 1$

$$(C(x) \neq (x^m + \dots + x^n))$$

$$(x^{k+m} + \text{something} + x^n)$$

we can never get just  $x^i$

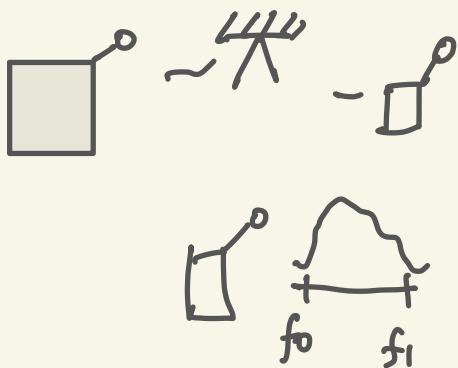
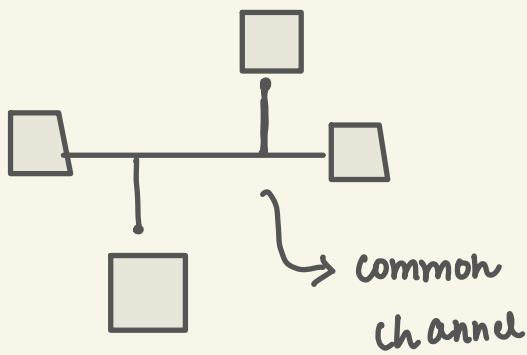
we can always detect one error



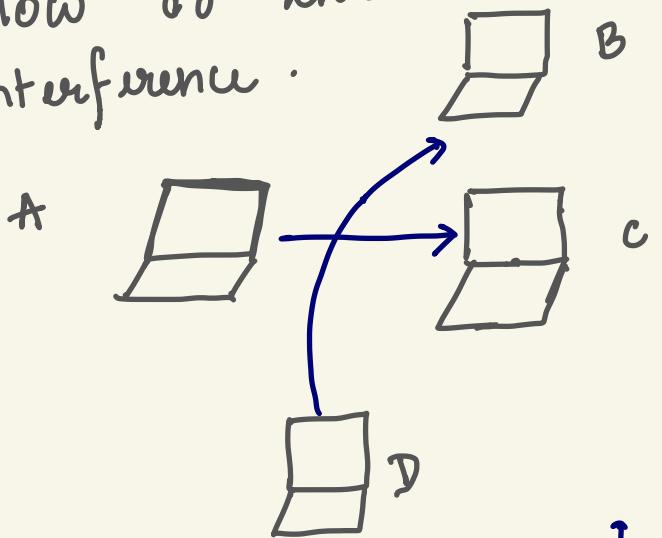




MAC layer  
many devices share the same channel

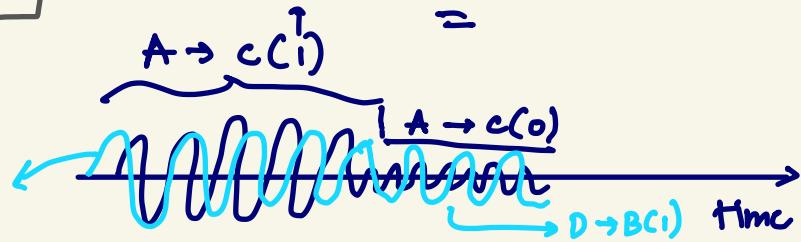


How do these communicate without interference:



A wants to send message to C  
D want to send message

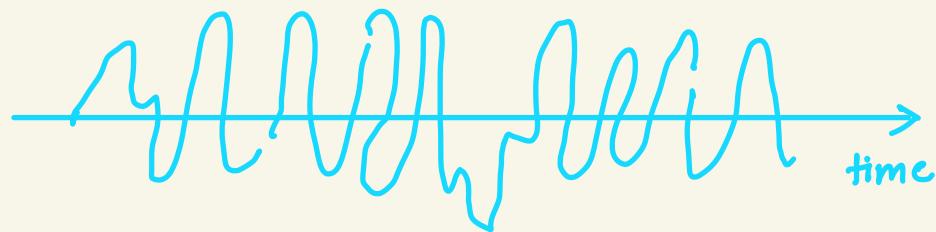
B =



At C :  
 $D \rightarrow B CI$

Now using same tech as above  
D starts communicating with B.  
Now C will also hear it.

How to make sense of this. So net signal received by C



Q> What is being sent, by whom for whom?

Q> We need some sort of identification  
So we know message is for whom.

Identifiers ?

This reduced bandwidth and throughput is reduced. Dividing band creates a performance issue

→ Performance: Throughput of A is  $\frac{1}{4}$  compared to if A used full band.

D

A

C

B

f<sub>0</sub>

f<sub>1</sub>

band

→ Divide the frequency band

into 4 parts

and each signifies who is sending.

→ What if we put a new device E!

→ Then what no way to make them a part of our network.

Addressing

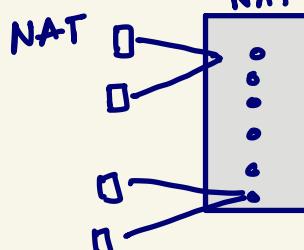
Assign bits to each node

$2^{32} \rightarrow$  problem (lack of assigning bits)

IPv4: 32-bits (layer-3)



$\{ 2^{32} \sim 4 \text{ billion} \}$  → We ran out these soon. So we use hacks  
many could be using same outer IP for world



↓ Solution increase number of bits!

IPv6: 128 bits

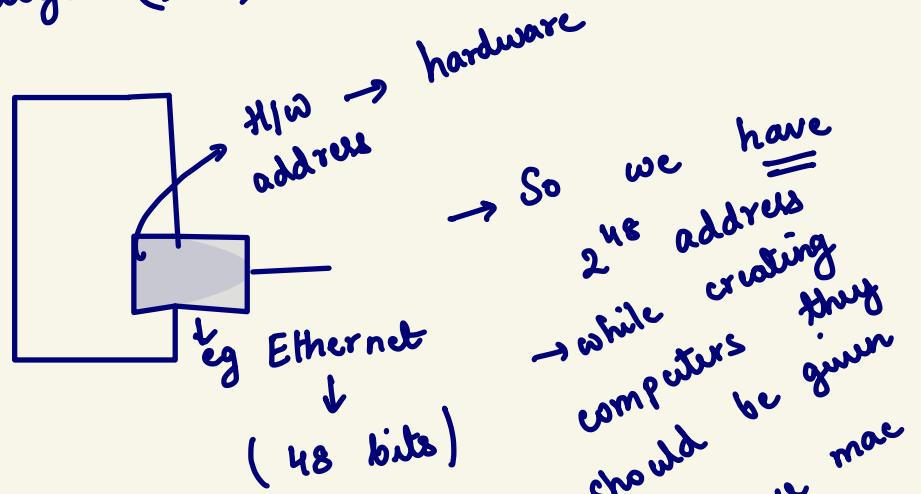
So for identification at MAC?

Q) How to assign address?

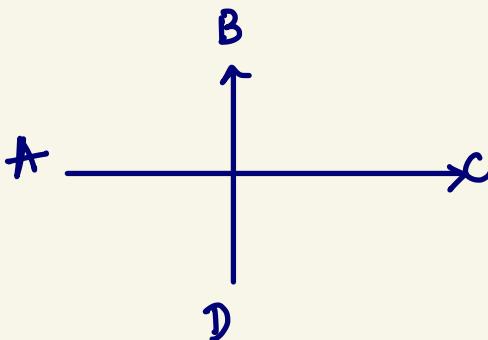
{ → manually  
→ automated }

} Some people could be entering the pool, some could be leaving.

We are at a low level, when phone boots up we want it communicate without manually given address  
M.A.C layer (DLL) → use hardware address

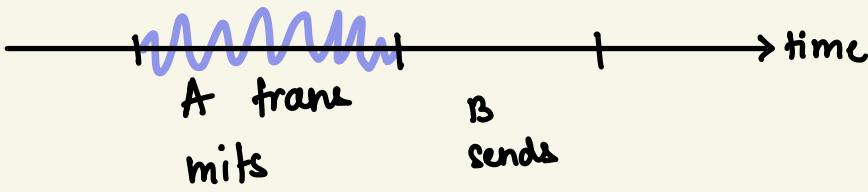


How to deal with interference  
→ Signals interfering with each other



→ One Solution  
diff freq band ?  
Who use what  
band !

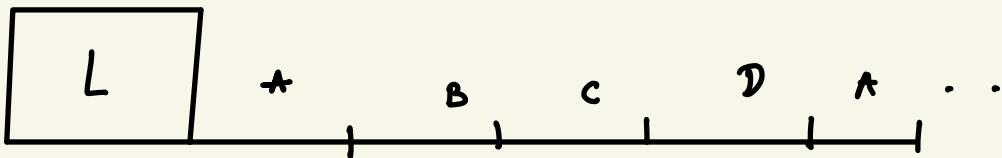
One Solution :- only one device sends at a time



→ many problems: Let us say 'E' was out and E did not hear A and E also sends message.  
→ Now when A is done who nexts send message.

→ How to decide who is allowed to send and when!  
→ One thing we can do is to have a leader!

Suppose leader is A.  
Leader could decide who is allowed to send message.  
Leader could decide who is allowed to send message. He could decide a temp schedule!  
and this to everyone. And we could keep renewing when a new node enters!



Leader sends schedule

So slots should not go waste.  
Also how should a new node communicate  
to leader to join.

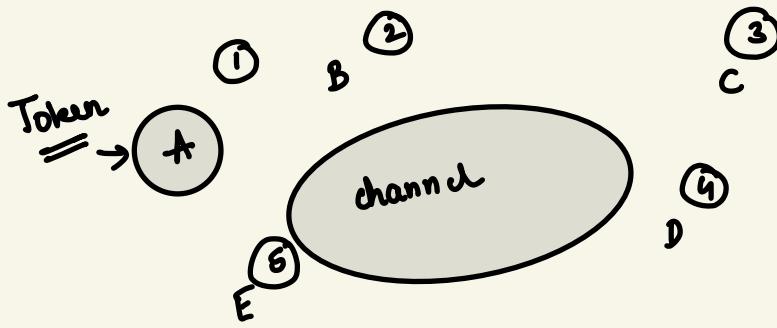
If we want to send request then  
I should have a different band.  
So we should be able to request leader  
for slots. We should schedule request rather  
than devices!

We should ask base station (leader). We will  
need multiple channels so the leader can  
get all possible information.

Who should be allowed to be send to  
leader. Here we face collisions.  
How to avoid collisions here while sending  
to leader.

Token Based → method

Whoever has token is leader



A is transmitting, when A is done he will release a special message I am willing

to release the token:-

Now let us say B has highest prio

next, so A will give token to B and

should communicate this over channel

So we should be able order devices and token should not be with one device for a long time.

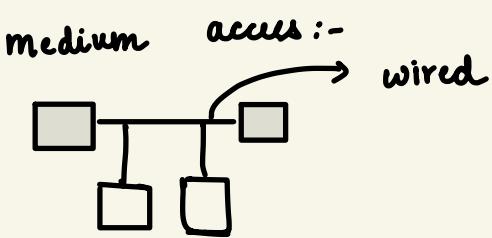
What if a particular node fails?

Some issue ↑

---

No leader, No mechanism →

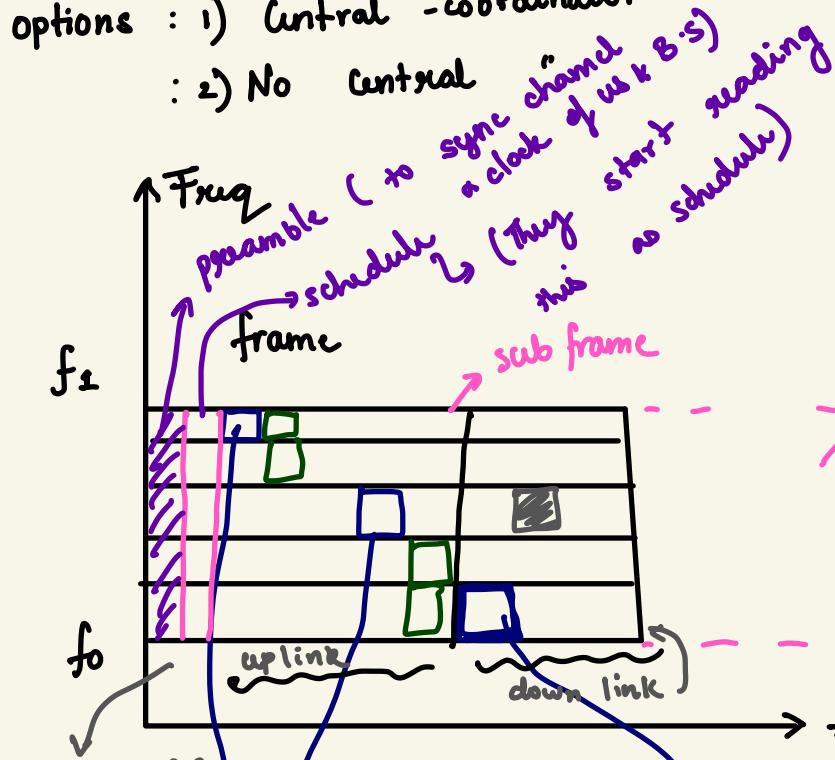
Wifi → allow collisions      } CSMA protocols  
Ethernet → collisions



The central co-ordinator decides when which devices can transmit and which devices can receive.

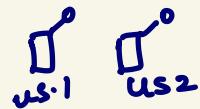
options : 1) Central - coordinator

: 2) No central



In these three tiles user one is allowed to transmit to B.S

user 1 receives data from B.S



same frame repeating again



So green is for us2

So how do users know which tile is theirs.  
A schedule is sent out before hand so  
each user knows which section of time &  
frequency belongs to them.

why not something simple by giving all  
freq at a particular time instead of giving  
files! Why?

User - Sending Data - UP-LINK

" Receiving Data - Down LINK =

Every user sees different kind of attenuation  
for the same file. So  $U_1$  can see  
more attenuation for grey shaded tile  
and  $U_2$  has less. So he can push  
more data to  $U_1$  at this file.

\* So we can optimize.

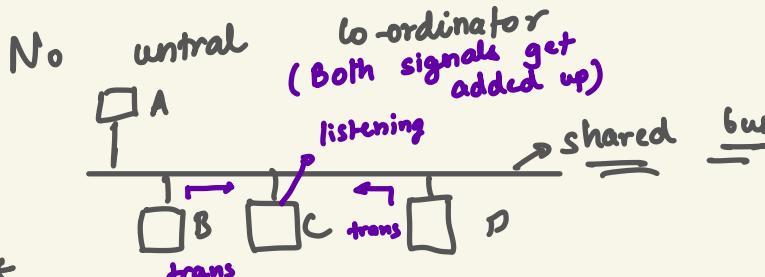
So users actually can tell B-S what tile is good.

Here in same frame we have up link and down link in the same frame.

In Ethernet only one frame can be from either side. In 4G both devices are transmitting. The up link ( $US \rightarrow BS$ ) and down-link ( $BS \rightarrow US$ ) is in the same frame. The measure from time preamble for their slot.

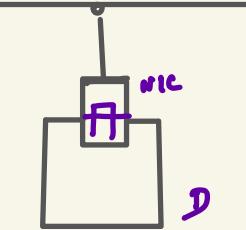
### Issue with central co-ordinator

- Single pt of Failure
- May not be possible in certain situations



\* If message not for them they discard. We want to connect and remove devices and use them for our advantage. We want play and pause. We have a single channel, no multiple channels.

- Want
- 1) Plug & play - connect & disconnect
  - 2) No central co-ordinator



- We could use token system
- 2) We could use a buffer.
- 3) We could have collisions because many people are transmitting

Idea: Suppose most of time 1 node transmitting only rarely we have collisions.

Collisions: Multiple nodes transmit simultaneously signals add up at receiver

↳ Neither signal can be deciphered.

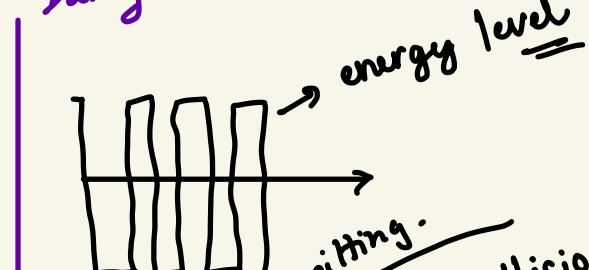
(On wire not much attenuation)

Normally



During collision amplitude increases and we can use energy level

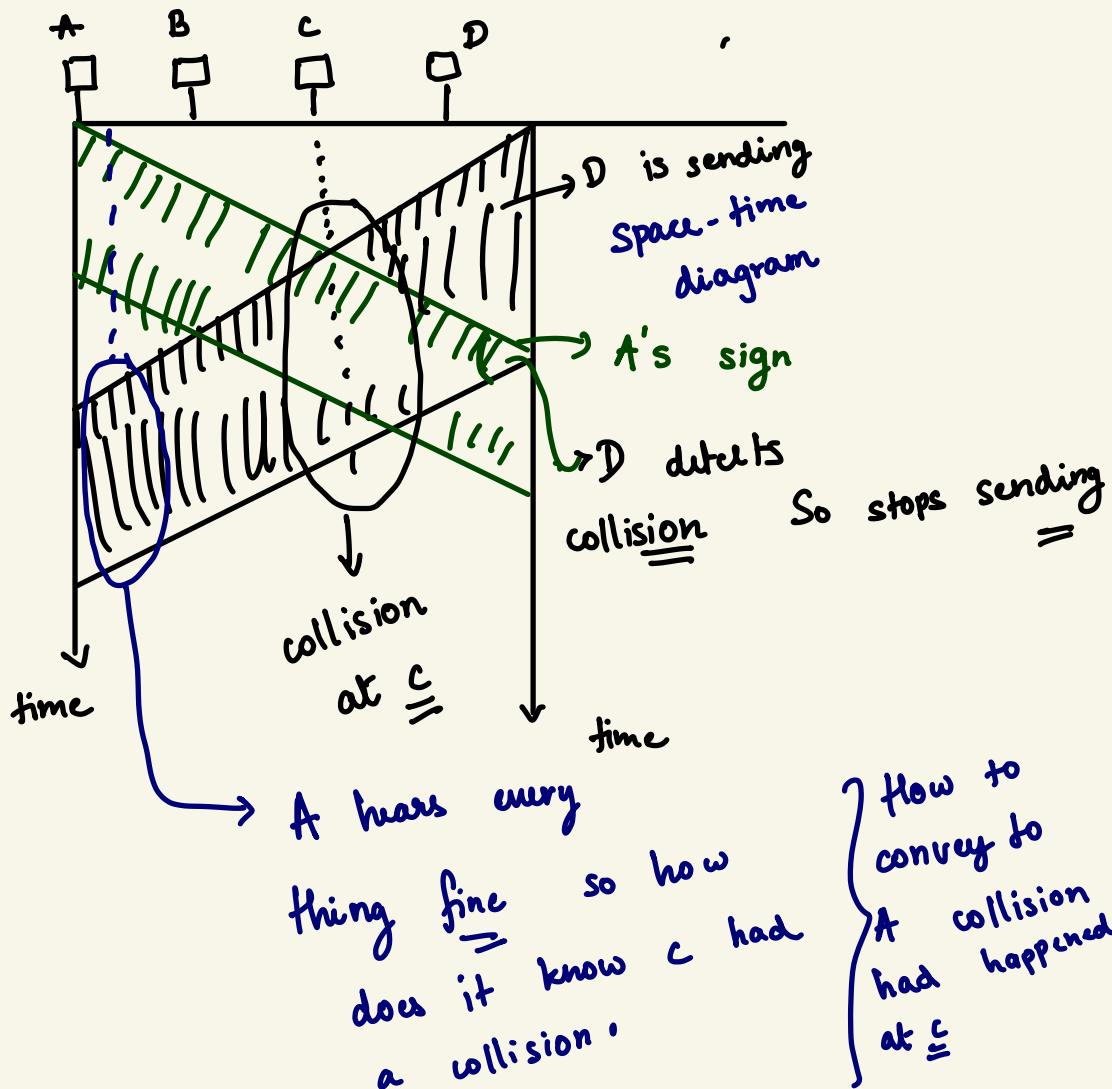
During collision



Once we see collision we should stop transmitting.

we should stop transmitting

If collisions occurs  $\rightarrow$  stop transmitting



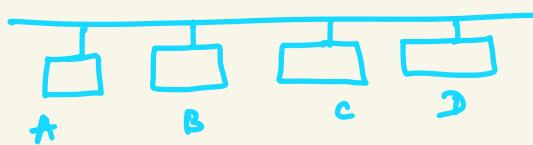
The ethernet does not allow to send small frames. There is a minimum size of frame. So this make sure every one hears the collision

# Ethernet

CSMA-CD → collision detection

carrier sense

→ CD: check energy levels.



→ if higher energy seen than expected then collision

Ethernet is created by

IEEE : 802.3

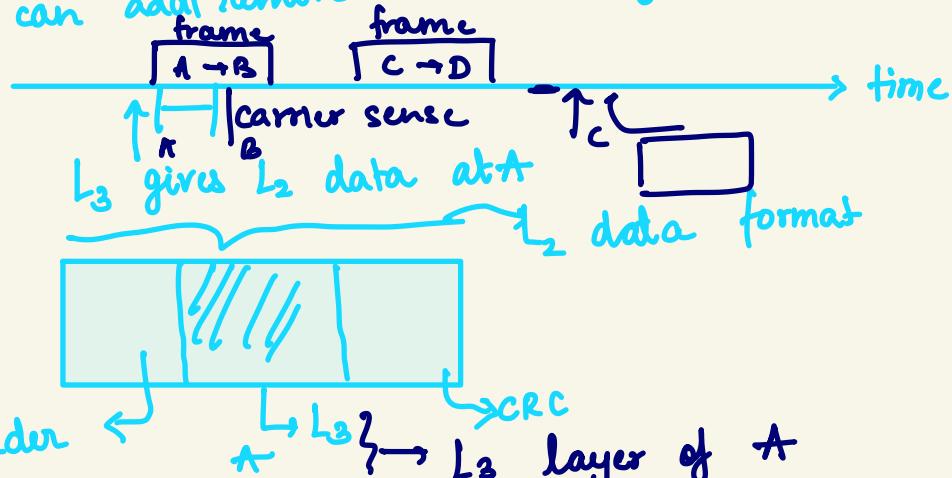
Wifi is a sister protocol

WIFI: 802.11



→ No centralized co-ordinator

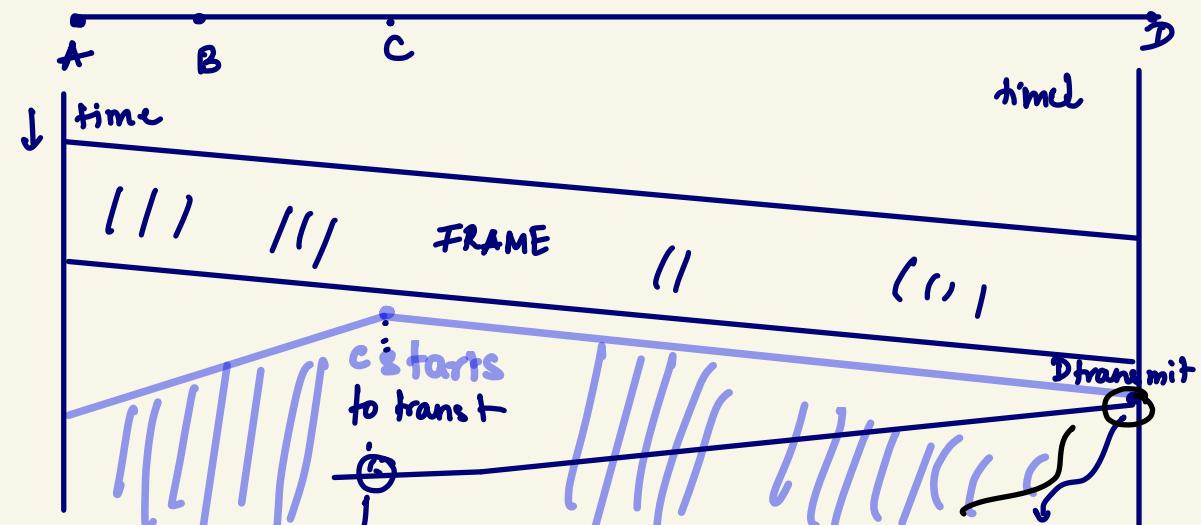
→ We can add remove nodes easily



When I am transmitting no other should transmit  
We use carrier sense to detect energy if  
some body is transmitting then carrier detects  
it and gives a "no" signal.  
Ethernet uses manchester encoding

- A is sending to B
- If no collision then assumed the other party got data perfectly
  - We want to pass data, it checks when there is silence on the channel.

space-time



~~||||| | | | | | | |~~

point in time  
where C  
detect collision

\* When stop collision is detected they  
transmitting.  
D stops  
immediately

If D detected and stop after a very small time and C misses this small signal sent by D then C has moved up.  
So D sends out jamming signal.  
The send one. So all hear collision.  
So all who detect collision we send a jamming signal.

Jamming signal has to be long enough!  
No how to decide who gets to transmit next!  
Because collision can happen with n senders where  $n \in \mathbb{N}$ .  
So we use random wait time

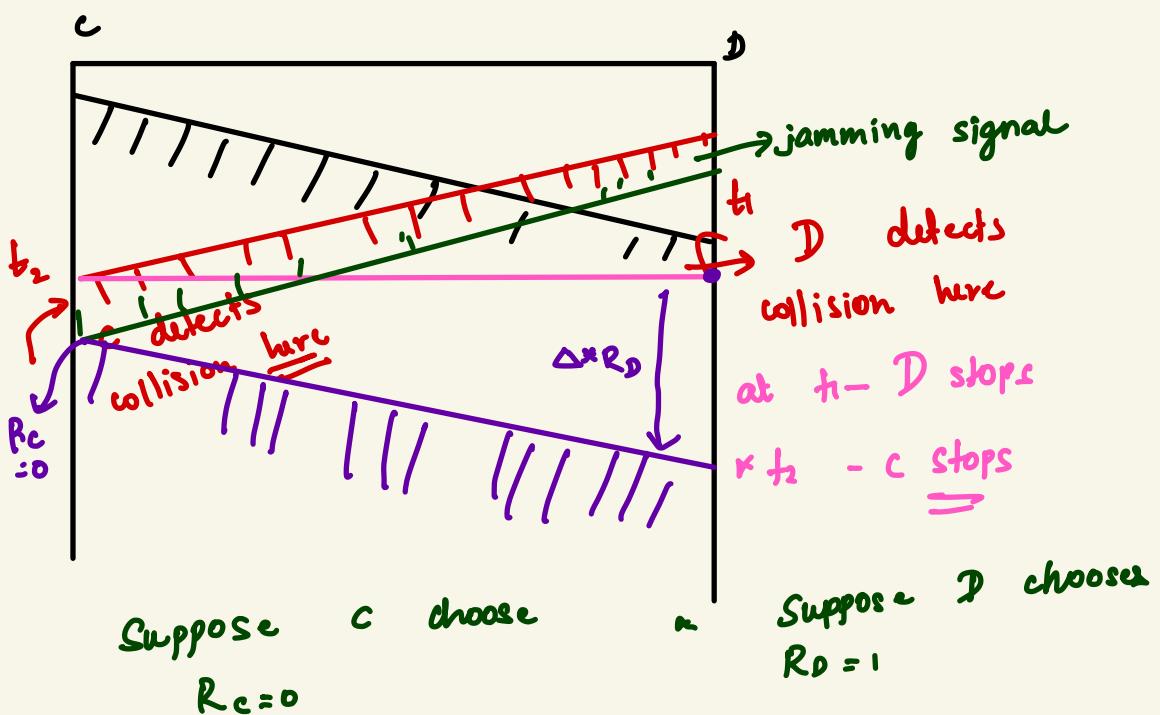
Each time node throws Random Number  $R$ ,  
 uniformly from some range wait for  $R \times \Delta$   
 before trying to send again:-

The one with unique min will win

Idea:

$R_c < R_D$  then  $c$  will transmit next  
 without colliding

We could again have collision if  $R_c = R_D$   
 and if  $R_c < R_D$  we want  $c$  to win so  
 what  $\Delta \leq \therefore$



we want  $c$  to transmit without collision

We D stops hearing signal he will wait  $R + \Delta$   
 so  $R_D + \Delta$  should put D in place where C  
 starts transmitting and C's signal should  
 reach D.

$\Delta$  should be large enough for D to  
 hear C before transmitting.  
 So what should this  $\Delta$  be.  
 If  $\Delta >$  Round trip time then we are safe  
 Can show, if  $\Delta > RTT$  then D will not transmit  
 before C

$RTT = 2x(\text{time for signal to go from one end}$   
 $\text{of network to another})$

Max cable length = 2500m

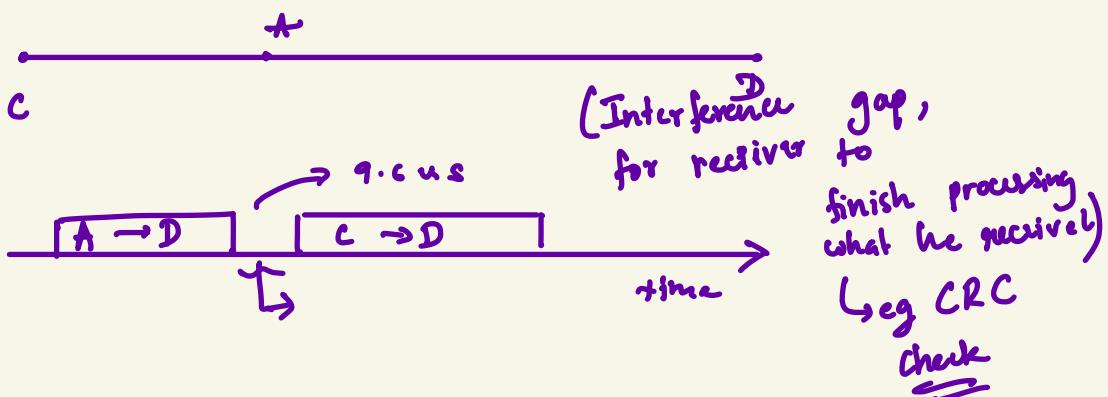
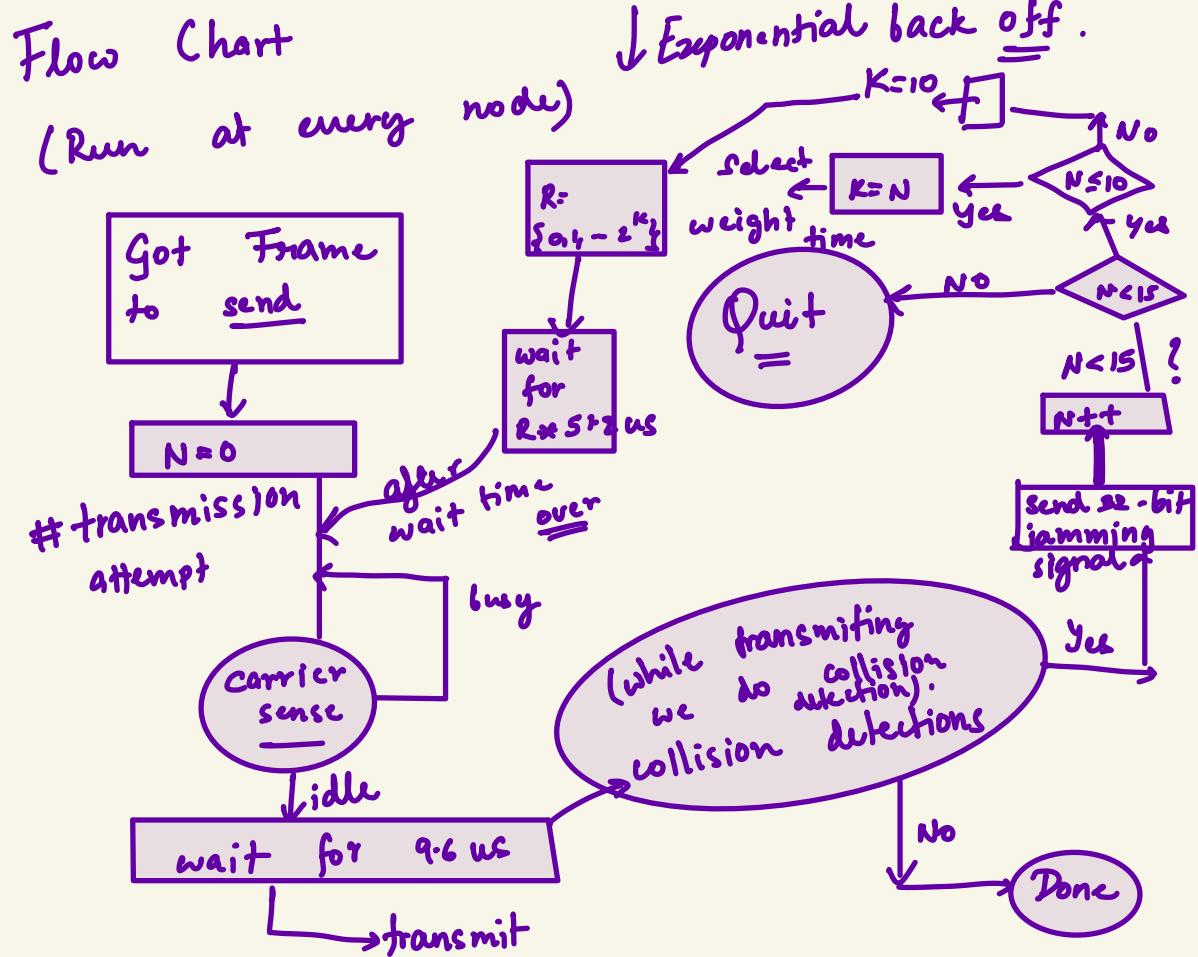
speed of signal =  $2 \times 10^8$  m/s

$$OWD = \frac{2500}{2 \times 10^8} = 12.5 \mu s$$

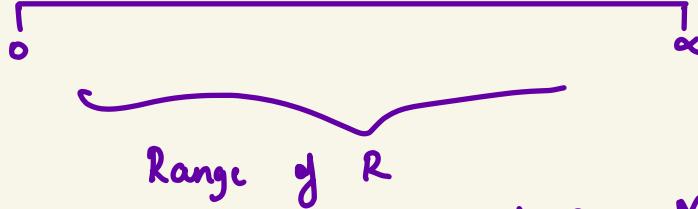
These repeaters are also there so we delay due to them



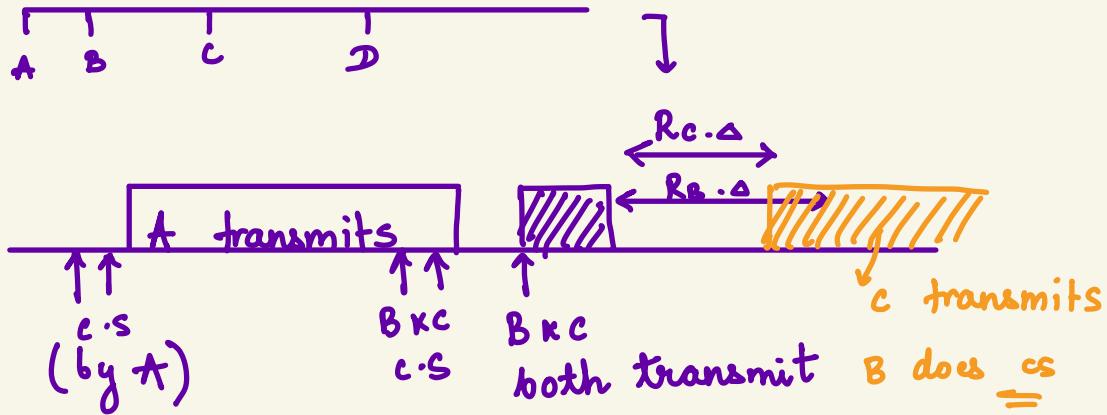
delay due to repeaters  $\rightarrow$  deal  $2 \times 12 = 24 \mu s$   
 otherwise  $51.2 \mu s$



Exponential Back Off  
 • Don't know how many people colliding  
 Suppose M colliding



this  $\alpha$  should depend on  $M$ .  
 There will some optimal  $\alpha$  we want  
 go over it so we keep increasing  $2^k$ .  
 CSMA - CD - Ethernet



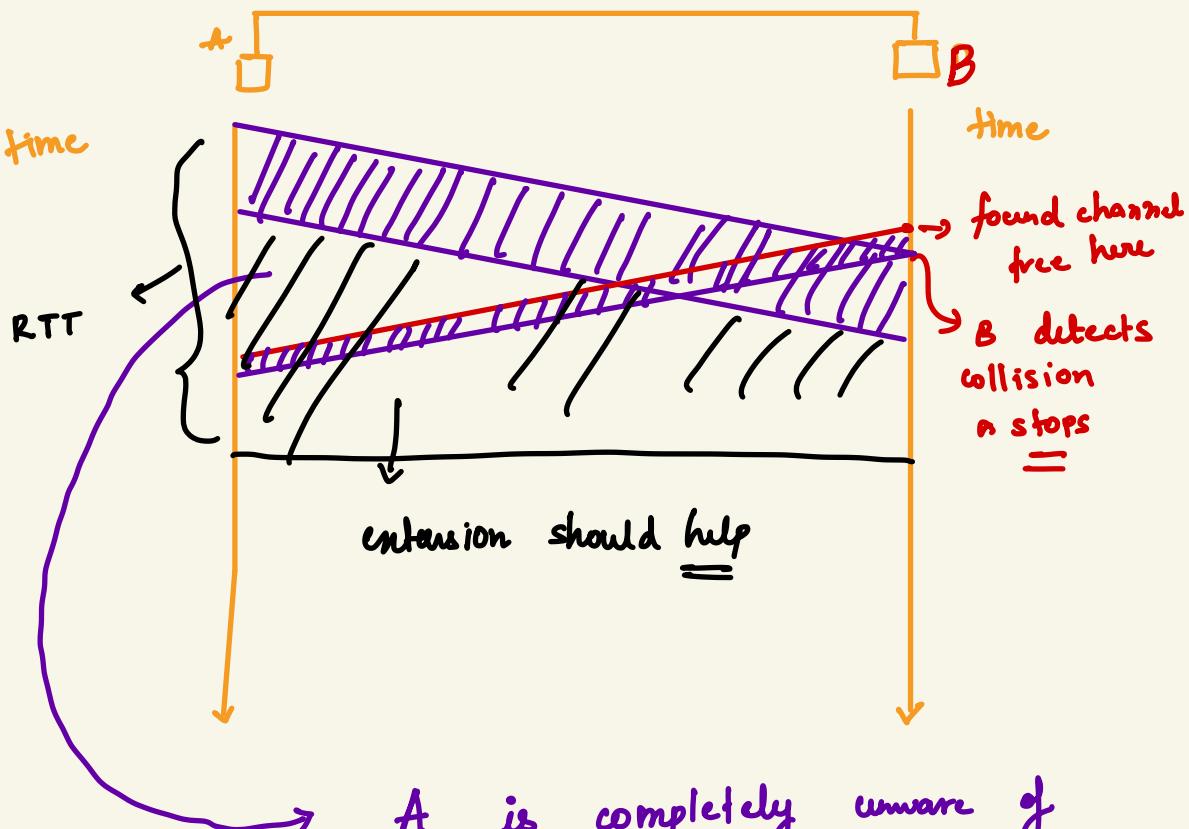
$$R \in \{0, 1, 2, \dots, 2^k\}$$

(we do this with hope that there is a unique min)

K increments upto 10 for every collision  
 (for same frame)

Simple, Robust  $\uparrow$   
↳ just connect machines to the bus  
↳ limits on frame size (Ethernet)  $\underline{\underline{802 \cdot 3}}$

## 1) min frame size



A is completely unaware of it because B's signal reached to it when A finished transmitting so all is fine for him

So we can do one thing we can make frame of large enough to detect collisions while detecting:-

$\Delta = \text{worst case RTT}$

$\Delta = 51.2 \mu\text{s}$

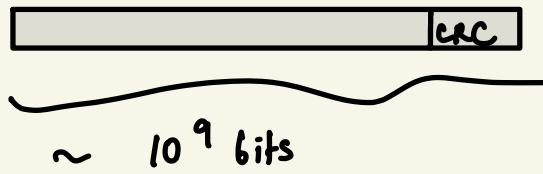
Original standard      10 Mbps  
Frame size 64 bytes = 512 bits

Time to transmit @ 10 Mbps =  $\frac{512}{10^7} = 51.2 \mu\text{s}$

If Mbps ↑ we reduced length by that factor so RTT ↓ as  
@ 100Mbps S.12 RTT

## 2) Maximum frame size

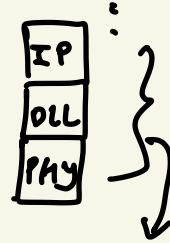
### (a) Bit errors



$\sim 10^9$  bits

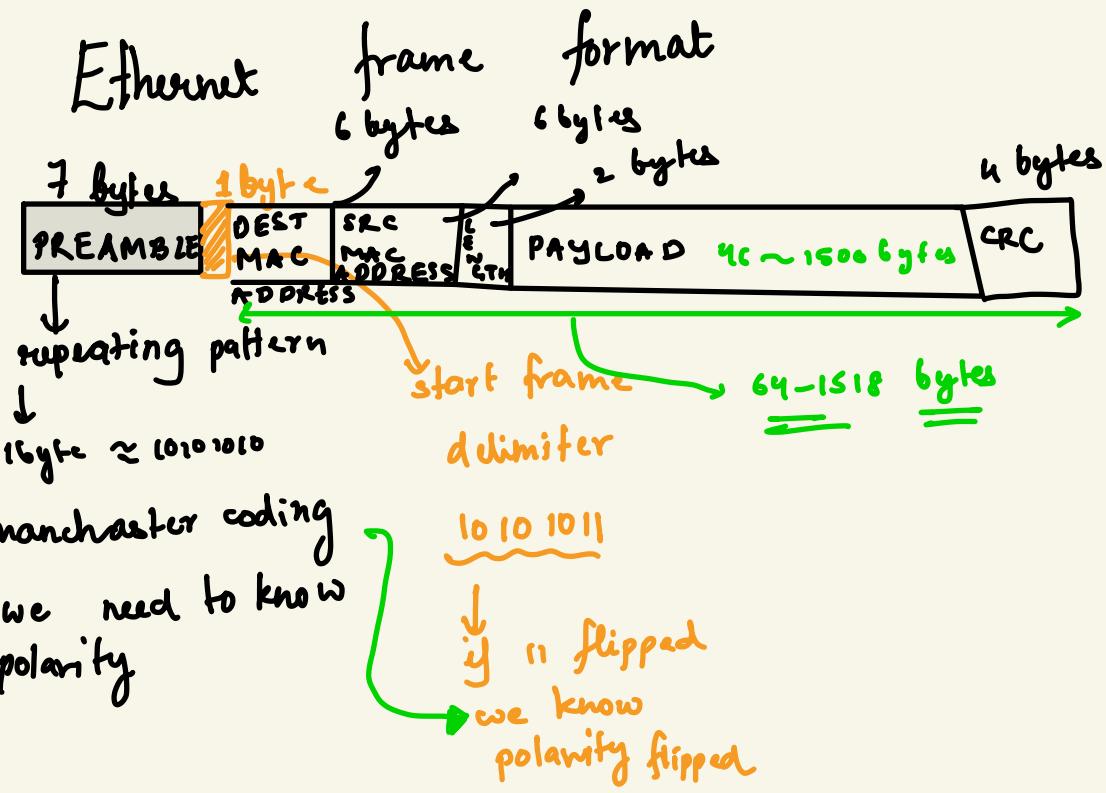
1 bit error in  $10^6$  bits  
so almost 1000 errors  
∴ High prob of at least 1 bit error

b) time of transport increase  
= if 10 Mbps then 100 seconds  
message!

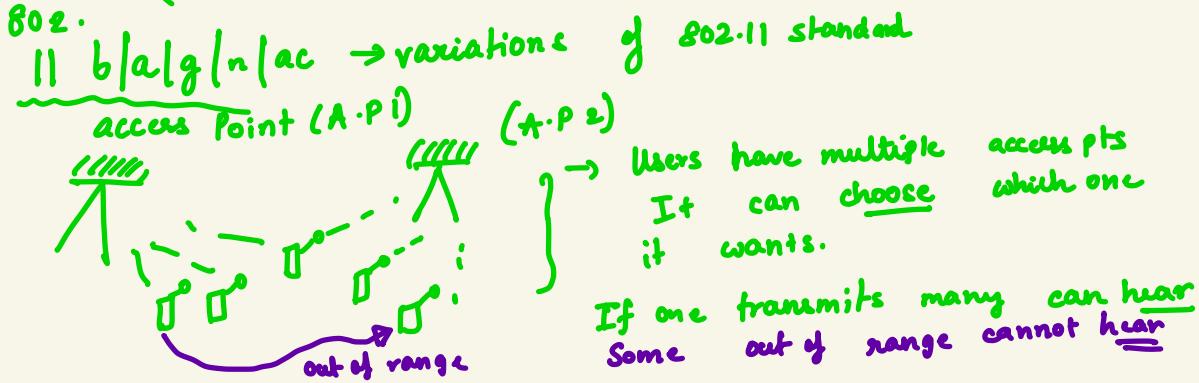


So while encapsulation we should keep in mind because for above max length

b) Don't want single user to monopolise the channel for very long.  
Max frame size = 1518 bytes



## WIFI (IEEE 802.11)



↓ Ethernet topology in wifi

will carrier sense for CSMA-CD work?



eg)

A → AP2  
B → AP2  
C ← AP2

A cannot sense B's (signal)

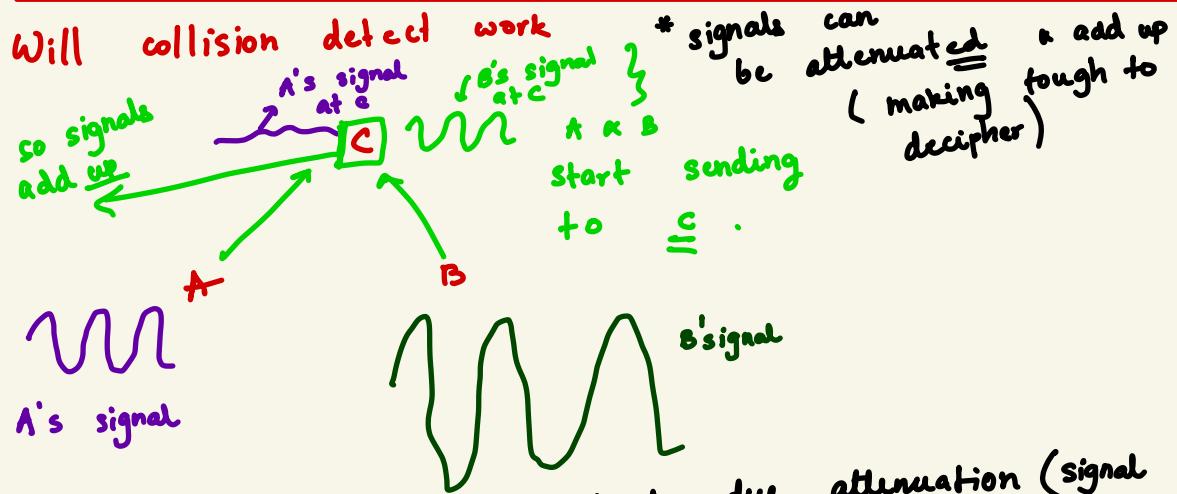
No because out of range issues!

So A is transmitting to AP2 and B is transmitting already to AP2 but out of range of A. So collision.

Hidden terminal problem

hidden.

C.S may not work



B's signal is like noise at A due to attenuation (signal decay with distance  $d \sim \frac{1}{d}$ )

When one person is transmitting he cannot hear anybody else.

Everybody is whispering compared to no.

While transmitting A cannot hear anyone else.

Collision detection does not work!!.

Here we cannot have same type of assumption as ethernet "if we hear no collision there is no collision".

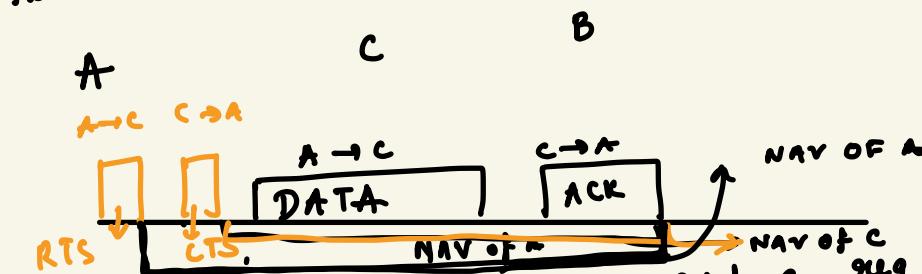
Fix:

Acknowledgment.

Receiver sends ACK

  
If ACK collided we send  $A \rightarrow C$  again!!

How to handle Hidden Terminal?



Before sending data we send a req to send to  $\leq$   $C$  will send a clear send. This to avoid the case where  $A$  sends while  $B$  is already transmitting to  $\leq$ :  $B$  will not hear RTF but will get CTS. This makes sure all potential sender to  $\leq$  get CTS meant allocation for one. CTS & RTS have a NAV (Network vector. (Time during).

Rule : All hearing RTS & CTS (other than ACK must remain silent for NAV duration)

VIRTUAL CARRIER SENSE TO B due to CTS OF C.

All senders to ACK are frozen during NAV.

• What if A send RTS but does not get CTS?

- A will assume collision  
and retransmit!

• What if RTS, CTS, DATA acknowledgement to A?

same over here  
are sent but no  
we again send RTS!

## WIFI MAC (IEEE 802.11)

### LAN

In wireless, it is so not straightforward

→ Cannot collision detect when we are transmitting our self.

→ Also carrier sense is difficult due to hidden terminal node

New protocol

CSMA - CA

↳ collision avoidance



virtual carrier sensing

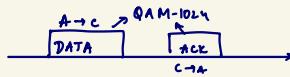


→ optional, can be disabled

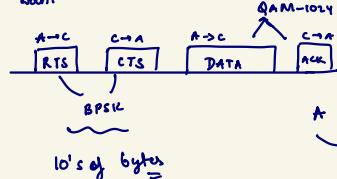
Why would one switch off RTS & CTS? This is overhead and it takes a lot of time..

Different modulation schemes

BPSK, QPSK, QAM-16, 64, 256, 1024



I want B to hear RTS & CTS of A to C



lots of bytes =

QAM-1024

CTS

RTS

DATA

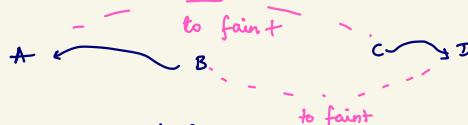
ACK

→ We use BPSK  
for RTS & CTS

→ Because we want to  
avoid bit errors

→ So we use less modulation  
because we detect with low  
probability of others.

exposed terminal problem



A can hear only B

B hears A, C

D only hears C

B sends to A

C sends to D

B → A, C → D possible together in theory

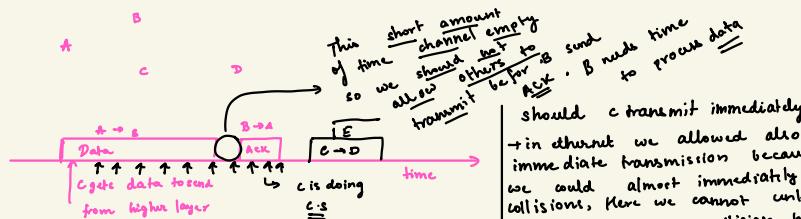
Case - i)

RTS/CTS enabled

If B starts first C remains silent (due to

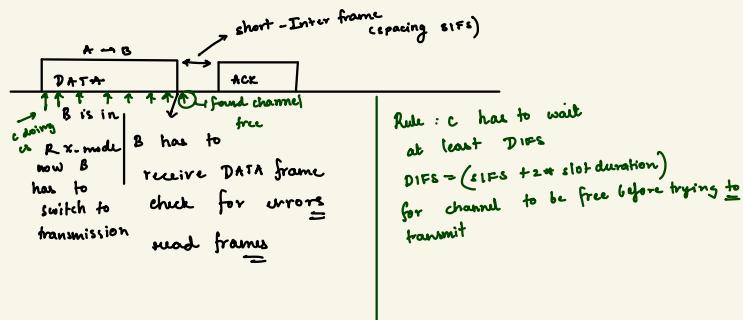
RTS).

When we do RTS/CTS we block C from sending data to D because C does C.S. But it could have sent to D because B's signal is too faint at D. So it would not affect C → D transmission. So due to WiFi MAC protocol could not allow such theoretical transmission.

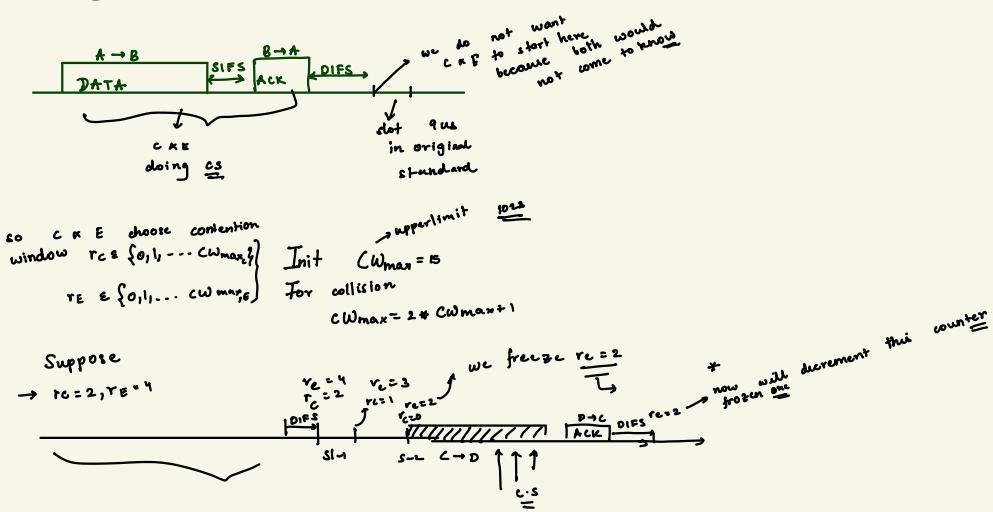


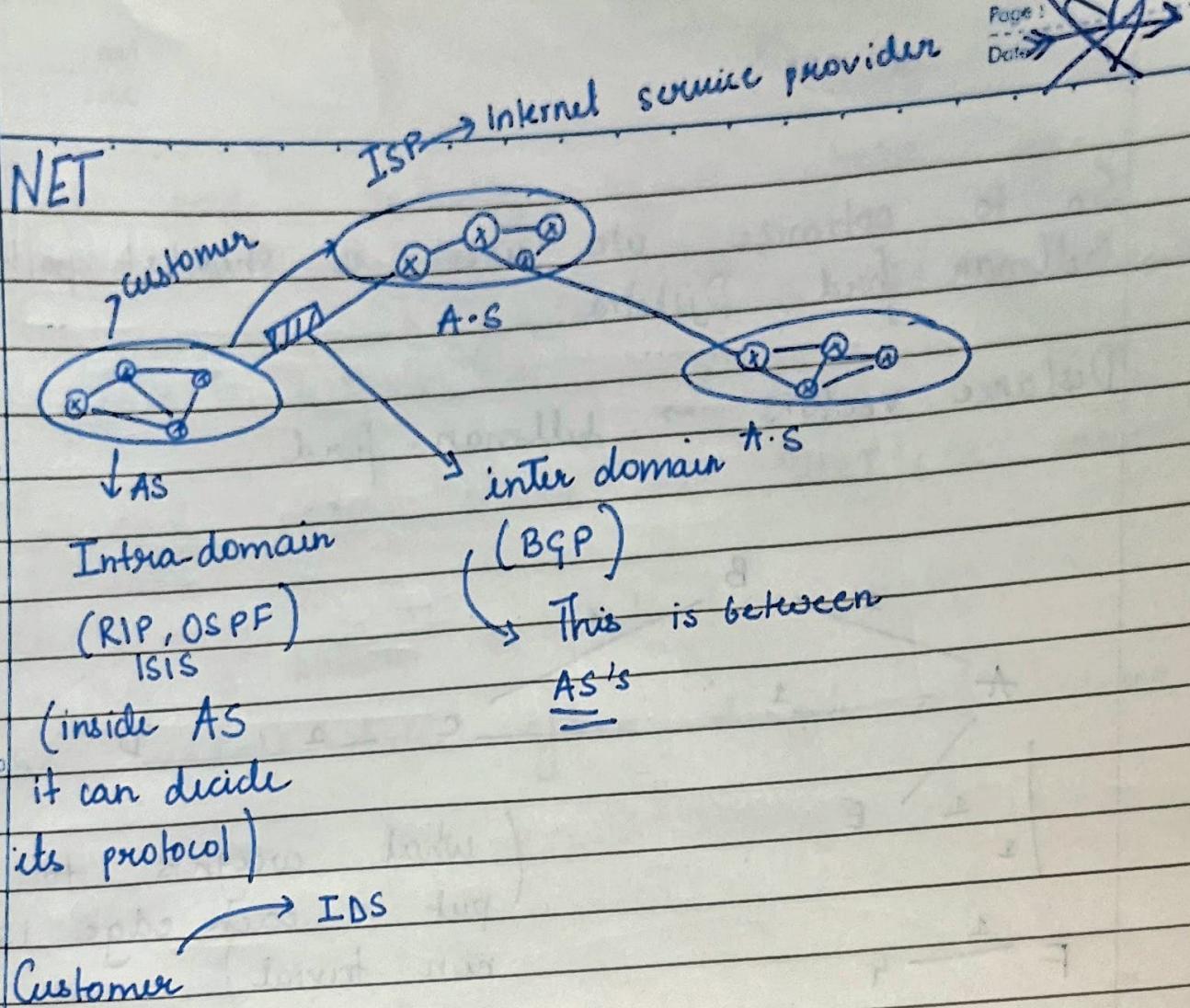
should C transmit immediately

in ethernet we allowed also immediate transmission because we could almost immediately detect collisions. Here we cannot unless we get an ACK, so collision here are expensive because if two simultaneously transmit we have to wait for ACK



Assuming if DIFS ends





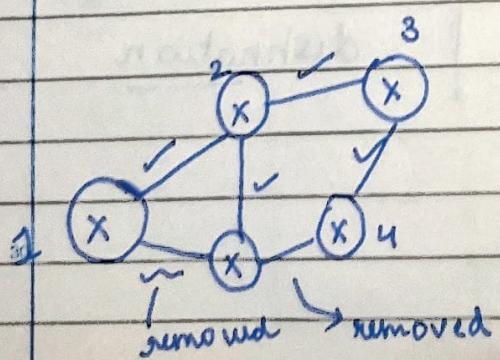
### service-level agreement (SLA)

Ex

- 100 Mbps
- > 99.9%
- UP-TIME
- 30 ms latency  
within own AS
- Packet-drop rate < 1%

} TCP is designed to create congestion

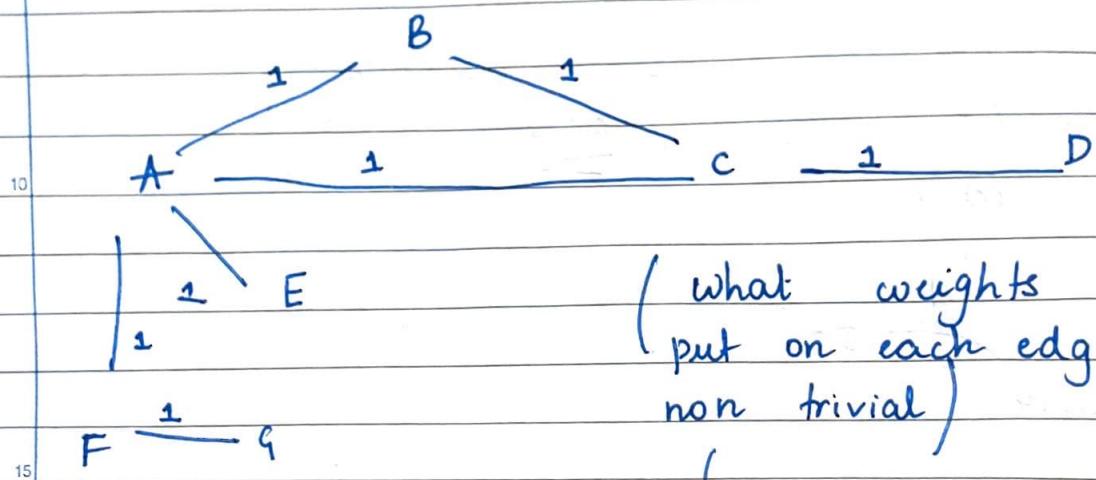
### intra-domain Routing (interior-gateway protocol) (IGP)



↓ spanning tree protocol disadvantage  
if 1 wants to send to 4  
he has to via 1 → 2 → 3 → 4  
if we use shortest

So to optimize we use a shortest path algo  
 → Bellman-ford, Djikstra.

5 Distance vectors  $\rightarrow$  bellman-ford



(what weights to put on each edge is non trivial)

At A (init):

| DEST | COST | NEXT-HOP |
|------|------|----------|
| A    | 0    | -        |
| B    | 1    | B        |
| C    | 1    | C        |
| E    | 1    | E        |
| F    | .    | F        |

\* min-cost known  
to various  
destination

At C (Init)

| DEST | COST | NEXT-HOP |
|------|------|----------|
| C    | 0    | -        |
| A    | 1    | -        |
| B    | 1    | -        |
| D    | 1    | -        |

- Each node initially knows about its immediate neighbors
- But we want whole network to be known to each network :-
- So each nodes shares its DEST | COST columns of table
  - ↗ this is known as destination vector.
- So update table after A gets info from C.

| DEST | COST | NEXT-HOP |
|------|------|----------|
| A    | 0    | -        |
| B    | 1    | B        |
| C    | 1    | C        |
| E    | 1    | E        |
| F    | 1    | F        |
| D    | 2    | (C)      |

↓ info from C

if A  $\xrightarrow{\text{edge}} B$  edge had 5 value, after getting info from (C)  $A \rightarrow B$  path can have value 2

| DEST | COST | NEXT |                                      |
|------|------|------|--------------------------------------|
| (A)  | 0    | -    | Because going from A-C reduces load. |
| B    | 2    | (C)  |                                      |

Finally at (A)

- Periodic exchange of information
- Convergence of tables :-

At A

| DEST | COST | NEXT-HOP |
|------|------|----------|
|------|------|----------|

|   |   |  |
|---|---|--|
| A | 0 |  |
|---|---|--|

|   |   |   |
|---|---|---|
| B | 2 | C |
|---|---|---|

|   |   |   |
|---|---|---|
| E | 1 | E |
|---|---|---|

|   |   |   |
|---|---|---|
| F | 1 | E |
|---|---|---|

|   |   |   |
|---|---|---|
| D | 2 | C |
|---|---|---|

|   |   |   |
|---|---|---|
| G | 2 | F |
|---|---|---|

D.V exchange periodically  
with neighbours

After few iterations, tables  
converge

15

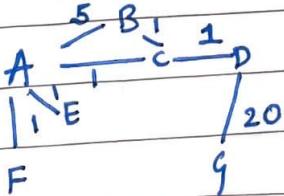
what F → G fails (This entry is no longer valid  
so F should tell all its neighbours if can  
reach G → )  
So F says my distance to G → ∞

So A should also send info to its neighbours it  
can reach G  
Everyone updates their table :-  
new graph

At A

| DEST | COST | NEXT-HOP |
|------|------|----------|
|------|------|----------|

|   |   |   |
|---|---|---|
| G | ∞ | X |
|---|---|---|



## count-to-infinity problem

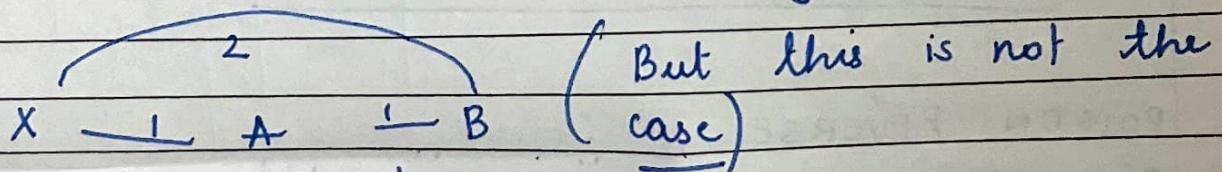
| X | A | B | DEST | COST | NXT   |
|---|---|---|------|------|-------|
| X | 1 | X | X    | 2    | A     |
| B | 1 | B |      |      |       |
|   |   |   |      |      | ↑ A+B |

→ A+B

suppose  $A \rightarrow X$  fails      A sends message to B  
 $(X, \infty)$  But B sends his table  $\begin{pmatrix} X, 2 \\ A, 1 \end{pmatrix}$

When A was sending B sends this

so A can think the topology to B is like



Now A makes table

| DEST | COST. | NXT |
|------|-------|-----|
| X    | 3     | B   |
| B    | 1     | B   |

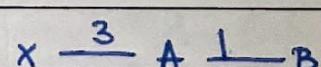
then sends to B  $(X, 3)$

B checked out  $\textcircled{X}$

But gets  $(X, 3)$  from A

| DEST | COST | NEXT |
|------|------|------|
| B    | 4    | A    |

Assumed network topology



Now B sends this to A

A makes  $X = 4$  then B makes  $X = 5$   
 this keeps happening till  $\infty$

- So we have loop  
→ So RIP (routing information protocol) say  $\infty = 16$   
→ So after cost = 16 we chuck it!  
So after few seconds/minutes this problem is caught! (During this time lot of packets in between A & B for x) Latency, pkt drop increased because we know only partial info of network thus this happens

10 SPLIT-HORIZON : DON'T share DV Entry with a neighbour, if that neighbour is next hop to concerned destination.

Ex : B does not tell A  $(x, 2)$  to A

15 since A is NXT-HOP for  $\infty$

→ Here we hide information

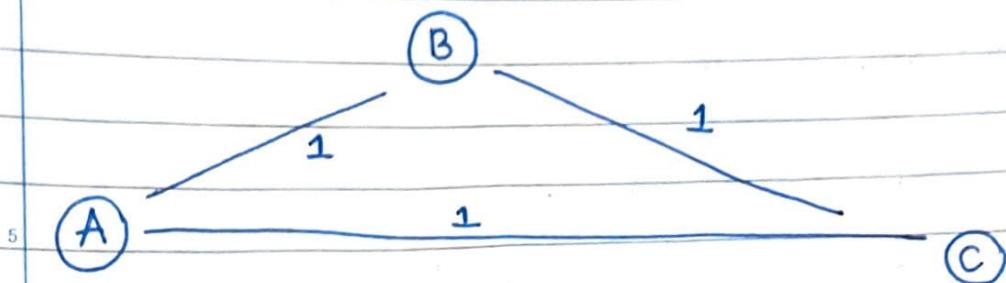
### POISON REVERSE

20 TELL NEXT HOP NEIGHBOUR , my distance is dist. (lost)  
to dest  $\infty$  is  $\infty$

Ex: B tells A  $(x, \infty)$  instead of saying  $(x, 2)$

They are not perfect } \* Some failures will too

## FAILURE OF SPLIT HORIZON



A + B

|    | DEST | COST | NXT |
|----|------|------|-----|
| 5  | A    | 1    | A   |
| 10 | C    | 1    | C   |
| 15 | E    | 2    | A   |

A + C

|    | DEST | COST | NXT | DEST | COST | NXT |
|----|------|------|-----|------|------|-----|
| 15 | B    | 1    | B   | A    | 1    | A   |
|    | C    | 1    | C   | B    | 1    | B   |
|    | E    | 1    | E   | E    | 2    | A   |

Suppose A-E fails!

$\left\{ \begin{array}{l} A \rightarrow B (E, \infty) \\ A \rightarrow C (E, \infty) \end{array} \right.$

→ A sends message to neighbours

→ Suppose  $A \rightarrow C (E, \infty)$  gets lost

→ So ~~sends~~  $c \rightarrow A$  no entry for (E)  
 "  $B \rightarrow A$  " "" ""

Now table of B AC

| DEST | COST     | NXT | DEST | COST | NXT |
|------|----------|-----|------|------|-----|
| A    | 1        | A   | A    | 1    | A   |
| C    | 1        | C   | B    | 1    | B   |
| E    | $\infty$ | +   | E    | 2    | A   |

\* So B sends this to  $\frac{C}{=}$  |  $(E, \infty)$ , C sends  $(E, 2)$  to C

Now B gets entry for E  $\rightarrow$   $(E, 3)$

$\rightarrow$  Because B thinks it can reach E by  $E$  not  $NXT$

Now B sends to  $\frac{A}{=}$  since E not  $NXT$

Hop for E =

$\rightarrow$  So A adds entry for  $(E, 4) \frac{B}{=}$

$\rightarrow$  Now " sends to  $C$  | NOT B E update E 5 +

$\rightarrow$  E sends  $(E, 5)$  to B, B makes  $(E, 6)$ , sends to A  $(E, 7)$  for entry of A ... till  $\infty$

### ADV. OF DV

very simple

PIS ADV. OF DV

creates count to  $\infty$  problems

### LINK STATE ROUTING (OSPF, ISIS) $\rightarrow$ DIJKSTRA

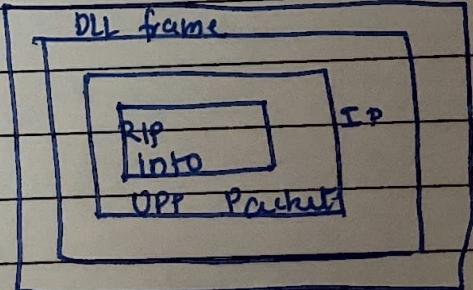
open shortest path first

RIP (sends distance vector as)

OSPF a UDP packet which has RIP info

PHY

This will be allowed to modify Routing tables



These modify local routing table

## INTRA-DOMAIN ROUTING

### D.V - BELLMAN-FORD

each node tells neighbours : distance to all others

+ simple

- count to  $\infty$
- convergence time high

we can use this for small network (10-20 nodes)

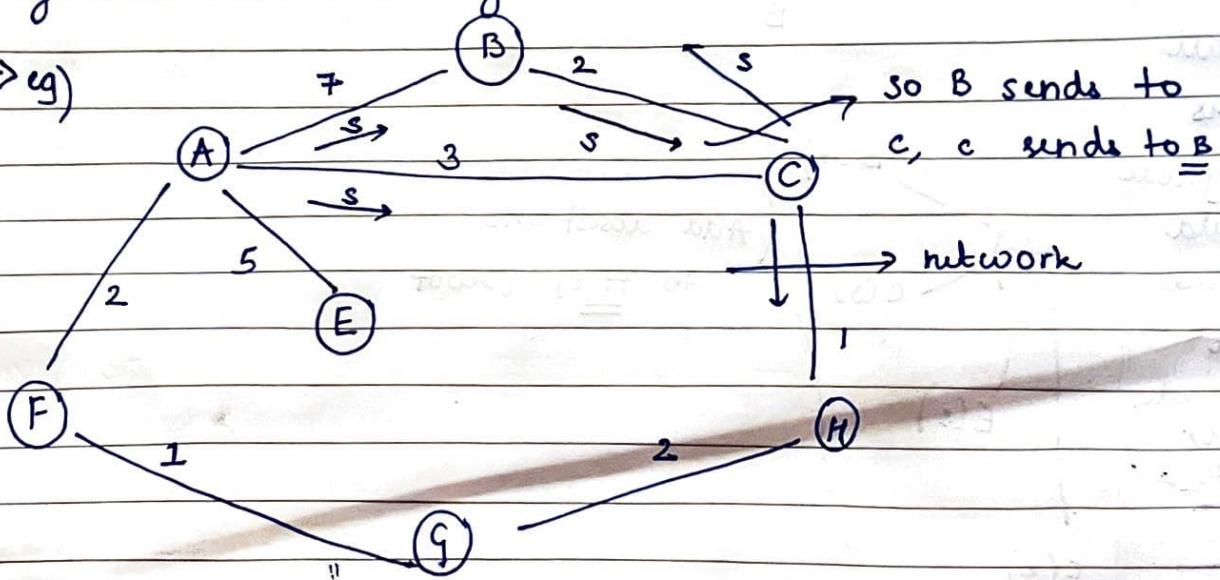
To avoid count to  $\infty$  we do RIP (where  $\infty = 16$ )

S.R: Link state routing  $\rightarrow$  djikstra :-

+ each node tells all others in the network :

My distance to neighbours  $\rightarrow$  BROADCAST

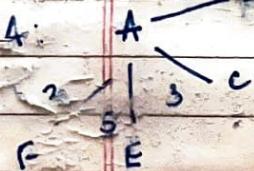
$\rightarrow$  eg)



A broadcast message to everyone that

+  $s \rightarrow$  A tells to everyone that is connected to everyone and forwarded throughout network, everyone gets the message  
Forward

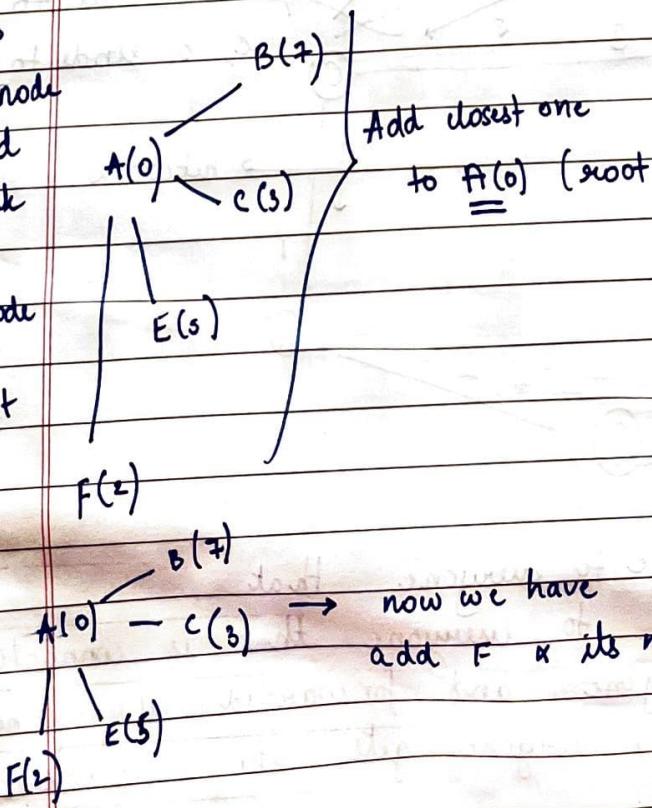
$A \rightarrow C \rightarrow D \rightarrow G$



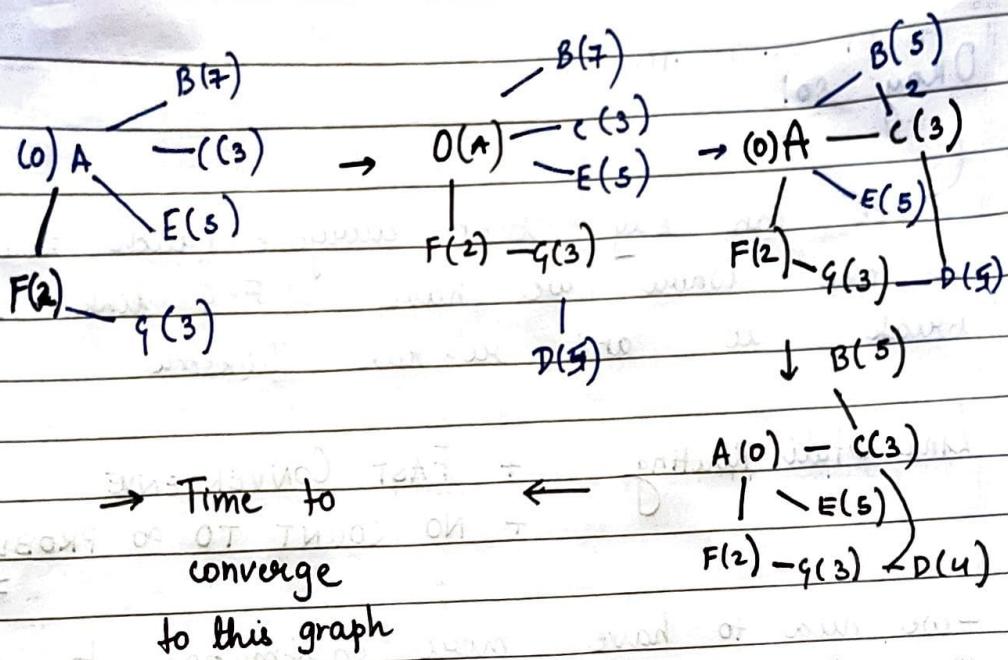
- So there are loops, we can one node receiving the message twice, so it should not forward it again more also
- We want to sent to all neighbours who have probably not got the message
- So  $\oplus$  of this everyone has information of the full graph, so they can run Dijkstra Day Run

\* shortest path tree from A to all others

Every node should the link that one node has for shortest



$A(0)$  → distance to itself = 0



Time for broadcast + Algorithm time complexity  
 $\hookrightarrow O(\text{diameter of network})$

| DEST | COST | NEXT |
|------|------|------|
| A    | 0    | -    |
| B    | 5    | C    |
| C    | 3    | C    |
| D    | 4    | C    |
| E    | 5    | E    |
| F    | 2    | F    |
| G    | 3    | F    |

This algorithm does not (not) run on one node as this keeps track of compound to distance vector

Suppose  $F \rightarrow G$  fail

$\rightarrow$  So F broadcast I cannot reach  $\underline{\underline{G}}$  } All re-run  
 $\rightarrow$  So G " I cannot reach  $\underline{\underline{F}}$  } Dijkstra

$\rightarrow$  Faster than D.V

May have loops during this time

Okay so

→ we can say that everyone (each router has a tree where we have F-G link) we break it and re-run - Dijkstra

Link State Routing + FAST CONVERGENCE

+ NO COUNT TO  $\infty$  PROBLEM

- we need to have more information to run Dijkstra

- need full topology info
- Full Broadcast for any change
- Dijkstra algo must be run

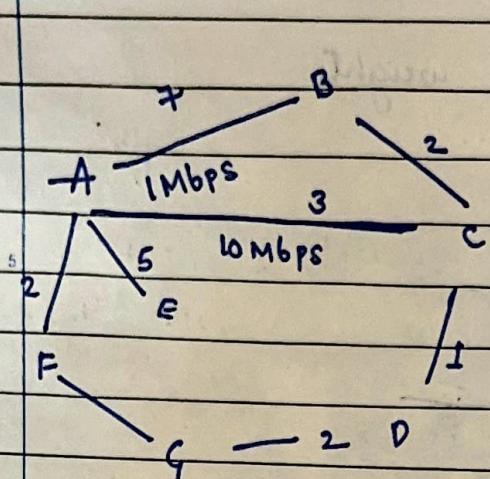
→ Count to  $\infty$  avoided as here we have full topology

→

Q: What should link weights be?

→ Links with extreme high weights would not be used!

→ So carefully choose that we can include!

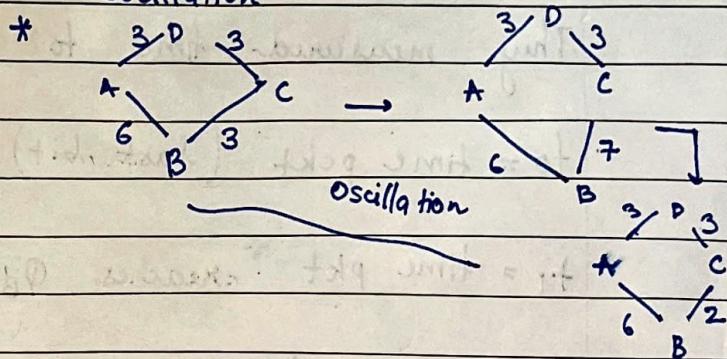


\* So what if A-C gets full then in that case actual speed  $< 10 \text{ Mbps}$

\* if we update weights it should broadcast that

\* then all re-run Dijksha

\* continuous changing diverts traffic \* these can cause oscillation

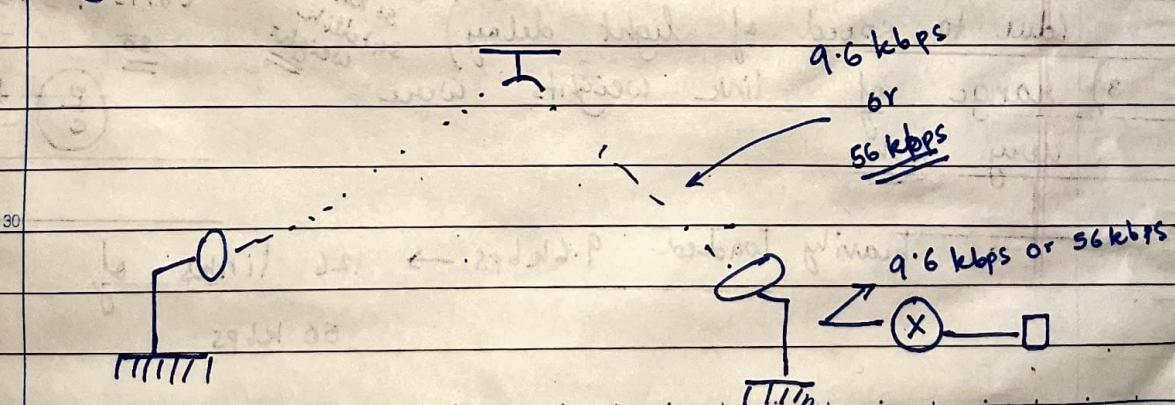


$\therefore$  RIP : default weight = 1 for every link (P.V)

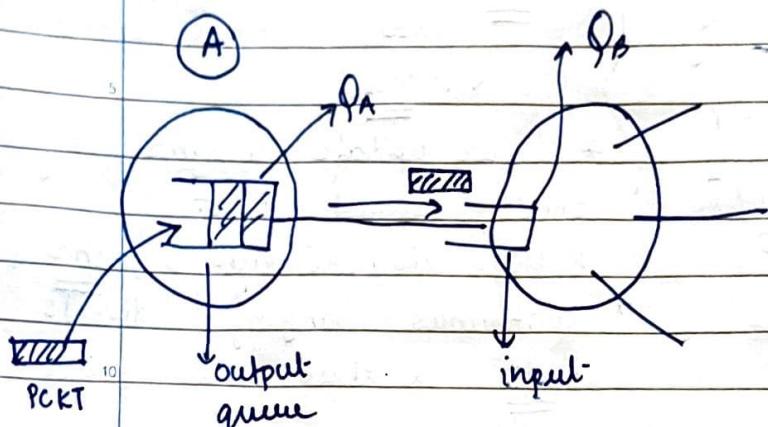
$$\text{OSPF (LSR)} : \text{weight} = \max \left( \frac{10^8}{\text{Link speed}}, 1 \right) \text{ in bits/sec}$$

Original Arpanet :- (1969 - )

They had low speed satellite link



So what should link weights



They measured time to reach other queue

$t_0$  = time pkt (last bit) enters  $\Phi_A(\pi)$

$t_1$  = time pkt reaches  $\Phi_B$  (last bit)

measure  $(t_1 - t_0)$  = Queuing delay + speed of light delay  
(at  $\Phi_A$ )

Avg these delays

over some time period

Set link to this avg of delay

Issues

} How they gave weights

+ transmission delay  
time empty

bits of pkt on wire c bits per sec

so  $\frac{p}{c} \text{ time}$

→ 1) OSCILLATIONS

2) satellite links were penalized a lot  
(due to speed of light delay)

9.6 kbps theoretical link had less weight than 56 kbps satellite weight

3) range of link weights were very high

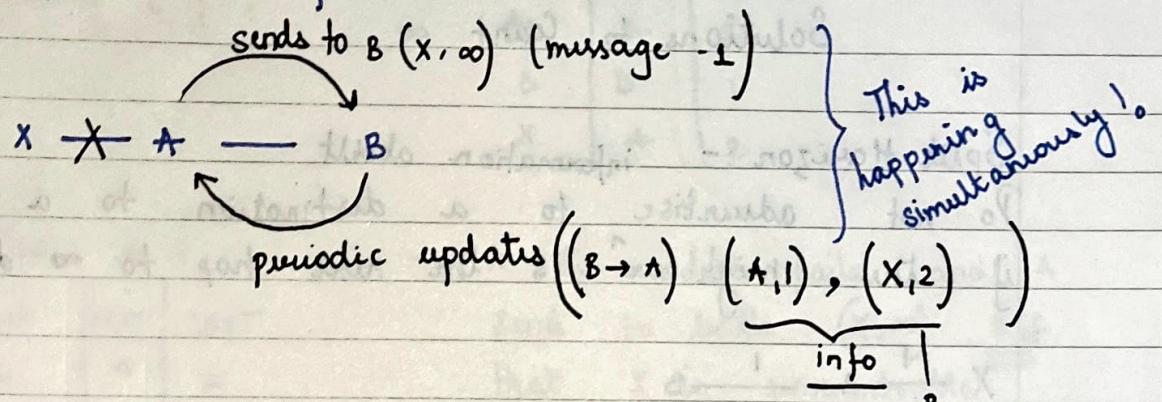
↳ 1 heavily loaded 9.6 kbps → 126 links of 56 kbps

Count to infinity problem :-

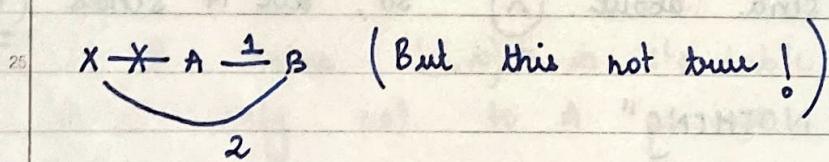
D-V  
Dest | Routing Table  
Next | cost

|      |      |     | $\frac{1}{A}$ | $\frac{1}{B}$ |      |      |     |
|------|------|-----|---------------|---------------|------|------|-----|
| Dest | Next | Hop |               |               | Dest | Next | Hop |
| A    | *    | 1   |               |               | X    | X    | 1   |
| B    | *    | 2   |               |               | B    | B    | 1   |
|      |      |     | at X          |               |      |      |     |
|      |      |     |               |               | X    | A    | 2   |
|      |      |     |               |               | A    | A    | 1   |

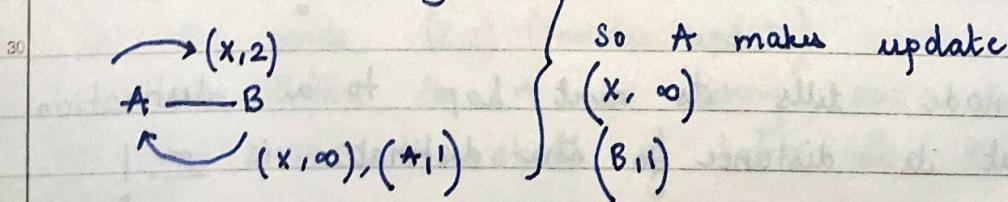
Suppose  $X \rightarrow A$  fails



B tells my distance to X is 2 then assume there is a path from  $\underbrace{B \rightarrow X}$



B changes its table and makes  $(X, \infty)$ . Now again both send messages



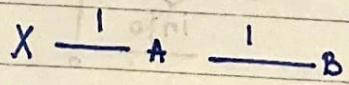
$\rightarrow$  B makes  
 $X \rightarrow A - B \rightarrow$  update  $X, 3$  val  
 $\rightarrow$  B thinks there is a path from A to  
X with weight 3  
then make an entry

|       |       |
|-------|-------|
| Table |       |
| X     | A : 4 |

- This keeps happening till  $= \infty$
- This is called count to  $\infty$
- So if A had packet for X it would send it to B, B would send it back A.
- So RIP allows a max distance of 16  
 $\text{cost} = 16 = \infty$  implies we cannot reach it

### Solutions to Count $\infty$

Split Horizon :- information about  
Do not advertise to a destination to a neighbour  
if the neighbour is the next hop to  $\infty$  destination



Now  $X - A$  fails

Now A adv  $(X, \infty)$ , simultaneously B sends  $(X, \infty)$  or does not send about  $(X)$  so, all it sends  $(A, 1)$

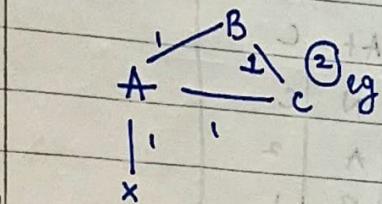
B TELLS "NOTHING"  
ABOUT X

Split Horizon with poison reverse!

A node tells its next hop to a destination that its distance to the destination is  $\infty$ !

So "B sends Routing ADV  $(x, \infty)$  to A"

Another Example



|      |     |           | At B |     |      |
|------|-----|-----------|------|-----|------|
| Dest | Nxt | Cost      | Dest | Nxt | Cost |
| A    | A   | 1         | A    | A   | 1    |
| C    | C   | 1         | C    | C   | 1    |
| X    | X   | 2         | X    | A   | 2    |
| At A |     |           | At B |     |      |
| Dest | Nxt | HOP(cost) | Dest | Nxt | Cost |
| X    | X   | 1         | A    | A   | 1    |
| B    | B   | 1         | B    | B   | 1    |
| C    | C   | 1         | X    | A   | 2    |

|      |     |      | At C |     |      |
|------|-----|------|------|-----|------|
| Dest | Nxt | Cost | Dest | Nxt | Cost |
| A    | A   | 1    | A    | A   | 1    |
| B    | B   | 1    | B    | B   | 1    |
| X    | A   | 2    | X    | A   | 2    |

| At X |     |      | At A |     |                                 |
|------|-----|------|------|-----|---------------------------------|
| Dest | Nxt | Cost | Dest | Nxt | Cost                            |
| A    | A   | 1    | A    | X   | link fails and A                |
| B    | A   | 2    | B    | X   | send to both $(X, \infty)$ B, C |
| C    | A   | 2    | C    | X   | that X is unreachable           |

→ But the message to is C  
is lost!

→ So B makes  $(X, \infty)$  in its table and sends this to C only not to A as we are using split horizon:-

→ But C can tell about its distance X to B as A is next hop not B.

→ So C sends  $(X, 2)$  (periodic update)

→ Does not tell this to A as using split horizon.  
So Now current state of table at B

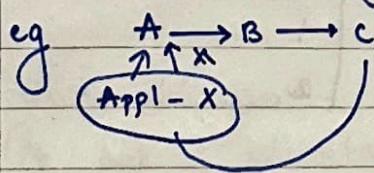
|   |   |   |
|---|---|---|
| D | N | C |
| A | A | 1 |
| X | C | S |

So now B thinks there is path from c of dist 3 to X

- Now B sends  $(X, 3), (C, 1)$  to A as A<sub>5</sub>  
no more next hop for X.
- Now finally table Status

| A + A |   |   | A + B |   |   | A + C |   |   |
|-------|---|---|-------|---|---|-------|---|---|
| D     | N | C | D     | N | C | D     | N | C |
| X     | B | 4 | X     | C | 3 | X     | A | 2 |
| 10 B  | B | 1 | B     | B | 0 | B     | B | 1 |
| c     | c | 1 | *     | A | 1 | *     | * | 1 |
|       |   |   | c     | c | , |       |   |   |

Now suppose any of them gets a message for X it will go in a loop for



So after four round trip travel = 16 ∵ not reachable

→ Do RIP stops it.

RIP: Routing information protocol :- (D.V based)

Cost of all links = 1

max cost allowed to a destination is 16 ( $16 = \infty$ )

Advantages of DV

→ Simple & easy to implement (As no complicated algo, no BFS, ...)

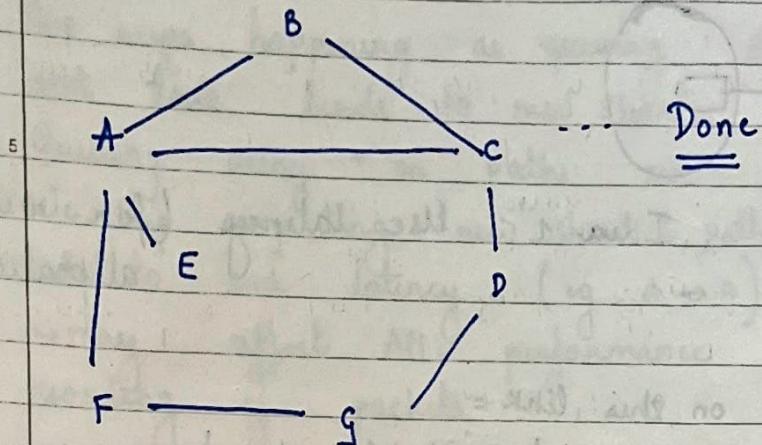
→ Disadvantages

- count to  $\infty$  & routing loops :-

- converge of routing table

→ So due this we do LSR

# Link State Routing!



**ADV:**

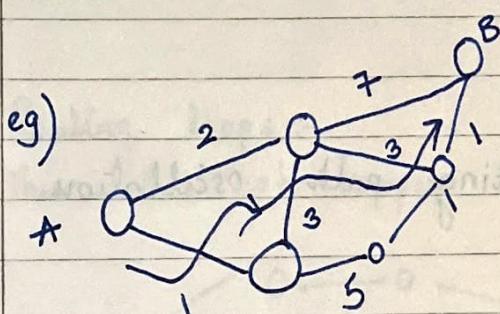
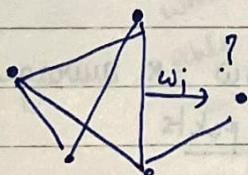
No routing loops, no count to  $\infty$  problems, convergence is faster of routing tables.

→ Disadv: Algo is more complex than D.V

How to choose LINK WEIGHTS

DV, LSR  
↓  
RIP

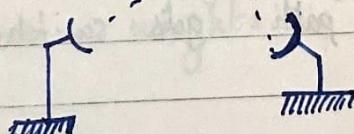
↳ OSPF, ISIS

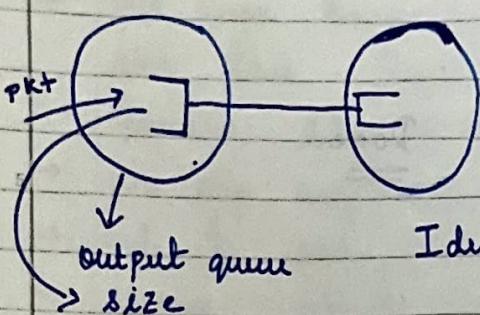


so how most packets will be sent through

→ Axipanit :- two types of links 9.6 kbps, 56 kbps  
U.S.A

Terrrestrial links, Satellite links





Idea 1 : Use latency (for weights allocation)

latency of pk on this link =

queue delaying + speed of light delay + transmission delay  
(when 1 first comes)  $\uparrow$  PLC

→ this same as time to put bits of wire

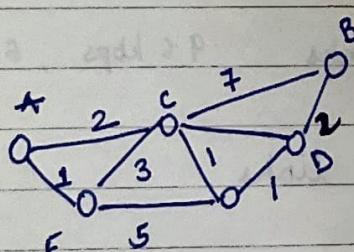
latency of packets keep changing! (due to queuing delay & transmission delay).

Take time window  $\times$  average out latency.

= latency of all pkts  
time window

#### ISSUES :-

① under heavy load, routing path oscillation



→ We are using path  
A - E - C - D - B

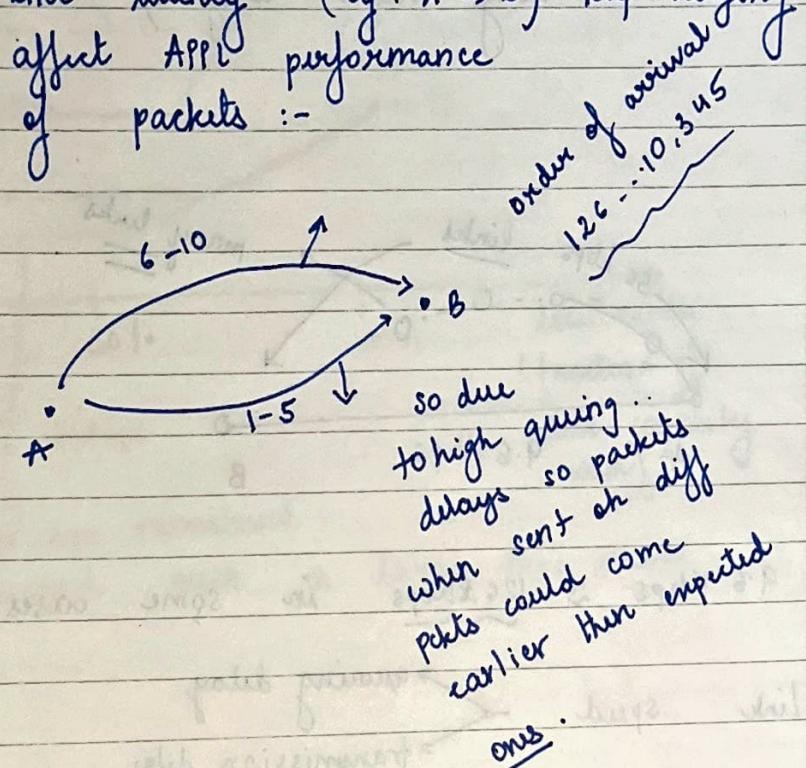
But since our load is high  
lot of traffic on this, we increase weight, so path gets switched to  
A - C - B.

- Then again since traffic shifts again we take another path ...
- This keeps happening as queuing delays keep changing and this leads to new shortest paths.
- Queuing delay ↑ on paths use
- We keep computing new shortest paths
- end to end latency (eg. A → B) keep varying
- May affect App performance
- reordering of packets :-

10

15

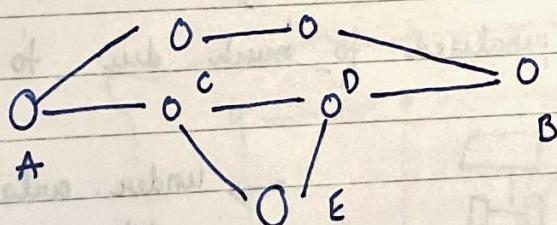
20



### Routing loops :-

- This does not happen when weights are fixed.

25

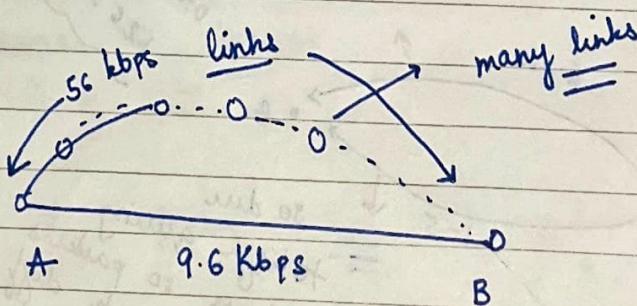


A thinks A-C-D-B is shortest path, C thinks it is even though what thinks, C does not have updated link weights.

A-C-E-D-B

→ So C sends it to E, E sends it back to some other node or back to C

(2) Range of link weights were very large because of this some links were penalized (due to very high weights).



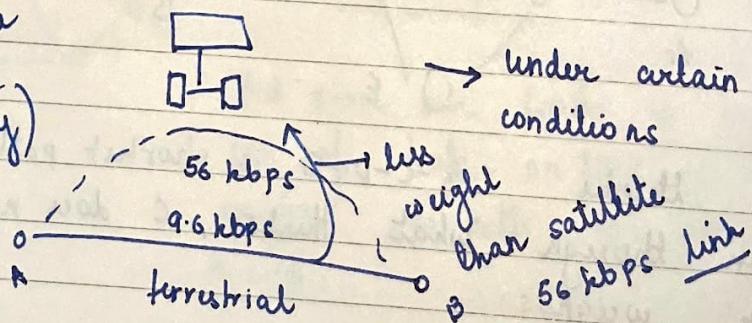
$9.6 \text{ Kbps} \approx \underline{126.5 \text{ Kbps}}$  in some cases!

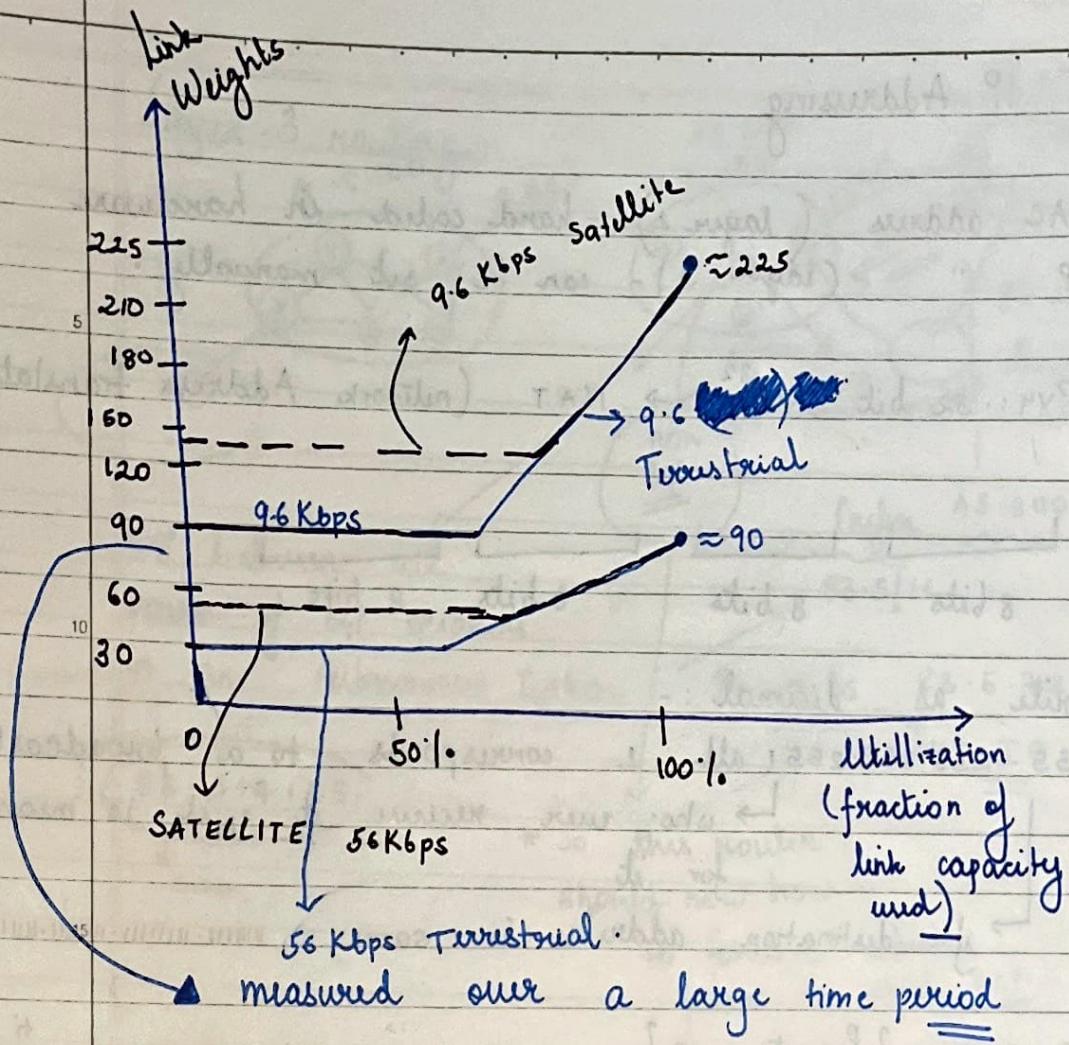
Link speed → queuing delay  
transmission delay  
transmission delay also affects queuing delay.

Is this logical?

3) Satellite links penalized to much due to transmission delay.

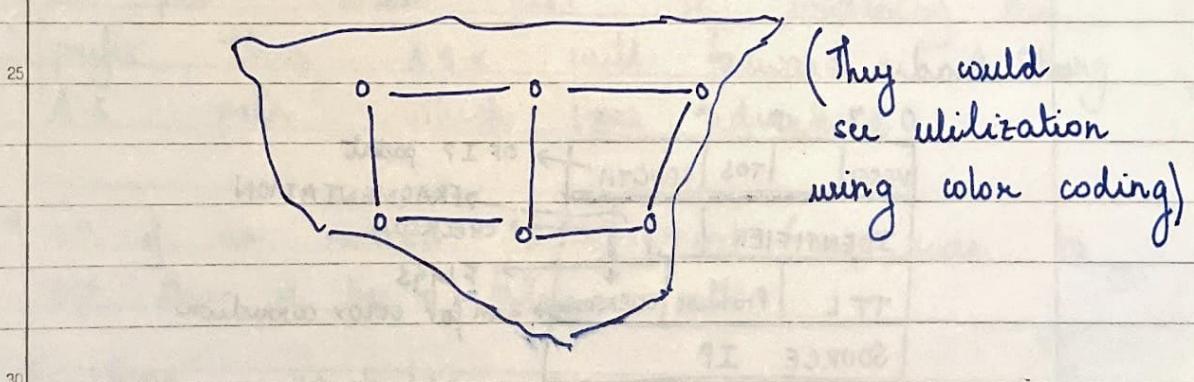
(No this due to speed of light delay)





OSPF :- wt of link =  $\max\left(\frac{108}{\text{link speed (bps)}}, 1\right)$

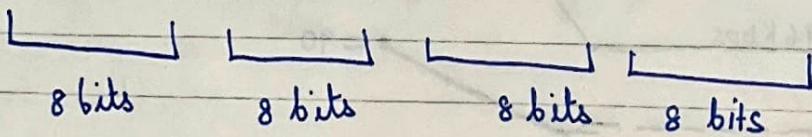
NOC : network operation centre (AT&T)



## IP Addressing

MAC address (layer 2) - hard coded in hardware  
 IP " (layer 3) - can be set manually.

IPv4: 32 bit  $\sim 2^{32} \rightarrow$  NAT (Network Address translation)



write as decimal :-

Ex 255.255.255.255 (all 1 corresponds to a broadcast)

↳ who ever receive it, it is meant for it

if destination address is same 1111.1111.1111.1111

10. \*.\*.\* { Private IP }  
 192. \*.\*.\* { Reusable } can be used in their own organizations

public IP → should be unique  
 in the internet

↳ only one host must use that IP address.

## IP Header

| 0          | 4        | 8        | 16                       | 31 |  |
|------------|----------|----------|--------------------------|----|--|
| VERS       | TOS      | LENGTH   | OF IP packet             |    |  |
| IDENTIFIER |          |          | FRAGMENTATION            |    |  |
| TTL        | Protocol | CHECKSUM | CHECKSUM                 |    |  |
| SOURCE IP  |          |          | FLAGS                    |    |  |
| DEST IP    |          |          | sum for error correction |    |  |
| DATA       |          |          | IP options / optional    |    |  |

Date \_\_\_\_\_

25.7.88



AS 300

### Layer-3 routing

A.S.100

eBGP

AS 200

eBGP

83.5/16

Belongs to  
AS 300

Prefix AS-300 BGP  
83.5/16 attributes

iBGP between all pairs of BGP speakers in an Autonomous System

R<sub>1</sub> gets 83.5.7.25

Dest IP

83.5.7.25

given

\* So this router should now know to forward which packet to which BGP

AS200 → AS100 : 83.5/16 AS200 - AS300

→ Golden rule : If an A.S. sends an advertisement for <prefix> x-y-z... to a neighbour, and the neighbour (sends a packet whose A.S. is a packet) sends this BGP speaker a pkt with dest IP matching the prefix, then A.S.x will forward it along A.S. path which was advertised.

→ So if we receive advertisement from node to BGP then it has to follow a path

→ Exceptions 83.5/16 AS200 - AS300

AS200  $\xrightarrow{\text{adv}}$  AS100 shorter prefix

what if

AS-200 → AS-100

longer prefix

83.5/16 AS200 - AS600 - AS300

Given

83.5.7.32

Date \_\_\_\_\_ / \_\_\_\_\_ / \_\_\_\_\_



→ Now we have two matches

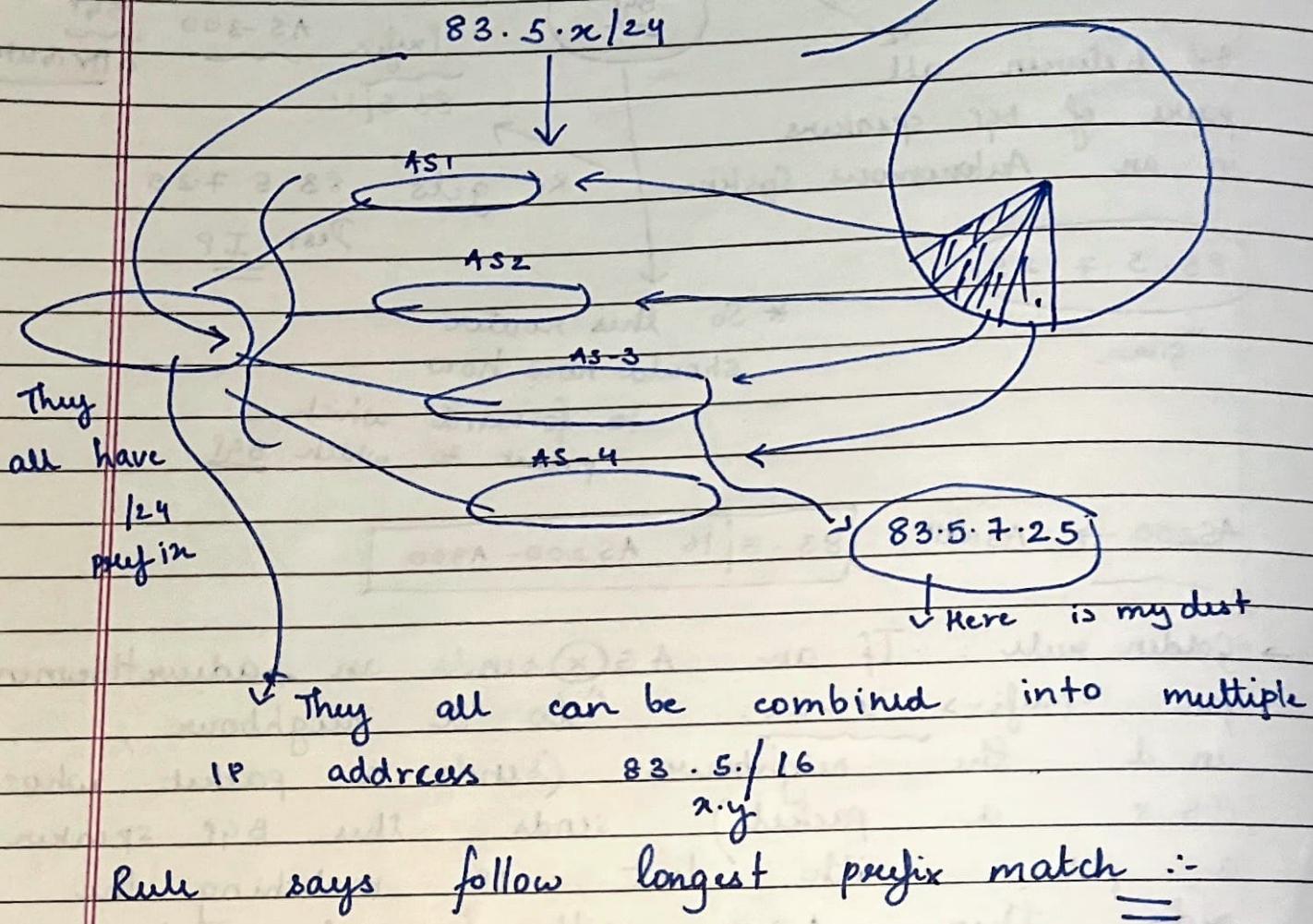
83.5.16 AS-200 - AS-300

83.5.7.24 ~ ~

SUPER  
NETTING

(COMBINING  
IP address)

So which one should we choose!



So AS-600 path will be used!

At destination we have the highest match

So if we try to take a one with shorter prefix match we can have a shorter path but we may flow away from the actual destination:-

Date \_\_\_\_\_

One entry for one prefix in routing table.

BGP- attributes

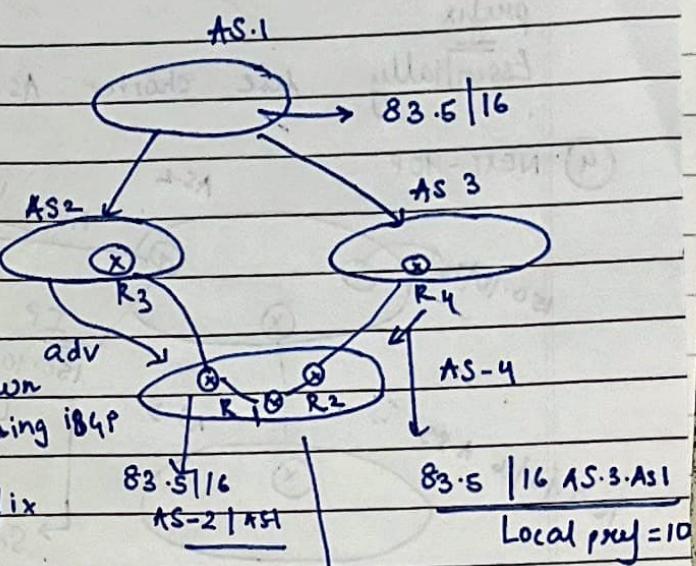
① Local-Pref

↓  
added locally

to advertisements

heard via eBGP.

This is then sent to other BGP routers in its own AS using iBGP.  
This makes sure the one with larger prefix is preferred:-



→ we do not forward local pref attributes. So any advertisement via eBGP we do not put our local pref, the receiver put his local pref

Local pref = 20

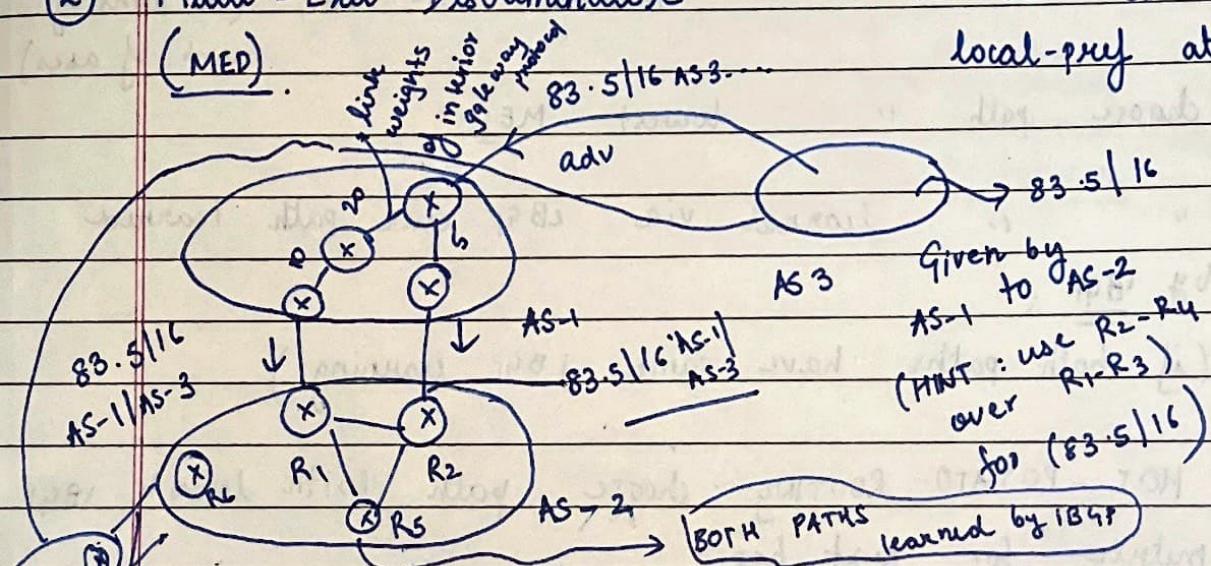
So now AS-4 can decide which path it wants to choose

\*  $R_1, R_2$  etc.

② Multi-Exit Discriminator (MED)

in AS-4 can add

local-pref attributes



So AS-1 can tell AS-2 that  $MED = 5$  at  $R_2 = 5$  and MED is 30 at  $R_1$

\* Rule lower MED values are

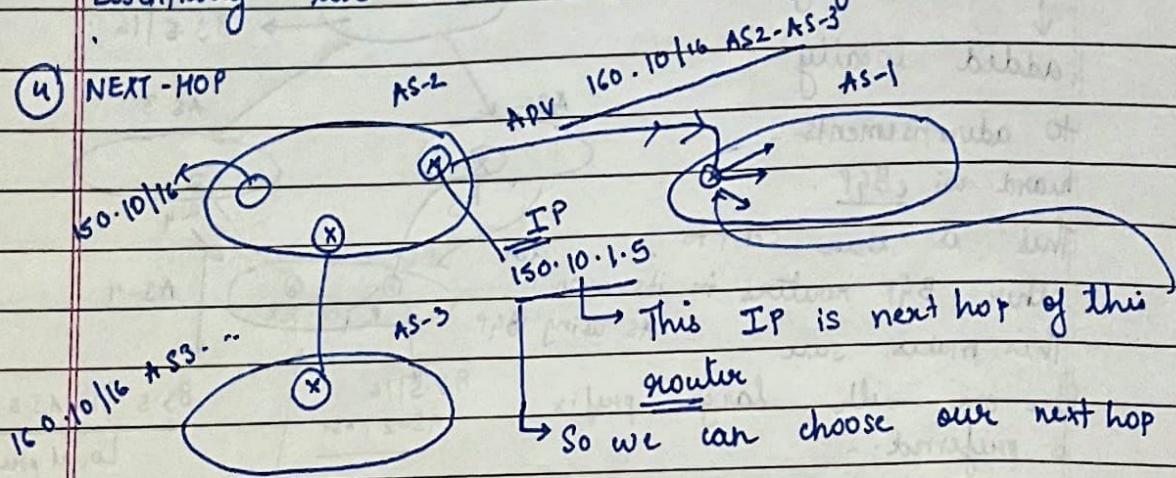
preferred (for same prefix)

→ It may be logical to compare MED of only 1 Autonomous system  
Because others should/could be different  
Date \_\_\_\_\_



(3) AS-PATH : list of ASes to destination as which has the prefix

Essentially use shorter AS-paths is preferred



Rules for choosing paths (used at each BGP speaker)  
(for each unique prefix)

- choose route with largest local pref
- " " " shortest AS-Path (in terms of num of ases)
- choose path w/ lowest MED

(a) " " learned via eBGP over path Learned by iBGP:

(if both paths have same iBGP learning)

(c) HOT-POTATO-ROUTING : choose path with lowest iBGP metric for next hop:-

(d) choose path whose exit router has the lowest router ID (in same AS) ~~Highest~~  
→ (= highest IP address on all interfaces) of router

IPv4 Header

bit

0      4      8      16      31

| VERSION | TYPE OF SERVICE       | LENGTH |
|---------|-----------------------|--------|
| 4       | (BEST-EFFORT SERVICE) | 80     |

| TTL | time to live | Checksum |
|-----|--------------|----------|
| 64  |              |          |

32 bits per row

This row skipped

$$\text{TTL} = \text{TLL} - 1$$

at every L3 router we decrement TTL & when TTL=0 =>  
 discard don't forward

So for TCP: 6, UDP: 17, ICMP: 1

This is used to check for errors

|                |                |
|----------------|----------------|
| SOURCE IP      | ADDR (32 BITS) |
| DESTINATION .. | ..             |

PAYLOAD

\* we could have TCP header / UDP header

IP packet

BGP forwarding rules :-

run at every BGP speaker (for each prefix)

1) LOCAL-PREF  $\rightarrow$  larger

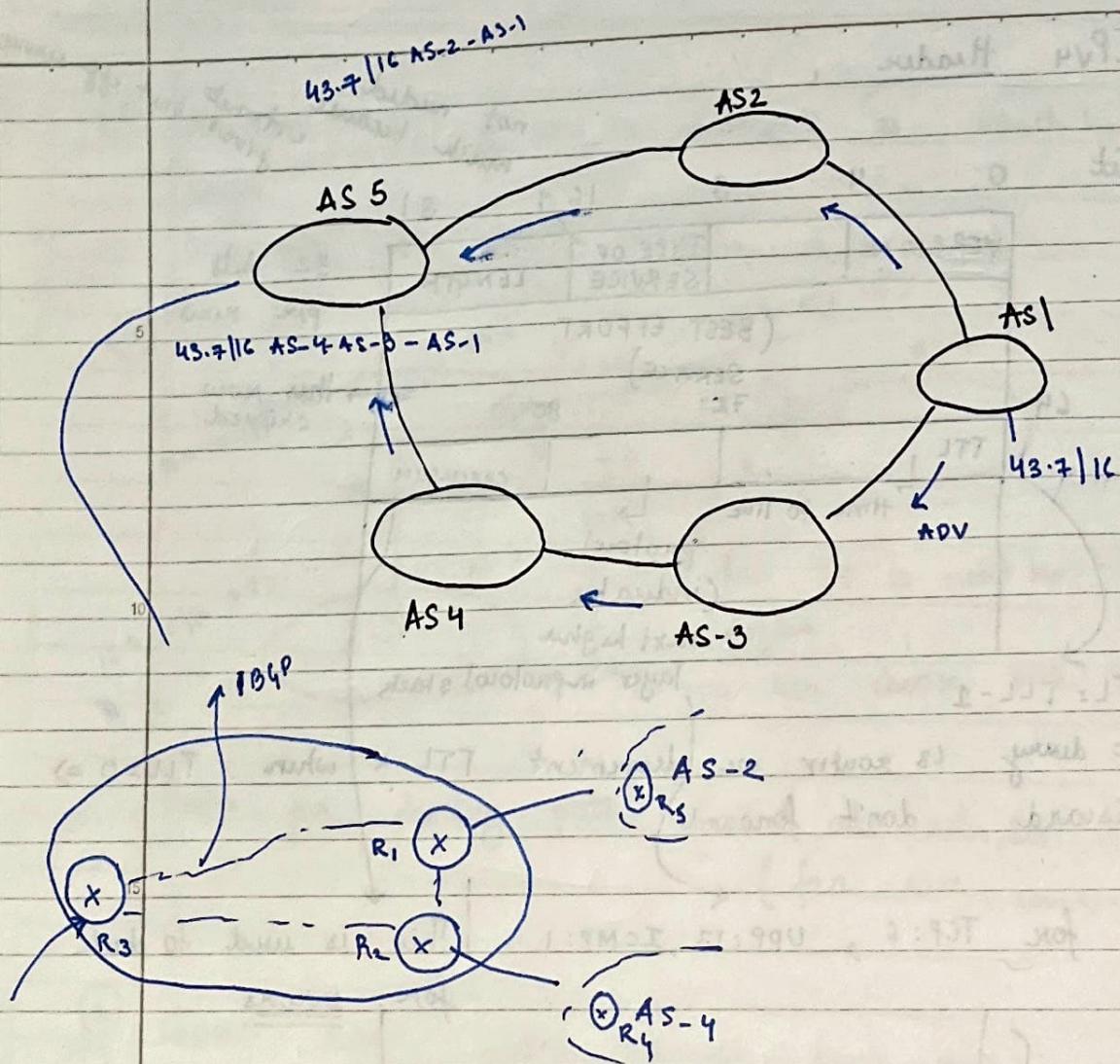
2) AS-PATH length - shorter

3) MED  $\rightarrow$  smaller

4) eBGP - over - iBGP

5) HOT-POTATO Router (shortest IGP to next hop)

6) lowest router ID :-



If the AS-5 decides that I want all traffic for IP 43.7/16 to go via  $R_2 \rightarrow AS-4$

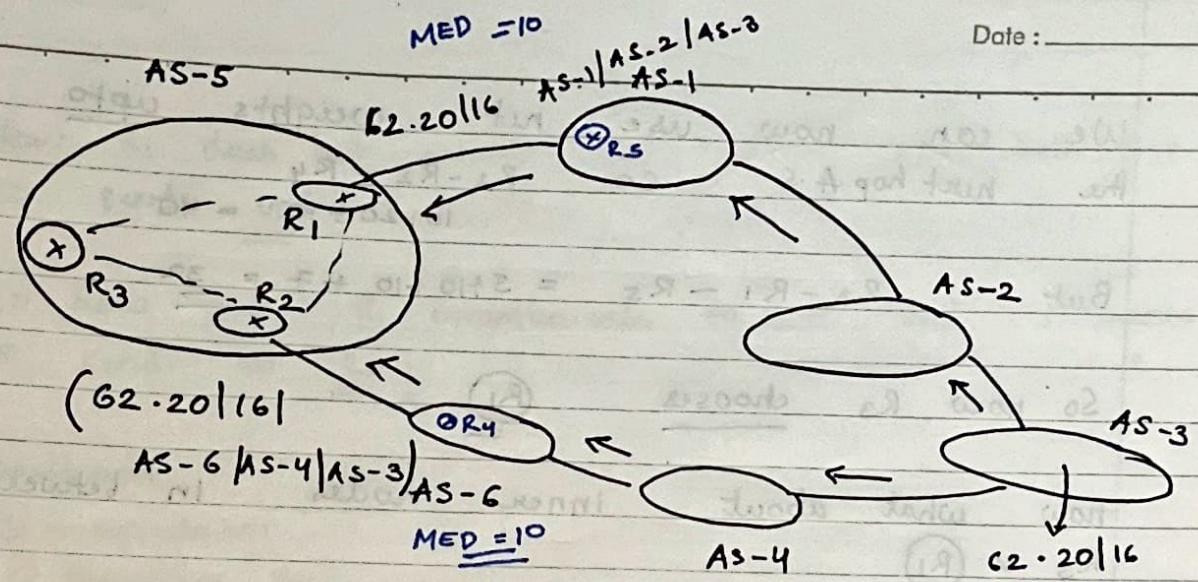
So what attributes we should use

① If no local pref given then will take path from  $R_1$  as shortest AS-path.

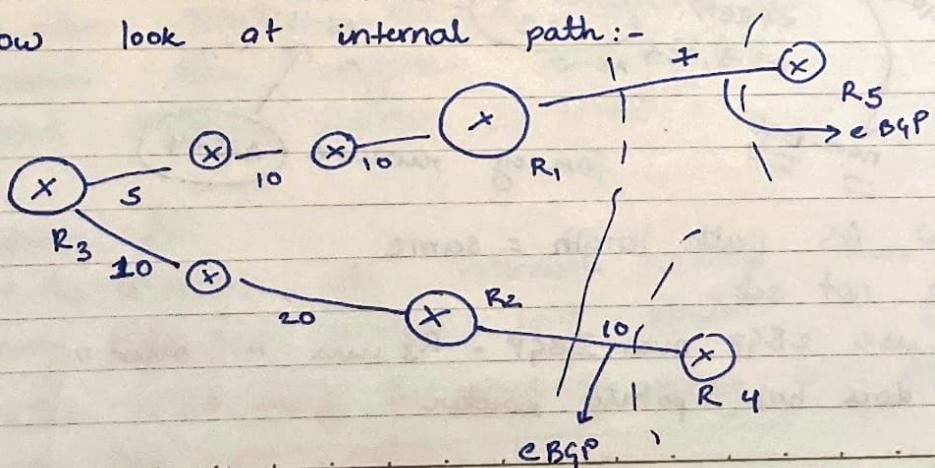
② So we should make  $R_2$  forward an ADV over BGP with higher local pref than  $R_1$  so traffic goes via  $R_2$ .

$\rightarrow R_2$  finds 43.7/16 AS-4 - AS-3 - AS-1 LOCAL-PREF = 10

$R_1$  " " AS-5 - AS-1 " " = 5



- NO LOCAL PREF SET
  - SAME LENGTH OF AS-PATH
  - Assuming we can compare MED across A.S we look at MED
  - 15 IF MED SAME
    - We can see R<sub>1</sub> learned about path (R<sub>2</sub>) over IBGP \* and it has its path learned over EBGP
    - 20 So R<sub>1</sub> got packet some how so would not send to R<sub>2</sub>, directly send to R<sub>5</sub>
    - Now what about (R<sub>3</sub>) both paths for this IP learned over IBGP..
      - ↳ Now look at internal path:-

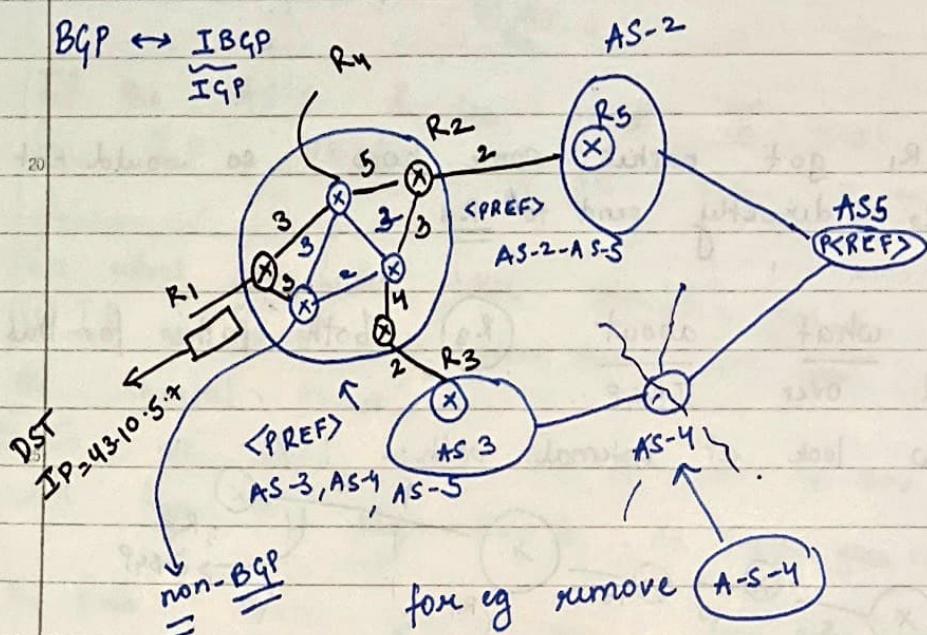
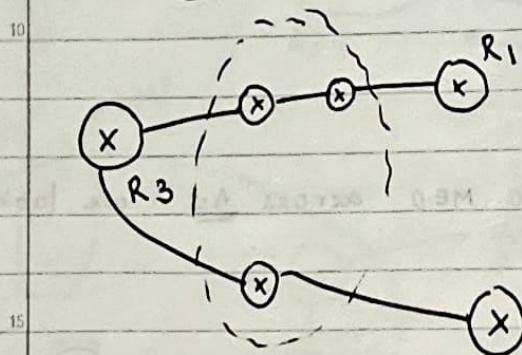


We can now use net weights upto the next hop A.S so  $R_3 - R_2 - R_4$   
 $10 + 20 + 10 = 40$

But  $R_3 - R_1 - R_5 = 5 + 10 + 10 + 7 = \underline{\underline{32}}$

So now  $R_3$  chooses  $\underline{\underline{R_1}}$

now what about inner nodes in between  
 $R_3 \times \underline{\underline{R_1}}$



- now AS path length = same
- MED not set
- R2 uses eBGP over IBGP & R3 uses " over it
- R1 does hot potato routing

→ Now  $R_1$  does hot potato as neither eBGP both iBGP so sends via =

So how does  $R_1$  communicate to  $R_4$  that it wants to send to  $R_2$ !

### Methods

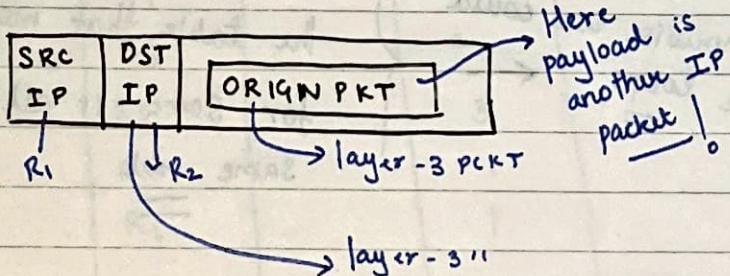
① encapsulation

② Pervasive BGP

③ Tagged IGP

① encapsulat<sup>n</sup> (Suppose  $R_1$  gets a packet for AS-5)  
( $\langle \text{PREFIX} \rangle \text{ IP-AS-5}$ )

So it encapsulates this by putting extra bits & making  
DST IP -  $R_2$



encapsulation of layer-3 into layer-3

→ So now  $R_1$  sends to  $R_4$  as runs Dijkstra

Now when reaches  $R_2$  it removes this outerlayer & forwards the original pckt to  $R_5$

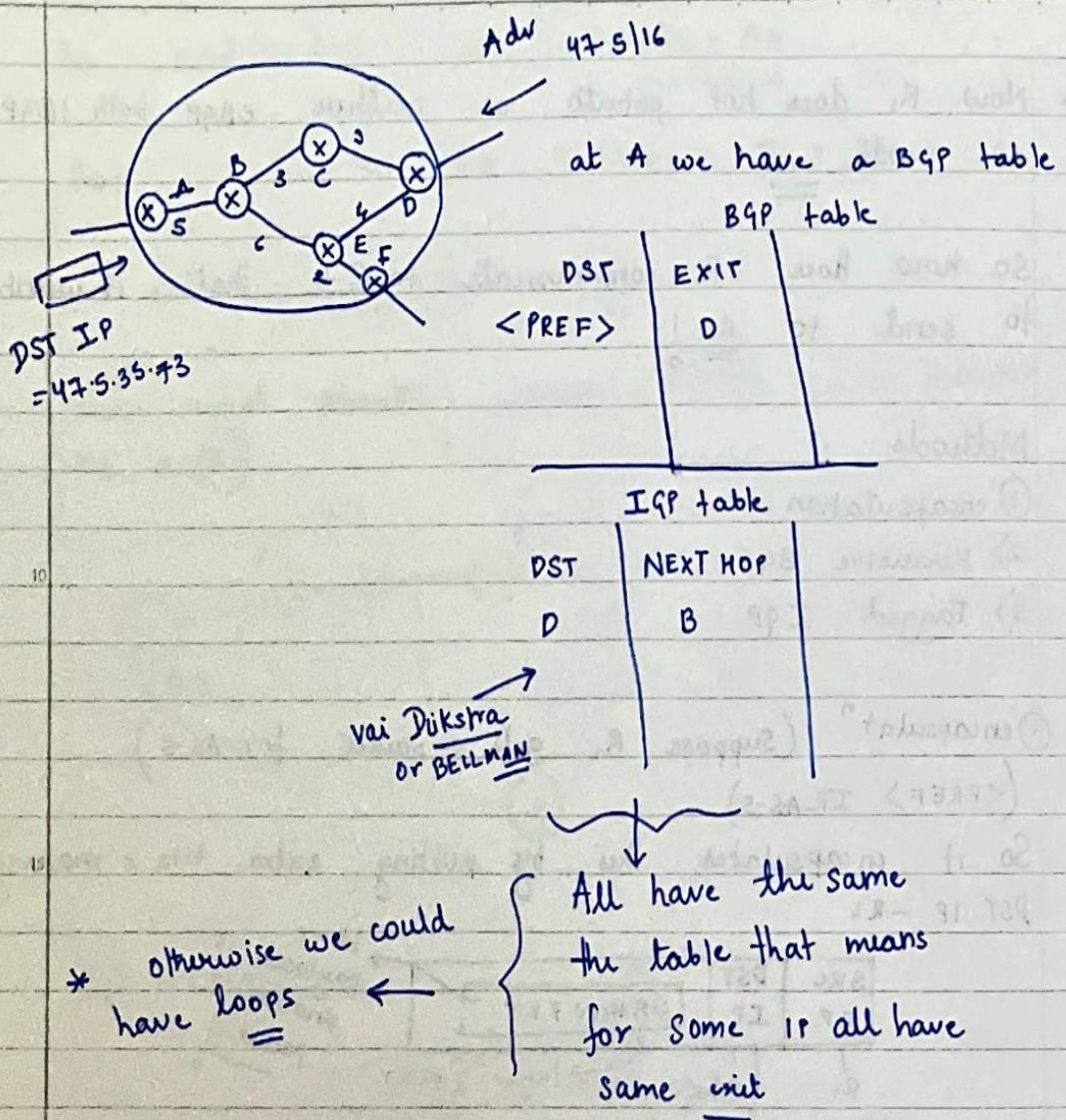
→ Payload itself is another IP packet:-

② Pervasive BGP

↪ Assumption : all internal routers run BGP

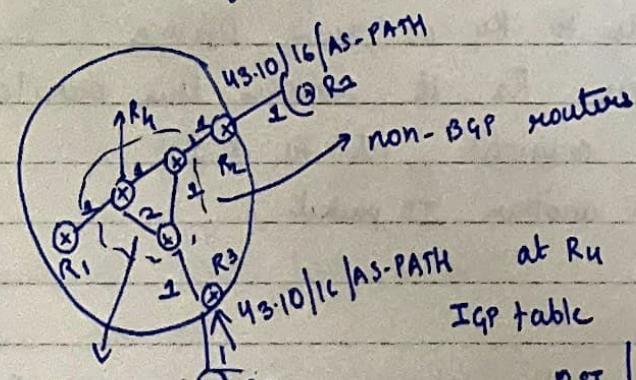
2) Every prefix has unique exit router in AS

↪ This could easily be solved by setting local pref for each node for each <prefix>



### → TAGGED IGP

BGP speaker (router) inserts some tagged information into IGP about prefixes learned via BGP



They are useless about these prefixes

| DST | NEXT | TAG |
|-----|------|-----|
| R1  | R1   |     |
| R2  | R5   |     |
| R3  | R6   |     |
| R5  | R5   |     |
| R6  | R6   |     |

R<sub>2</sub> can send information ~~43.10|16~~ TAG = R<sub>2</sub>

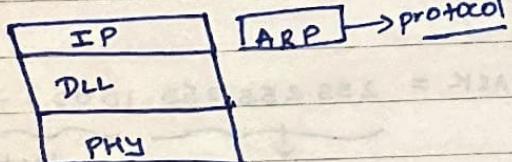
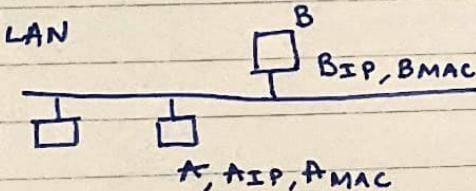
And suppose network (AS) runs OSPF (LSR-Dijkstra)

So R<sub>2</sub> puts 43.10|16 (R<sub>2</sub>) → TAG so in table for prefix we get a entry so now we can treat this as a node in table & create an Dijkstra table

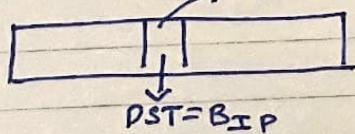
Now table at R<sub>4</sub>

| DST                                   | NEXT                             | TAG                        | COST |
|---------------------------------------|----------------------------------|----------------------------|------|
| "                                     | "                                | -                          | 7    |
| 43.10 16                              | R <sub>7</sub>                   | R <sub>7</sub>             |      |
| 43.10 16                              | R <sub>8</sub>                   | R <sub>8</sub>             |      |
| <u>tagged info</u>                    |                                  | → no lost for <u>these</u> |      |
| R <sub>1</sub>                        | R <sub>1</sub>                   |                            | 1    |
| R <sub>2</sub>                        | R <sub>5</sub>                   |                            | 2    |
| R <sub>3</sub>                        | R <sub>C</sub>                   |                            | 3    |
| R <sub>5</sub>                        | R <sub>5</sub>                   |                            | 1    |
| R <sub>C</sub>                        | R <sub>C</sub>                   |                            | 1    |
| R <sub>7</sub><br>R <sub>8</sub>      | R <sub>5</sub><br>R <sub>6</sub> |                            | 3    |
| So now we have match for 2 tags so we |                                  |                            |      |

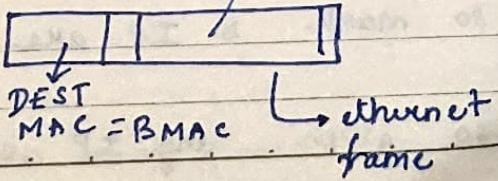
L3-L2 interaction



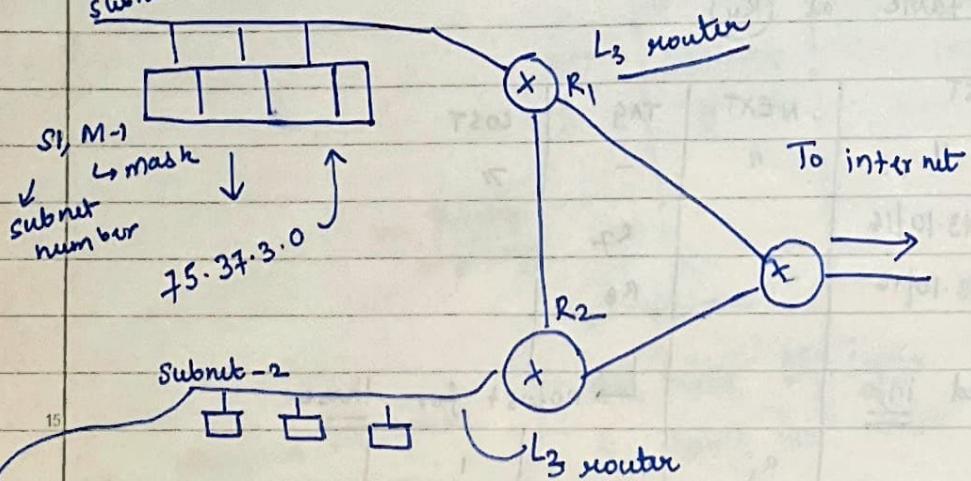
IP packet



IP packet



- If we do not have MAC address of B what to put  
DEST MAC - ADDR
- We use ARP protocol at layer -2 × make broadcast & ask who has BIP reply •
- So before this SUBNETTING  
NETWORK  $75.37.3 | 24$  prefix



Last 8 bits for all can be set by setting last 8 bits at L3 router :-

- So for all on S-1 put 0000 0000 → last 8 bits
- " " " S-2 .. 10 - - - → " "

Subnet num, Subnet MASK → identifies which bits of IP addr to consider

$$\text{MASK} = 255.255.255.1000 \dots$$

all 1's      last 8 bits

for this S-2  $75.37.3.128$

as 100 - 0 → 128

so mask & IP are stored.

so given any IP, say x.

Date : \_\_\_\_\_

If  $(X \& M_1) == S_1$  then  $X$  is in subnet-1  
"  $(X \& M_1) == S_2$  .. " ...  $\equiv$  -2

5 Table at Subnet Num  
R<sub>1</sub>

|                | MASK           | NEXT           |
|----------------|----------------|----------------|
| S <sub>1</sub> | M <sub>1</sub> | -              |
| S <sub>2</sub> | M <sub>1</sub> | R <sub>2</sub> |
| 0 other        |                | R <sub>3</sub> |

Protocol in IP header is next protocol :-

6: TCP, 17: UDP, 1: ICMP } TCP is above IP

TCP: time to live :-

documented at each IP router :-

0  $\Rightarrow$  drop packet :

Routing Table :-

Destination

Next hop

10

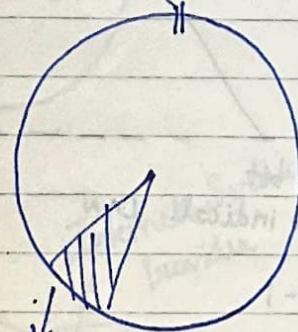
$\rightarrow$  small in number

IPv4  $\approx 2^{32}$  we need to give some kind  
on encoded IP address

15

a.b.c.d  
 $\hookrightarrow$  octet

20



this slice  
given to 111 B

25

730.52.30.\*  $\rightarrow$  address space

$\hookrightarrow$  common prefix

$2^8$  choices

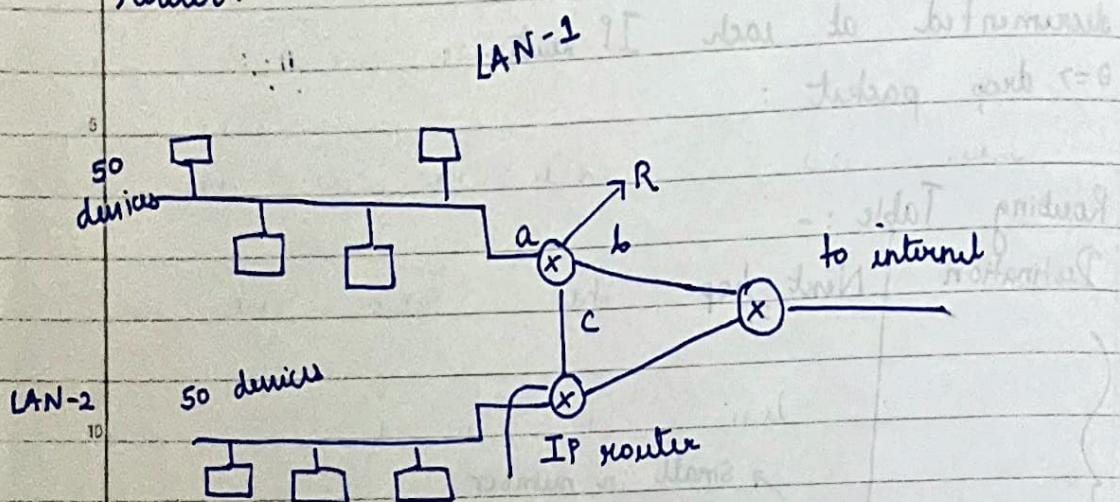
class A : 8 bits      24 bits  $\rightarrow$   $2^8$  hosts  
 Network address      hosts  
 (common prefix)

" B : 16 NW      16 bits host

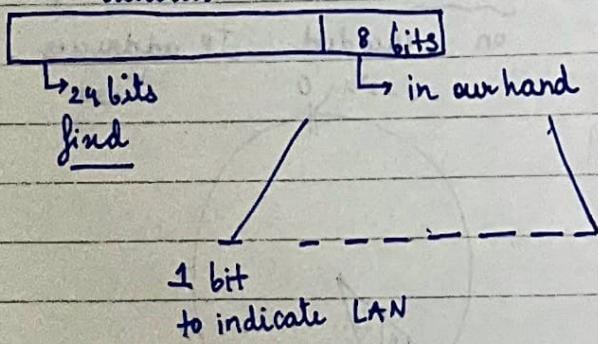
" C : 8+16 (NW)      8 " "

Subnetting: Given a slice of IP addresses how to divide among LAN's, setup/config, internal router

router:-



class C: address



0 - means LAN-1

1 - " " - 2

Subnet mask: Says which bits in IP address to use  
~~~~~ to decide which LAN to route

(eg) 25th bit is subnet mask:-

so how to choose mask

EX 11 - 11 00 - 0
25 bits 7 bits

Subnet address

S₁ for LAN-1

S₂ " " - 2

Router

Suppose R gets a packet and dest IP is 'D'

$$\text{Is } (\underbrace{D \text{ and } M_1}_{\text{bitwise}}) == S_1 ? \quad [S_1 = 730 \cdot 52 \cdot 30 \cdot 0]$$

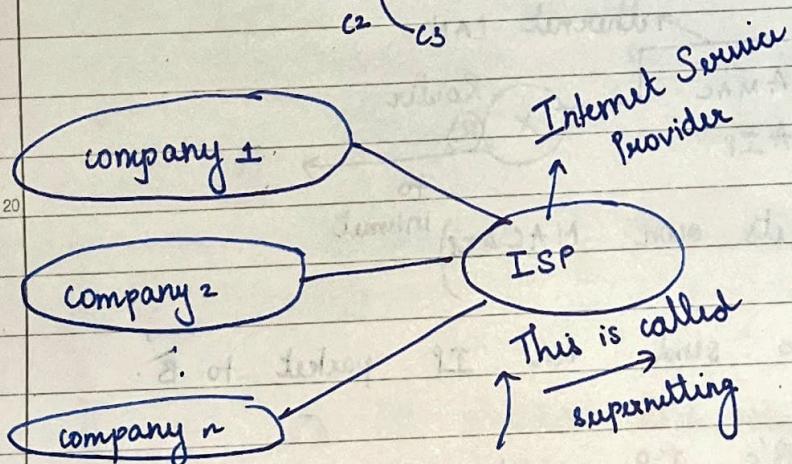
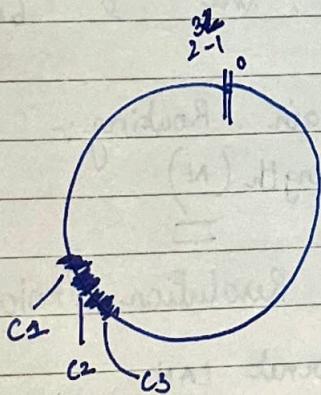
$$[S_2 = 780 \cdot 52 \cdot 30 \cdot \underline{\underline{128}}]$$

So packet was on interface @ do nothing

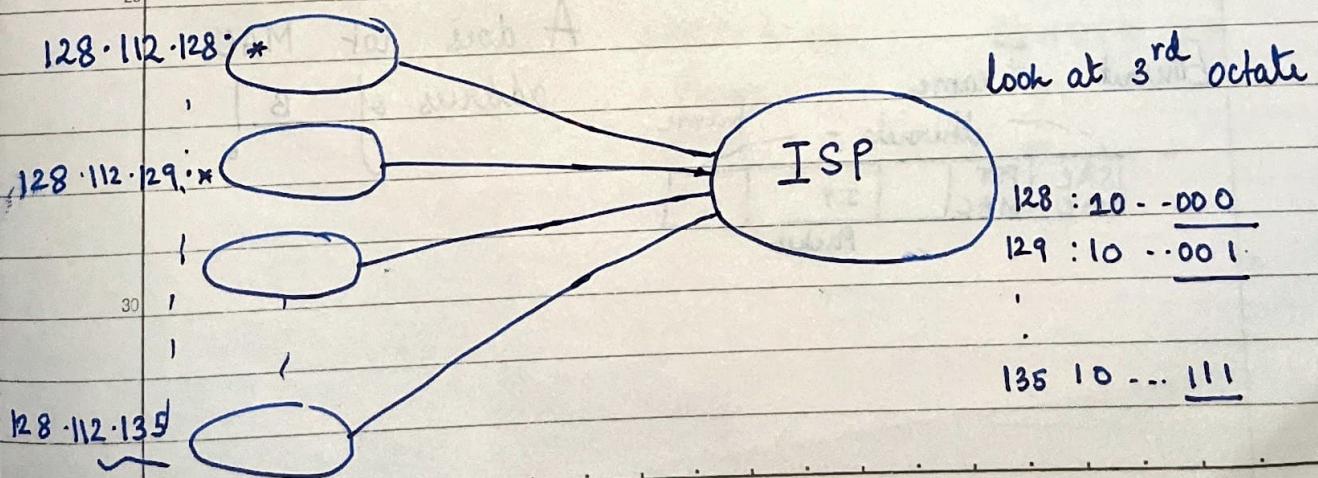
Is $(D \text{ and } M_2) == S_2$ then forward on interface

If not S_2 then forward on B.

Ex >



I want to combine them all in one if possible then how to do it !



Last 3 bits only change.
 → So we can write all these as a common IP prefix

$128 \cdot 112 \cdot 128 \cdot 0 / 21$

5 a.b.c.d / N → consider the N leading to get IP prefix
 Because these do not change for these company.

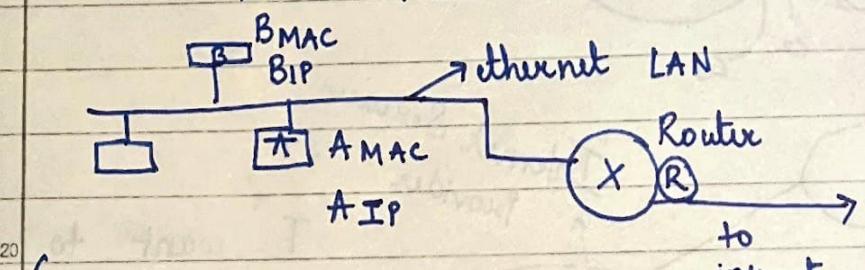
If given dest. IP address 'D'

10 If first N bits of 'D' match with the first N bits of a.b.c.d, the 'D' belongs to that prefix.

CIDR: Classless Inter Domain Routing :-

→ specify any prefix length (N)

15 ARP: Address Resolution Protocol



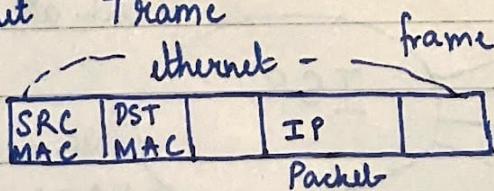
20 (R itself has its own MAC & IP)

A wants to send an IP packet to B.

25 A knows B's IP address :-

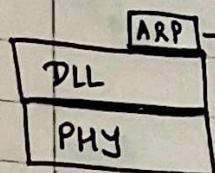
A does not know MAC address of B!

Ethernet Frame



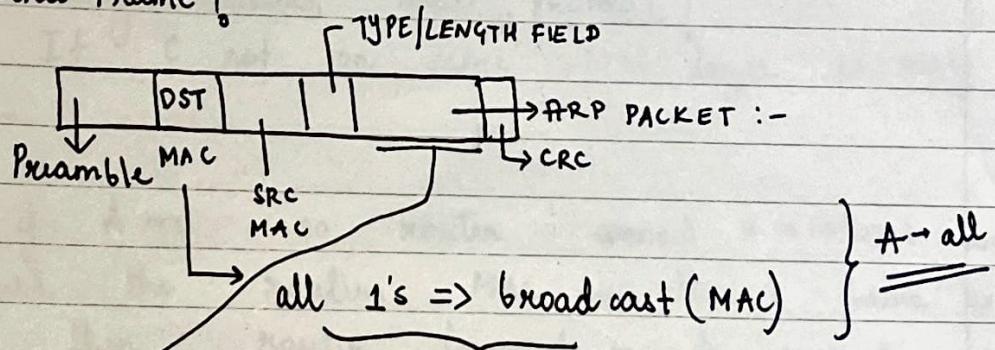
Now does A create ethernet header without knowing B's MAC address.

ARP to the rescue!



In practice we do not know exactly what is above what, by this I mean it varies with networks:-

Ethernet Frame!



ARP Request:

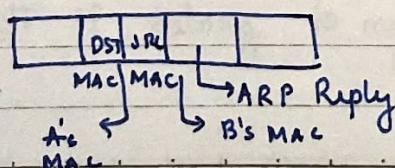
Request : SENDER MAC ; SENDER IP
 (A_{MAC}) (A_{IP})

TARGET MAC ; TARGET IP
 $(\text{all } 0\text{'s})$ (B_{IP})

So if Target IP matches then only packet useful!

ARP Reply $B \rightarrow A$

| | | |
|---------------|-------------------------|-----------------------|
| UNICAST FOR A | SENDER MAC
B_{MAC} | SENDER IP
B_{IP} |
| | TARGET "
A_{MAC} | TARGET "
A_{IP} |



- Information (Bmac) stored in ARP cache at A
 A maintains IP: MAC binding, to reduce every time asking
- This is not a permanent binding.
 Suppose some one changes card on device, or reassigned to someone else!
 → this has a time-out (order of minutes)
- What if B is not trying to send it someone on same ethernet
 10 ARP only handles local protocol!
 So If C not on same MAC layer no reply of ARP.

So if A → c so router connect A → internet we should the router MAC in ethernet frame as DST, then router forwards to its respective as

So A has to intelligent because what to put as MAC, whether to do ARP or not!

20 ① How does A know if DST IP belongs own network or not

② If DST IP on other network, how to know R_{IP, MAC}

1) A_{IP}; a₁, a₂, a₃, a₄
 eg)

Subnet mask Ex : 255 : 255 . 255 . 0

all 1's

If A's Target IP × subnet mask = A_{IP} × MASK

then same subnet so ARP

This DST IP belong to my address.

If not then we try to find RMAC

→ Assume we know RIP

→ We ARP the target address of R_{IP} & we get MAC of R.

→ How to get IP address?

i) Limited broadcast

DST IP = all 1's

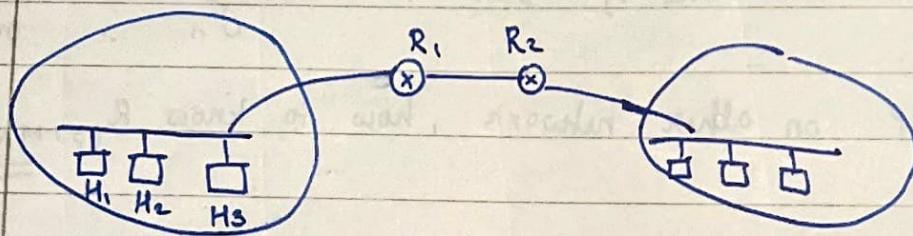
In this case all on same-subnet get message
no other device because

$$\text{DST IP} \times \text{MASK} = \text{SUBNET } \underline{\underline{\text{IP}}}$$

Now what is the issue, if we want to send broadcast out will this work no!

Why → The L3 router comes to know that this is a Limited Broadcast & stops it from leaving

ii) Directed Broadcast



10.1.1/24

subnet

Suppose we want to broadcast on 10.3.3/24 from 10.1.1/24.

→ If all IP dst bits = 1 then limited broadcast we send this back to subnet (I mean we do not send from R₁ → R₂)

DHCP:-

I want to my IP address. Default Router's IP address

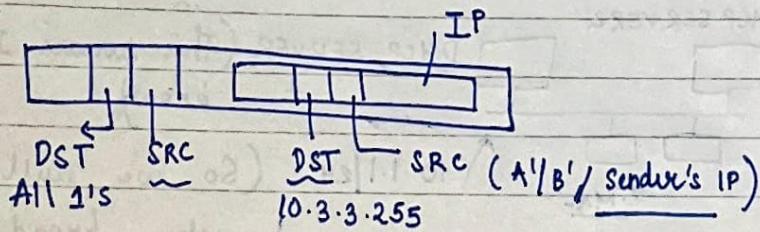
Date: _____

So In order to do this we need to know the subnet's IP or the LAN's subnet's mask!

Because IP PKT DST = subnet's IP!

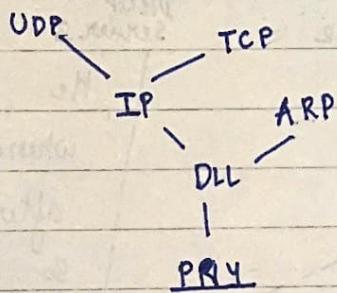
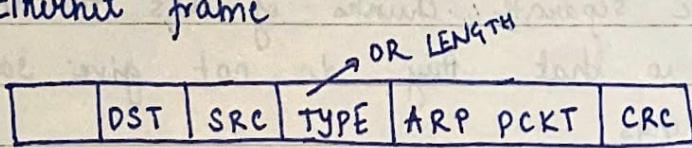
* all mac layer BITS = 1

10.3.3. 1111111 — HOST
24 bits 8 bits



To answer the previous Q we have DHCP
Dynamic Host Configuration Protocol.

Ethernet frame



We had the Type in DLL frame because sometimes we may want to do ARP and IP.

So thus we need a type/ length:

if value of that chuck < 1536 then this length if value > 1536 then this type field.
(eg : 0x0806) then ARP packet!

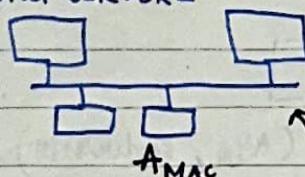
ARP: Know IP

Don't know Mac

DHCP: Node joins network
wants IP.

→ Suppose a Node wants to join a network
how to get an IP

DHCP SERVER 2



DHCP SERVER (This server's IP we do not know)

$10 \cdot 1 \cdot 1 / 24$ (So we will have to make broadcast in order to get to DHCP server:-)

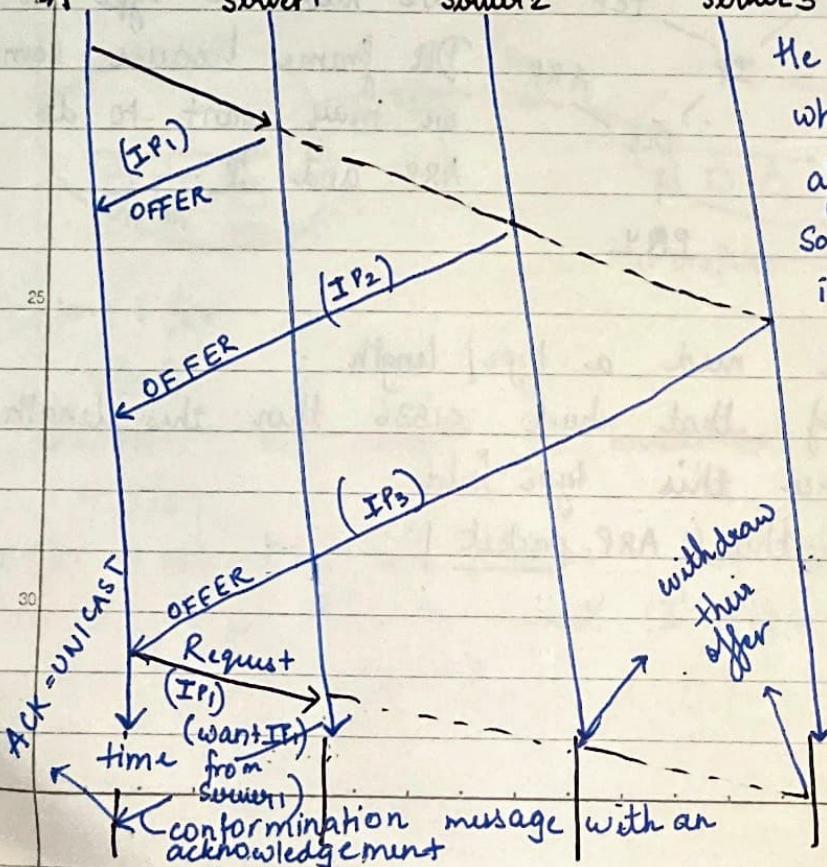
Sync time in servers

→ We usually give separate chunks of IP's to these servers and so that they do not give same IP to diff things

→ If given same pool then they have to co-ordinate

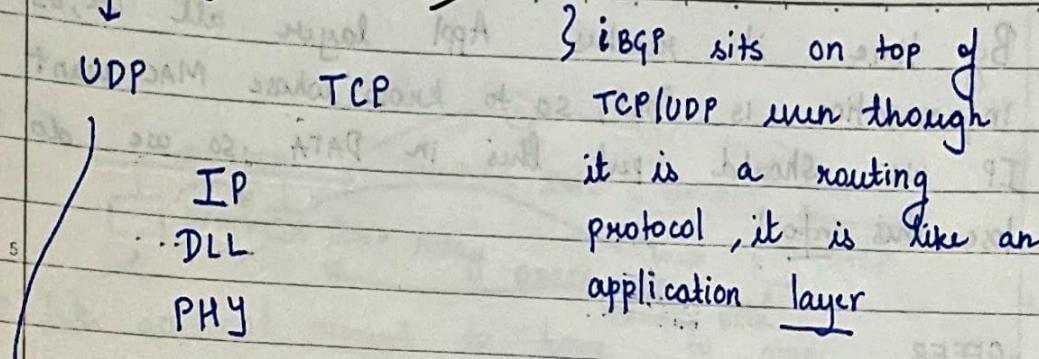
A $\xrightarrow{\text{DHCP}} \text{SERVER 1}$ $\xrightarrow{\text{DHCP}} \text{SERVER 2}$ $\xrightarrow{\text{DHCP}} \text{SERVER 3}$

He can accept whenever he wants after acknowledgement So after Δt time accepts it.

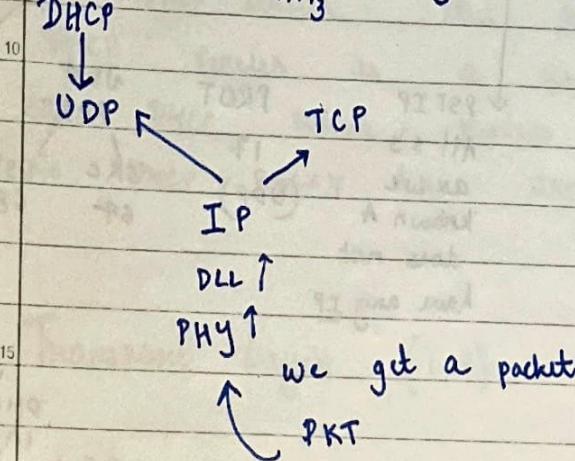


DHCP (port) = 68
for client

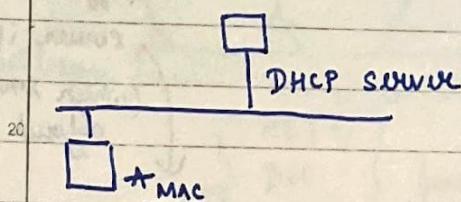
Date: _____



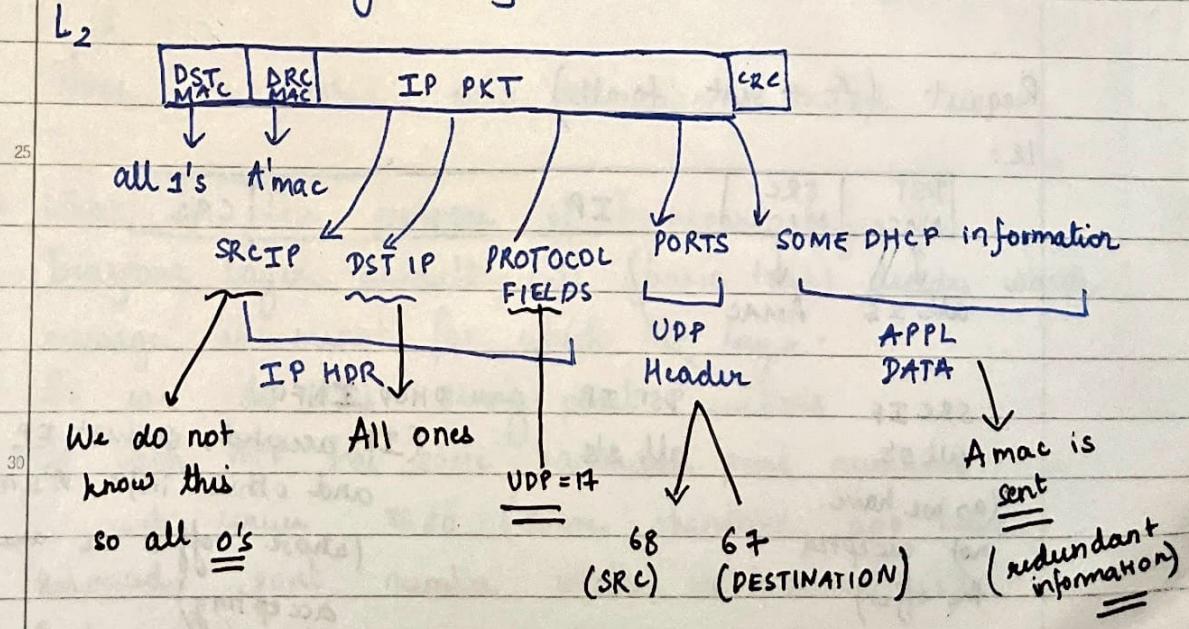
We want reply quickly so use UDP protocol



We have reached the IP field and IP decides next layer protocol after that which application it is meant, so we use PORT number to decide which application will get this pkt.



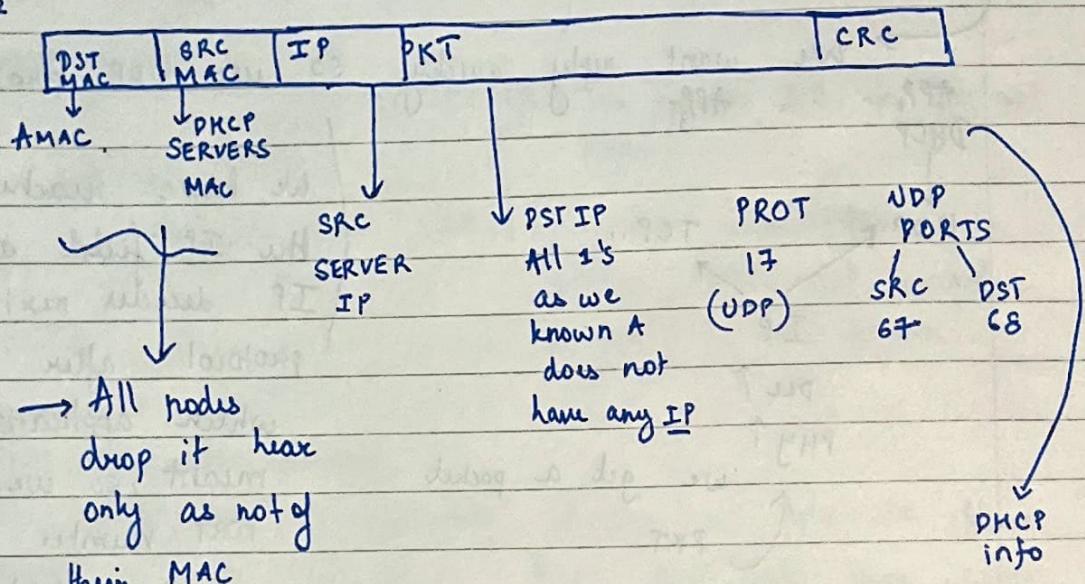
A sends discovery message



By time it reaches Appl layer all L2, L3 information is lost, so to know whose MAC want IP we should put this in DATA, so we do not lose this info!

OFFER

L2



Request ($A \rightarrow$ sent to all)

L2 :

| | | | |
|---------|---------|----|-----|
| DST MAC | SRC MAC | IP | CRC |
|---------|---------|----|-----|

↓
all 1's
AMAC

SRC IP

all 0's

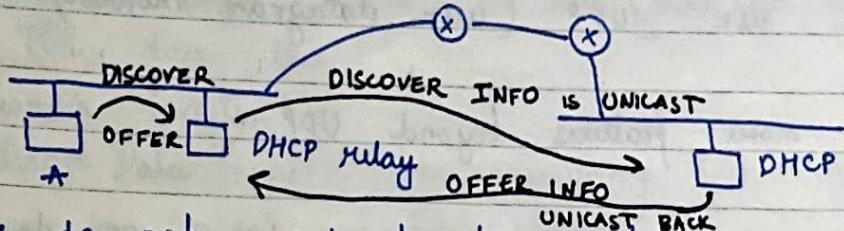
(as we have
not accepted
the offer)

DST IP ... DHCP INFO

all 1's

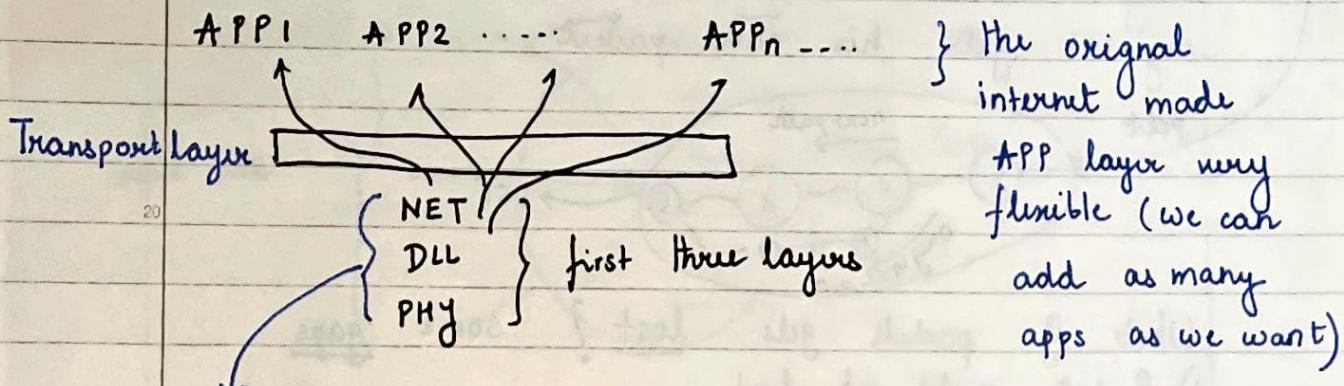
(IP accepted, server's IP
and other info A's mac)

(whose offer we are
accepting)



- We do not want to have too many sources on LAN, so only few subnets have DHCP
- So we have a DHCP relay which acts as a DHCP server, this forwards A's message to DHCP server as a unicast to Destination DHCP server and DHCP server replies via unicast
- So DHCP RELAY has SRC IP has RTP (GADDR)
↓
Gateway

Transport Layer (L-4)



These are pushed into OS, network drivers

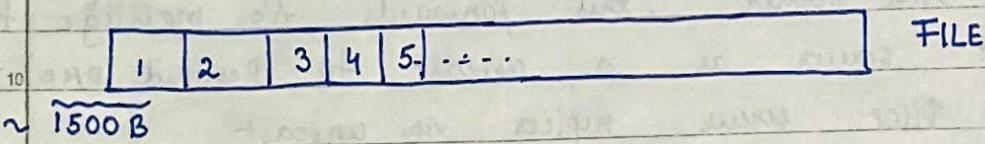
- * What is the purpose of Transport layer?
Transport layer demultiplexes (basic task) decides which message is meant for which APP layer.
- So we do this using port numbers
- So each APP has some particular port numbers.
- " web server #80 (some standard app have reserved port number which we cannot use)
- Port number in transport header

If we want a protocol doing only this we should use UDP (User Datagram Protocol).

Some more features beyond UDP :-

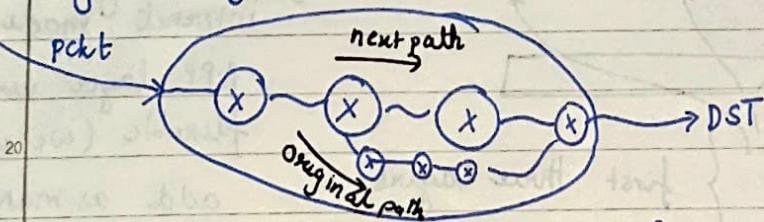
So we have some work where we do file transfer

FILE TRANSFER



We know while transferring we will encounter Ethernet somewhere (at some L-2) so max size of Data should be ~1500B

So we make sure the file is segmented in such a way that its data + size header ~ 1500B. So each segment gets his own packet.



What if packets gets lost? some gaps

1) Packets could get lost

so 1 2 4 5 6 7 ---

2) The weights of links could be changed, some link fails so a new path is taken

so packets could be faster on new one & there could be reordering

1, 2, 4, 5, 7, 8, 6, 9, 10

Even though Queuing implement FIFO, still reordering could happen.

Now files look bull shit
 So transport layer should do this (take care of this).
 So TCP does this

: Reorder Packets

: Reliable Data Transfer : retransmitting lost segments

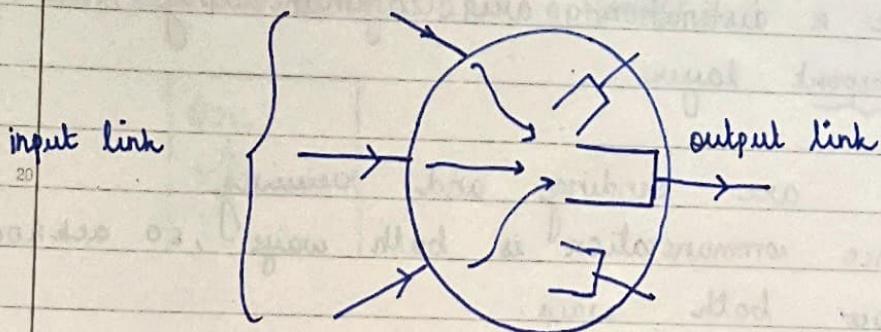
: Congestion Control and Flow control

So to this the source will have to do retransmission

→ So TCP will not allow any other packet to be received after 12, until 3 comes, so source retransmits ③ segment.

- So now receiver has to do some buffering.
- Routers are "Dumb" they cannot do this!

So TCP has to infer packet loss



So if input rate > output rate

Queues start filling up and packets are buffered.

- So Queues have drop tail mechanism.

- If this happens congestion occurs:-

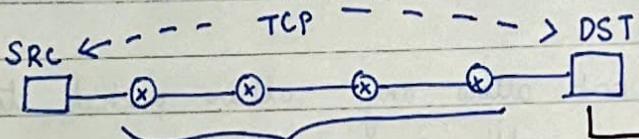
→ Queues full, packets are dropped.

- So if apps are using UDP a lot they just keep sending data & congestion prolongs for a bit.

So we need TCP to do congestion control.
"reducing segment sending rate"

So if many are using TCP so many input wires reduce input rate, so Queue starts emptying

FLOW CONTROL



So if congestion in router along SRC-DST we call it "congestion control"

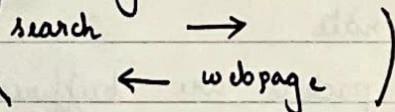
* if congestion occurs here handling it here is called flow control.

Why does flow control differ from congestion control.

→ The src & destination are communicating directly at Transport layer

→ So we are sending and receiving. Now since communication is both ways, so acknowledgement goes both ways.

e.g. Google search :-



→ So the receiver has a TCP buffer and it could full, this is because receiver could take time to process received data.

→ So receiver can directly say without informing (some mechanism of congestion control) that there is congestion.

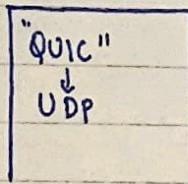
→ So can communicate this to src easily!

- When to use UDP?
- DHCP servers use UDP and avoid TCP
- When the latency is very low on the path so
 $321 \rightarrow 123$
 - kind of all will same transport time
- When data is so small that it fits in one packet (no reordering needed) (still message loss?)

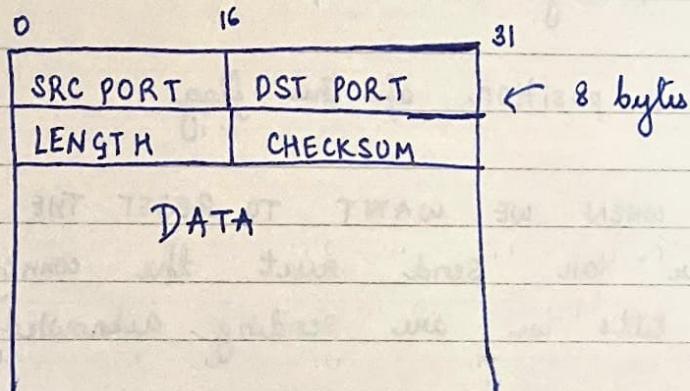
If high latency we can use TCP!

Because TCP somehow does not allow wrong order and if we sent $321 \rightarrow 1-3$ (second) 2^{nd} has not arrived, (got lost in buffer) it will only send 1 and withhold queue from 2-3-4-5-6....
Too much buffer!

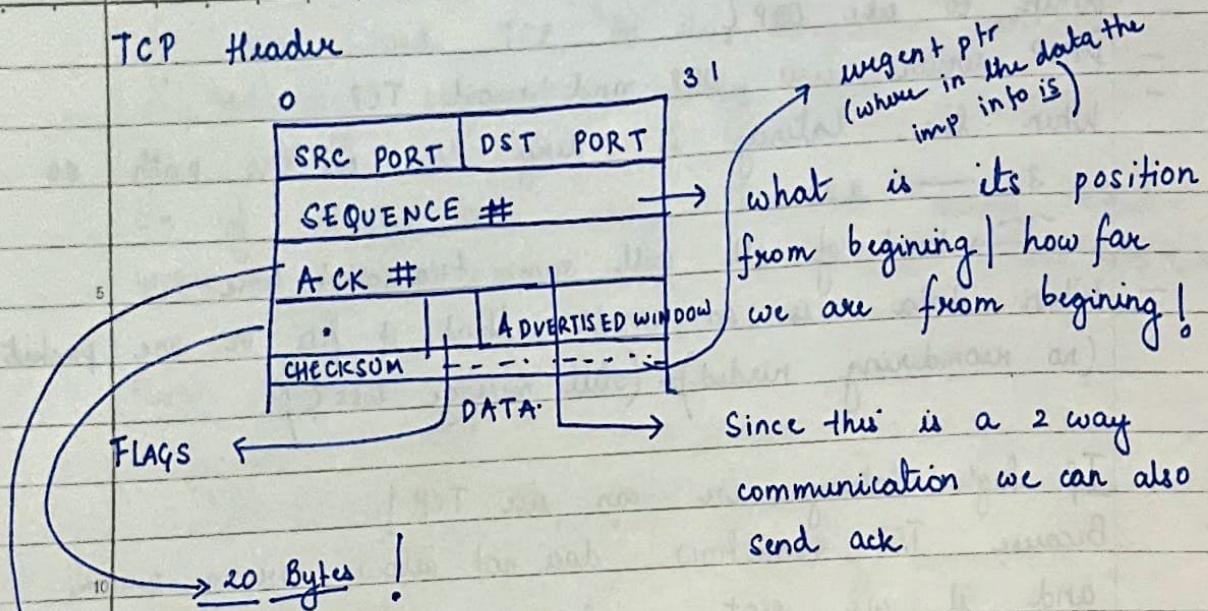
Google Chrome uses something called "QUIC"
It sits on top of UDP optimization!



UDP Header



TCP Header



20 bytes of our head so each size of data.
 → what byte # we are expecting next
 We have got everything till prev byte)

SEQUENCE NO

BYTE NO of 1ST byte in segment

FLAGS

↳ SYN | FIN | RESET | PUSH | URGENT | ACK

Initially we start a connection we send a syn pckt fin packet (end of connection)

→ So 0/1 in position of this flag.

- RESET IS WHEN WE WANT TO RESET THE CONNECTION
- So Attacker can send reset the connection!
- ACK bit tells we are sending acknowledgement.

Date : _____

ADVERTISEMENT WINDOW IS USED for telling how much
of TCP buffer is full/ so if too less space avail
the server receiving ^{empty} can slow down sending.

5

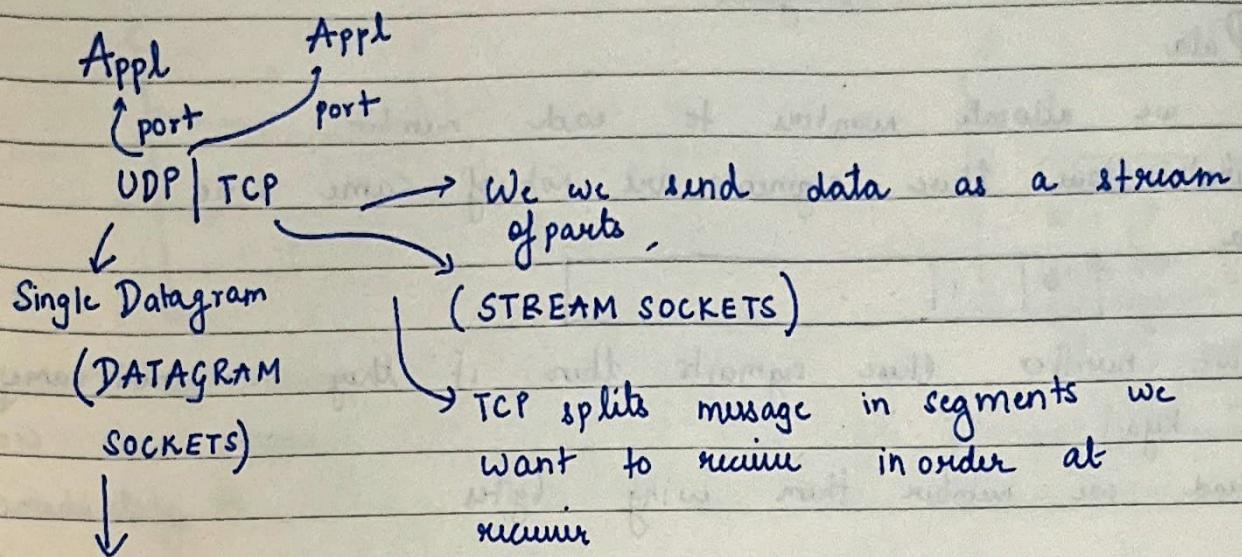
10

15

20

25

30



UDP:
Packet received ← give immediately

Here we need to remember a lot of state information

CONNECTION:: Sender $\xleftrightarrow{\text{TCP}}$ Receiver

* This connection is responsible for Congestion control, flow control

(Listens at port (e.g. 80))

TCP connection Establishment

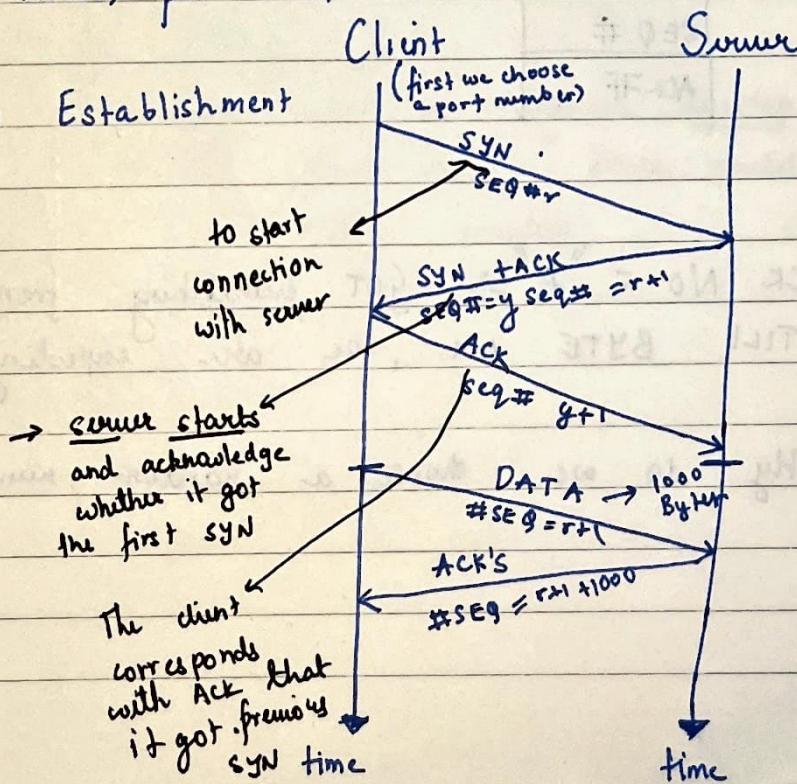
3-WAY HANDSHAKE

→ Only done for TCP

→ FLAG field of TCP Header

SYN | FIN | ACK

We set these bits



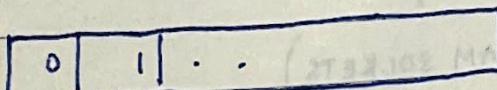
How is TCP designed

Data

→ So we allocate numbers to each number

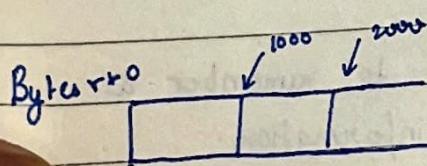
→ First issue these segments are not of same size

Data



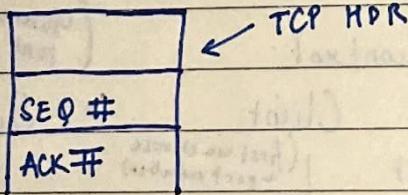
If we number these segments then if they are not same then kya?

Instead we number them using bytes



We use to send the starting number of byte in message and then using data length field we estimate length of packet

also we do not use 0, we use a random number
∴ r = random number

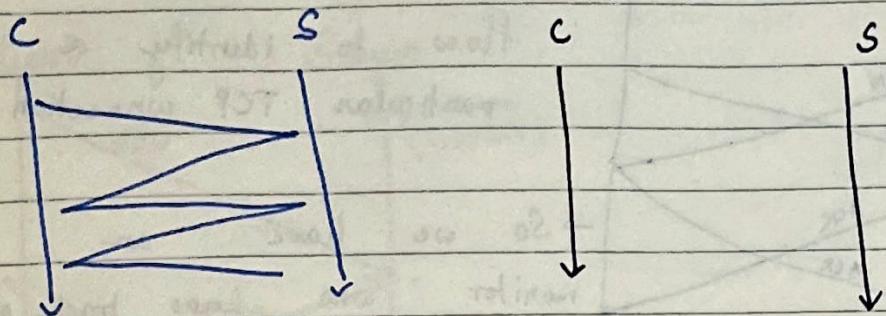


ACK NO = "A" \Rightarrow GOT everything from start (SYN SEQ #)
TILL BYTE A-1, we are expecting byte A next:-

Why - do we choose a random number?

If segments number from 0

new connection



1 2 3 4
data from old segment + connection (lost in network now)

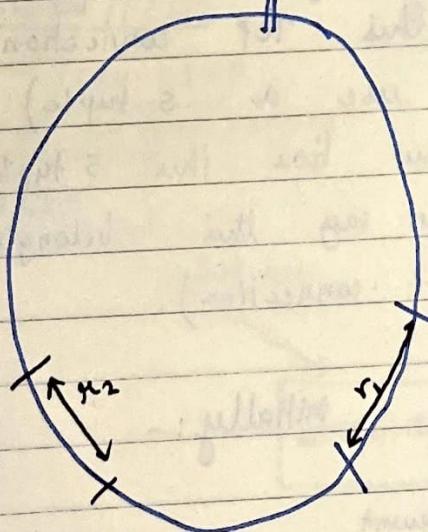
↓ → we close the connection!

Now we already (I mean we open an connection) immediately to server! (client port = same, server port = same)

↓
So these could be so close in time, some data by first old connection (was stuck in some router) and we have started receiving info from new packets.

So if it occurs it we put in some place with its new data. So to avoid this we need sequence numbers

SEQ #: 32 bit $0 \dots 2^{32}-1$



→ so for different connections we try to have different segment range for

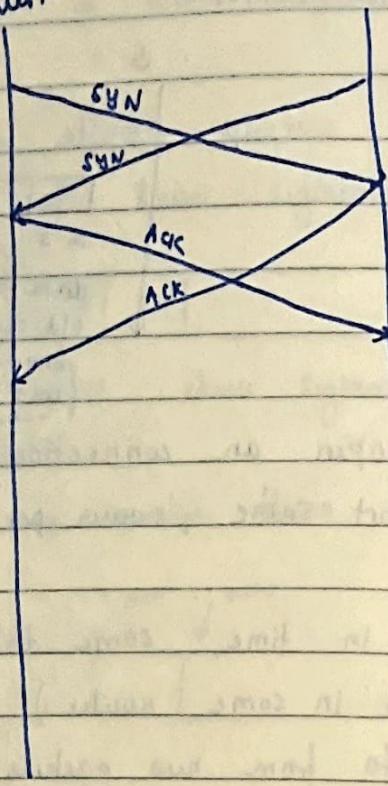
→ (we want to minimize overlap)
→ If files large there could be some overlap.

4-WAY hand shake :-

NFT pg

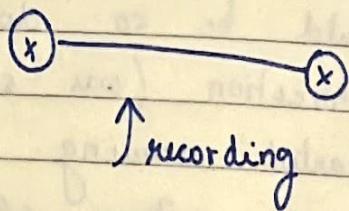
Client

Server



How to identify a particular TCP connection,

→ So we have an monitor who keeps track of TCP connections & dubs data if TCP sends it to suspicious place



So

IP : SRC IP, DST IP

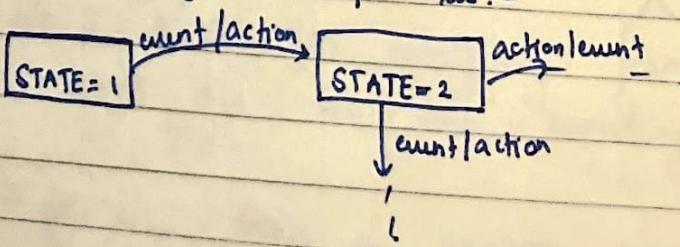
PORTS : SRC PORT, DST PORT

PROTOCOL : (INDICATES TCP)

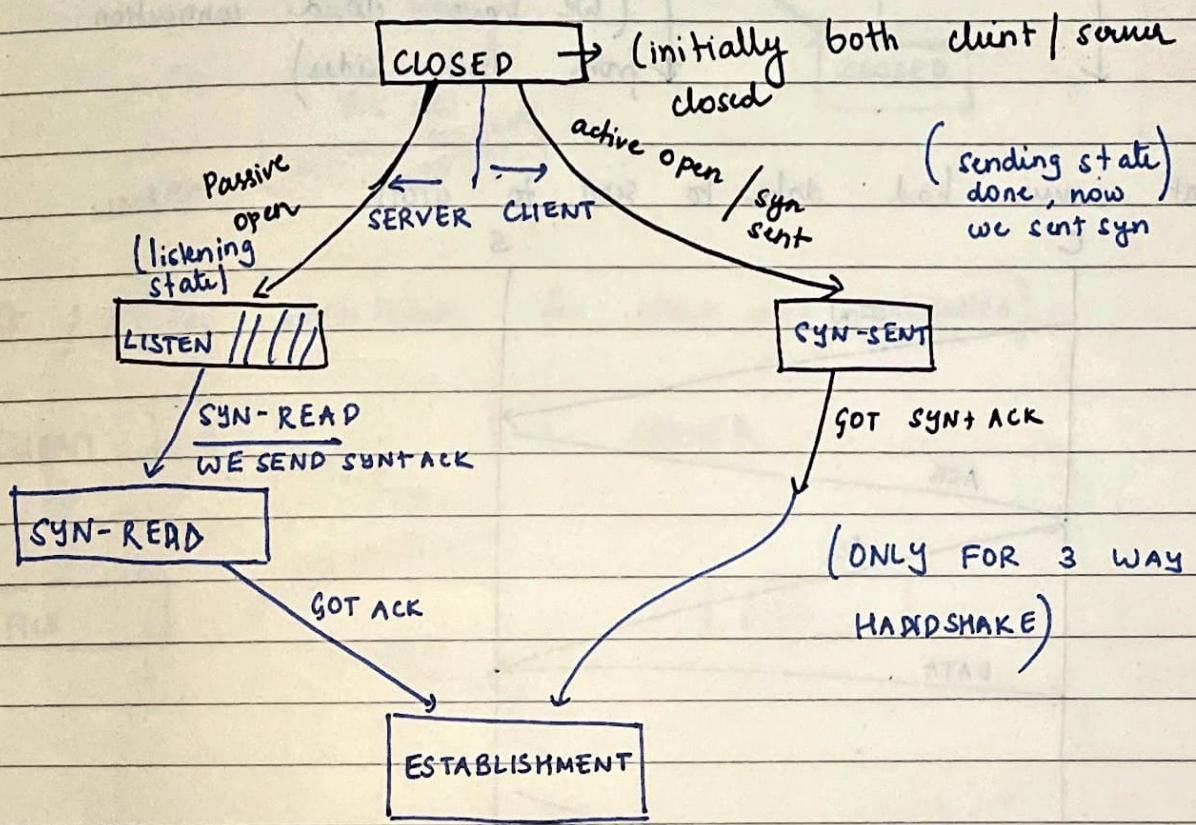
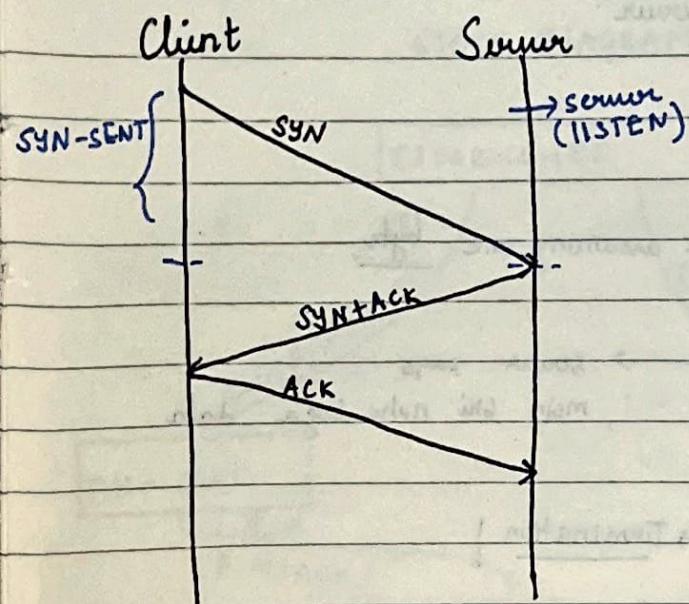
(IP HDR).

If we want to identify if this belongs to this TCP connection:-
(So we use a 5 tuple)
(so whoever has this 5 tuple we can say this belonging to this TCP connection).

So server is always listening initially:-
So state diagrams helps us:-

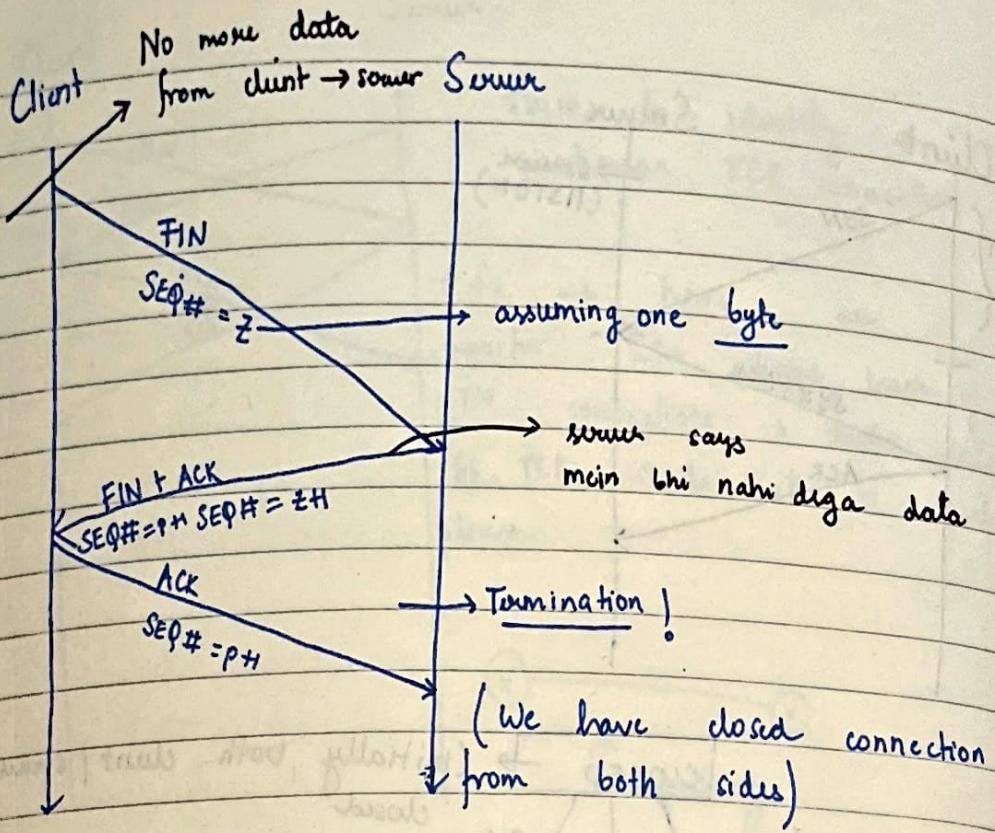


So state diagram for 3-WAY HANDSHAKE

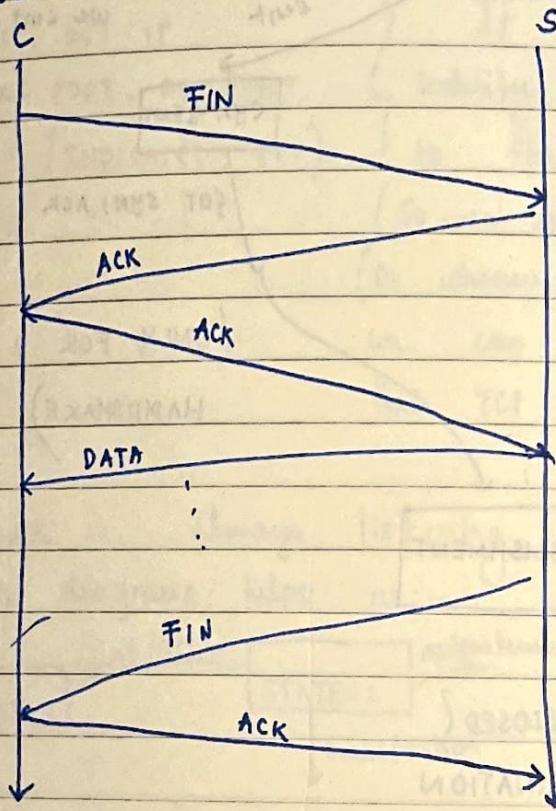


HOW CONNECTIONS IS CLOSED?
CONNECTION TERMINATION

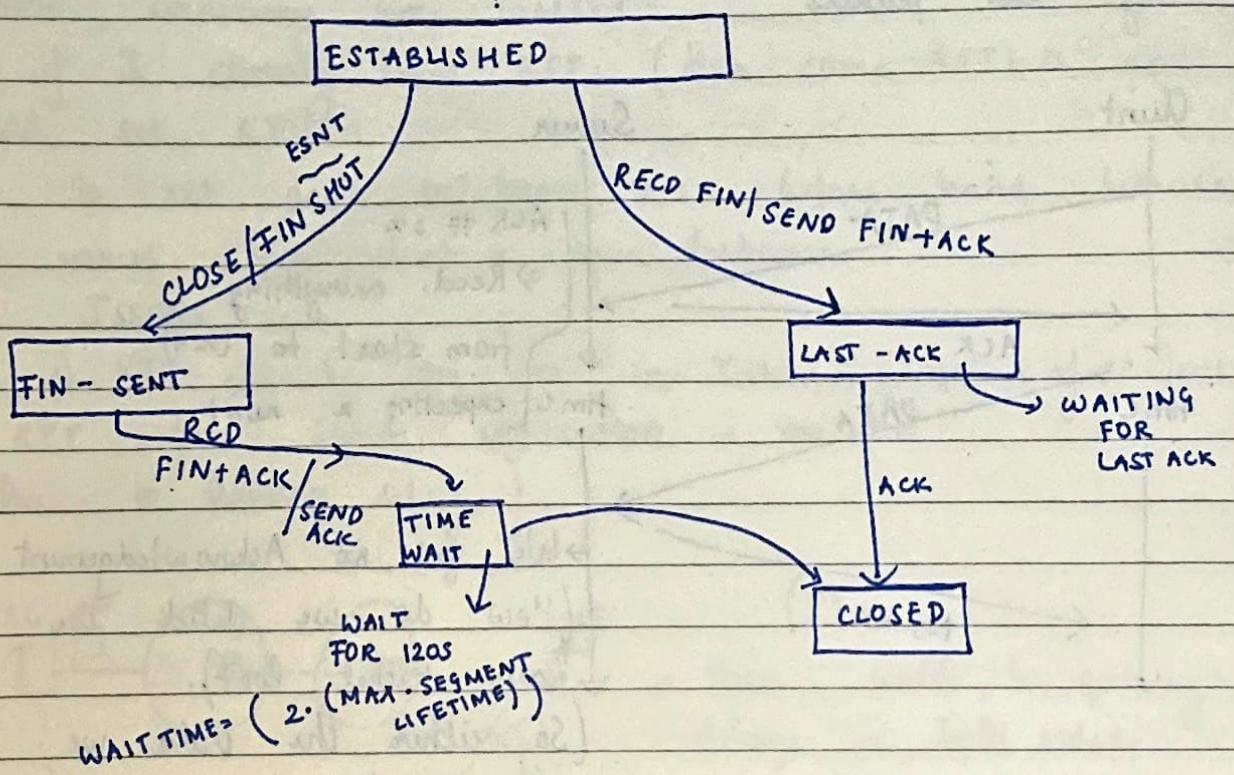
Nxt Pg



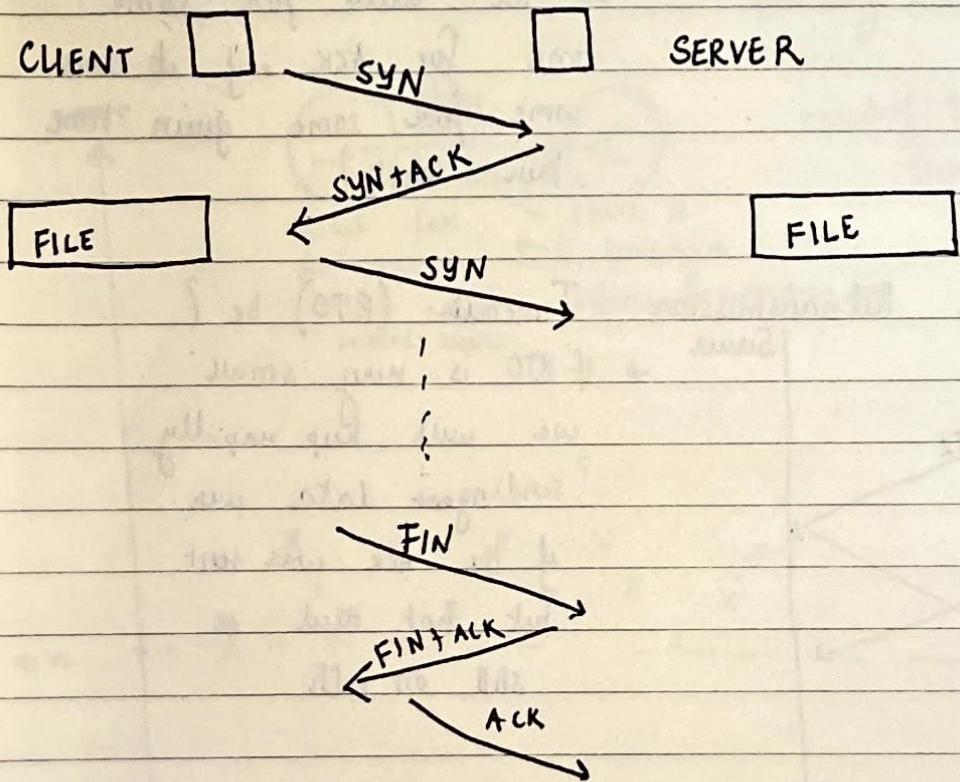
What server had data to send to client.



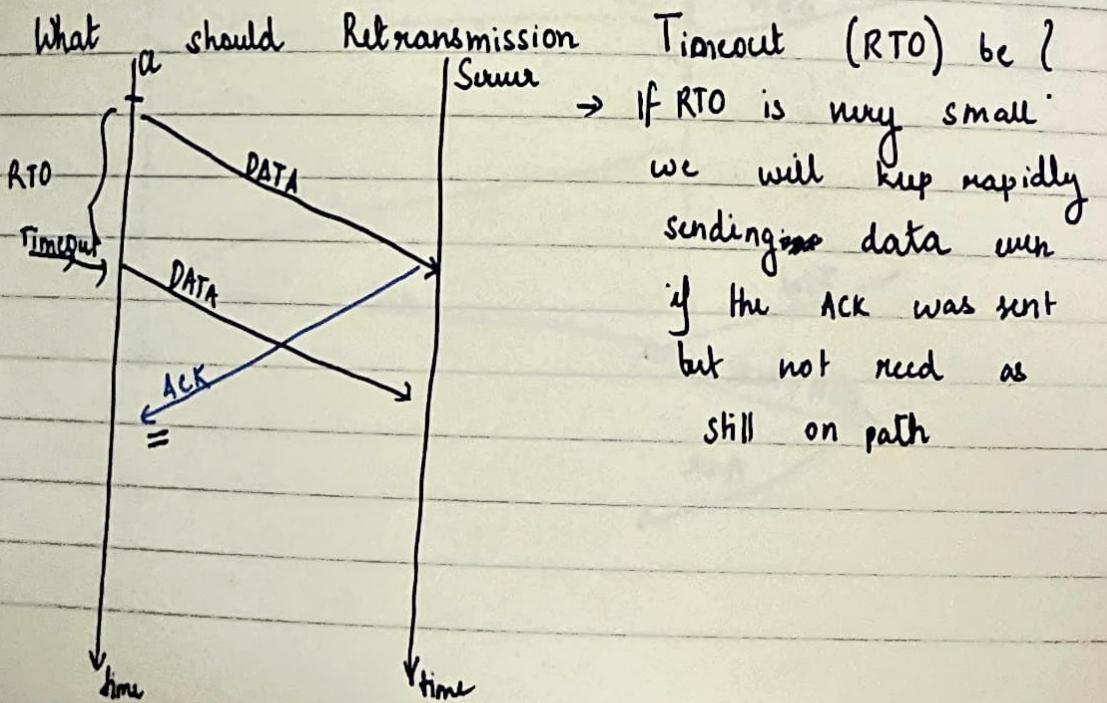
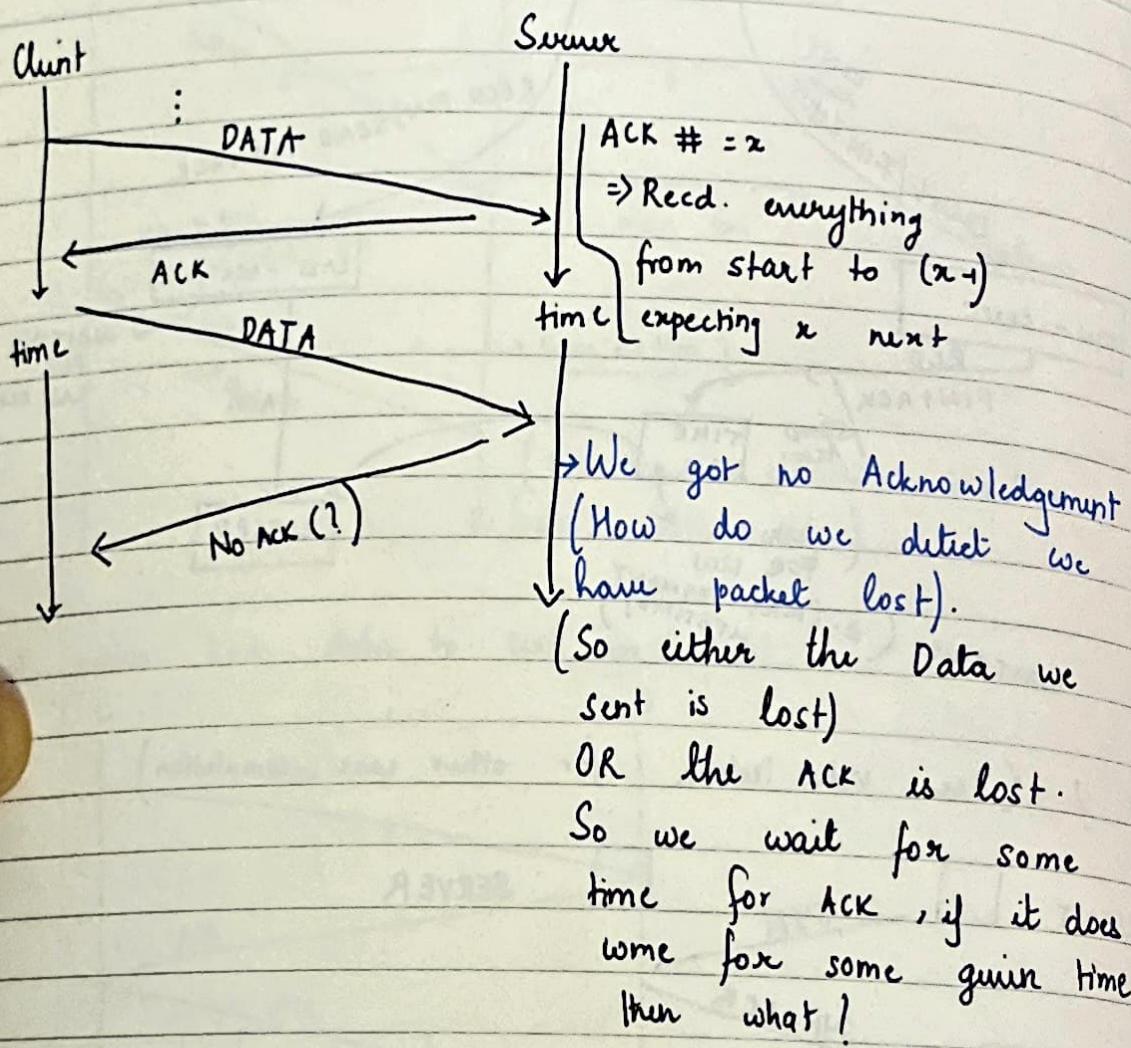
STATE DIAGRAM



IF ? (See video lecture for other case - completion)



① TCP has to do retransmission!
of lost packets



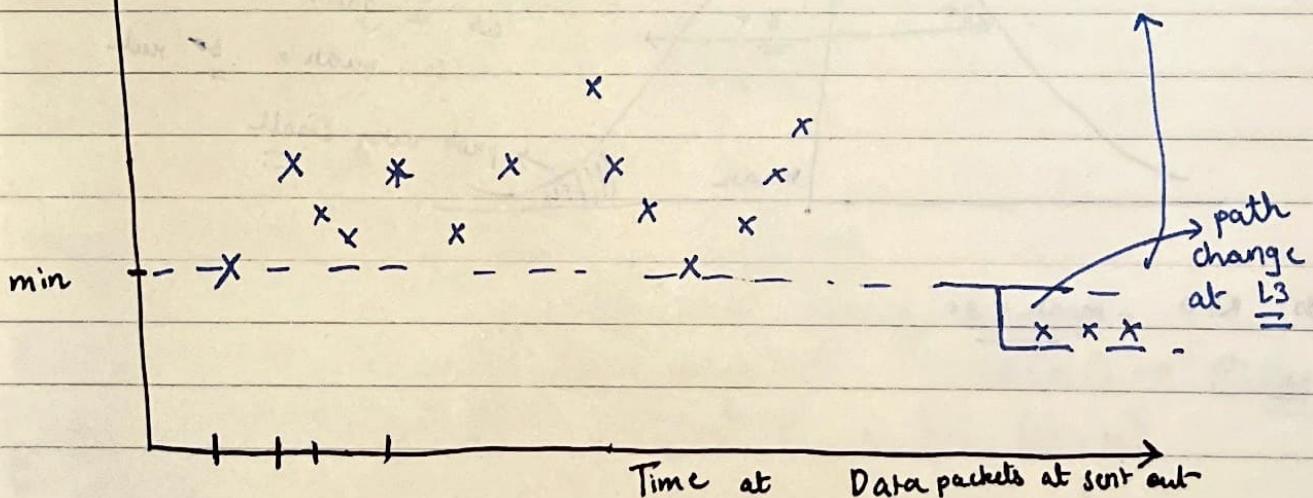
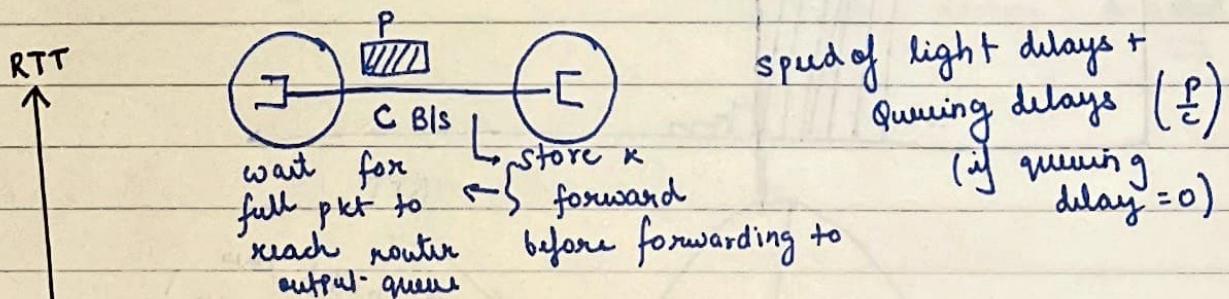
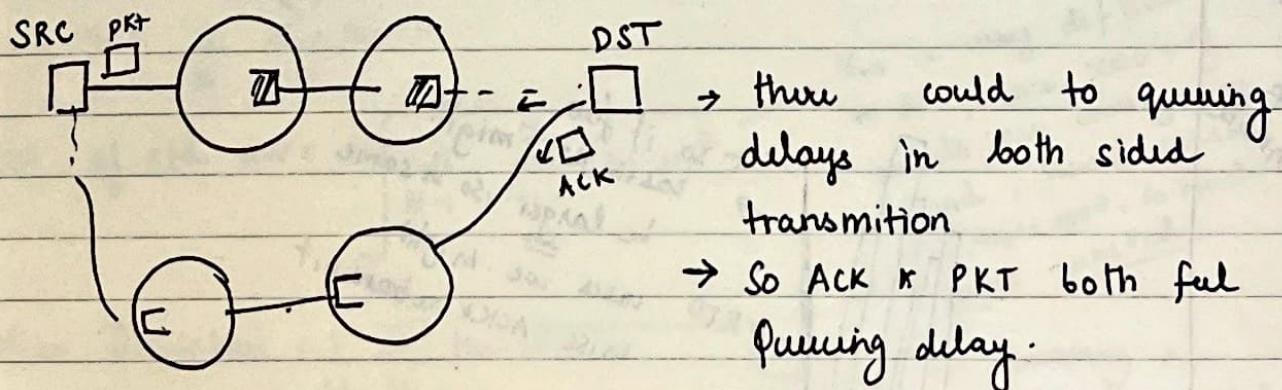
→ If RTO is very large, we will wait for a lot of time, unnecessary time wasted.

So if I clearly know RTT (then some $RTT + \Delta$ we could set RTO).

It is not good to know RTT before hand because it varies for client a server distance.

ISSUES

- So, 1) RTT can be ms - sec in Internet (location of SRC/DST)
- 2) RTT for same connection is variable.
(Due to Queuing delay)



So in store & forward we use checksum to check whether packet is corrupt or not

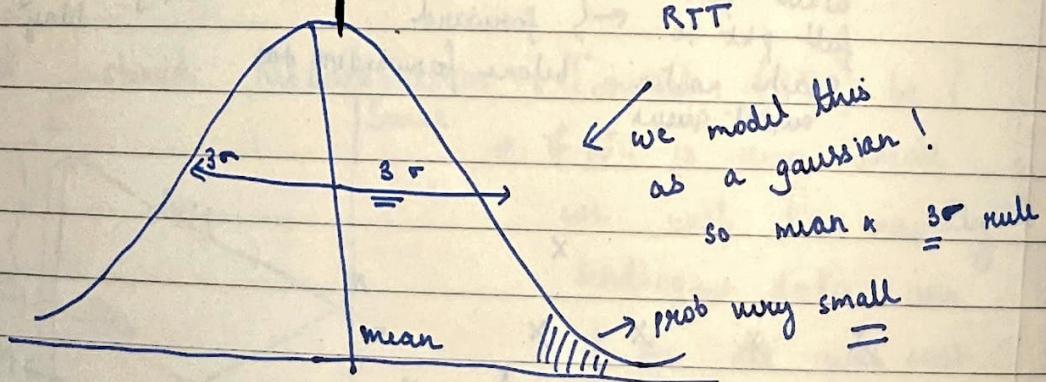
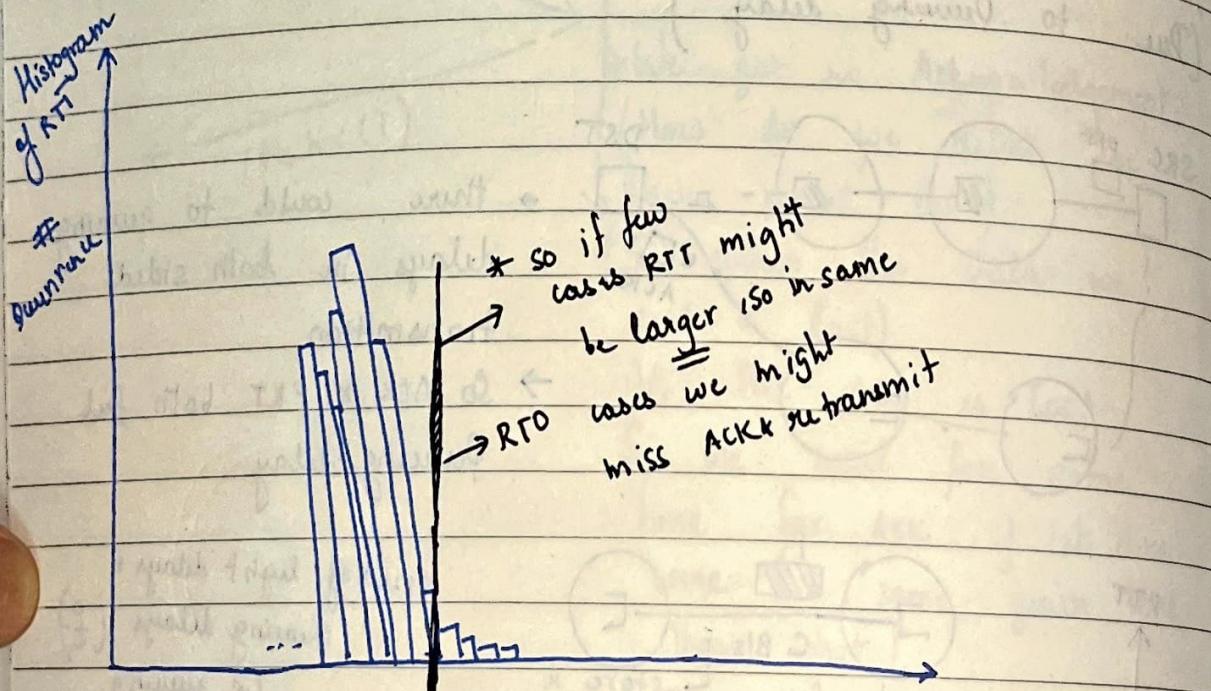
So CUT-THROUGH delay we immediately throw it in output queue the moment we start getting so we can reduce transmission delay.

We won't be able to check whether packets we have received is correct or corrupted, but in CUT-THROUGH we have speed.

So how should I choose RTO :-

So people plotted a histogram of RTT

Because we do not queue size, link speed, buffer capacity, queues ----- at all



$$\text{So } \text{RTO} = \underline{\text{mean}} + 3\sigma$$

Idea: Measure RTT over time, set $RTO = \text{mean} + \underbrace{\text{const}_n}_{\text{r}}$

→ Intuition:-

→ Also I want to look back in time quite less.
So that I can avoid errors due to path changes,
suddenly queues filling up, (congestion) ..
→ so limit look back period:-

∴ Random Numbers : x_1, x_2, \dots, x_N

$$M = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\text{Est of std dev} = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - M)^2}$$

this is very difficult
to do, since 1000's of
packets are coming every
second, how do square,
square root, so much
time wasted.

$$\text{Mean Deviation} = \frac{1}{N} \sum_{i=1}^N |x_i - M|$$

squaring &
sqrt rooting avoided
=

ALGO USED FOR TCP (for RTO)

RFC

(request
for
comments)

Sample RTT ;
↳ lastest RTT estimate ;
↳ current

current estimate of mean actin RTT

→ exponential moving average.

$$\text{EstimRTT} = (1-\alpha) \text{EstimRTT} + \alpha \text{Sample RTT} \quad \alpha \in (0,1) \rightarrow \textcircled{2} \underline{\text{step}}$$

① Difference = Sample RTT - EstimRTT = like $\underline{(x_i - M)}$
first step

$$\text{Deviation} = (1-\beta) \text{ deviation} + \beta \underline{\text{difference}}$$

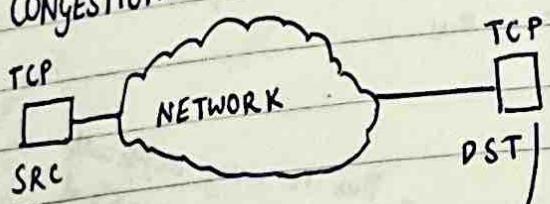
$$\text{Timeout} = \alpha \times \text{EstimRTT} + \phi \times \text{Deviation}$$

$\alpha = \frac{1}{8}$

$\beta = \frac{1}{4}$ } Default

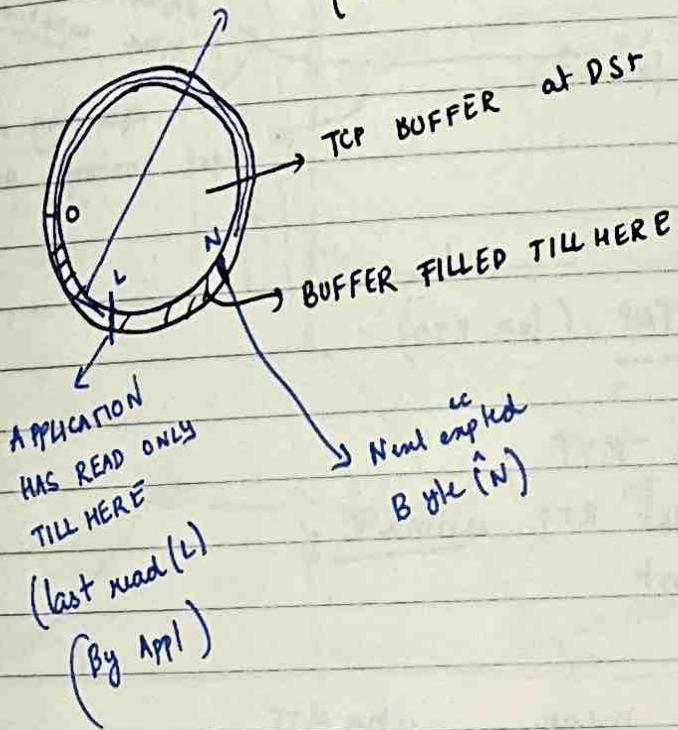
TCP Vegas:- (Some Class)

CONGESTION & FLOW CONTROL



→ (Network does not touch the packet just pass it to DST)

What if congestion here?
(Too much load)



So the major arc of LN is available as some delta can be overwritten as read and processed.

So if APPL is slow buffer could easily get filled.

So Field in TCP hdr = Adv. Window = Free Buffer available

So if this buffer space is reducing at a high rate we want to communicate to APPL at src to slow down.

Buffer size = M

so summing space = $M - (\underline{N-L-1})$

↳ (Because N is next expected byte)

so if $N < L$

: $(L-N+1)$

Because N has wrapped around.

Adv. Window = $(M - (N-L-1)) \bmod \underline{M}$.

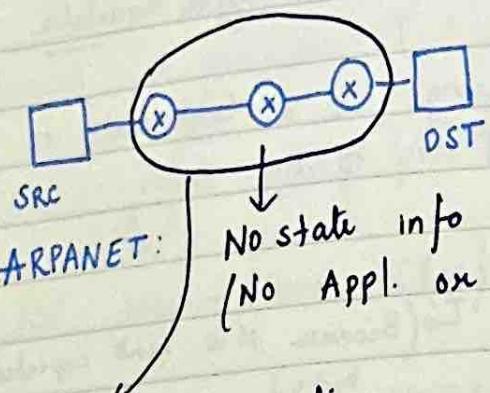
SRC using window how much max unacknowledged bytes I could have send out. ↳ Definition of window

WINDOW = min (congestion window, Advertised Window)

↳ Method to control Data rate.

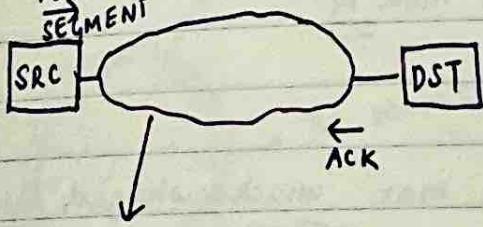
Depends on congestion window!

TCP CONGESTION CONTROL

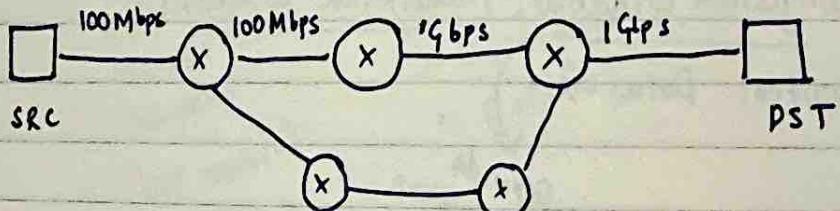


only pkt forwarding.
If network goes down, since we had no state, so no loss of huge amount of information.

So they put the complexity at end
So end nodes handle packet loss (packet forwarding) mech.
So for $\text{TCP} \times \text{SRC} \times \text{DST}$ Network is a black box



So these do not know Data rates in network
but somehow they have to estimate this.

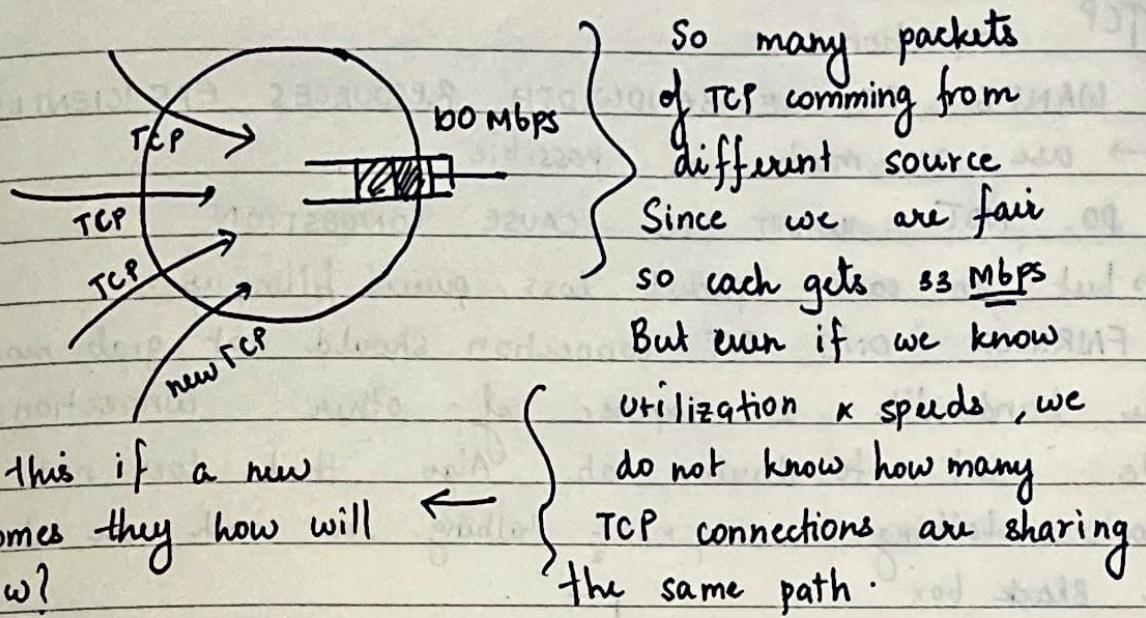


So they do not know bottleneck link speed & they do not know at what rate data should be sent.

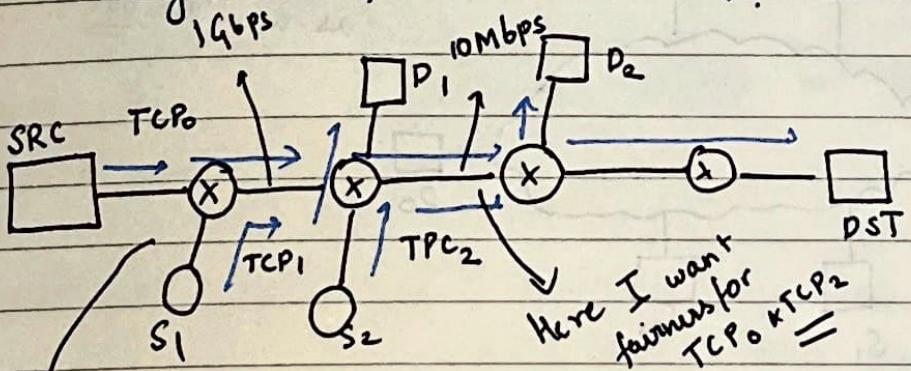
- End-Hosts running TCP

- DON'T KNOW LINK SPEEDS

- DON'T KNOW UTILIZATIONS



Also knowing LINK speed does not help!



So I do not know how many connections I share the path with, I can say about particular link.

$TCP_0 \rightarrow$ Data Rate : 10 Mbps
 $TCP_1 \rightarrow$ " " = 990 "

$TCP_2 \rightarrow$ " " = 5 "

TCP_i is not hurting TCP₀ because TCP₀ cannot ~~at max~~ get more than 10 Mbps since 10 Mbps is bottle neck

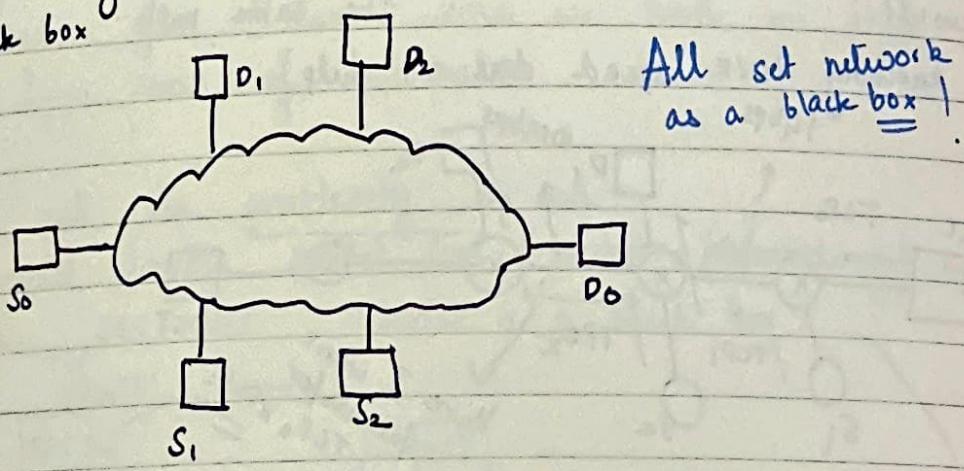
So $TCP_0 = 5 \text{ Mbps}$ } so always there not fair distribution
 $TCP_1 = 995 \text{ "}$
 $TCP_2 = 5 \text{ "}$

So Can these source and DST talk?

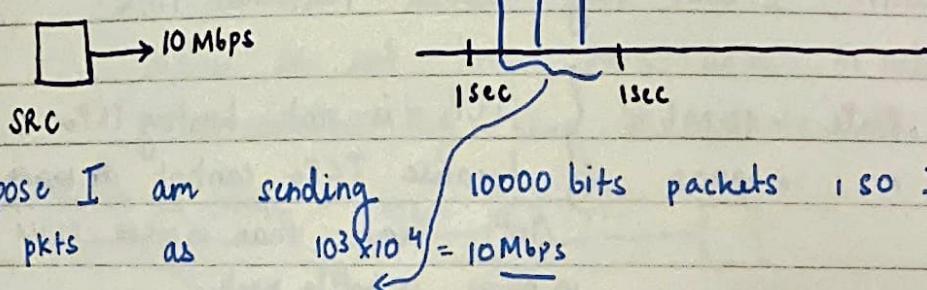
So Congestion Control Issues

TCP Connection

- WANTS TO USE BANDWIDTH RESOURCES EFFICIENTLY
 - use as much as possible
 - DO NOT WANT TO CAUSE CONGESTION
 - but do cause packet loss, queues filling up.
- FAIRNESS - ONE TCP connection should not grab most of the bandwidth at expense of other connections.
We want to develop an Algo that does not involve routers talking * $src_1 \times src_2$ talking, so each see network as a black box



Q) How to set data rate of TCP



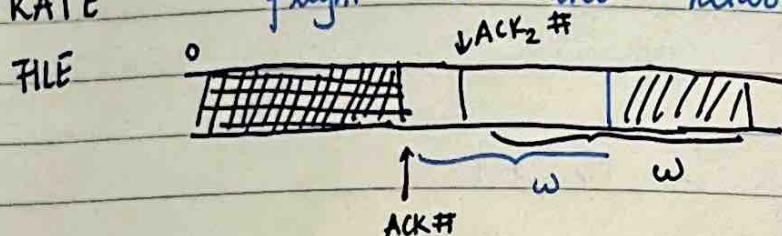
Suppose I am sending 10^3 pkts as $10^3 \times 10^4 = 10^7$ bits packets, so I can send

So we can send data for 1/2 sec or some part and then stay silent. So basically I burst into network.

But what if 1Gbps was speed then I would burst into network by sending 10^6 packets, so we can reduce window size to 10 ms and then we can stay silent, we have counter that maintains this.

WINDOW

BASED DATA RATE FILE W : max amount of unacked data in flight in the network.

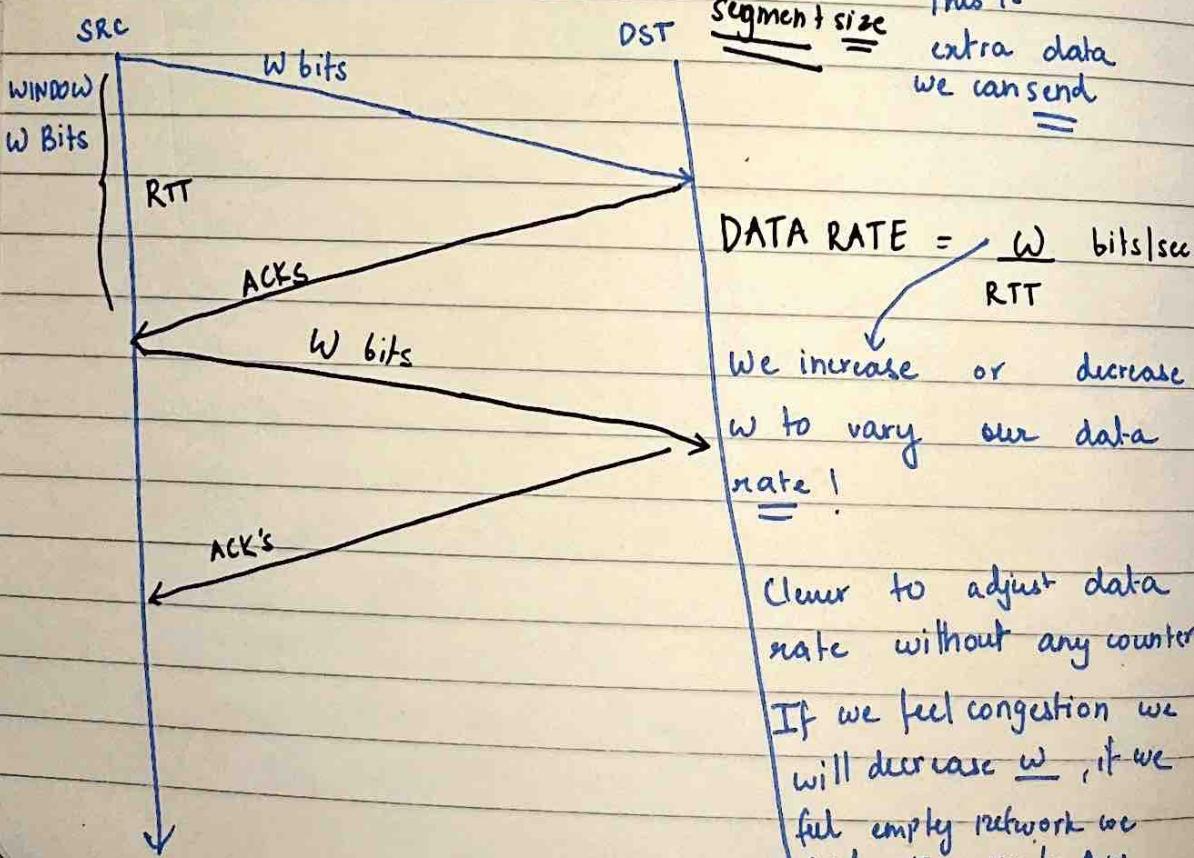
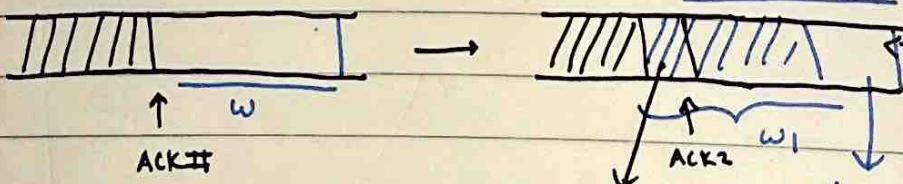


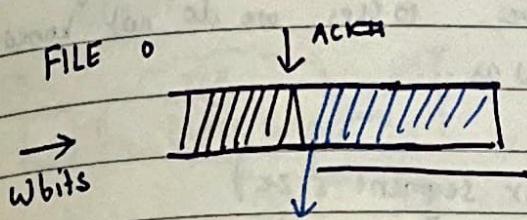
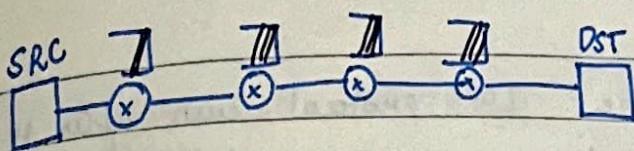
(the byte no server expecting)

So since window size w , we can send data from ACK no w size of.

But we get a new Ack

So I am allowed to send new data that is not in network after sliding window to right!

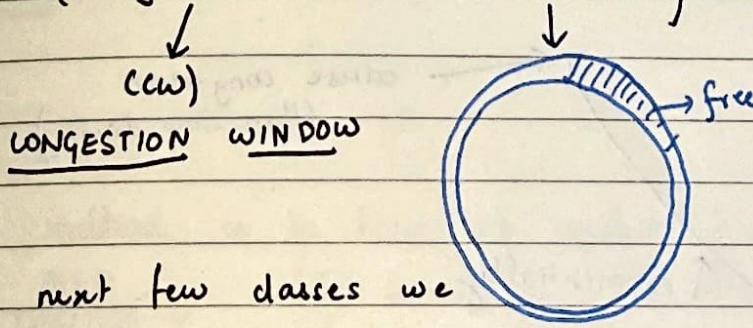




These w are in flight

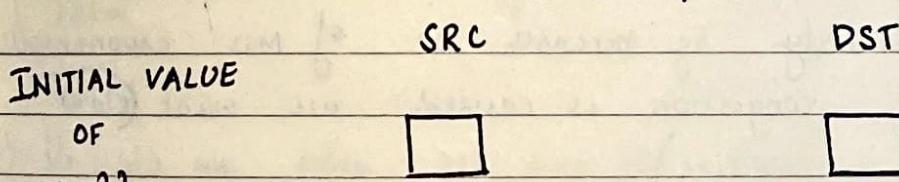
- so we can see that at max we would have occupied w bits in queue all over the path, so we come how heuristically control congestion!
- So may even have reached the DST, so would be queues. but worst case could be all/most queues across path!

$$w = \min (\text{long. Window}, \text{Adv. Window})$$



So far next few classes we will assume $\text{Adv. window} \geq \text{congestion window}$.

We $w = cw$ thus we do this!



Suppose $w = 10^3$ pkts $\rightarrow 10^7$ bits | RTT: 1 ms =
 $1 \text{pkt} = 10^4$ bits

So what is my data rate $= \frac{w}{RTT} = \frac{10^7}{10^{-3}} = \underline{\underline{10^4 \text{ bps}}}$
very high!

We want to use the same TCP protocol even after 10 years, so we can all have 10 Gbps, we do not know anything, what should be set w as.

IDEA: SET $w = 1 \text{ MSS}$ (max segment size)

$1 \text{ MSS} \sim 10^4 \text{ bits}$

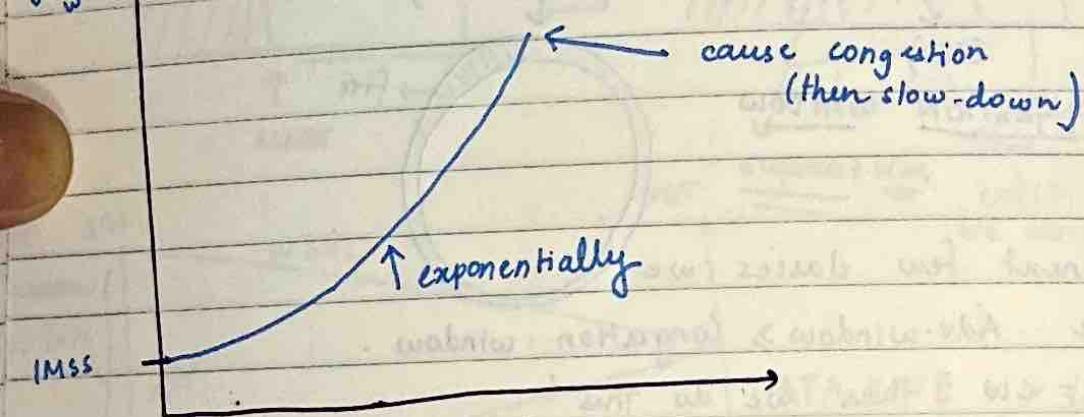
RFC 5681

less than $1 \text{ MTU size of Ethernet}$
so we send only one packet

So we are not using bandwidth efficiently
So this is called slow start: start with $w = 1 \text{ MSS}$

(we want
increase
 w rapidly
to get
appropriate)

Want w
Assuming
 $w = w_{\text{tiny}}$
 w_{tiny}

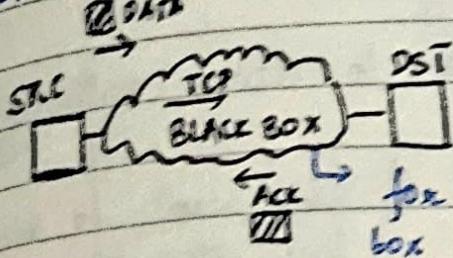


We start aggressively by increasing no of MSS exponentially and then when congestion is caused we want (slow down) and (lower wsize)

These two questions ↓
are unanswered

| |
|------------|
| TCP Tahoe |
| TCP Reno |
| TCP Vegas |
| TCP AFRICA |
| TCP CUBIC |

TCP CONGESTION CONTROL



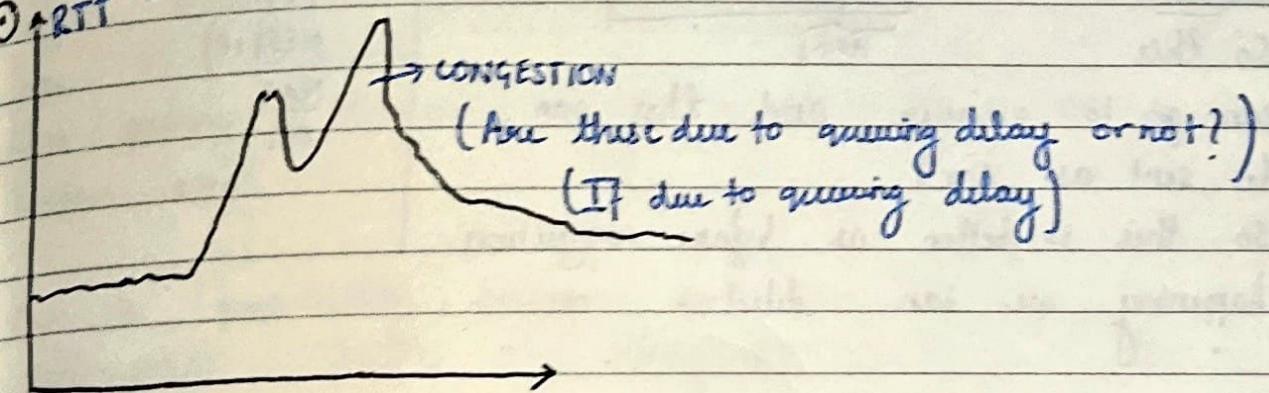
for source + dest. network is a black box

SIGNALS OF CONGESTION?

How do we figure congestion is there?

RTT may increase (due to increase in queuing delay)

RTT



② Other method is to look at packet loss!

FIFO



(If queue full we have packet drops)

(So we can somehow use this to figure out packet loss).

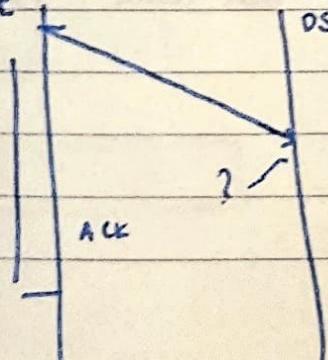
Data

③ How do we infer pkt drop has occurred?

(We use timer (for receiving ACK) so we infer pkt drop)

TIMEOUT

src

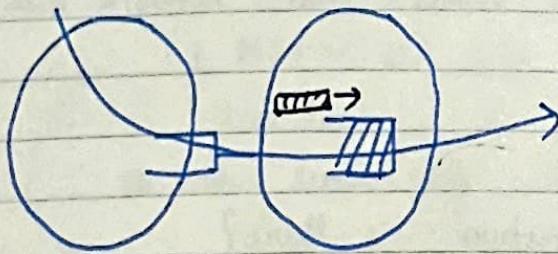


DATA →

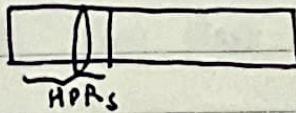
(We have to figure how to detect early packet loss because timeout could be very large)

③ ECN: Explicit Congestion Notification!

The routers should give feedback to source without changing protocol



So we set some bits in TCP header!



So this can go to receiver and this can be sent as ACK!
So this is better as before congestion happening we can detect.

(We want to keep fact intact that TCP does not happen at router level but we want some feedback)

TCP
NET(IP)
MAC.
DLL

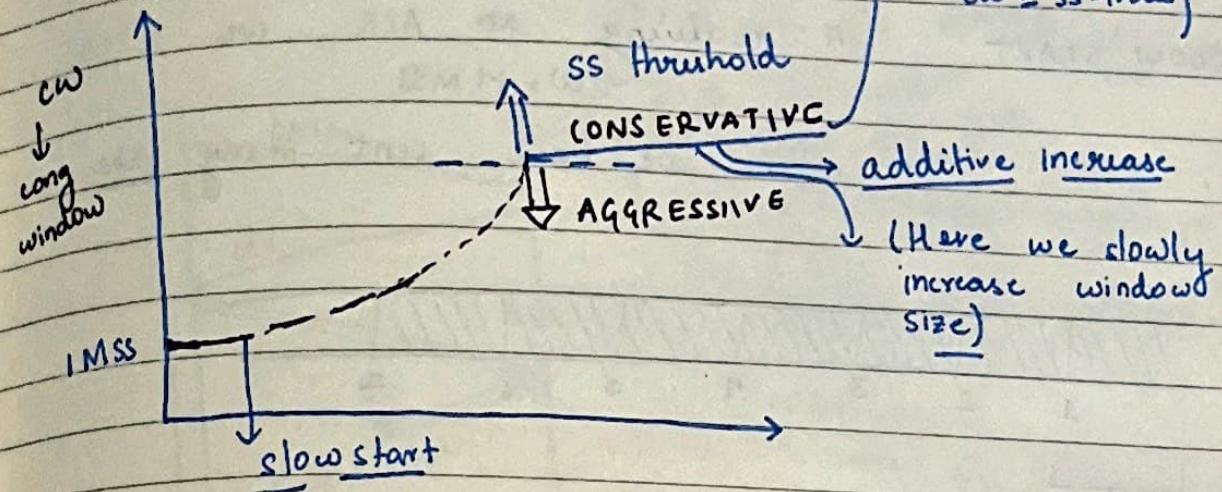
Before pkt is put in queue it modifies the TCP header bits, does not completely participate (some what of a hack)

ISSUE: Range of band width present (BW) is High
 $10 \text{ Kbps} - 10 \text{ Gbps}$

How should TCP figure out at what rate to send?

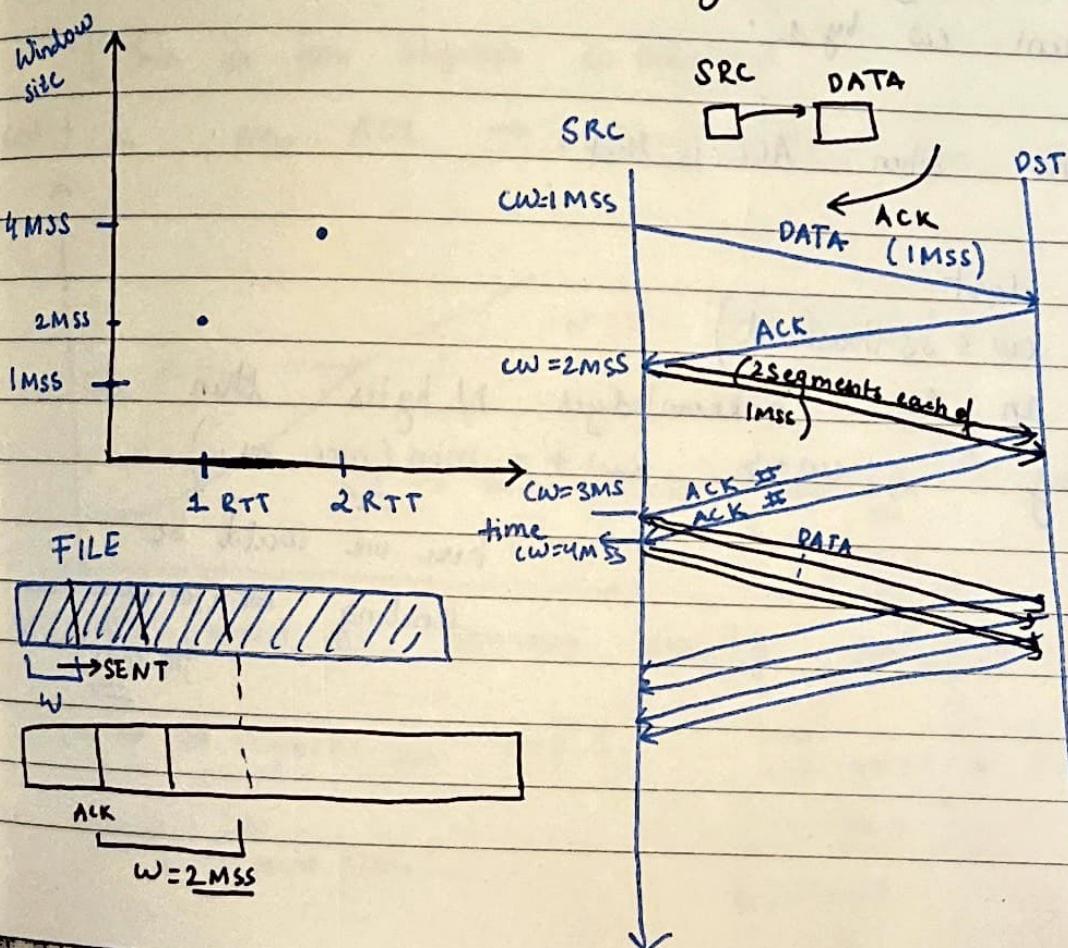
So we looked a slow start, we look CW and start it with smallest size and rapidly ↑ CW size till we hit danger mark, we can learn this overtime!

this phase is called
CA (congestion avoidance)
 $cw \geq ss\text{-thres}$)

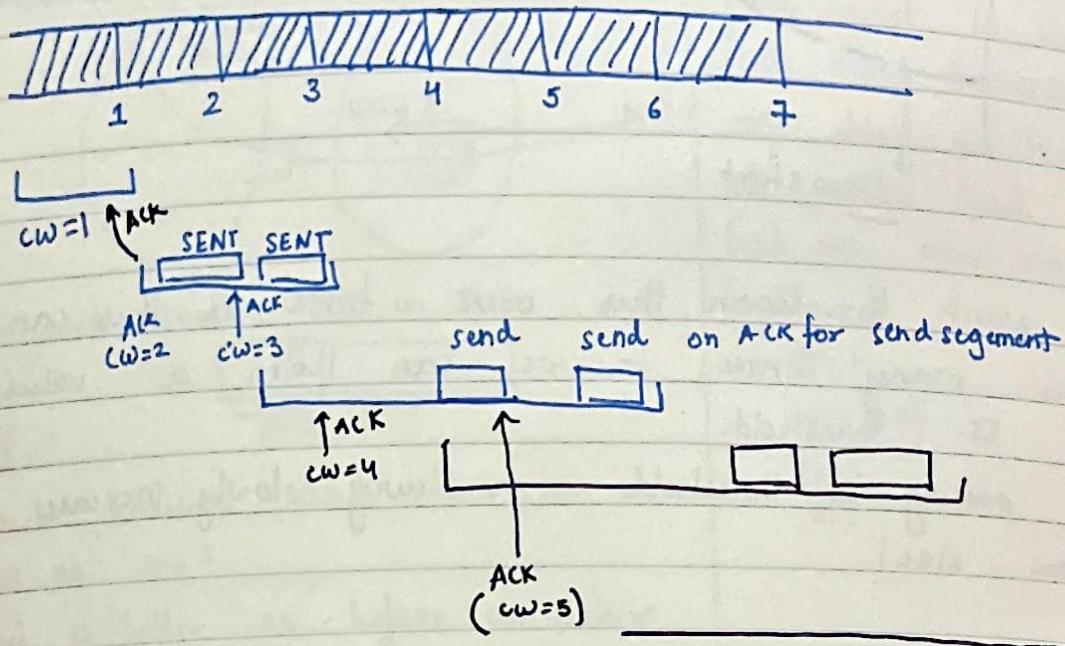


We want to learn this over time, as this can happen many times so we can learn a value for ss threshold
After passing ss threshold we very slowly increase window size!

How to practically increase cw !
SLOW-STRAT Double cw every RTT.



(IN SLOW START) : On receiving an ACK we
 $CW = +1 \text{ mss}$
 when we have sent many pkts



So everytime I get ACK I shift a window by eight and increment cw by 1.

What happens when ACK is lost?

[RFC 5681]

In slow start

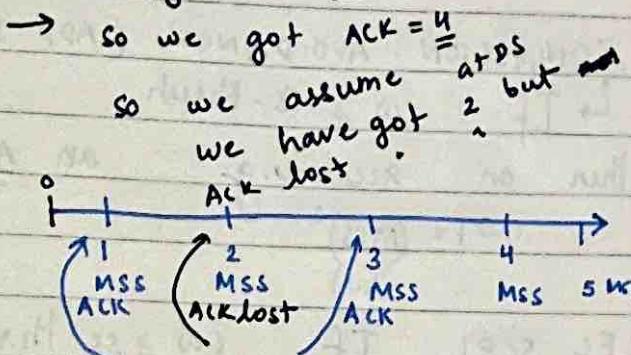
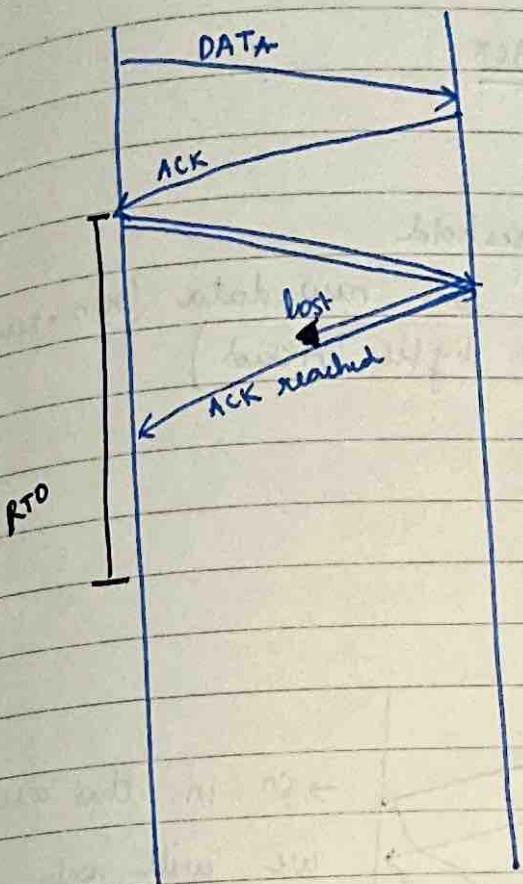
(ie $CW \leq ss\text{-threshold}$)

then if an ACK acknowledges N bytes then cw increases by $CW \Rightarrow CW + = \min(mss, N)$

here we could be sending less (for cases when file ends)

→ We retransmit if RTO expired and we have not received ACK for this or higher number!

CASE 1



→ This acknowledges for 2 bytes from 1-2 & 3 but for second lost ACK

→ (so we cannot increase by 2)

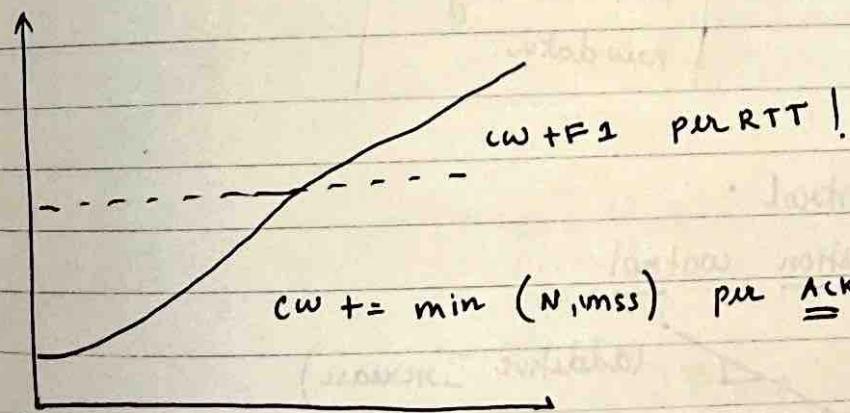
→ we should only increase by 1

Now window size = 3

and 3-4, 4-5, 5-6 sent

↓ This is for beyond ss. threshold!

$CW + 2 \text{ per ACK} \rightarrow CW = 1 \frac{\text{MSS}}{\text{per RTT}}$



So we want to increase CW by per RTT

($W \rightarrow \# \text{ segment per RTT}$)

so $\frac{Cw}{1 \text{ segment size}}$

I want to have

$$\frac{n}{n} = \frac{1 \text{ MSS}}{\frac{(MSS)^2}{(Cw)}} = \frac{(MSS)^2}{(Cw)}$$

CONGESTION AVOIDANCE (ADD. INC)

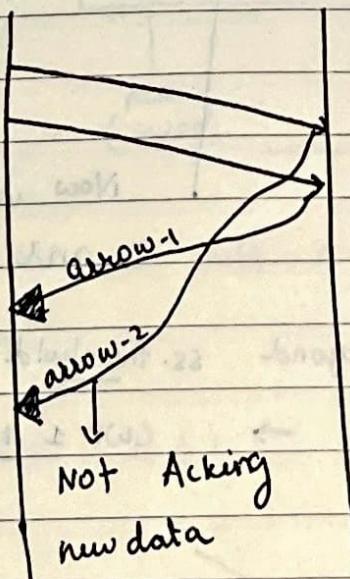
↳ If $cw \geq ss\text{-thresh}$
 then on receiving an ACK
 $cw+ = \frac{(mss)^2}{(cw)}$

RFC 5681 : If $cw \geq ss$ threshold
 on getting ACK for new data (non-zero
 number of new bytes Acked)

$$cw+ = \frac{(mss)^2}{(cw)}$$



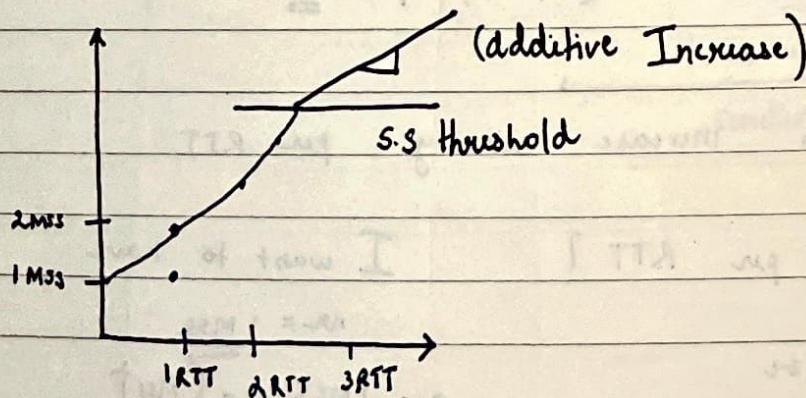
This case



→ so in this case
 we will only
 increase for arrow-1
 not arrow-2

TCP Congestion Control

slow start, congestion control



slow start

$CW = 1 \text{ MSS per ACK}$

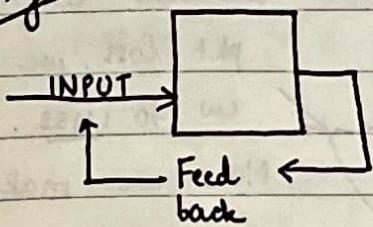
for congestion avoidance

$$CW += \frac{(\text{MSS})^2}{CW} \text{ per ACK}$$

How to detect and react to congestion?

STABILITY OF NETWORK

No prolong condition



ACK ~ Feedback } This helps
input } us control
(data) } congestion

How to deal
with feedback
=

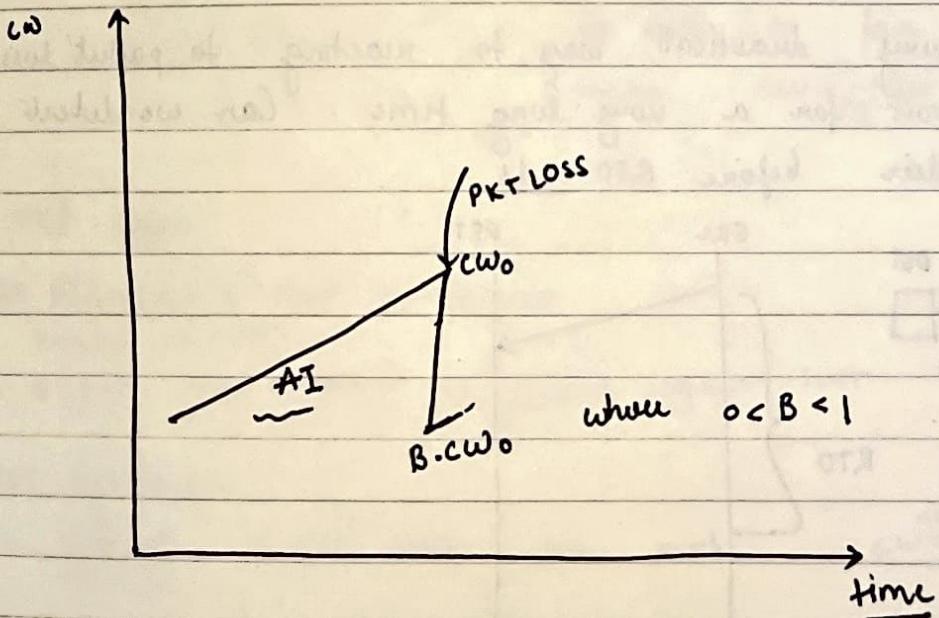
AI-MD:

Additive increase - multiplicative Decrease

↓
increase
conservatively

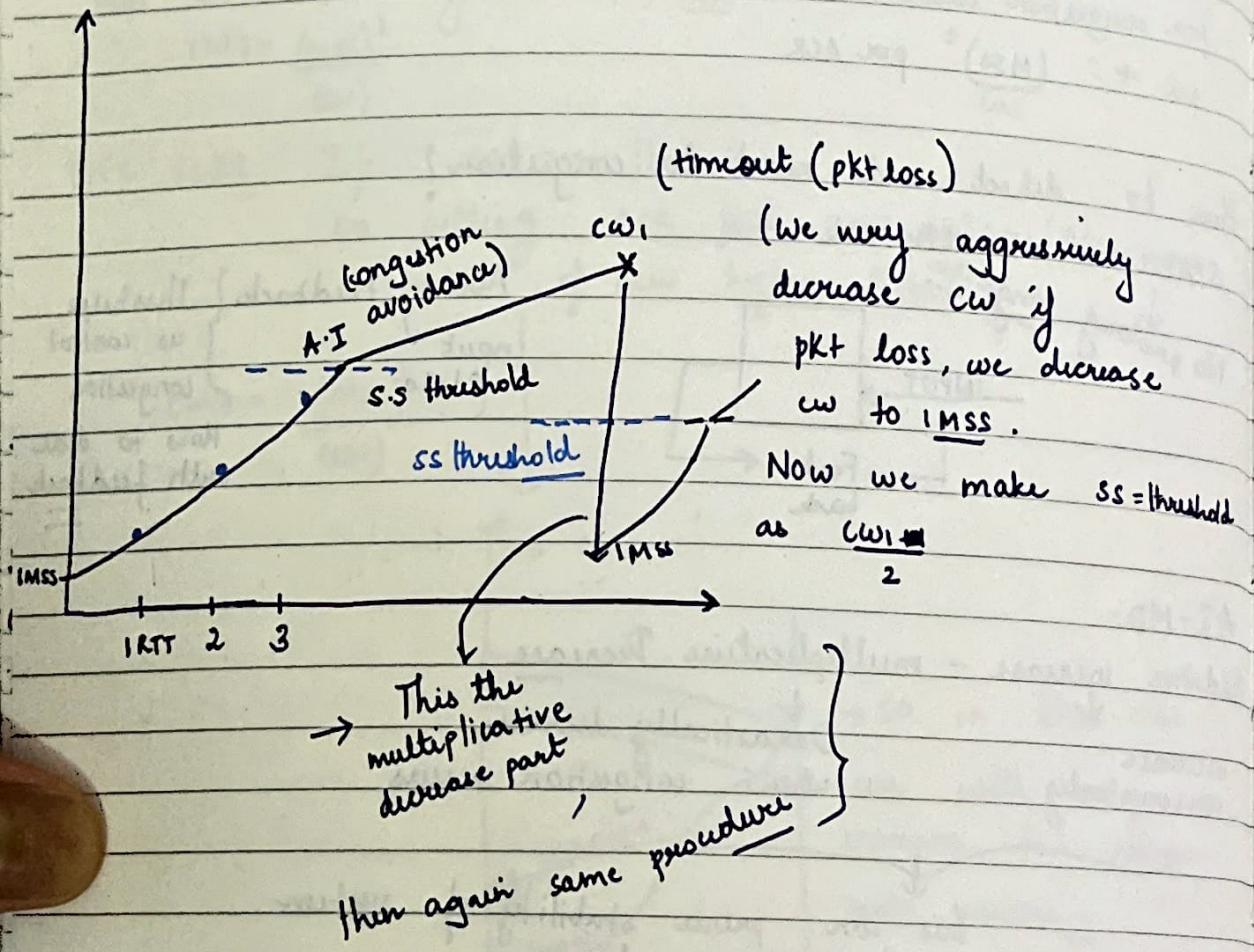
↓
Darastically decrease
when congestion occurs

This can prove stability of system
under some conditions.



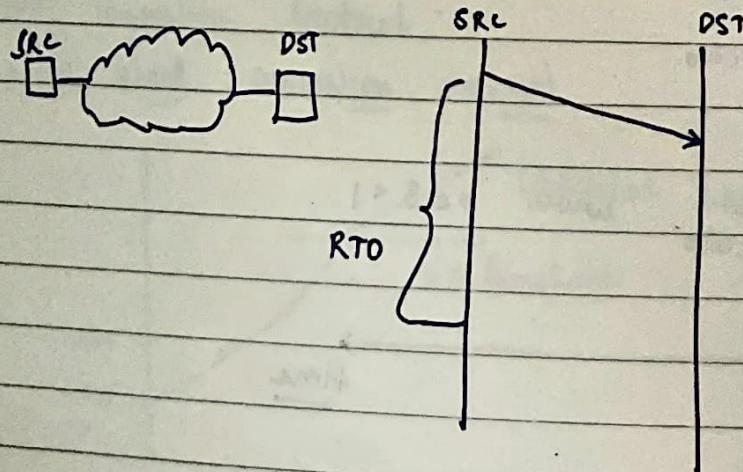
TCP TAHOE (CALIFORNIA LAKE)

→ One of early versions:-



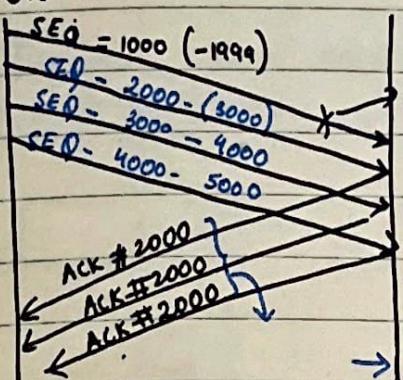
TCP RENO

Timeout is a very drastic way to reacting to packet loss. We have to wait for a very long time. Can we detect packet loss earlier before RTO ends.



SRC

DST



lost

Does getting multiple acknowledge, it means pkt loss?

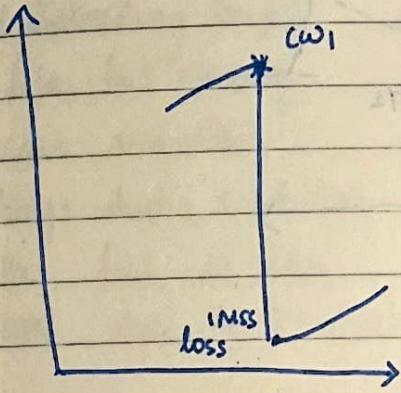
No, it could be packet reordering.

→ If we have lots of multiple ACK for same number the packet loss

(RULE)

→ 3 DUP ACKs, assume pkt loss

Should we reduce cw to (MSS for every pkt loss)?



let us differentiate between

T.O loss
 (Very drastic congestion)
 (No duplicate ACK getting back, no pkts getting through)

T.D (triple duplicate acknowledge)

→ This means some packets are getting through
 → CW not very drastic

→ So we should react differently

to both, so for TD ACKS we make $CW = \frac{CW'}{2} \rightarrow CW' = \underline{CW_{old}}$

TCP RENO

FAST RECOVERY & FAST RETRANSMIT

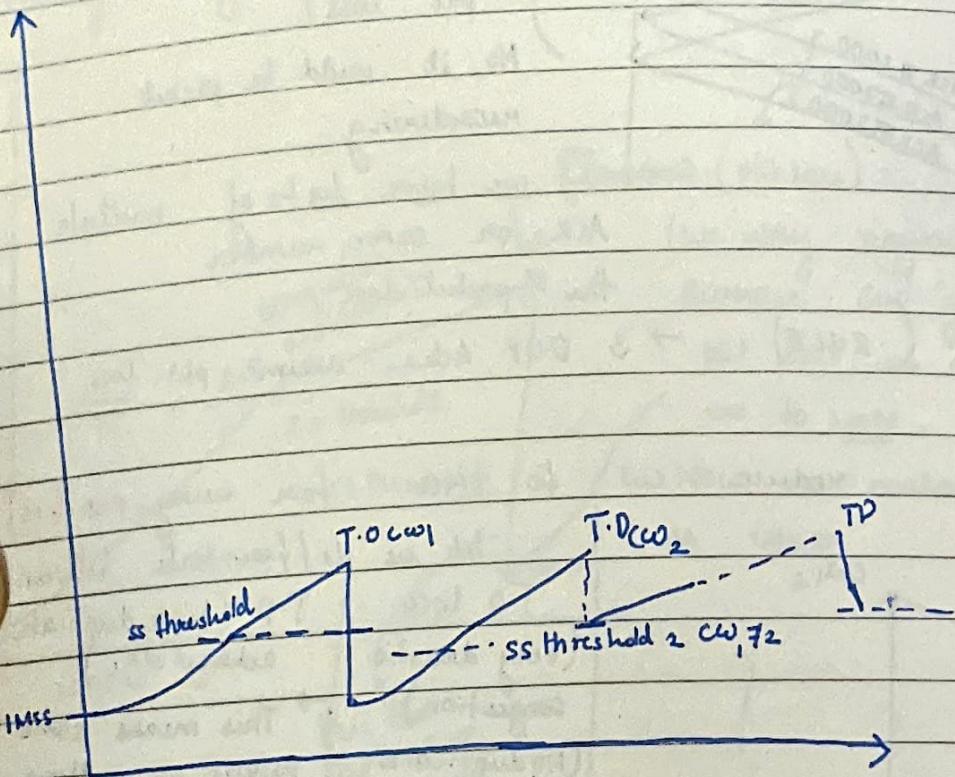
If 3 DUP ACKS received → assume segment lost × retransmit it

FAST RECOVERY

On getting 3 DUP ACK'S we make

$$CW^{\text{new}} = \frac{1}{2} CW^{\text{old}}$$

TCP RENO



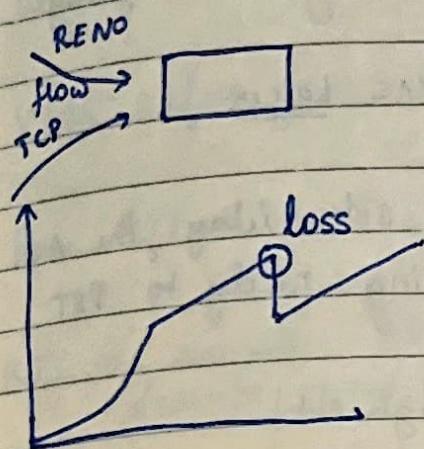
TCP VEGAS

↳ RTT (and vary much by this)

TCP RENO, TAHOE

→ triple duplicate
T.D → $CW = \frac{CW}{2}$
→ Fast retransmit
lost segment
→ Fast retransmit

ISSUE: Designed to cause congestion



→ We are actually causing congestion by increasing window.

So they are designed to create loss. The best thing is they use available bandwidth.

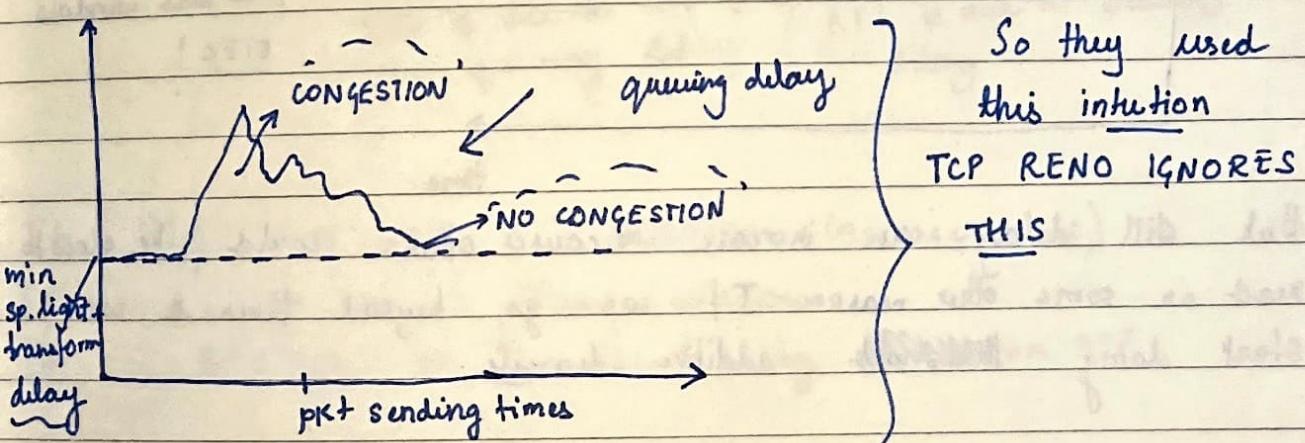
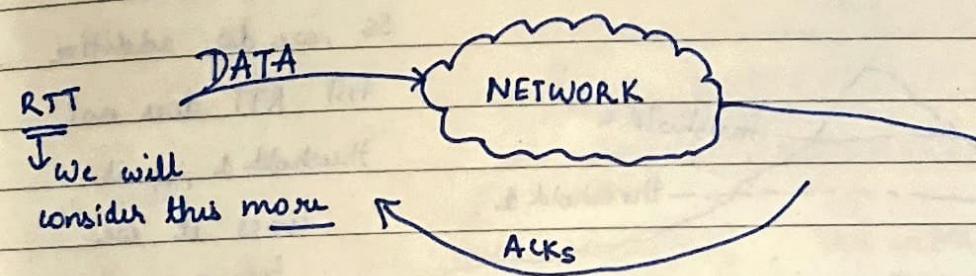
So they actually ~~try~~ fill up the queue, so if we are sending UDP packets
↓ what could happen?

→ Due to queues getting filled up UDP or loss sensitive can suffer.

Also we get increase in JITTER
Also overall $OWD \times RTT$ increase.

Only good for downloading large files, but playing video and other things which needs less jitter fail in some cases due to TCP RENO.

TCP VEGAS



(So looking at this graph we can take early action!
So is this intra delay due to Queuing while going
to DST?

- ① This delay could be due to MAC Layer!
(some protocol failure)
- ② Since this delay is due to both sides delay, the ACK
could be reason. The data could be going nicely to DST
but ACK's could be queued.
So why should we reduce SRC throughput?

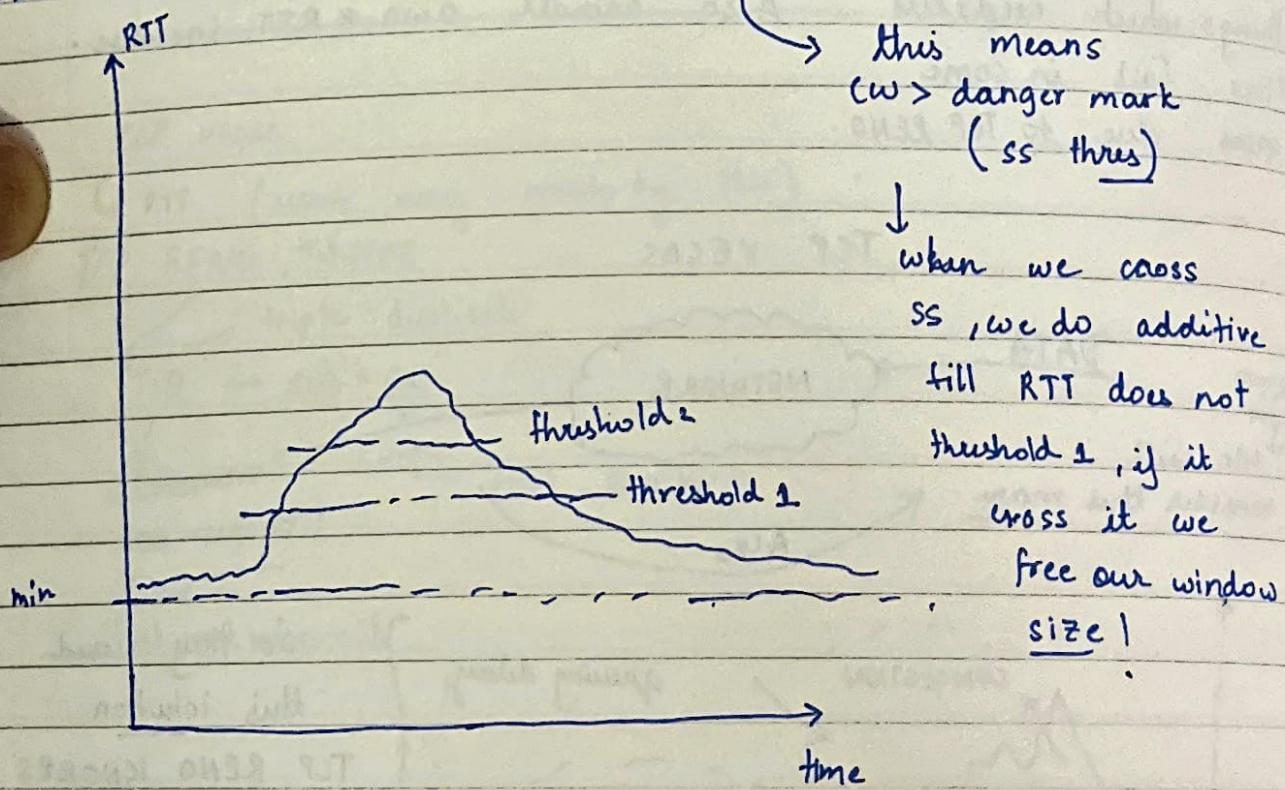
TCP VEGAS

S.S \rightarrow as in RENO

T.O \rightarrow " " "

T.D \rightarrow " " "

But we can change congestion avoidance



But still delay can increase because others could be slow to
react on some other reason. If we go beyond thresh 2 we
start doing ~~multiple~~ additive decrease.

→ The problem is thus-1 (α) & thus-2 (β)
 what values should we choose?

VEGAS RULES (FOR CONGESTION AVOIDANCE)

Base RTT — min observed RTT in some recent time window.

RTT → smooth estimate of the RTT
 ↗ current smooth

Suppose no congestion then

$$RTT = \frac{\text{Base RTT}}{w: \text{window}}$$

$$\text{Expected rate} = \frac{w}{\text{Base RTT}}$$

$$\text{Actual rate} = \frac{w}{RTT}$$

$$\text{Diff} = \text{Expected Rate} - \text{Actual Rate} = w \left[\frac{1}{\text{Base RTT}} - \frac{1}{RTT} \right]$$

$$= w \frac{(RTT - \text{Base RTT})}{(RTT) * (\text{Base RTT})} \geq 0$$

We are assuming this extra value in
 this is due to the same queuing delay } RTT is due to queuing delay

If $\text{Diff} < \alpha \Rightarrow$ A.I (additive increase in w_s)

If $\alpha \leq \text{Diff} < \beta \Rightarrow$ Freeze window

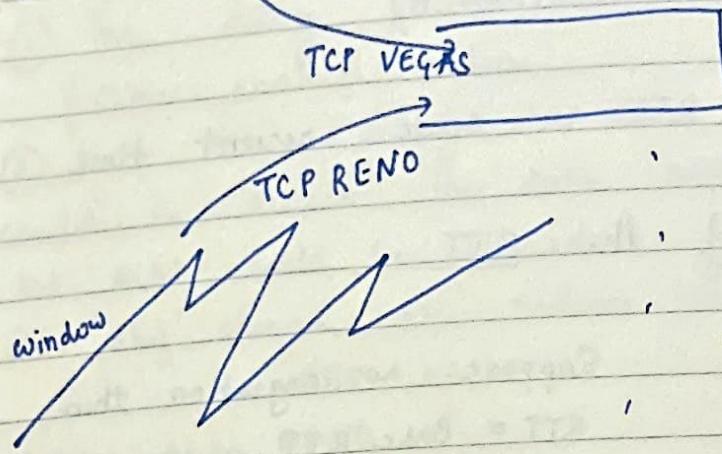
If $\beta \leq \text{Diff} \Rightarrow$ Decrease w by 1 MSS per RTT

Suggested

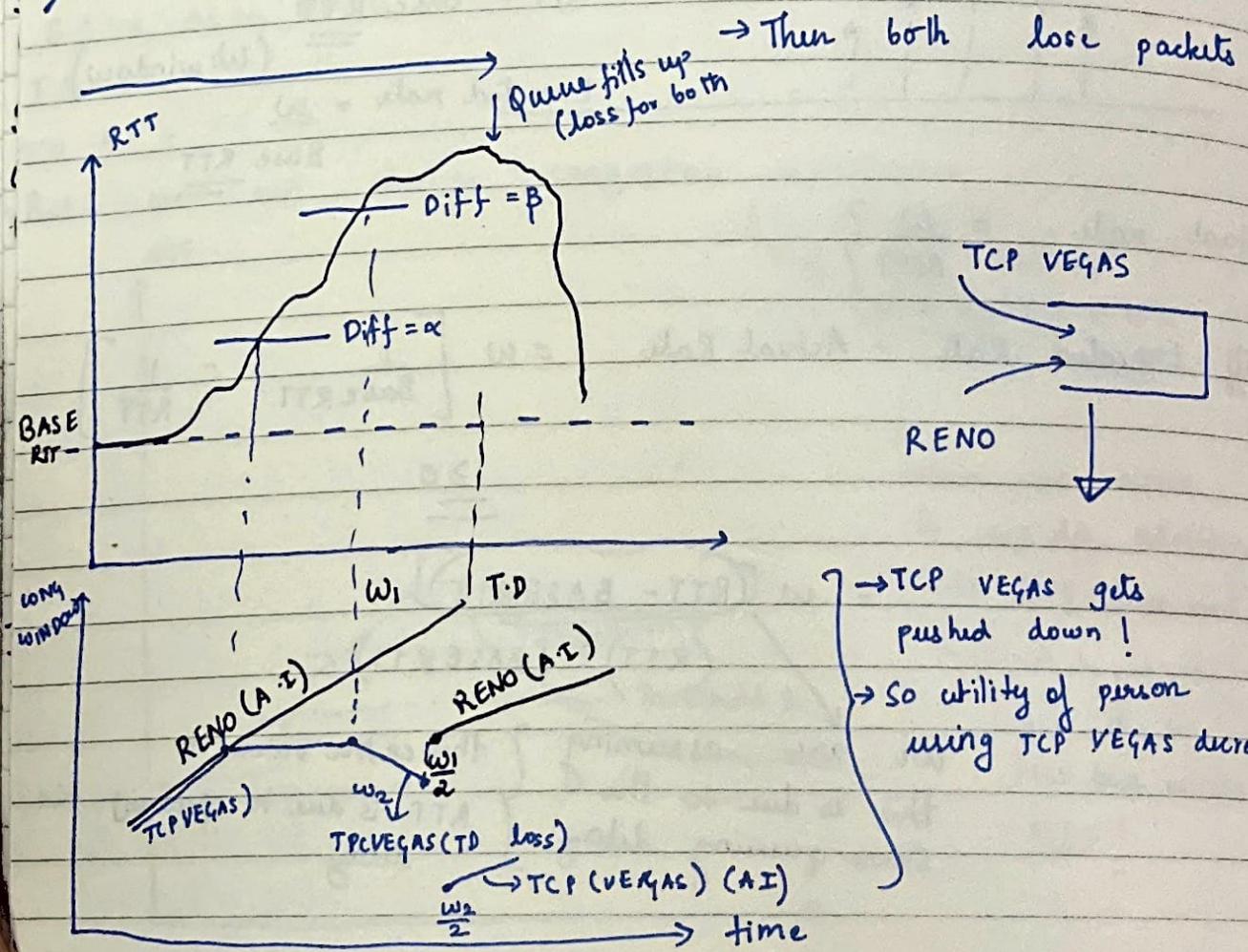
$$\underline{\alpha = 30 \text{ kbps}}$$

$$\underline{\beta = 60 \text{ kbps}}$$

CAN TCP VEGAS be compatible with other versions of TCP?



TCP starts of same way TCP RENO,
→ Queues started filling
→ TCP ~~RENO~~ ^{VEGAS} Freezes window size
→ TCP VEGAS Then starts reducing window

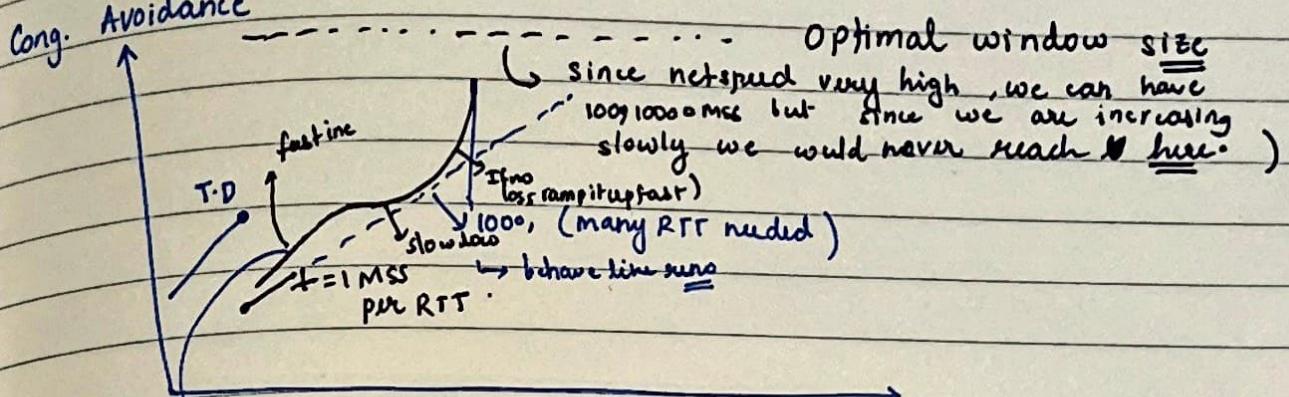


TCP RENO Default TCP for a very long time
 TCP RENO \rightarrow 1988 - 2012 (default)
 by many OSes!

TCP CUBIC \rightarrow (now default) (Linux / Mac OS X)

Why wasn't RENO good enough?
 Now we have very high speed networks

Cong. Avoidance

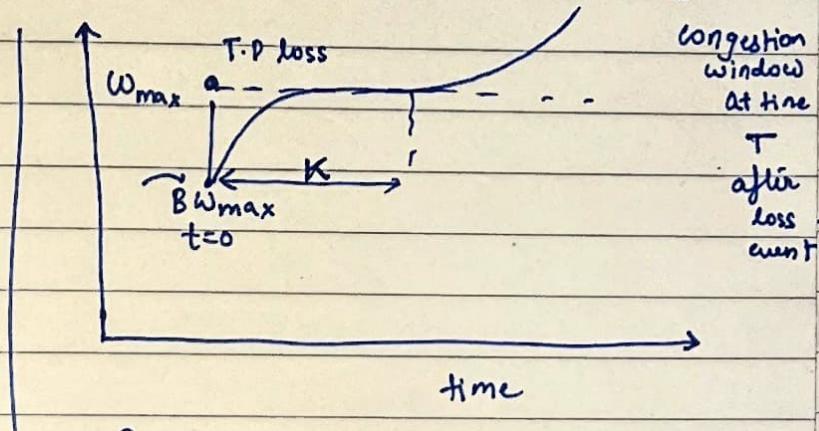


\downarrow This increase function looks like a cubic function
 So additive increase was very slow.

\hookrightarrow So what if others are using TCP Reno and we are being aggressive, we could kill TCP Reno.

So we want fast inc, but we want slow down (without killing TCP Reno).

$$x^3 = y$$



$$B = \underline{0.7}$$

$$cwnd = c(T-K)^3 + w_{max} \quad 113$$

$$\text{let } c=0.4, K = \left[\frac{w_{max}(1-\beta)}{c} \right]$$

$$\text{at } T=0 \quad cwn = -w_{max}(1-\beta) + w_{max} = w_{max}\beta$$

$$\text{at } T=K \quad \text{window size} = \underline{w_{max}}$$

\rightarrow K does not depend on RTT

\rightarrow

Application Layer

DATE

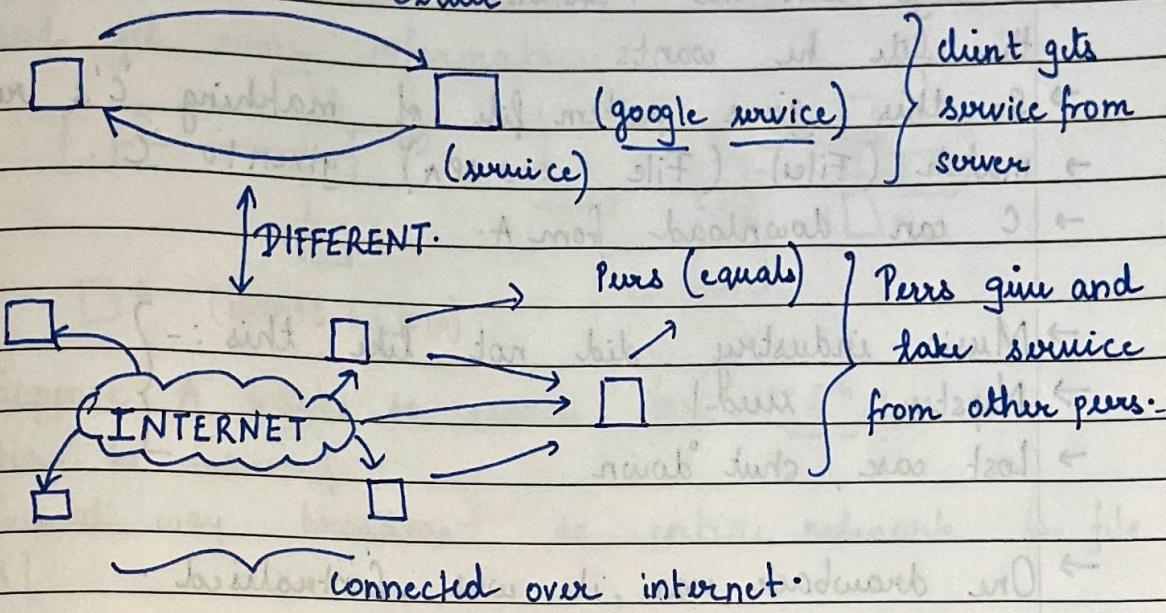
| | | | | | |
|--|--|--|--|--|--|
| | | | | | |
|--|--|--|--|--|--|

P2P networks.

Peer to Peer Networks.

Client

Server



- Bootstrapping

- 1990's, why don't we share digital content with others

- | | | |
|-------|-------|------|
| peer | movie | peer |
| music | music | peer |

 } file: Michael Jackson's Music
hosts

NAPSTER} (originated peer to peer network)

(anyone can join this network)
(we could connect to servers)

Peers



NAPSTER|

SERVER



- A logs in
- " gives detail of the files he wants to share.

| | |
|--------|------|
| File 1 | IP A |
| File 2 | IP B |

- C logs into the Napster server and searches for the file he wants.
- So this gives pattern file of matching 'C's request.
- match (File) (File 1 ↔ IP A) given to 'C'.
- C can download from A.

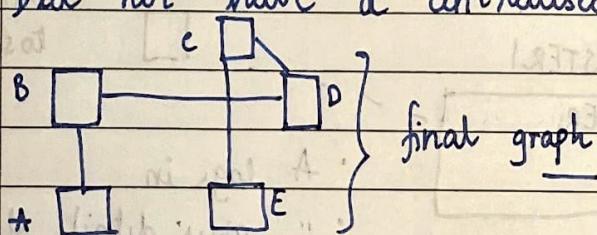
- Music industry did not like this :- }
- Napster sued }
- Lost case, shut down }

- One drawback was, it was Centralised.
- It makes an easier target.
- Single point of failure.
- (not fault tolerant)
- Somebody brought this down, 'ded'.

- We need a decentralised server.
- legally easy to take down.

GNUTELLA (Next generation)

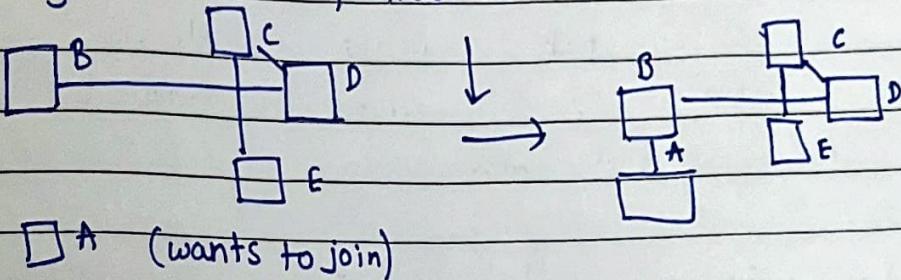
- Did not have a centralised system.



So how this network is created!

Bootstrapping

- (1) Application may have IP's of some other peers
- (2) So we have some website which we can contact and get some information.



Suppose A wants to search for "f"!

→ How to search?

- Trivial way - broadcast to entire network for file 'f' b/w A?
- Over kill, we do a limited broadcast
- LIMITED BROADCAST ($TTL = n$, TTL -- (at every hop), $TTL == 0 \Rightarrow$ don't fwd query).

Suppose $n = 2$

$\left\{ \begin{array}{l} A - B - D \\ TTL=2 \quad TTL=1 \quad TTL=0 \end{array} \right\} \rightarrow$ So let D have the file.
 \rightarrow So how do they transfer?

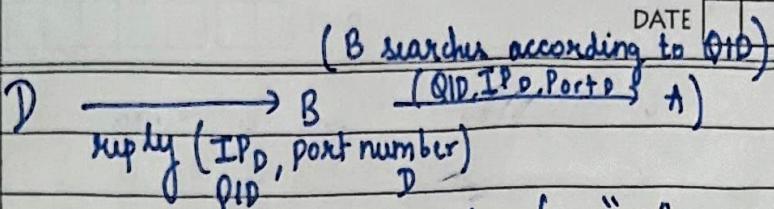
First Meth : QUERY had "f" , QID (unique ID of query)

A $\xrightarrow{(f, QID)}$ B $\xrightarrow{(f, QID)}$ D (file NAME)

(So D has some info about "f")

B caches (QID, A) (we save "QID" and person from whom we heard request from).

D " (QID, B).



→ A can now request info from "D".

- Now B can store some information about "f".
- So B caches reply. (QID, "f", IP_D, Port_D).
- So now G joins and asks for "f", B directly replies to G.
- If we do not file we can search again by increasing ttl.
- We keep increasing ttl till max value reached.
- So in future GNUTELLA put A's info in query.
- So A & D directly communicate. (But now B cannot cache all this information). (The reply information)

→ Second Method:-

Query has IP_A, so directly send to A.

- So the QID helps to prevent multiple forwarding of a query!
- So each person makes sure to forward to all neighbours only once.

$$\begin{array}{l} \text{V.0.4} \\ \text{V.0.6} \end{array} \left\{ \begin{array}{l} \text{max TTL = 7} \\ \text{max TTL = 4} \end{array} \right\} \left\{ \begin{array}{l} \text{Why do we dec max TTL} \\ \text{allowed?} \end{array} \right\}$$

When we have large networks!
 (Graphs become dense)

TTL = 4 could make Query reach $\sim 10^4$ nodes
TTL = 7 " " $\sim 10^7$ "

→ So we decreased TTL overtime as GNUTELLA became popular. Overkill hai!

- Good → No centralised node!
- We have many broadcasts! (Bad!!)
- (Is there a clever way to avoid these broadcasts?)
- (Where info is stored?)
- Suppose there are N-nodes!
- can win by contacting only $\log N$ nodes
- limited broadcast does not work.
- How do we store info.

Suppose we have identifier of nodes, let's use IP addresses.

IP_A f₁ Files } suppose we could
IP_B f₂ map these file
⋮ ⋮ name to IP address.
IP_E ⋮ } (And we store this
IP_F ⋮ info somewhere else).

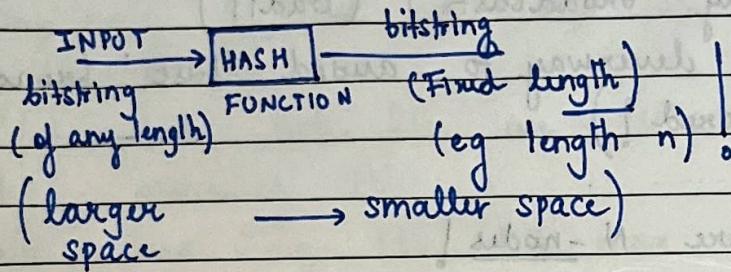
IP_G - has f₁ } Can A in some clever fashion ask
IP_S (has f₂) { Q where is f₁ or some other node
→ So Q tells E I have "f₁" }
So A asks E " " " "

So all nodes having "f₁" tell E. So A gets info and downloads this from "Q".

S stores it has "f₂" in "B". So if 'f' node wants to search for f₂ it contacts "B".

{ So f₂ mapped to B }
 { So f₁ .. " E }

One method is to use a hash function.



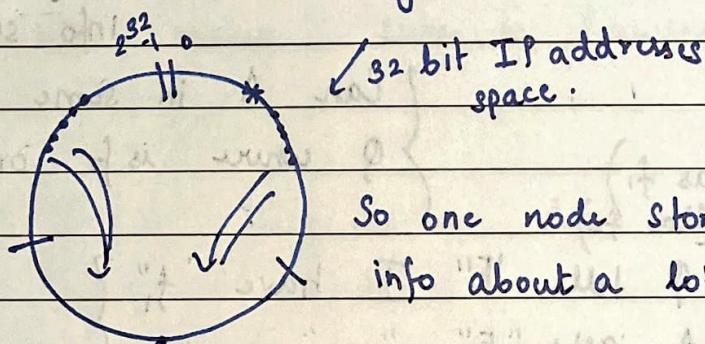
{ f₁ → H(f₁) } we can use
 { f₂ → H(f₂) } deterministic this is hash function

H(f) are in domain file of file to IP Address.

→ So we choose closest IP to H(f).

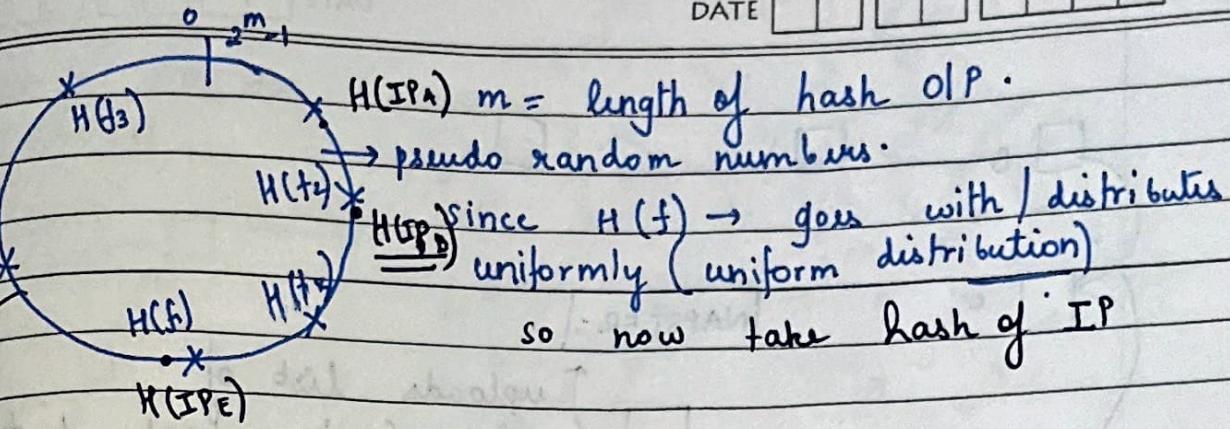
But the issue is the networks has no uniform distribution of IP addresses.

→ So some nodes do a lot of work.



So one node storing info about a lot of files.

So they suggested to hash of IP addresses



(so f_1 is closest to E node)

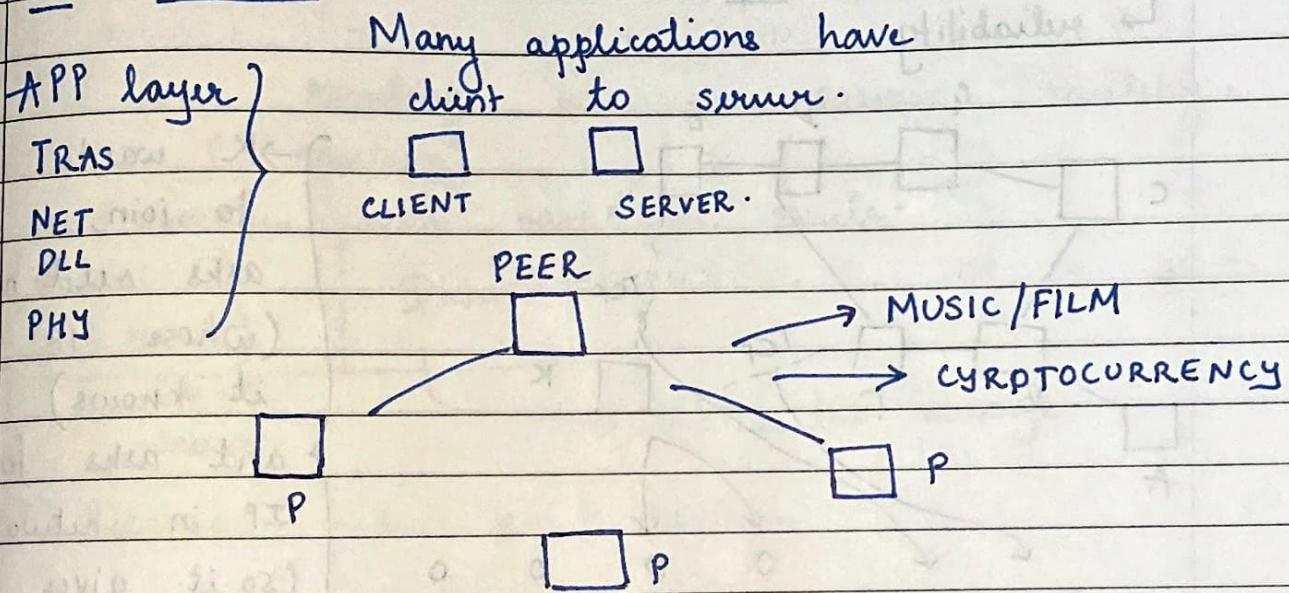
(so E must store information about where f_1 is stored.)

→ In this abstract space we are embedding information.

→ So we can zoom in and find where files are.

BIT TORRENT } many files can have many nodes.

P2P Networks :-

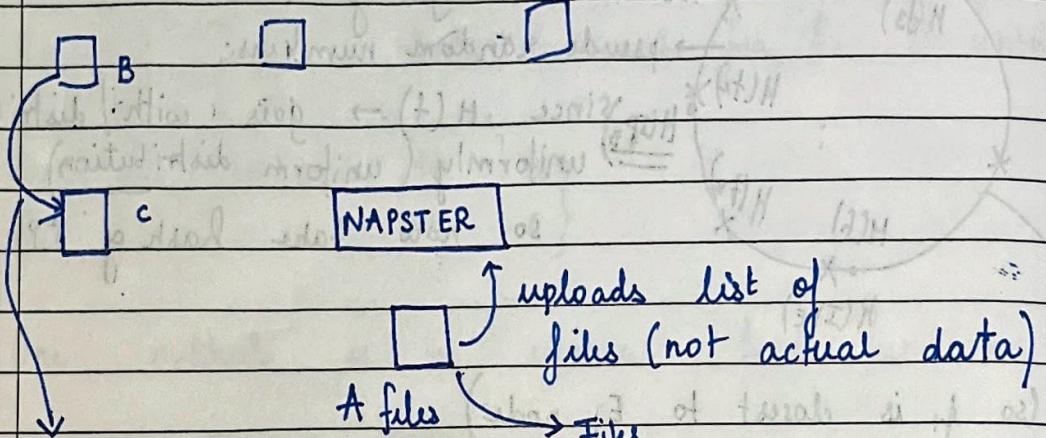


1ST GEN:-

→ NAPSTER } → Ads people do not want info to be always be on their laptop.

→ So a main server is useful.

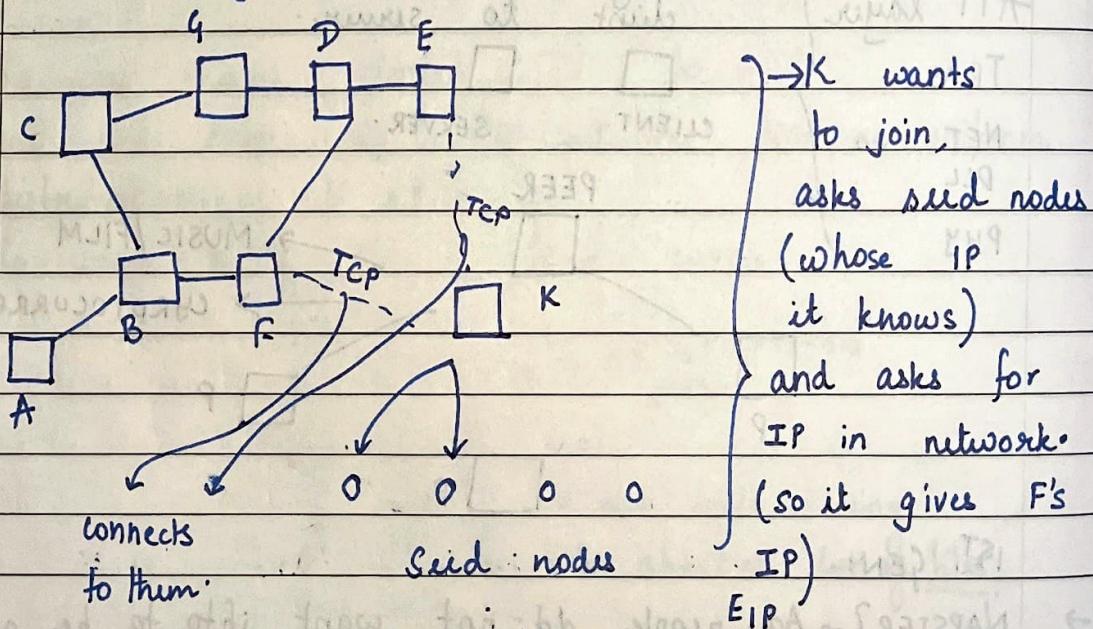
→ But centralised >



→ Suppose
C has a
file for B
and NAPSTER told B

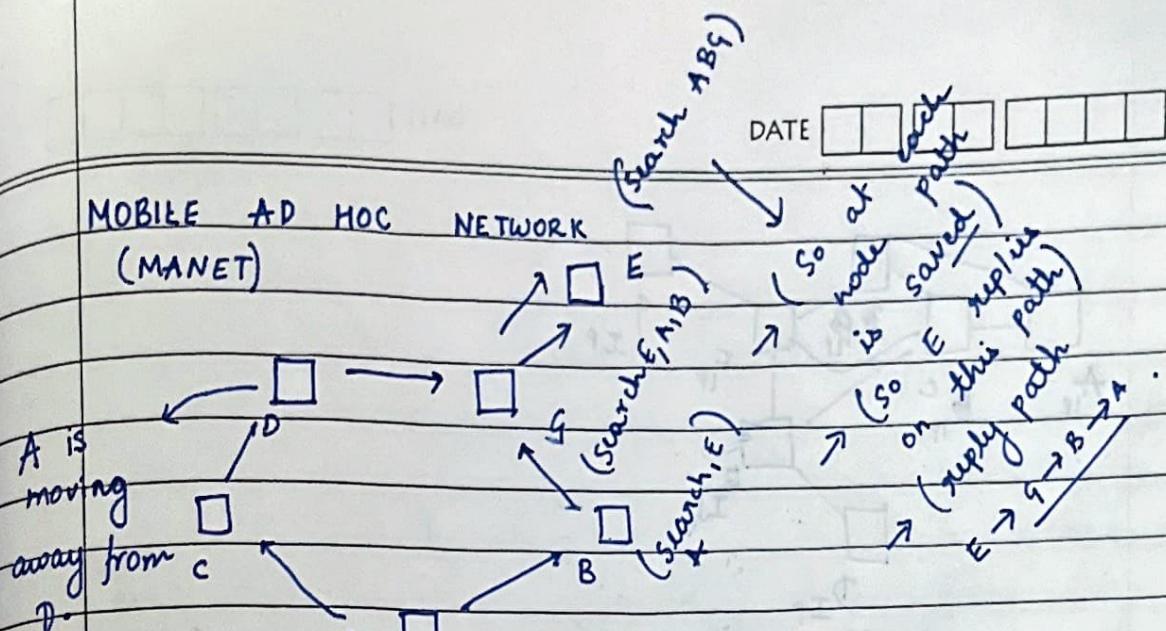
Now B & C connect

via TCP and share info
→ reliability



- As people start downloading files they make new friends.
- The network becomes more dense!

MOBILE AD HOC NETWORK (MANET)



tanks (moving around in a battle field | cellular networks are down).

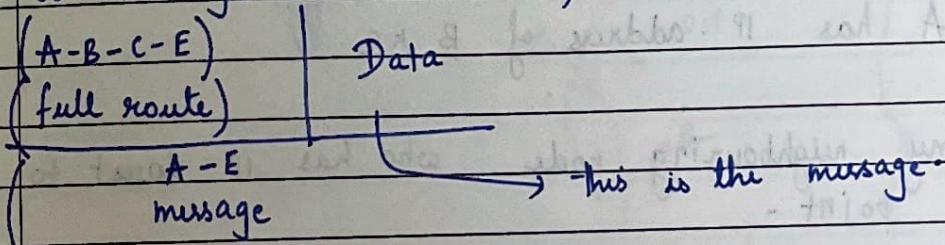
How to do routing.

(Searching for E)

(Suppose D links breaks down)

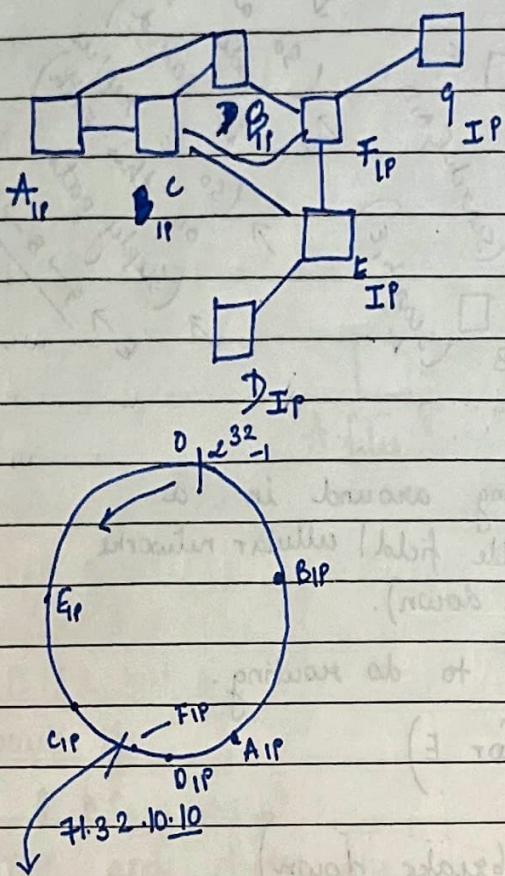
So for routing we use source routing.
(Not IP routing)

→ So message will have full route.



Because nodes are moving.

3RD GEN P2P : Distributed HASH TABLES
(DHT)



So A wants to know who has IP closest to point $71.32.10.10$

This a subproblem for P2P networks efficiency.

→ So A has IP address of B & C.

→ In my neighbouring node who has IP closest to this point -

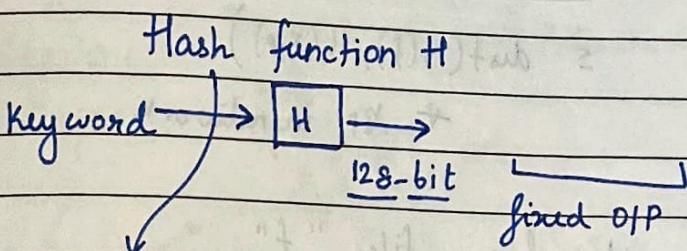
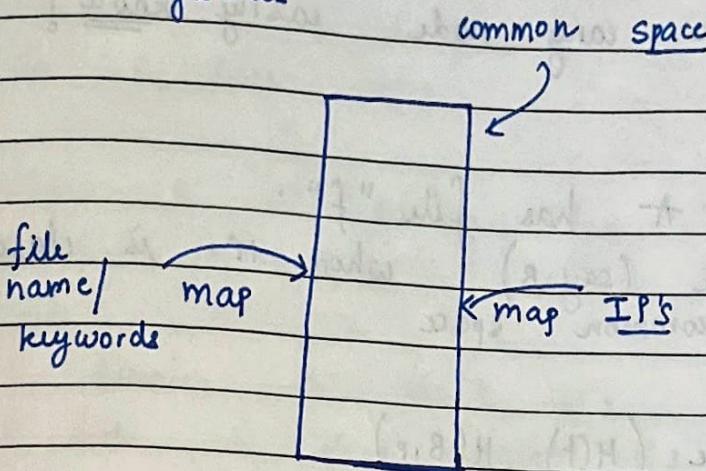
→ So C gets my IP request

→ Now C checks in its neighbour who is closer,

→ So we send it E_IP but if we had F_IP as neighbour

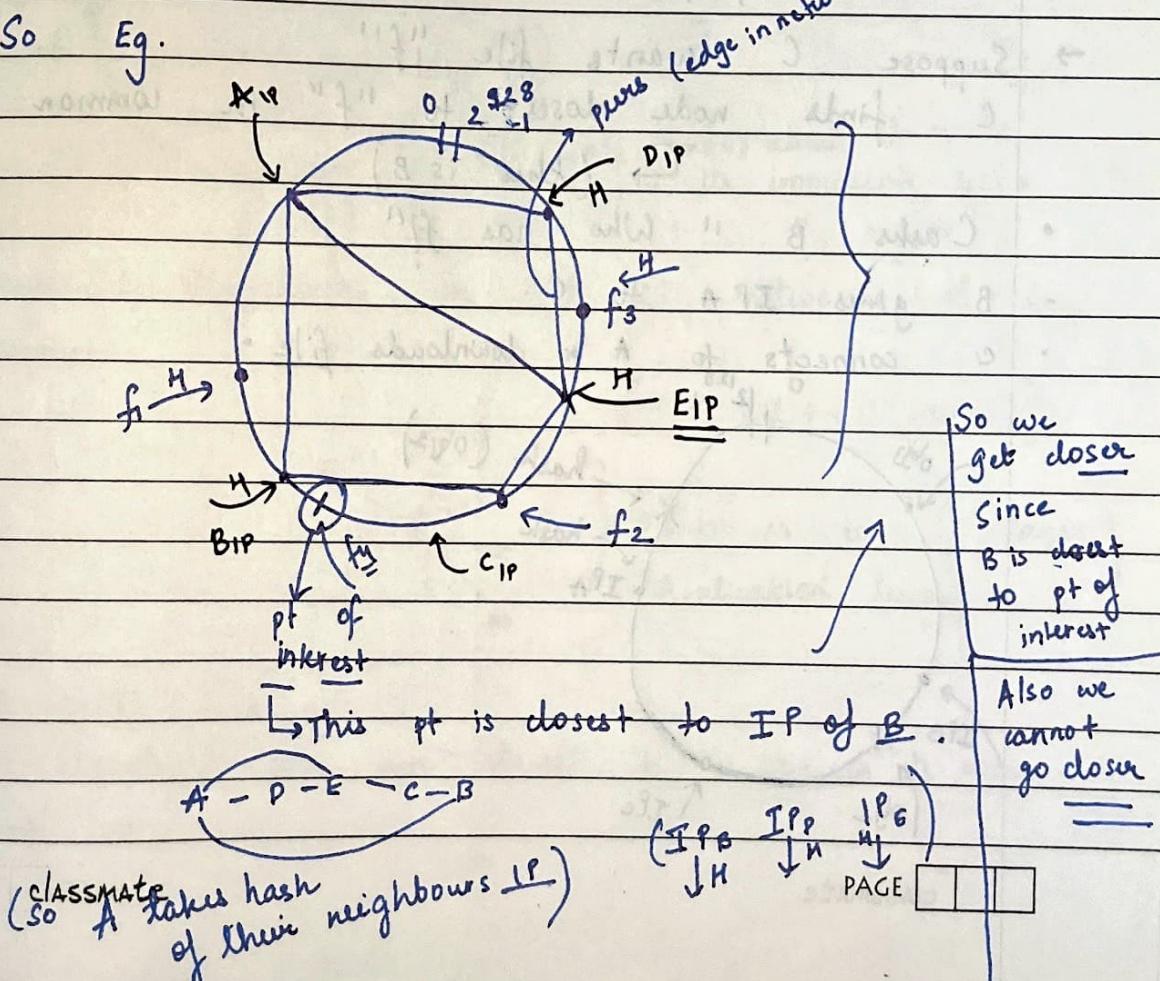
→ So we would send E_IP to C.

IDEA: file names $\xrightarrow{\text{map}}$ IP
keywords



→ for even a single bit change the O/P is very random

So Eg.



So we want to who has f, so we give information:-
we can reach any node easily know!

→ HIGH LEVEL:-

Suppose node A has file "f".

A finds node (say B) where IP is closest to f in common space

Means distance ($H(f)$, $H(B_{IP})$)

$$= \text{dist}(H(f), H(X_{IP}))$$

+ XIP in network

→ A tells B it has file "f"

A tells B it has file "f"

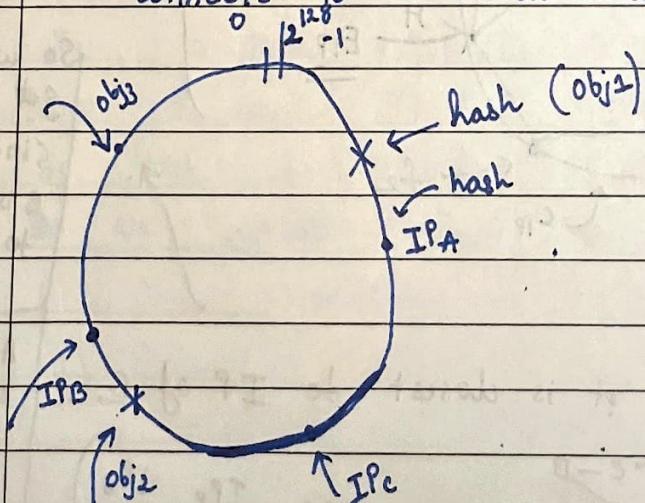
| File | IP |
|------|-----|
| "f" | IPA |

→ Suppose C wants file "f"

C finds node closest to "f" in common space.

→ (this is B)

- Asks B "Who has f?"
- B gives IPA to C.
- C connects to A & downloads file.



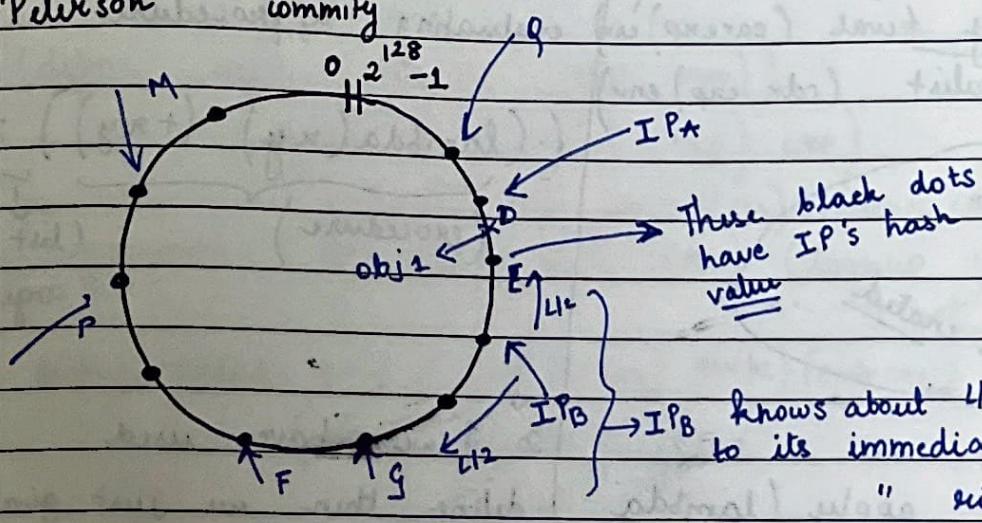
Node C has obj₁

- Now hash of obj₂ is closest to hash of value IP_A it has.
- So C tells obj₁

→ Node B is looking for obj₁.

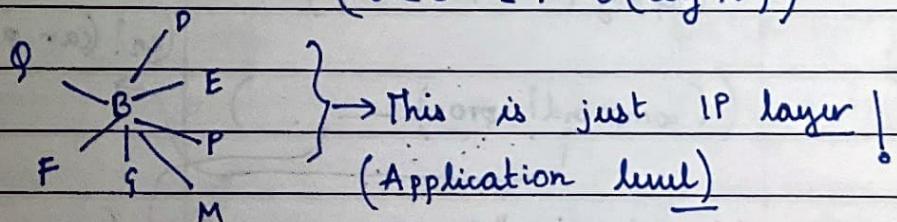
- B finds A, which is closest to obj₁ in common space.
- B queries A about obj₁ | • B downloads
- B tells A that C has obj₁ | obj₁ from C.

PASTRY - Dashed Advisor Peterson
Peterson This community



→ Leaf Set: Neighbours in P2P networks of any node

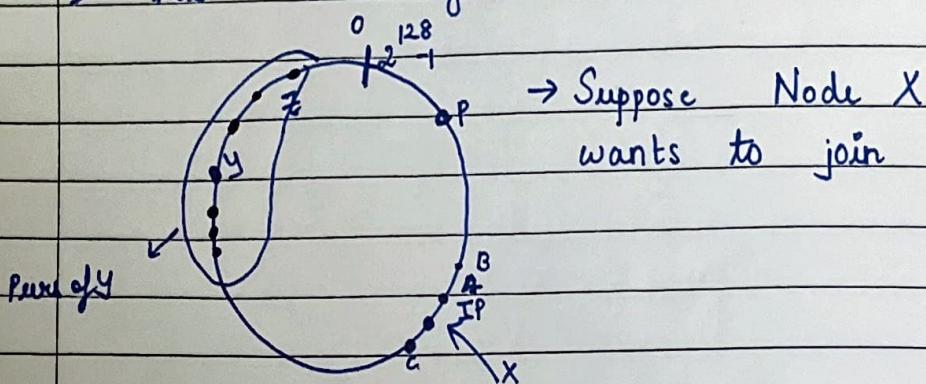
(size $L + O(\log N)$)



→ M wants obj₁

→ M should find peer closest to obj₁ in its own leaf set.

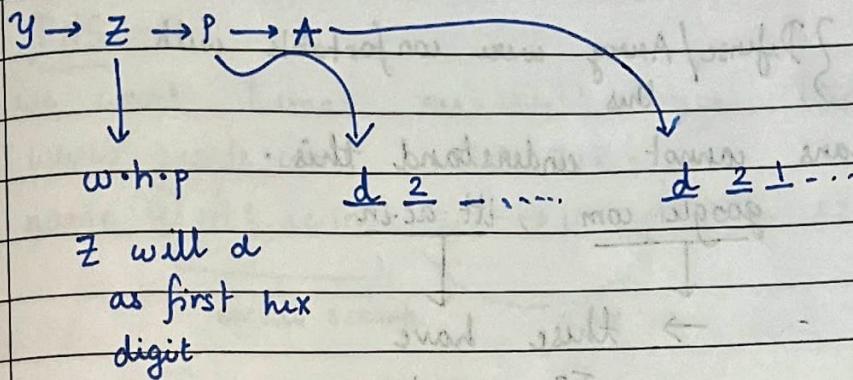
- So M will find B to be closest, it will forward message to B.
- So now B now D is closest in its leaf set.
- It forwards to D.
- D has no neighbour closer, so D tells M.



- So X does not know its next IP's ($\frac{1}{2}$ walk)
- So X asks Y to help it to join.
- So Y finds in its leaf set Z is closest
- Z finds P, P finds A.
- A finds itself to be closest.
- Now A can tell $\frac{1}{2}$ people on its right and $\frac{1}{2}$ people of left to X.

- X contacts some other node already in network (Y)
- "Y" routes message to node closest to hash (IP X).
- Each node on path (Y, Z, P, A) become part of leaf set of A.
- A carries info about $\frac{1}{2}$ on left & right
 - ↳ join leaf set of X.

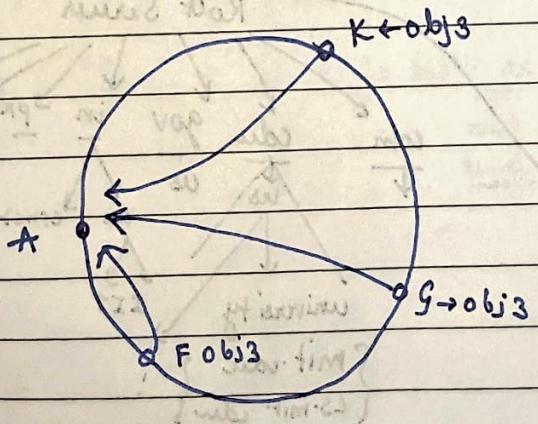
d z f e 3 7 \downarrow 128 bits \rightarrow No of digits here $O(n)$ ($O(\log n)$)
hex looking for



hash IP_z: d \dots
 (we are looking for longer prefix match)

- Now what if a node fails
- So this can be done by communication of periodic messages at application layer.
- So we can take care of this by asking other neighbours of their L & R.
- Also we store copy of some information at some nodes.

BIT-TORRENT



Now files are with multiple user, so we can download from all of them by allocating which chunks of data we want.

DNS, HTTP

→ DNS → domain name system.

→ IP { Defense / Army were comfortable with
32.75.5.9 this

→ But humans cannot understand this.

→ so for " google.com, itc.ac.in

→ these have

IP mapping

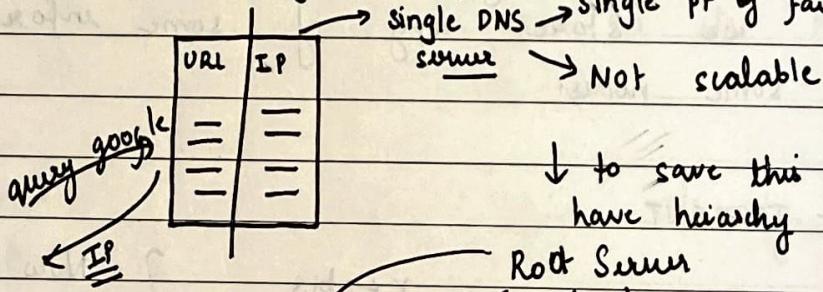
→ www.google.com - IP ?

→ So server should know who to send message for this.

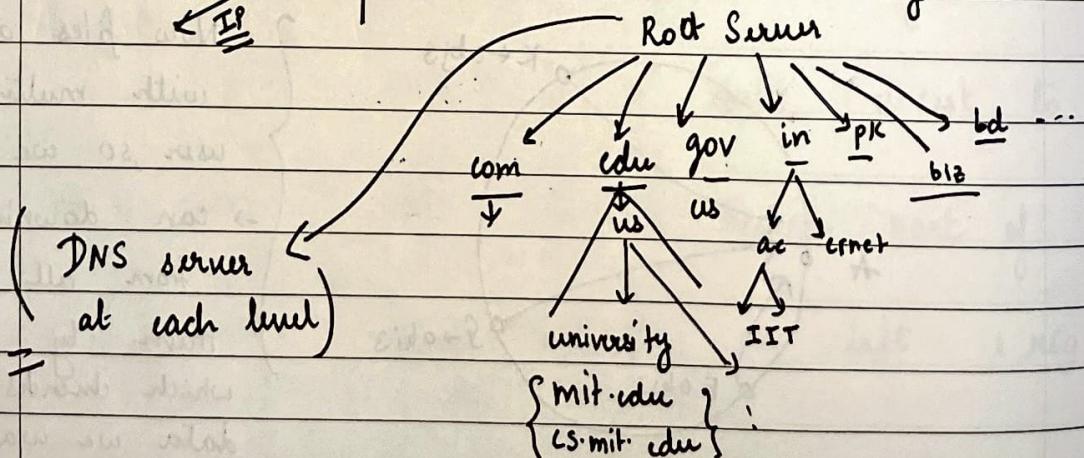
→ But network protocols do not understand this.

→ Also one more ad. any url can change its IP.

→ But if we have a single server it is not scalable, because if it fails then what to do.



↓ to save this we have hierarchy



→ so far IN → we have a DNS server and it should have IP address of DNS server of all lower in hierarchy at its node.

DNS :-

we want human readable strings like

www.google.com 8.32.7.9.

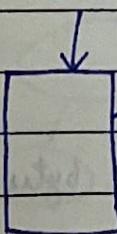
name @ itb.ac.in → IP address :- 32.23.5.22.

email server

DNS

DNS

takes us from name to IP address.



NAIVE

→ single pt of failure

← many people bombarding it.

→ Hierarchical Solution.

:- Root Servers

Root Server - Top Level

Top level domain

in

.com

.edu

.gov

.bd

ac

.cn

.nic

itb itt k

cs ee

eee

eee

CS. itb. ac.in

--- next top level

level

- So at Root Domain we will have DNS servers.
- and at node we have DNS servers.
- Also point to be noted :- DNS server of parent domain knows the DNS servers of immediate children. ↓
- name of DNS server, IP address of DNS server.

Sanya.cse.itb.ac.in

So I have to change its name, IP address corresponding to it by just telling its local friend.

ROOT SERVERS.

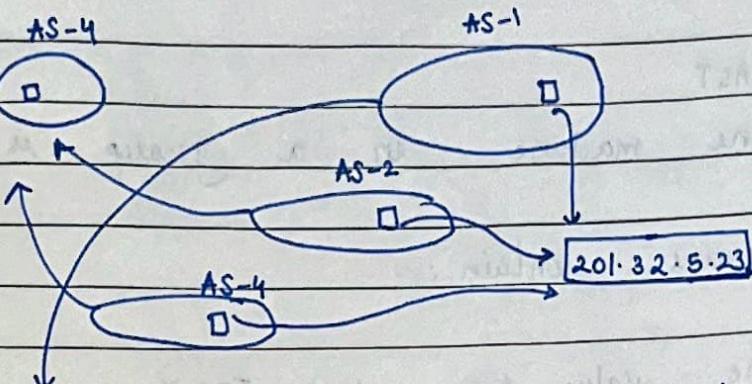
DNS pkt size : 512 bytes:-

13 ROOT servers → info fits in 512 bytes.

| Names | IP |
|---------------------|----|
| A. root-servers.net | |
| B. - | |
| . | |
| . | |
| . | |

A.root-server.net → multiple machines which correspond to a single name.
201.32.5.23

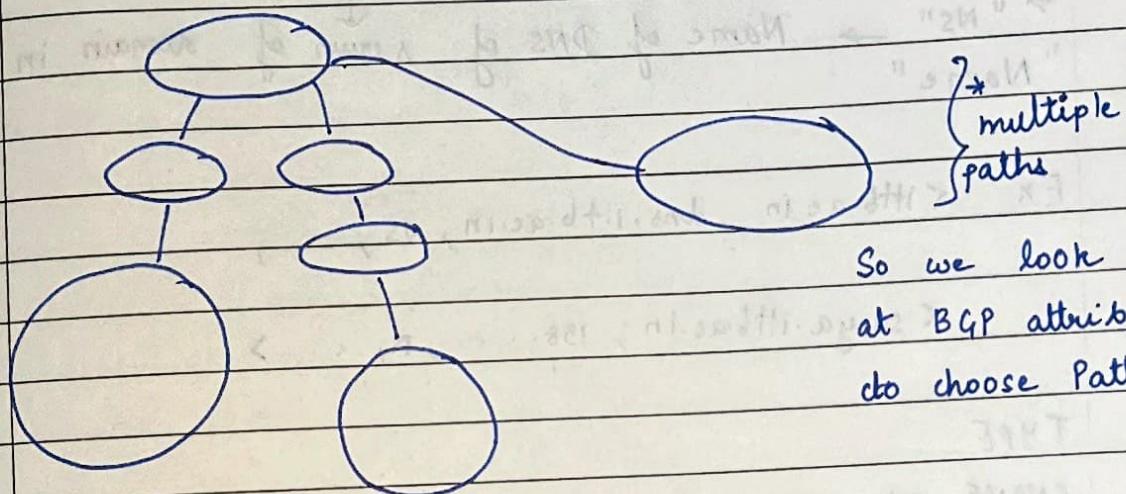
So I have many machines same IP address, so how does layer 3 react to it



→ BGP advertisements of this IP coming from multiple as path.

→ so it seems a server connected to many but diff machines.

→ so what happens in case of multiple AS.



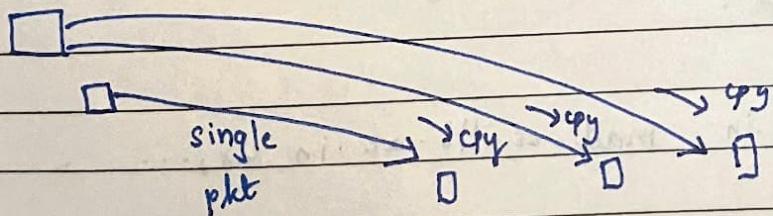
So we look
at BGP attributes
to choose Path.

UNICAST - Single destination.

BROADCAST - All nodes are destination.

Multicast → Subset of nodes in internet where they are destination.

Video Server



* one packet
is split into
many packets
at network.

ANY-CAST

- Any one machine in a group is destination.
- DNS servers contain :-

<name , value, type , class , TTL >

TYPE: "A"
→ value corresponds to a IP address.
→ "NS" → Name of DNS of server of domain in "Name"
(in: internet) ↓
eg * how long this record is valid.

Ex <itb.ac.in , dns.itb.ac.in , NS , ... >

< surya.itb.ac.in , 198... , A, ... >

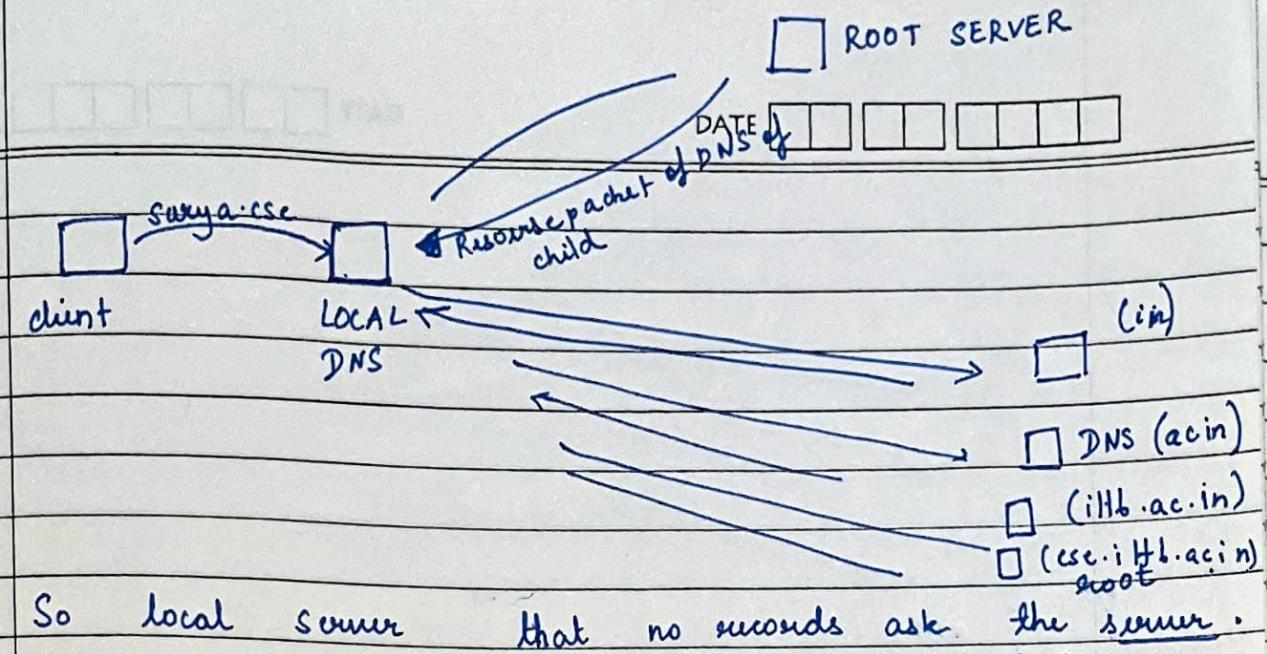
TYPE

CNAME → canonical name (alias)

cs.mit.edu
ce.mit.edu } → both have same IP

Mx → email server of domain specified in name field.

<cs.itb.ac.in , mail.cs.itb.ac.in , MX, ... >

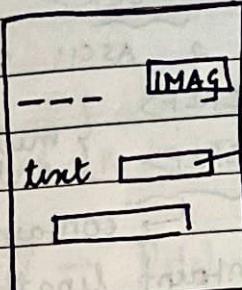


Application Layer :-

P2P, DNS :-

see detailed example in text book.
HTTP: Hyper - text transfer protocol.

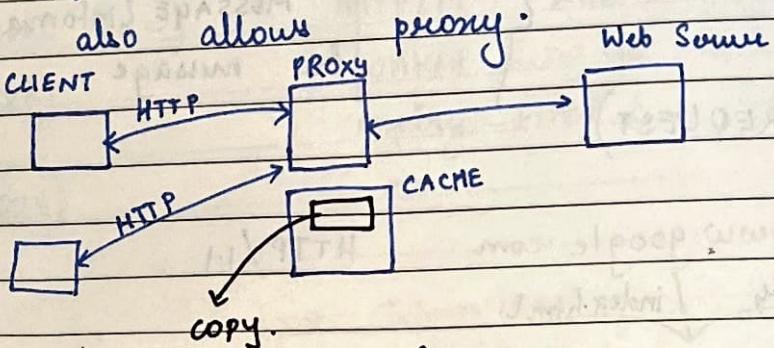
→ Help us download and submit information to webpages.



So from one page we can navigate to another web page.

So DOWNLOAD PAGES & SUBMIT INFORMATION:-

It also allows proxy.

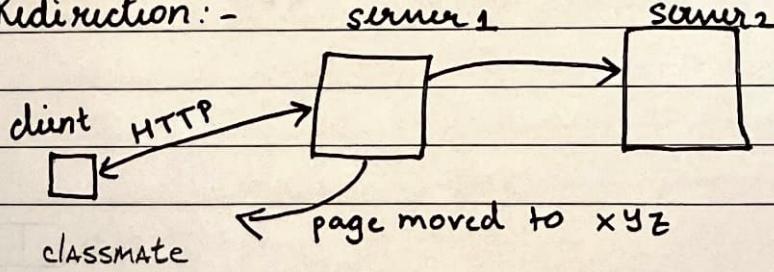


(we will have some expiry time).

* So let us say people want download same page.

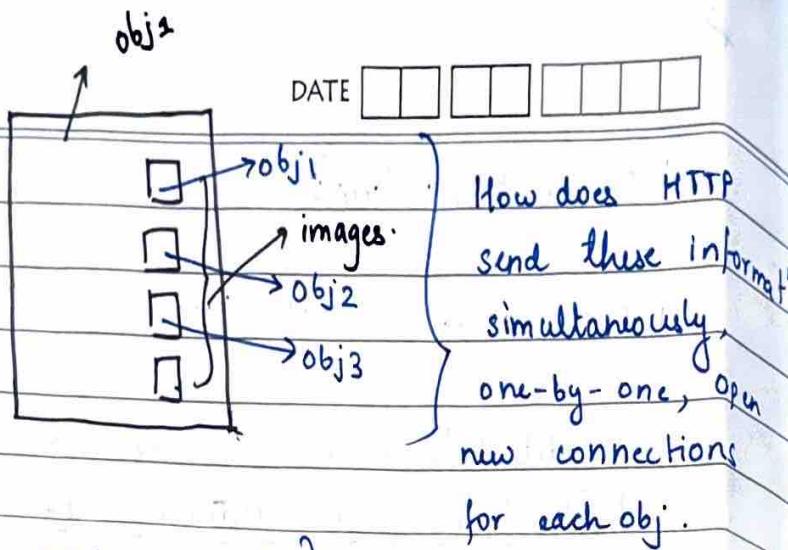
So at proxy we save it in cache and we distribute to users.

Redirection :-



} we can shift information from one url to different url.

HTTP
TCP



So HTTP uses UDP quickly (i.e.)

HTTP MESSAGE format
(Start line)

(MSG HDR)

Header ends at <CRLF>
empty line just having CRLF.

Message body

Example:- (OF REQUEST)
GET <URL>

http://www.google.com
Alternatively index.html.

<CRLF>
2 ASCII CHAR
<CRLF>
<CRLF> } multiple lines ending with crl.
<CRLF> contains]
content length : 1256 bytes
Expiry time : -- ?
ENCODING USED FOR
MESSAGE { information about
the message .

This information can also be put in the header.

IN HDR: HOST : www.google.com.

REQUEST OPERATIONS

GET → RETRIEVE DOCUMENT.

HEAD → RETRIEVE ONLY META INFO (expiry time, last modified time, length, to know if we want download new copy of file or not.)

OPTIONS → AVAILABLE OPTIONS (HTTP version at server).

POST → GIVE NEW INFORMATION TO SERVER.

PUT → WE MODIFY SOME INFORMATION AT URL OF THE SERVER / store.

DELETE → Delete specified URL.

Responses : like

Example

| | | | | |
|--------------|---------------------------|-----|----------|--------|
| (Start line) | HTTP/1.1 | 202 | ACCEPTED | <CRLF> |
| (HDR) | content-length :- | | | |
| (BODY) | expiry-time : (CRLF)
— | | | |

RESULT CODES :-

1XX Informational

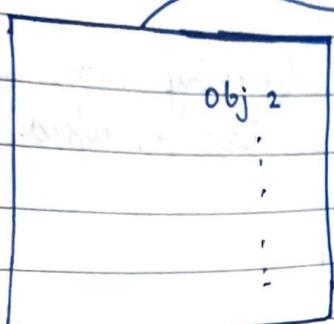
2XX SUCCESS

3XX REDIRECTIONS

4XX CLIENT ERROR

5XX SERVER ERROR

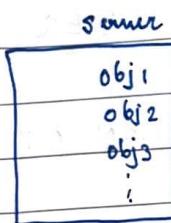
HTTP & TCP



obj 1 (index.html)

when browser download the
first one, it asks for all
objects within.

BROWSER
↓ REQUESTS
HTTP
↓
TCP



tricky thing to understand how to send all
Objs, obj 1 > obj 2 !
Then if we keep transferring one obj
then we can be stuck while sending
a large object, while others could
be quickly sent.

Q. SEND all request at once or sequentially
→ one request per TCP connections!

1.0

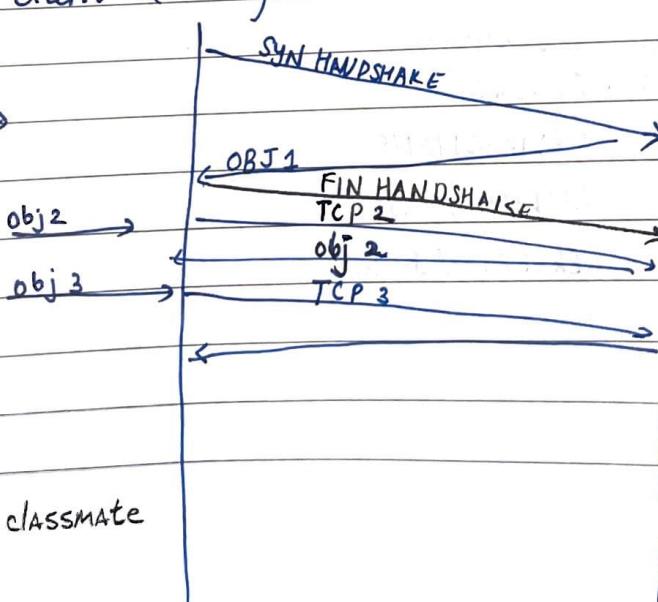
→ Send obj by opening new TCP connection for each
obj request.

client (HTTP)

Server

BROWSER

obj 1 (REQ) →



→ For each request
opening a new
TCP connection.

classmate

ISSUES:-

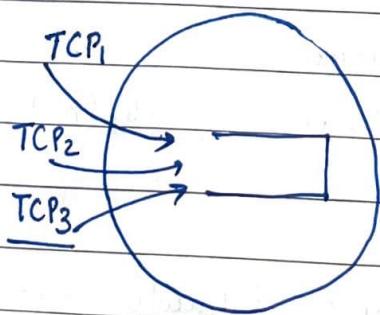
- 1) EACH TCP connection takes time to learn optimal CW
- 2) overheads (Handshake)
- 3) Server state is large due to multiple TCP connections being open.

Server:

(port : 80) https: 43? 443?
 (8080) =

http

4) FAIRNESS



* By opening our own many TCP connection, we get an unfair advantage.

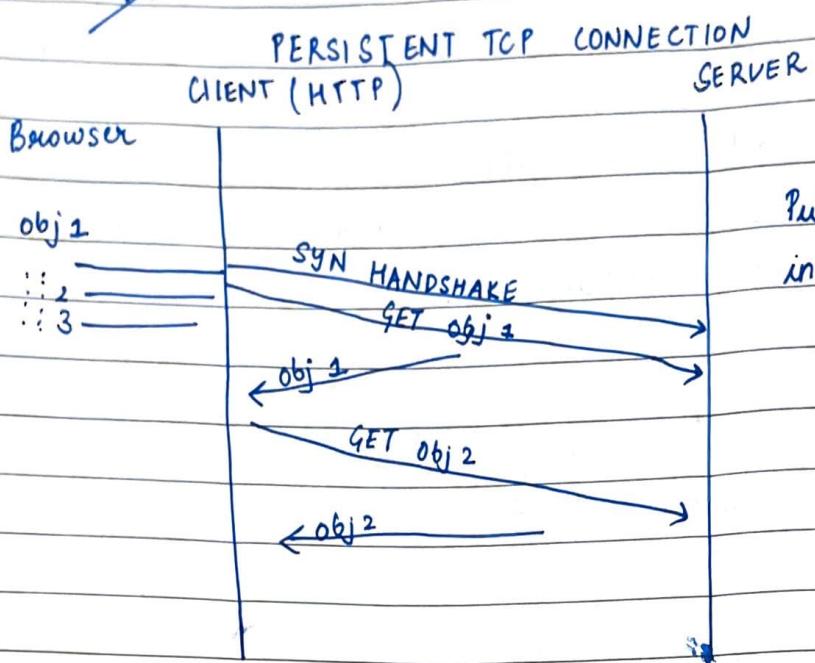


each link gets $\frac{c}{m+1}$ of capacity.

$$\therefore \text{Client 1} = \frac{mc}{m+1}$$

$$\text{Client 2} = \frac{c}{m+1}$$

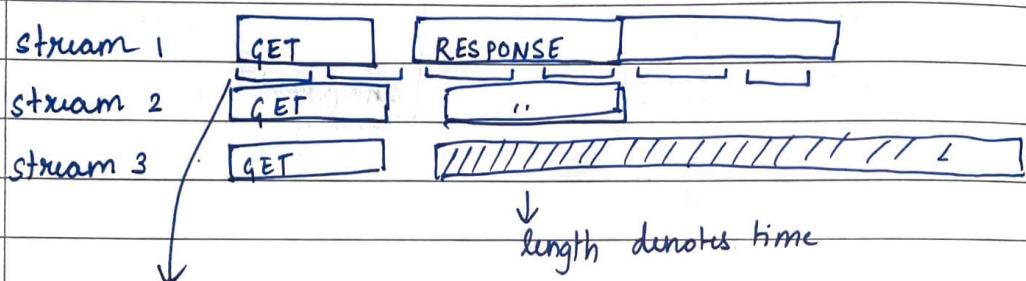
HTTP ~~1.0~~ rejected 1.1



Head of line blocking issue

→ Obj 'n' gets delayed because obj(n+1) is taking time to download.

HTTP 2 → single TCP, streams divide into streams



divide them into frames.

→ $F_{1,1}, F_{1,2}, F_{1,3}, F_{1,4} \}$ frames of first obj

Over the same TCP connection

