

COMPUTERS : NETWORKS

Protocol: Rules for transmitting & formatting of data.
set of rules.

- * internet: inspired from ARPANET

1960s

data → bits → signal.

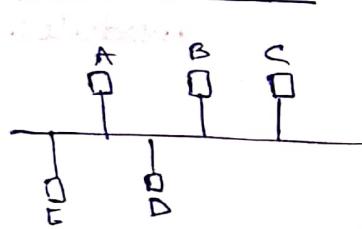
(wire/wireless).



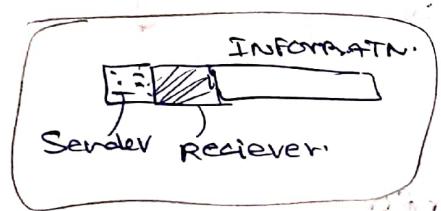
1) Physical Layer:

Some headers, to "denoise".

2) Data Link Layer:



medium access.



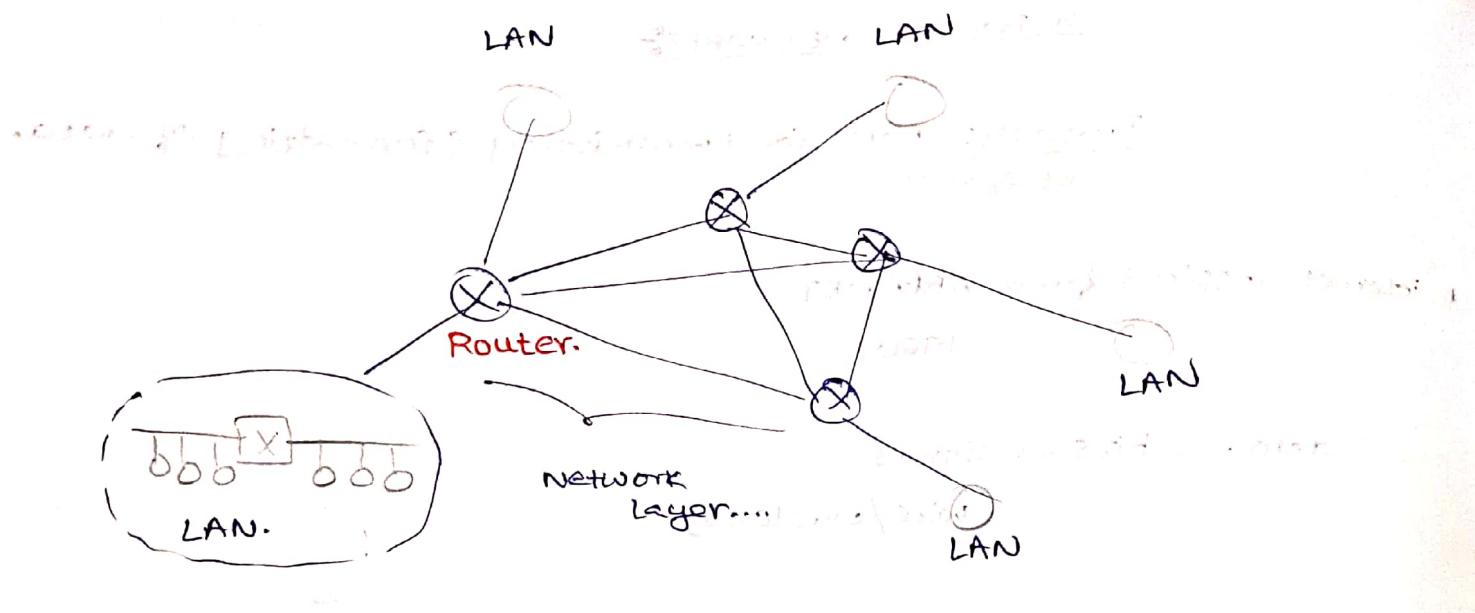
PACKET.

- Also gotta resolve interference.

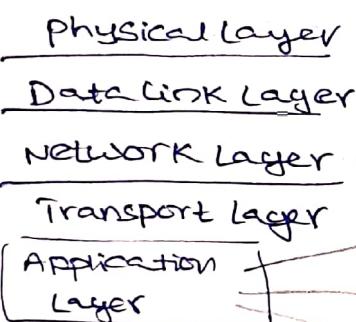
- * Level-2 switch: **repeater.**

(repeats only those signals which are necessary).

(routers are Level-3 switches).



*

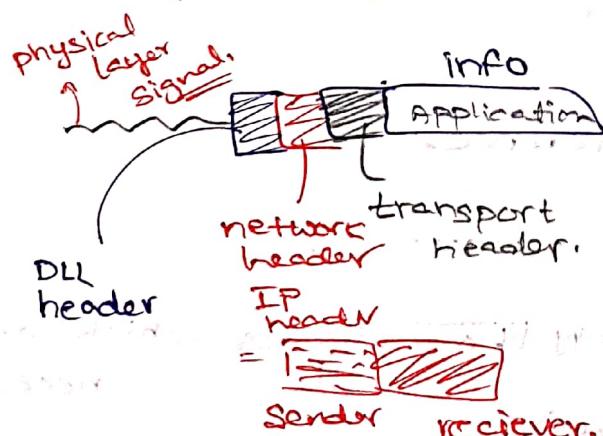


Protocol stack

Protocol
Stack.

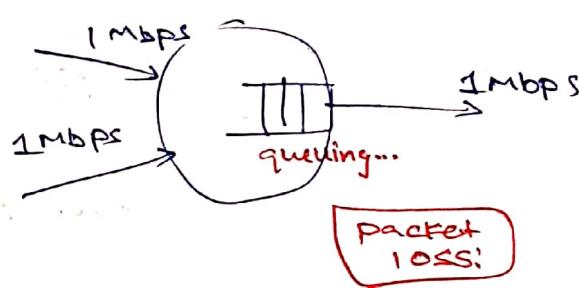
Single thing is
broken down into
modules.

message packet:-

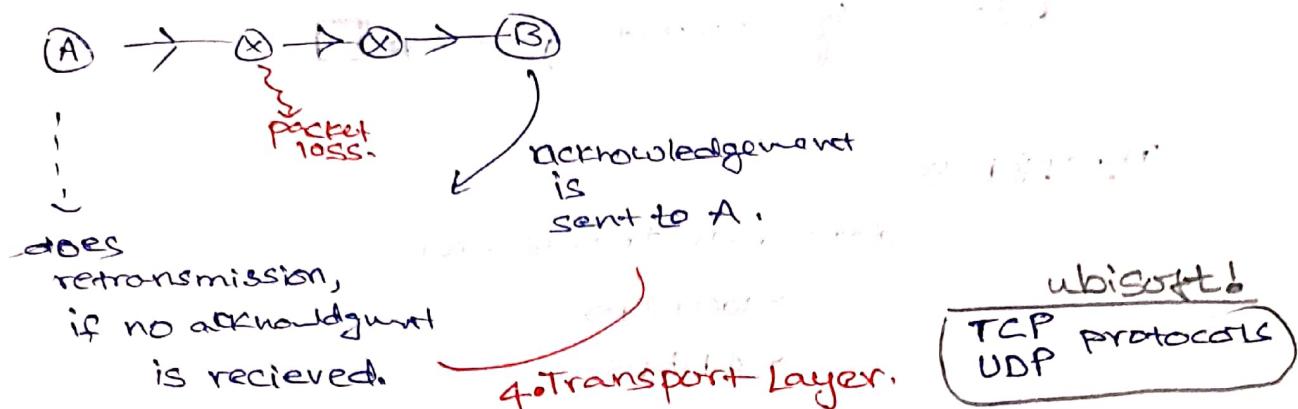


*

router:



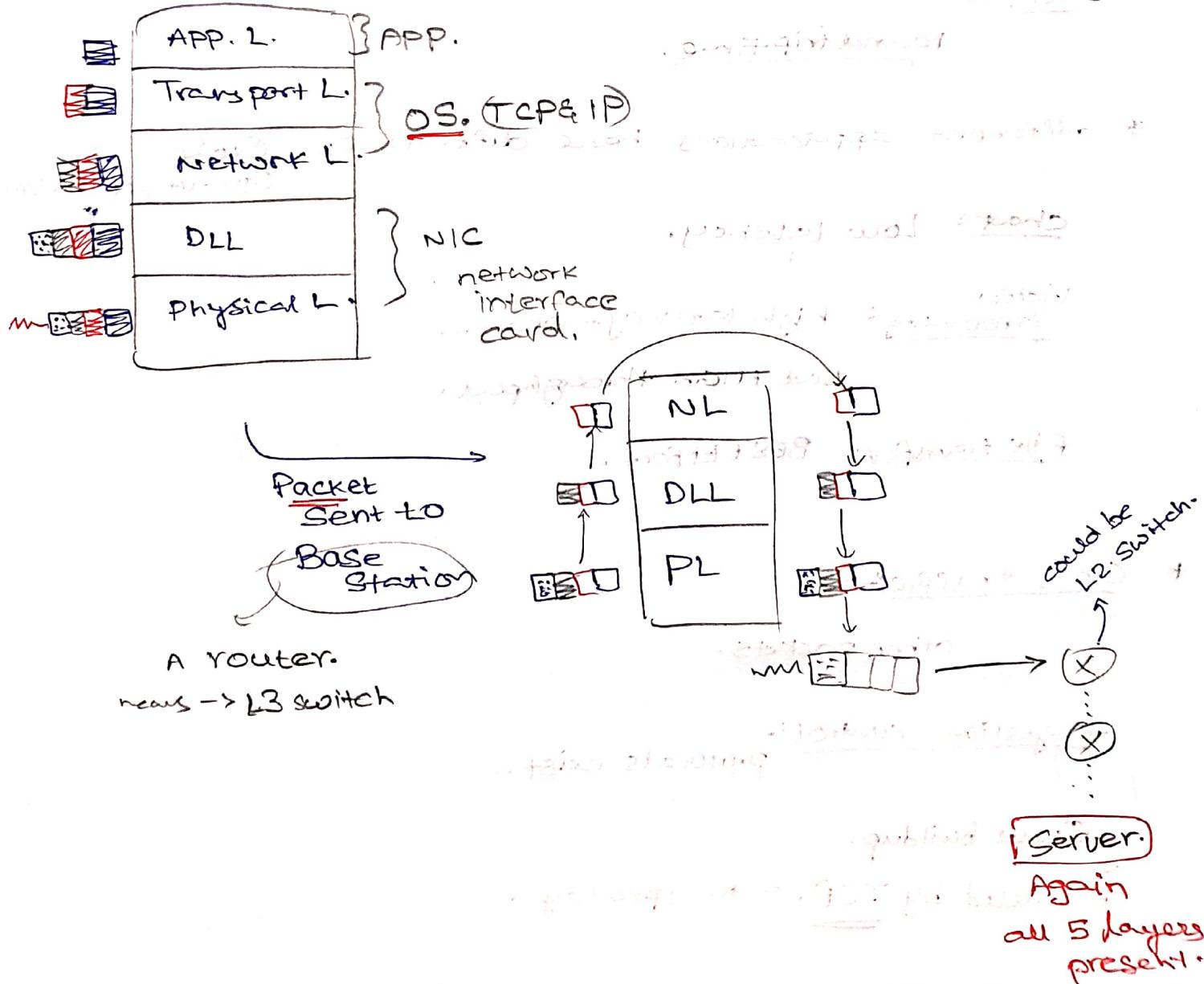
- * When A & B talk.
(end devices; with transport layer intact).



this design, retransmission by source; rather than intermediate is by choice.

(Some military history).

* My mobile:-



* Performance:-



Throughput:-

data rate transferred

> 100 Kbps

> 1Mbps

Latency:-

Delay of a message across network.

~ 100ms for one way latency.

> time for reaching.

RTT:-

round trip time.

* different applications have difference QoS (quality of service)

chat: Low latency.

video streaming: high latency... okay...

But HIGH throughput.

File transfer: Best Effort.

* Cross-traffic:

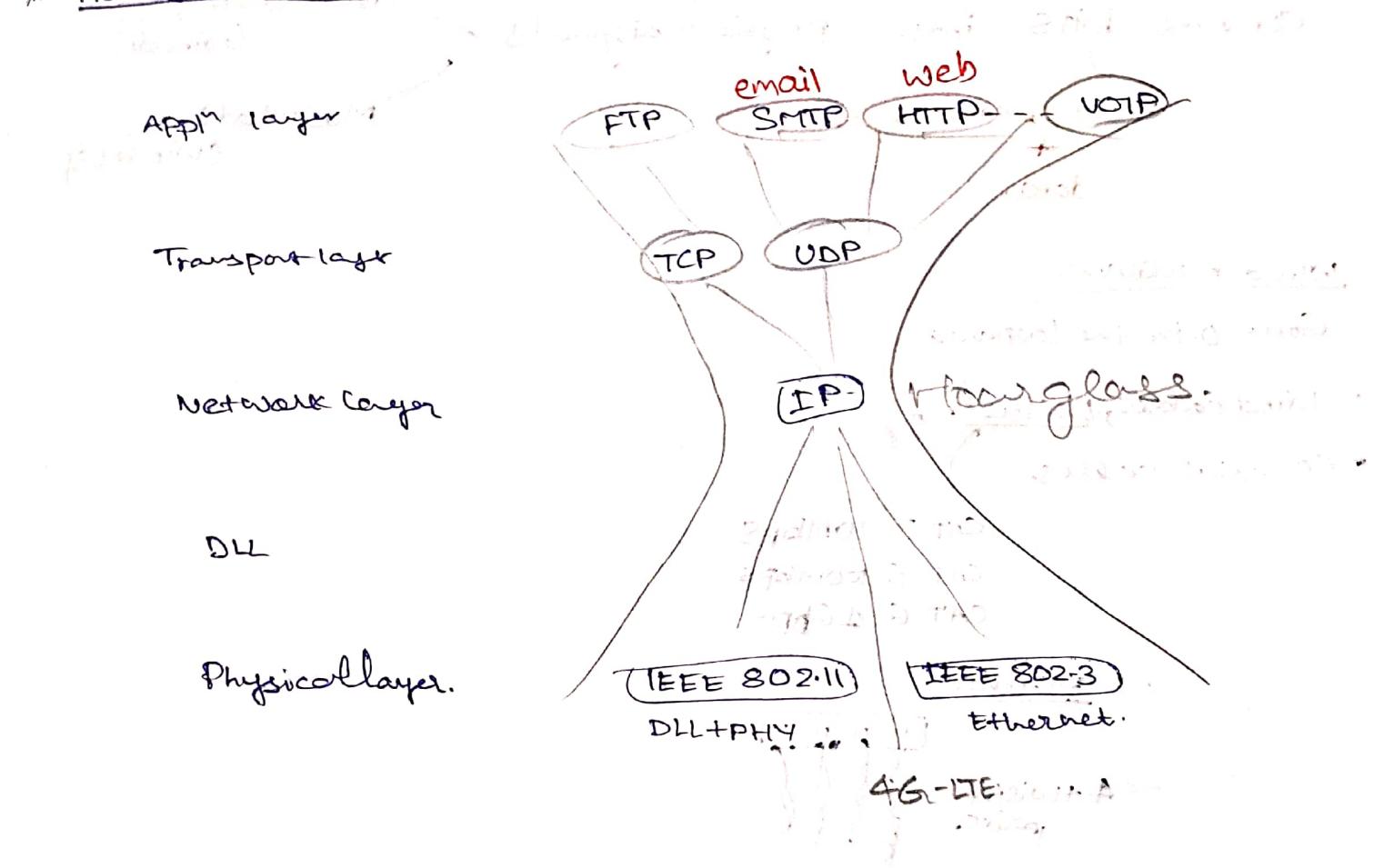
other packets.

Congestion control:- protocols exist...

Queue buildup.

Handled by TCP, or transport layer

* Hourglass model:



Protocol stack

Layered structure

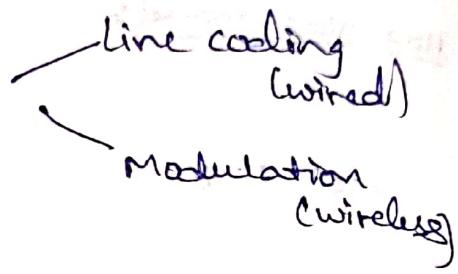
Protocol stack (bottom to top)

Layered structure

Protocol stack (bottom to top)

I) Physical layer:-

Converts bits into physical signals
+
Headers of other layers + data.

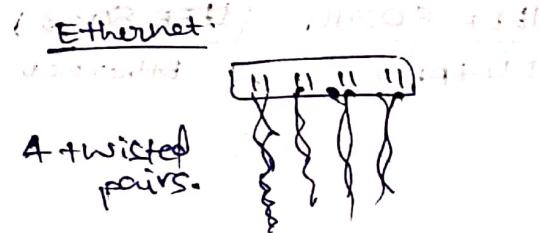


Wired medium:-

Gotta size the loop area.

- wired cables, twisted.
- co-axial cables.

CAT 3 10Mbps
CAT 5 100Mbps
CAT 6 1Gbps

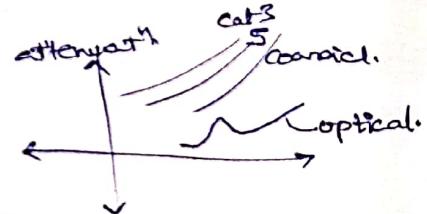


- optic fiber.

→ Attenuation:-

$$\begin{aligned} \text{Attenuation in dB} &= 10 \cdot \log_{10} \left(\frac{P_{in}}{P_{out}} \right) \text{ dB} \\ &= 20 \cdot \log_{10} \left(\frac{A_{in}}{A_{out}} \right) \text{ amplitude} \end{aligned}$$

Usually: $\text{dB}/1000 \text{ ft.}$



$3 \text{ dB} \rightarrow \frac{1}{2}^{\text{th}}$ power transmitted

* As distance \uparrow ; attenuation \uparrow .

→ Power in dBm, dBW:-

$$P_{dBm} = 10 \cdot \log_{10} \left(\frac{\text{Power}}{1 \text{ mW}} \right)$$

$\therefore 3 \text{ dBm} \rightarrow 2 \text{ mW.}$

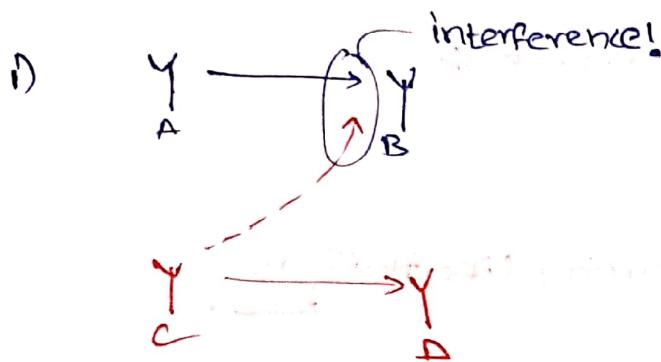
$0 \text{ dBm} \rightarrow 1 \text{ mW.}$

wireless



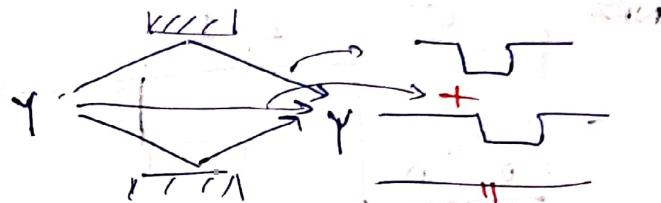
- directional antennas.
- MIMO - massive MIMOS.

problems in wireless:-



2) Obstructions

3) Multipath effect.

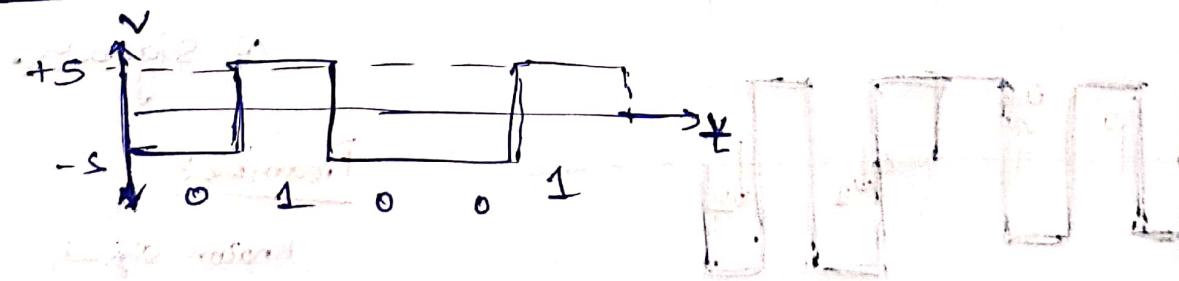


noise too!

signalling :- bits to SIGNALS.

1) Wired :- aka Line Coding.

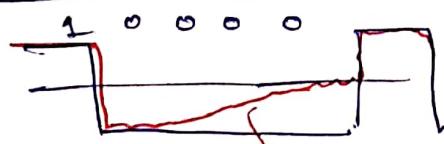
1) NRZ: non return to zero!



issues:-

1) clock recovery (receiver's clock may not match sender's)

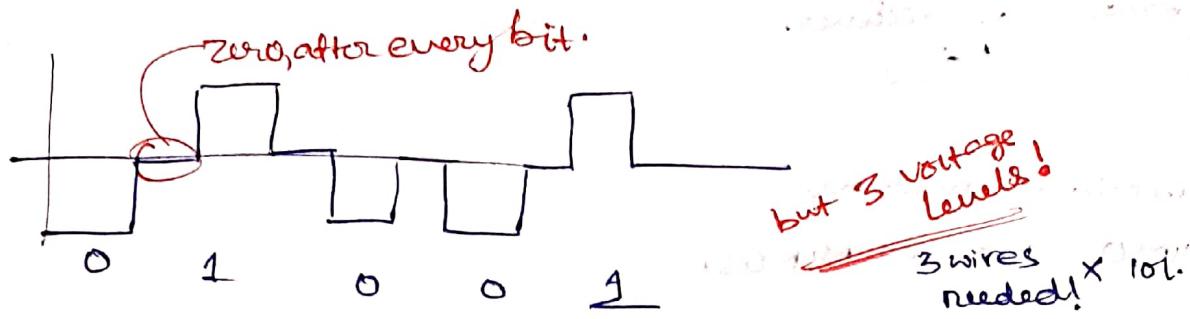
2) Baseline wander-



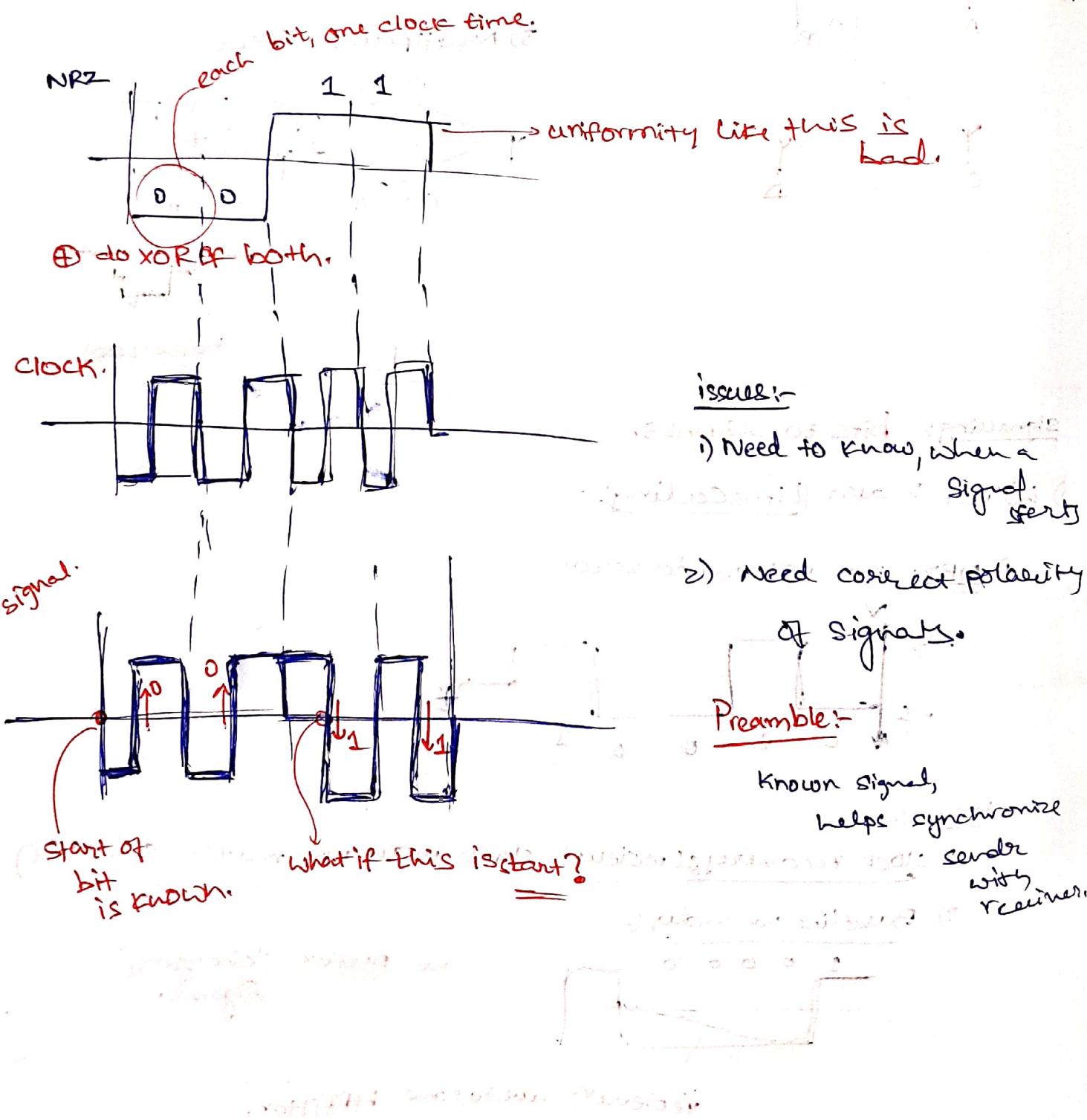
∴ we prefer "changing" Signal.

receiver; due to one HPfilter.

2) Return to Zero (RZ) :-



3) Manchester Coding (IEEE 802.3, ethernet) :-



A) Differential Manchester coding.

3) 4B15B coding:-

4 Bits } written & sent as
2⁴ } 5 Bits;

so that we don't have

more than 3 consecutive zeros.

Ex: 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000

0000 0000 0000 0000 0000 0000 0000 0000 0000 0000

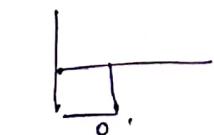
0000 0000 0000 0000 0000 0000 0000 0000 0000 0000

→ BAUD RATE:-

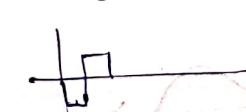
bitrate: no. of bits transferred per sec.

baudrate: no. of distinct symbol changes, made in the medium, per sec.

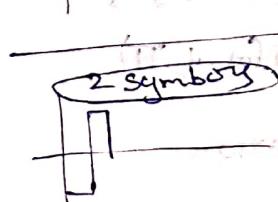
Ex:- In manchester coding:-



$\therefore \text{Baud rate} = 2 \times \text{Bit rate.}$



relates to frequency in the wires.



$\therefore \text{Baud rate more relevant to physical medium.}$

- its limitations for attenuation

→ bandwidth

No need of much of diff. bits of info.

most freq. are used up

→ MODULATIONS



each signal carrier is allotted a bandwidth.

(After F.T. of original wave)

* we can do FM, AM, PM

Amplitude phase modulation.

)

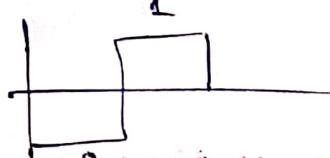
now!

for high speed

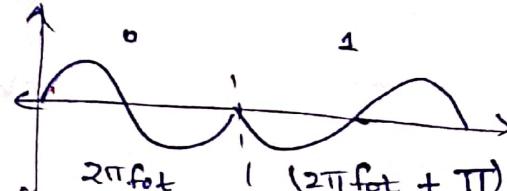
data transfer

$$\checkmark \cdot \cos(2\pi f_0 t + \theta) \checkmark$$

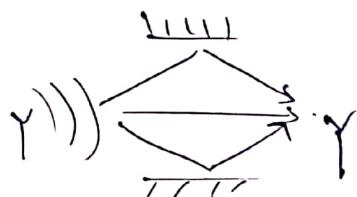
*



→



- * when $s(t)$ is sent out;



$$r(t) = \sum_i a_i \cdot s(t - t_i) + n(t)$$

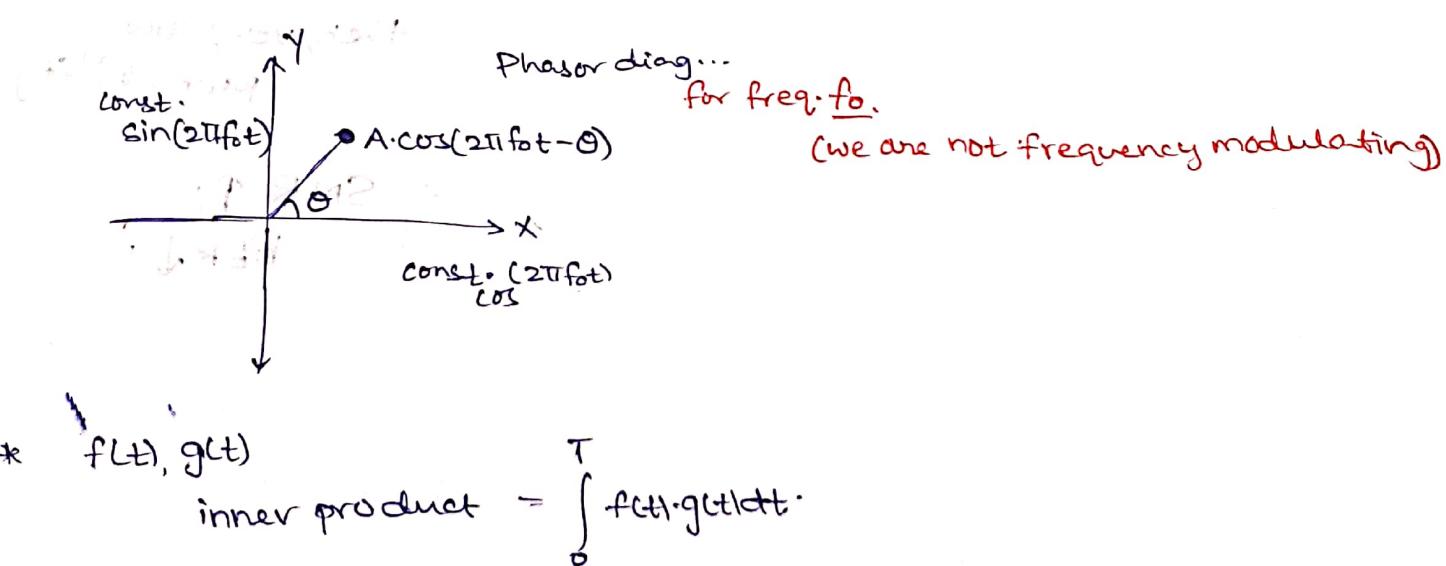
multi path
 phase differences
 due to multipath

white noise.
 from source.
 (Gaussian)

BER: bit error rate.

SNR: signal to noise.

→ Signals in vector spaces constellation diagram.



$$\star f(t), g(t)$$

inner product = $\int_0^T f(t) \cdot g(t) dt$.

$$\therefore \hat{e}_x = K \cdot \cos(2\pi f_0 t)$$

$$\langle \hat{e}_n, \hat{e}_x \rangle = 1$$

$$\therefore K^2 \int_0^{N \cdot T} \cos^2(2\pi f_0 t) dt = 1$$

$$\therefore K^2 \frac{\frac{N}{2} \cdot N}{f_0} \cdot \frac{1}{2} = 1$$

$$\therefore K = \sqrt{\frac{2f_0}{N}}$$

$$\hat{e}_x = \sqrt{\frac{2f_0}{N}} \cdot \cos(2\pi f_0 t)$$

$$\hat{e}_y = \sqrt{\frac{2f_0}{N}} \sin(2\pi f_0 t)$$

- * we have $r(t)$. we implement in circuit itself, that it "computes" $\langle r, e_x \rangle$ & $\langle r, e_y \rangle$

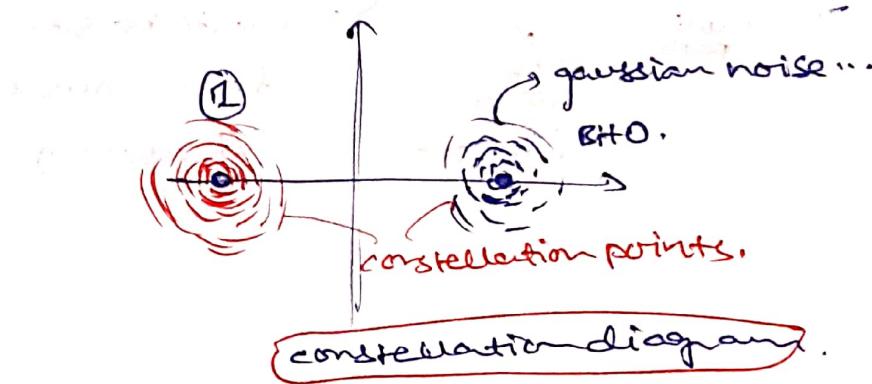
$$\therefore r_x = a \cdot \langle s(t), e_x \rangle + \langle n(t), e_x \rangle$$

noise in x
attenuation

gaussian
white noise

* if $s(t) = A \cos(2\pi f_a t)$ Bit 0.

$$= A \sin(2\pi f_a t) \cdot \text{Bit 1}$$



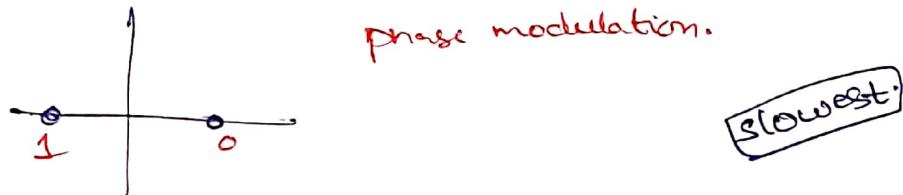
low probability that
bit 0's signal
goes over to
1's side
bit error

if $s(t) \gg n(t)$
then, very low
probability

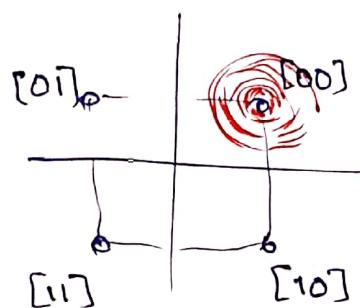
- SNR ↑;
BER ↓

modulation eggs:-

1) Binary phase shift keying (BPSK):

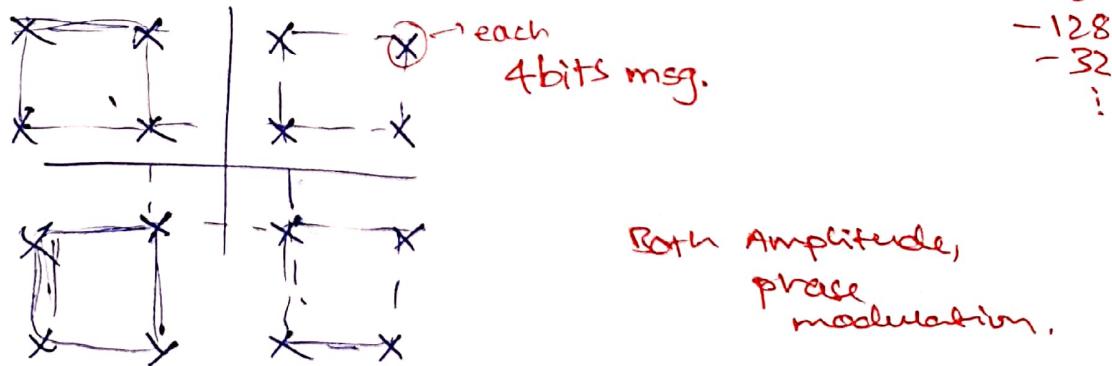


2) Quadrature Phase shift keying (QPSK):



∴ Bitrate = $2 \times \underline{\text{Baudrate}}$.

3) QAM - 16: Quadrature Amplitude Modulation. — family of modulations



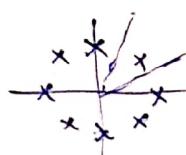
QAM-16, 32 ... offer high speeds;

but require good SNR values.

for detecting (•: constellations points will be far away).

a) 8-PSK:

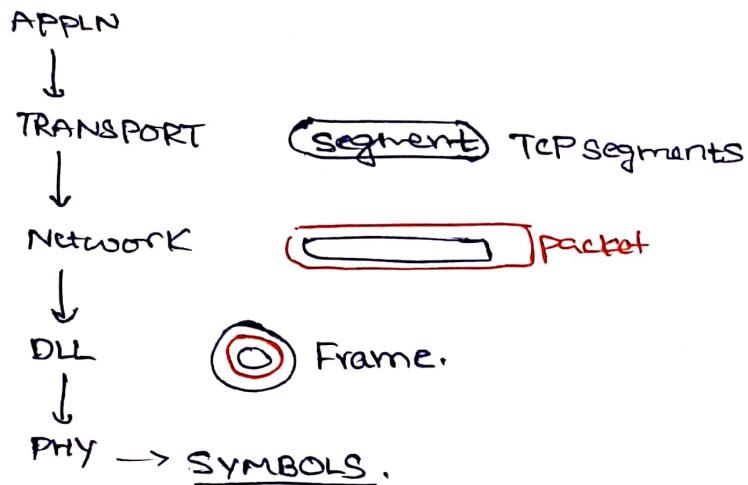
(only phase modulation.



(2)

Data Link Layer:-

Protocols for when & when not to send information.

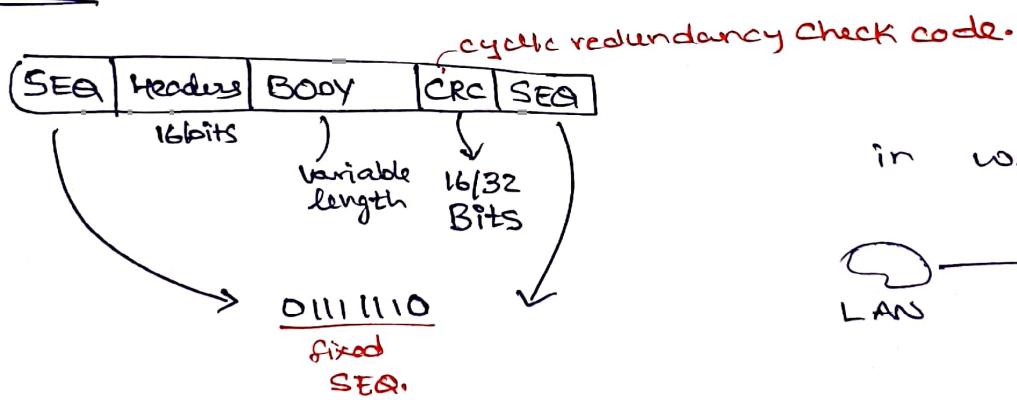


→ HDLC:- Highlevel Data Link Control. used as layer-2 in

WAN.

(wide area network).

FRAME:-



in WAN...



* when ideal; SEQ is sent:-

seq seq seq... (useful for clock synchronization).

now;

seq seq seq... | HEADER | BODY | CRC | seq...
 the SEQ "MAY" come here or not
 do Bit STUFFING:-

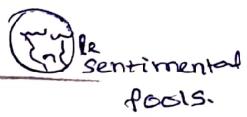
Algorithm for every 5 consecutive 1's;
 (encoding)
 place a zero.

Algorithm:-

(decoding) if we encounter 5 consecutive 1's;

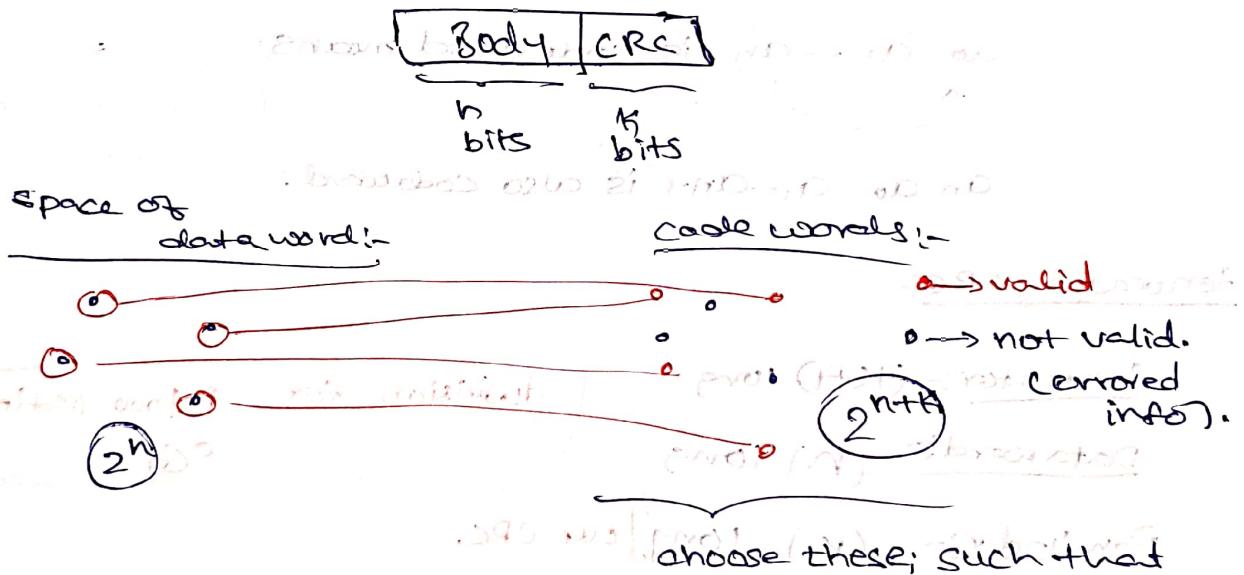
011111 | 0 → remove this & thus decoded.
 101111 | 10 → SEQ ✓
 111111 | 11 → No chance!

111111 | ∵ Bit error! discard the entire frame.



→ cyclic redundancy check (CRC) :-

By coding theory, we try to somehow "code" the body & give a codeword (= BODY + CRC).



→ Hamming Distance:-

→ How many bits different b/w two strings?

0	1	1	0
1	0	1	0

Hamming distance = 4 ("far")

Hence; very less probable,

by bit errors.

they are "far" apart.

So; if bit error comes, then we shouldn't be able to reach one codeword from another

We'll See

Galois Field GF(2) :-

2 elements $0, 1$.

Add", Subtr":-

do XOR.

Multiplication:-

$$(1011) \times (1)$$

$$\begin{array}{r} \text{Multiplication} \\ \hline 1011 \\ \times 1 \\ \hline 1011 \end{array}$$

$$0 \oplus 1 \times$$

$$11101$$

int "1011" is 5 bits of $2^3 + 2^2 + 2^1 + 2^0$ \Rightarrow xor not.

* CRC; such that $D \oplus R = 0$ becomes a valid word

as $a_1 \dots a_n$ is code word means,

$$\begin{array}{r} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_n \end{array}$$

(try to verify).

an $a_1 \dots a_m$ is also code word.

Generating CRC:-

divide D by G

Generator:-($K+1$) long

Dataword:- (n) long

division for D \oplus R (K times)

E_G

Reminder:- (K) long] our CRC.

code word = Data + CRC (append!)

2 methods exist for it

Eg: $G = 1101$ and $D = 110110$

$$D = 110110$$

not $D \oplus G$.

its $D + Q(K \text{ times}) E_G$

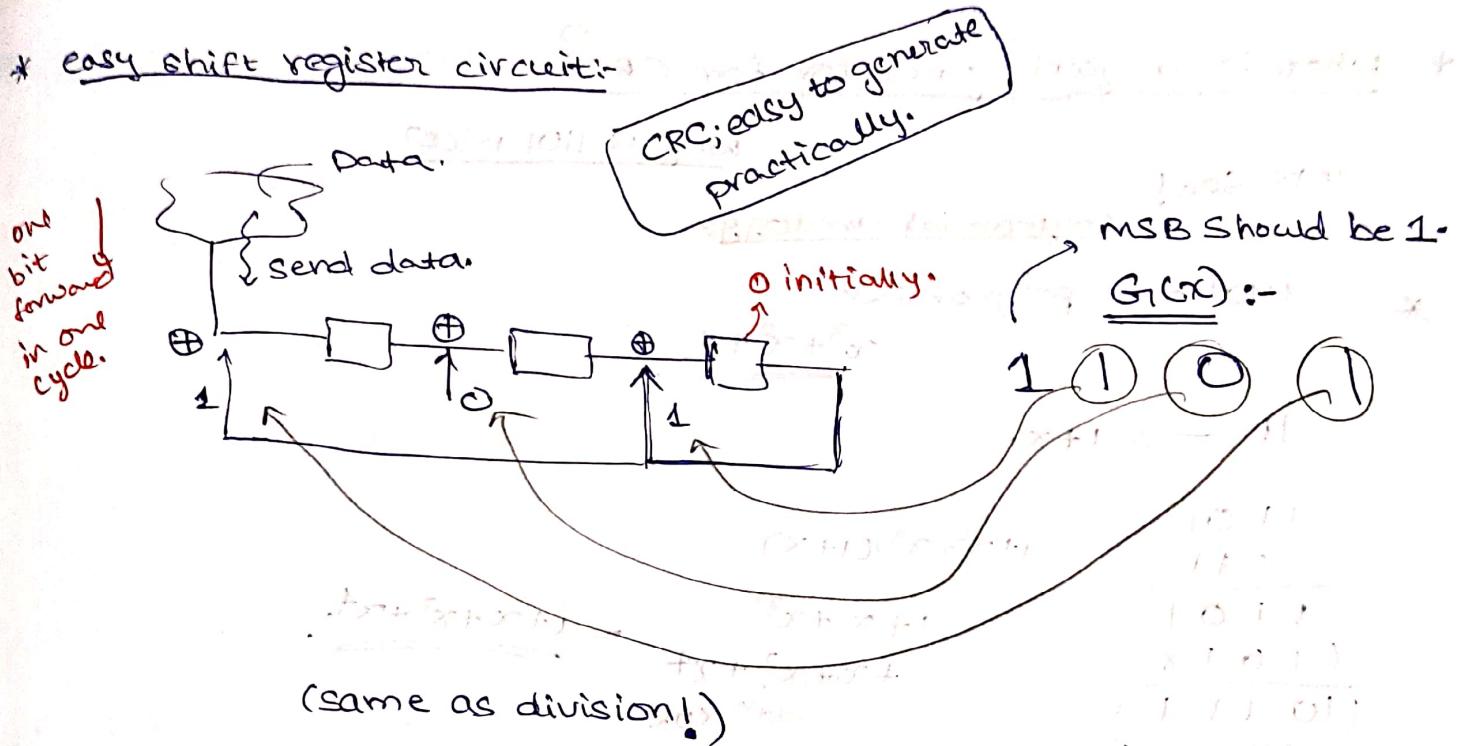
Quiz-1
error!

$$\begin{array}{r} 100011 \\ \hline 1101 \quad | \quad 110110000 \\ \oplus 1101 \\ \hline 0001 \\ \oplus 0000 \\ \hline 0010 \\ \oplus 0000 \\ \hline 0100 \\ \oplus 0000 \\ \hline 1000 \\ \oplus 1101 \\ \hline 01010 \end{array}$$

$$\begin{array}{r} 01010 \\ \oplus 1101 \\ \hline 01010 \end{array}$$

$$\begin{array}{r} 01010 \\ \oplus 1101 \\ \hline 01010 \\ \oplus 1101 \\ \hline 0111 \end{array}$$

our CRC.



Last 3 bits will be remainder! (from right \rightarrow MSB to left \rightarrow LSB).

* So, now the receiver receives the code word.

How to decrypt?

Data | CRC.

Data | 000

↓ CRC compute

\equiv CRC sent

① recalculate
CRC for data.

Data | CRC

↓ CRC compute

\equiv 000

② $(000 + 000) \mod (x^3 + x + 1)$

will get 0;

when divided data | CRC.

(Data + CRC) mod $x^3 + x + 1$

* What is a good generator for CRC?

lets see!

why is 1101 nice?

at at time we have polynomial analogy:

* $1101 \rightarrow \text{polynomial}$:

$$x^3 + x^2 + 1$$

$$11 \rightarrow 1+x.$$

$$\begin{array}{r} 1101 \\ \times 11 \\ \hline 1101 \\ 1101 \times \\ \hline 10111 \end{array}$$

$$(1+x^2+x^3)(1+x)$$

$$\begin{aligned} &= 1+x^2+x^3 \\ &\quad + x+x^3+x^4 \\ &\quad \text{add XOR} = 1+x+x^2+x^4. \end{aligned}$$

* $P(x)$ for 1001 is x^3+1 . $\frac{P(x)}{G(x)}$ = 0

Code word.

let error in transmission be $E(x) = 00\cdots010\cdots$

show that $G(x) \mid P(x)+E(x)$ for 1-bit error.

$$\text{Received} = P(x) + E(x)$$

\downarrow
in $G(x)$.

NOW; $G(x) \rightarrow \text{generator}$.

Note! $\frac{P(x)+E(x)}{G(x)} = 0$; then we say, no error! even if $E(x) \neq 0$.

∴ for 1-bit error detection,

$$E(x) \neq 0; \text{ then } \frac{P(x)+E(x)}{G(x)} \neq 0.$$

$$\frac{P(x)}{G(x)} = 0 \quad (\because P(x) \text{ is codeword})$$

$$\therefore \frac{E(x)}{G(x)} \neq 0.$$

Say $G(x)$ is of form $1(1101-011)$.

$$\text{then } G(x) = 1+x^4+()$$

never divides (x^4)

* How to say $1+x^{k+1} + \dots$ never divides x^i ?

this is GF(2)!

By contradiction

Say

$$(x^{k+1} + \dots + 1)(x^p + \dots + x^q) = x^i$$

not intuitive...

Easy! ~~It's not~~

XOR Swallows terms.

$$\underbrace{(x^{k+1} + \dots + 1)}_{\text{terms}} \underbrace{(x^p + \dots + x^q)}_{\text{terms}} = x^i$$

$$x^{k+1+p} + x^q + \dots$$

(odd number) at least 2 terms.

Can't divide x^i . why! because x^i must be equal

∴ 1-bit errors are always detected, for $G(x) = 1 \dots 1$

→ 2-bit errors → ~~can't be divided into smaller parts~~ model.

→ 2-bit errors

$$E(x) = x^i + x^j \quad (j > i)$$

$$= x^i (x^{j-i} + 1)$$

$$\frac{E(x)}{G(x)} = \frac{g_1(x) g_2(x) \dots}{f_1(x) f_2(x) \dots} \quad [\text{can't be divided into smaller parts.}]$$

Let $G(x) = 1 \dots 1$.

then $x^i \rightarrow$ not divisible.

∴ $x^i + x^{j-i} \rightarrow$ can be divisible but ...

$(1 + \dots + x^{k+1})$ divides $((+x^n))$ for large n .

* Order

Defn:

small 'i' s.t. $G(x)$ divides $1+x^i$ is

Called order of $G(x)$.

∴ $G(x) = 1 + \dots + x^k$; then we can have

$G(x)$ with order $2^k - 1$.

∴ \exists

∴ if we have 16-Bit-CRC;

then min. length

b/w two erroneous bits

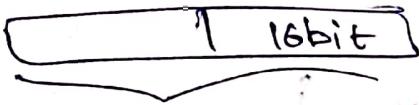
(maths exist)

crazy good!

$= 2^{16} - 1$. for non-detection.

Codewords :-

Now if the frame length is



length $\leq 2^k - 1$ then; we can always detect

2-error frames.

(with $G(x)$)

→ handling odd no. of errors: - (with $G(x)$) being discussed as before).

$$E(x) = 1001 \cdot 1 \cdot 0 \dots \text{ (odd no. of 1s)}$$

say $G(x)$ has $(1+x)$ as factor.

then $G(x)$ doesn't divide $E(x)$.

since,

Thus error detected!

in $GF(2)$:

$$(1+x)(\underline{\quad}) = \frac{\text{error with odd}}{\text{error with even no. of 1s.}}$$

eg:

$$\begin{array}{r}
 10110 \\
 \times \quad 11 \\
 \hline
 1011010 \\
 \underline{+} \quad 1011010 \\
 \hline
 (11)(1011011)
 \end{array}$$

even!

converse:

if $G(x)$ has even 1s
then $(1+x)$ is a factor.

so; if $G(x)$ has even 1s; then we can detect

odd-error frames

(and discard them)

* HDLC:-

16 bit CRC:

$$G(x) = x^{16} + x^{15} + x + 1.$$

32 bit CRC:

$$G(x) = x^{32} + x^{26} + x^{23} + x^{22} + x^{16} + x^{12} + x^{11}$$

$$+ x^{10} + x^8 + x^7 + x^5 + x^4 + x^2 + x + 1.$$

(Lord!)

HDLC over. Lets see ARQ:-
(b/w WANs)

damaging

burst of errors

$$ECD = 00 \dots 11100\dots$$

l bits.

$$= x^l(1 + x + x^2 + \dots + x^{k-1})$$

$$G(x) = x^k + \dots + 1$$

$\therefore G(x)$ won't divide x^l

$G(x)$ won't divide $1 + x + \dots + x^{k-1}$ if $k < l$

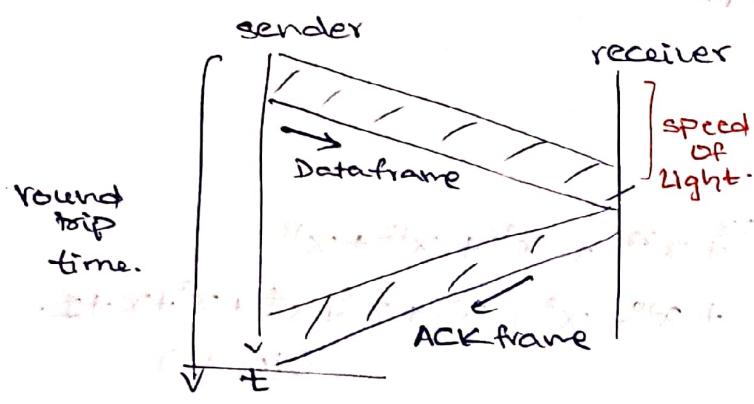
\therefore bursts of length $l < k+1$ are detected.

Protocol:-

i) ARQ (Automatic Repeat Request).

in DLL

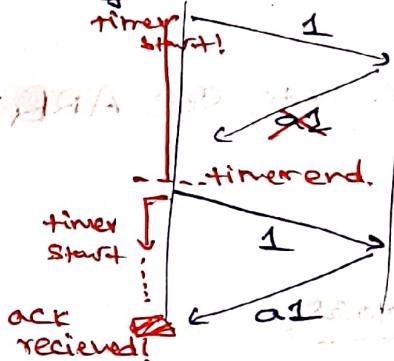
Handles reliability.



∴ what if acks don't come?

- set a time out!

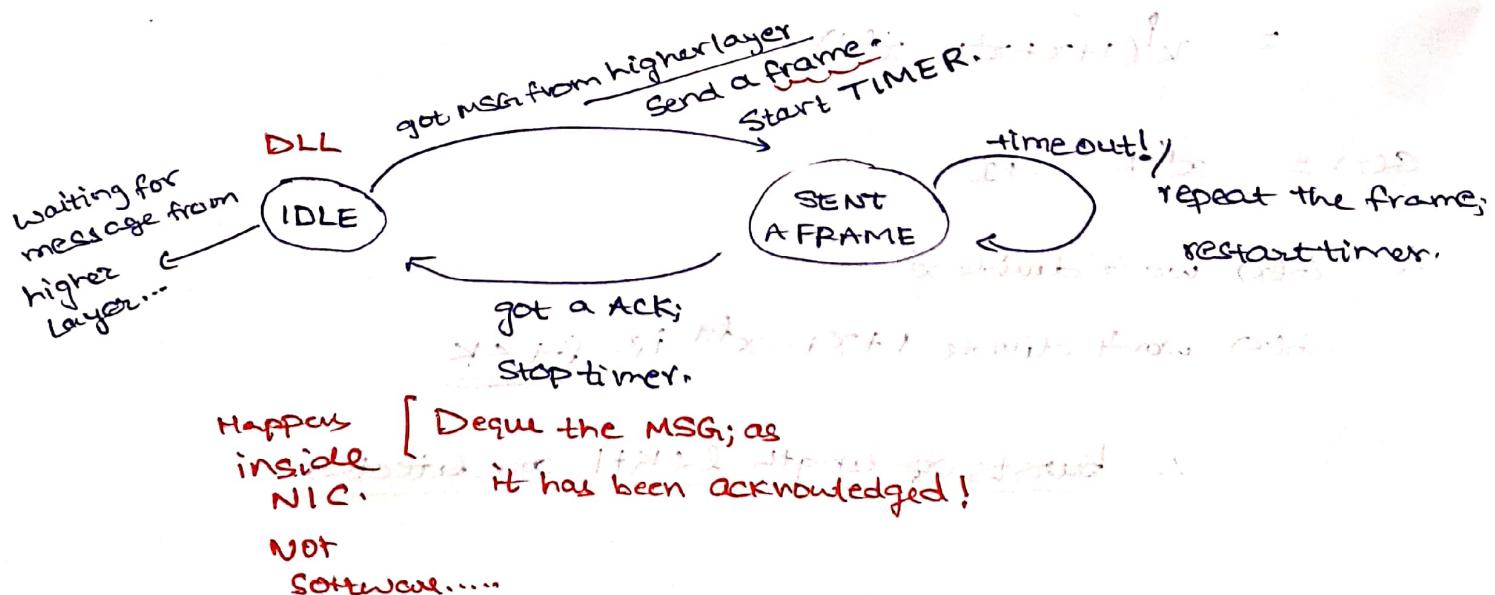
Eg:



We did a socket coding for this nai!

State diagram at sender; for ARQ

(∴ receiver; just sends ACK).



Medium Access:

Both wired & wireless!

- * what will we do, to prevent interference of ≥ 2 messages?

OR how to proceed, after we've two msgs at once?
we can't!

we drop that msg.

Our design is to make sure, that at any instant,
only one msg is in the medium. (medium Access)

- if, by chance, we send a packet & detect collision
2: Don't receive ACK's
we simply need to re-transmit that packet!

< Looks very elementary, crude! But this gives you the 150Mbps !
wow!

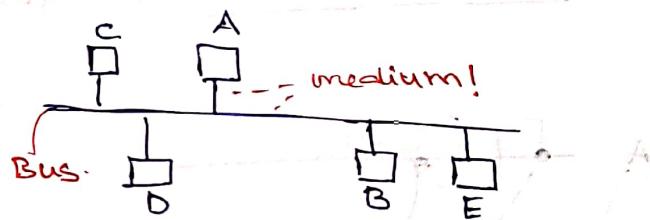
→ wireless protocols:

Ethernet LANs

(most used wireless part)

CSMA :-

carrier sensing Multiple Access.



- * These are random access protocols: schedule for transmission; not decided before hand.

- unicast:- single node is destination

- broadcast:- Every node is destination.

Basic Idea:-

Create a frame, send it! (But do carrier sensing first!)

- If we get a collision; Stop transmission; Wait for RANDOM time & repeat!
- How to detect?
- RANDOM BACKOFF.

IEEE 802.3: Ethernet

CSMA-CD: No Acknowledgements.

collision detection. (not possible in wireless; heavy attenuation).

- carrier sensing: (is possible in wireless; but "HIDDEN TERMINAL ISSUE")
Is the medium free, or busy?

> Detect energy.

Not busy: noise

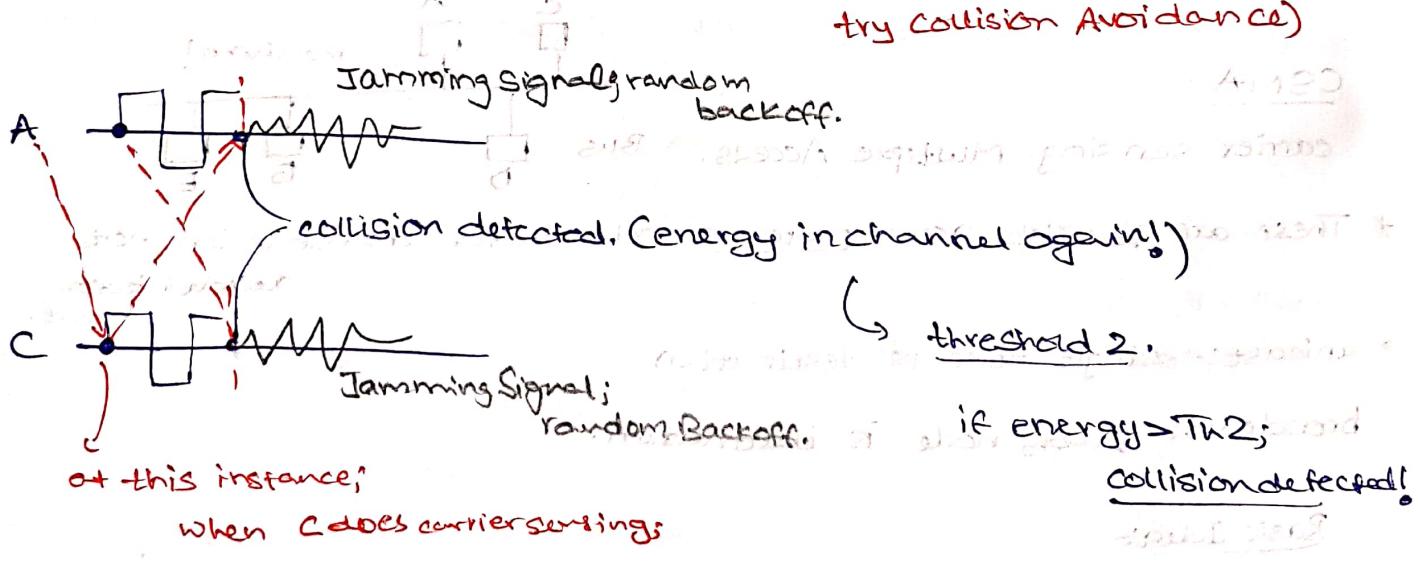
Busy: noise

if energy > threshold 1; channel is busy.

just above
noise energy.

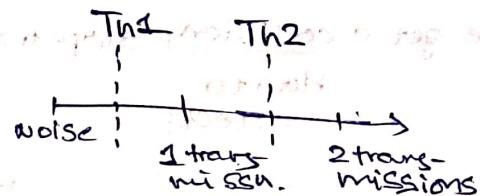
- Collision Detection: (not possible in wireless;

try Collision Avoidance)

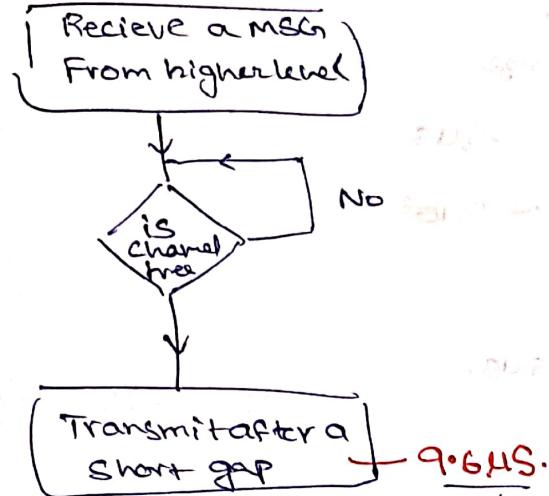


Channel is free.

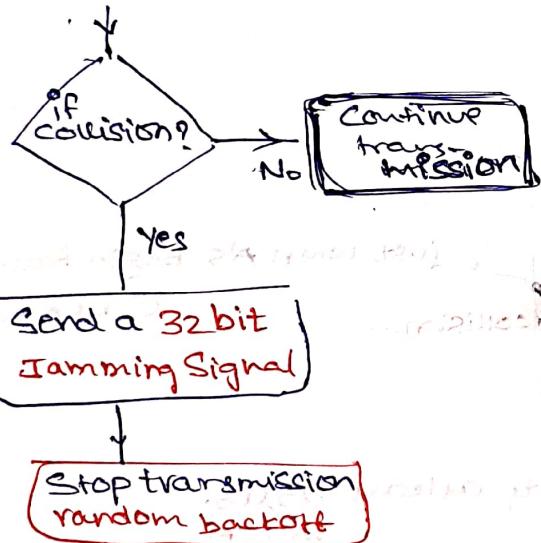
So transmit!



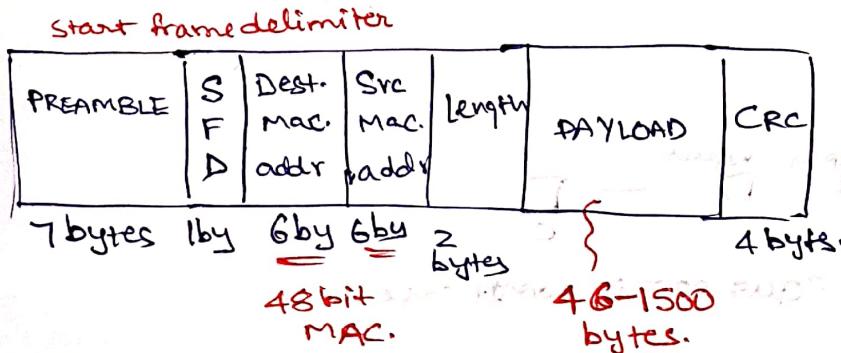
Flow charts:-



Collision detection:-



Frame details in 802.3:-



why this min. limit?

Why this max. limit?

memory in NIC goes up!.

→ to give others also, access to channel
 even if 1 bit error; we need to discard frame.

→ why the min-frame size?

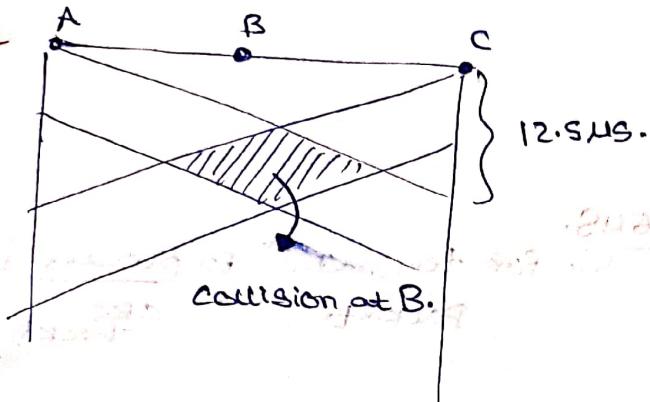
Ethernet:

2.5 Km max. range.

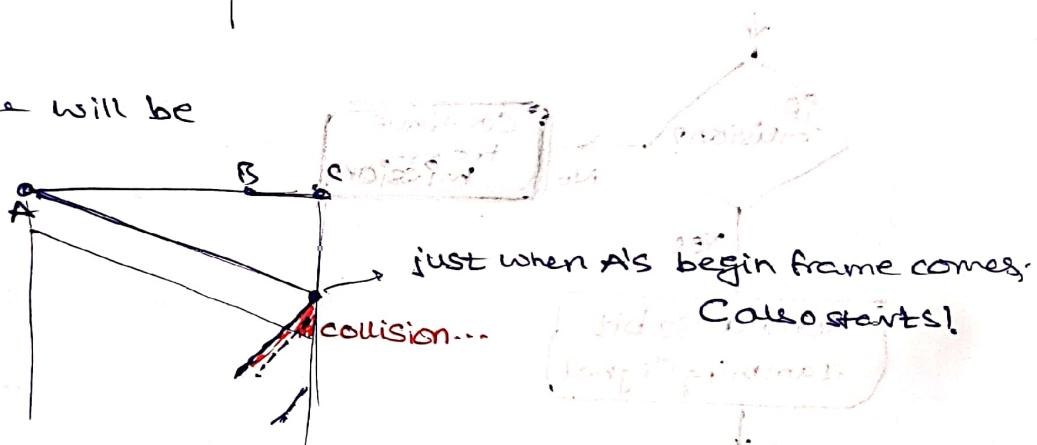
$$RTT = \frac{5000}{2 \times 10^8} = 25 \mu s.$$

✓ not $3 \cdot 10^8$

A must detect collision; before his frame ends. (since there is no ACK in CSMA-CD)



∴ worst case will be



∴ A's frame must last, atleast $25 \mu s$.

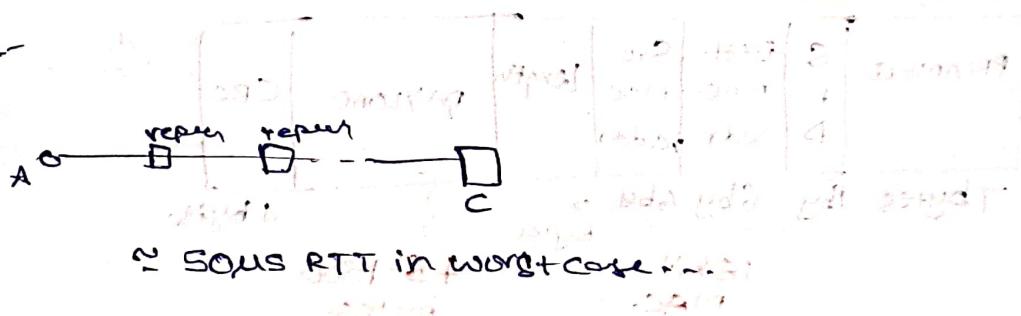
Hence,

minimum size is imposed.

2.5 KM

is range!

But, Scenario -

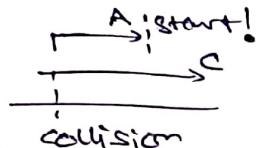


& 64 bytes; 10Mbps → $51.2 \mu s > 50 \mu s$.

the min. imposed on payload!

so that A can detect collision before his frame ($51.2 \mu s$) ends.

random Backoff:-



$\Delta = \text{min. wait time}; \leq 2 \text{ ms. give enough time}$

to detect
col.

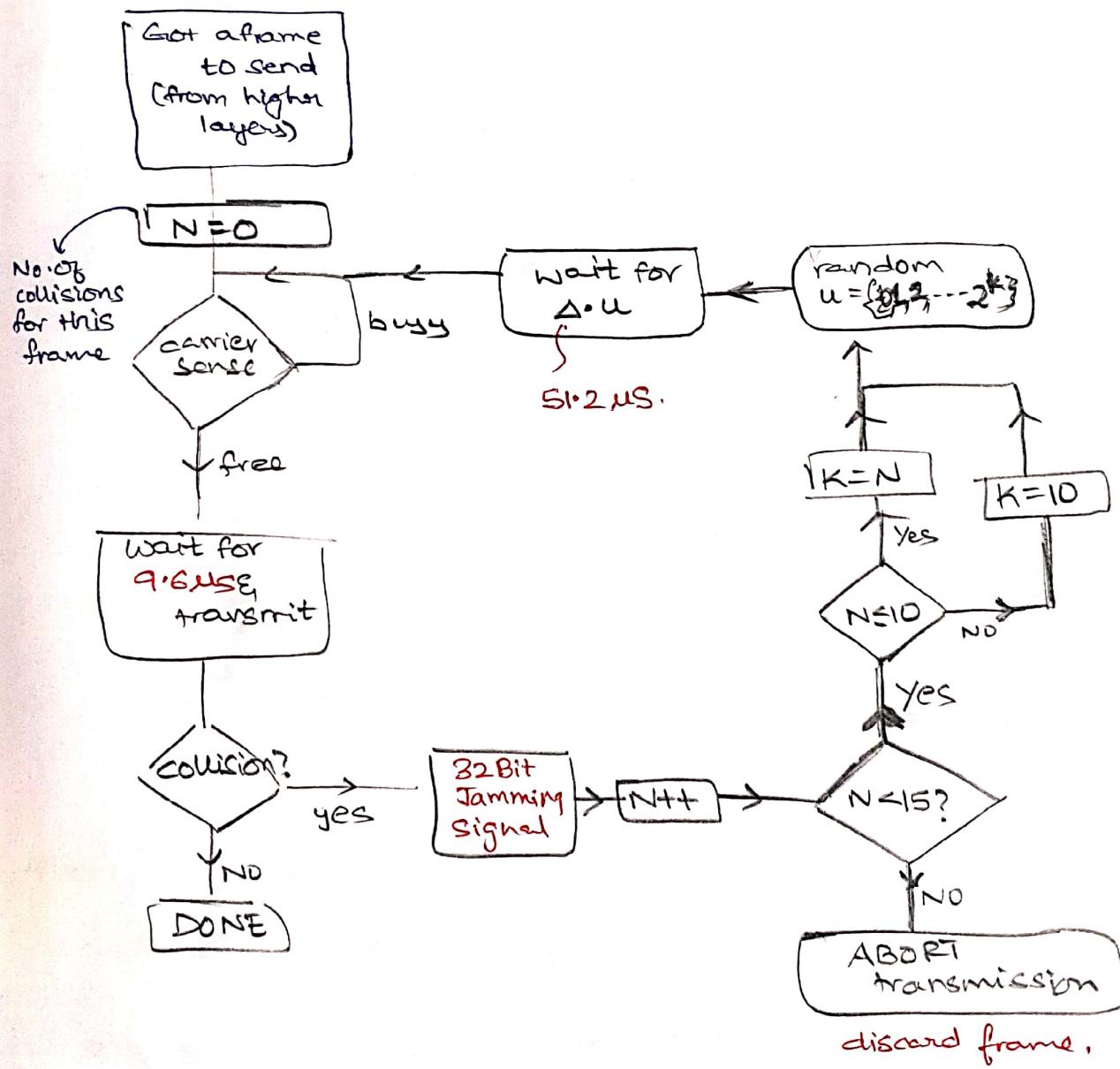
Wait time $\in \Delta \cdot u$

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

$K = \text{no. of collisions (consecutive)}$

exponential
backoff.

Final, CSMA-CD flowchart:-



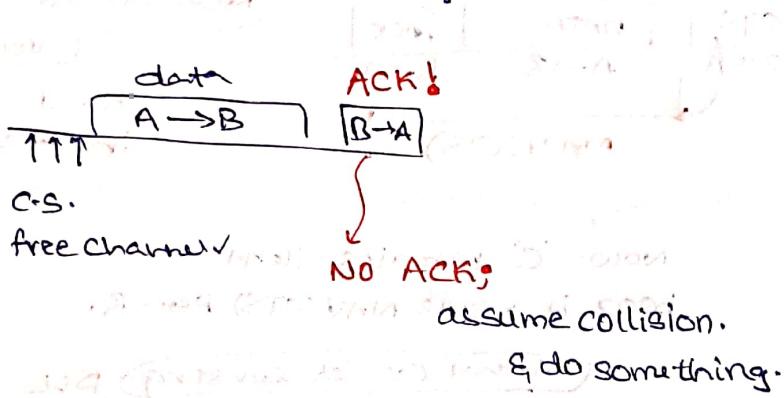
卷之三

WIRELESS (WiFi):

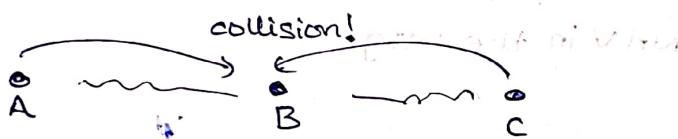
high attenuation! maybe CS; but no hopes for CD.

$$\frac{1}{d^2}, 2 \leq d \leq 5$$

so, how TO detect collisions?



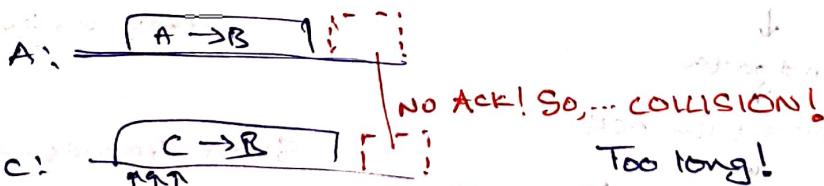
→ Hidden Terminal problem:



B can hear A & C.

A, C can't hear each other.

So; consider following scenarios:-



Too long! in ethernet;
we would stop
DATA in middle.

But here; we're
waiting for full
transmissⁿ.

DO something!

* lets do virtual carrier sensing!

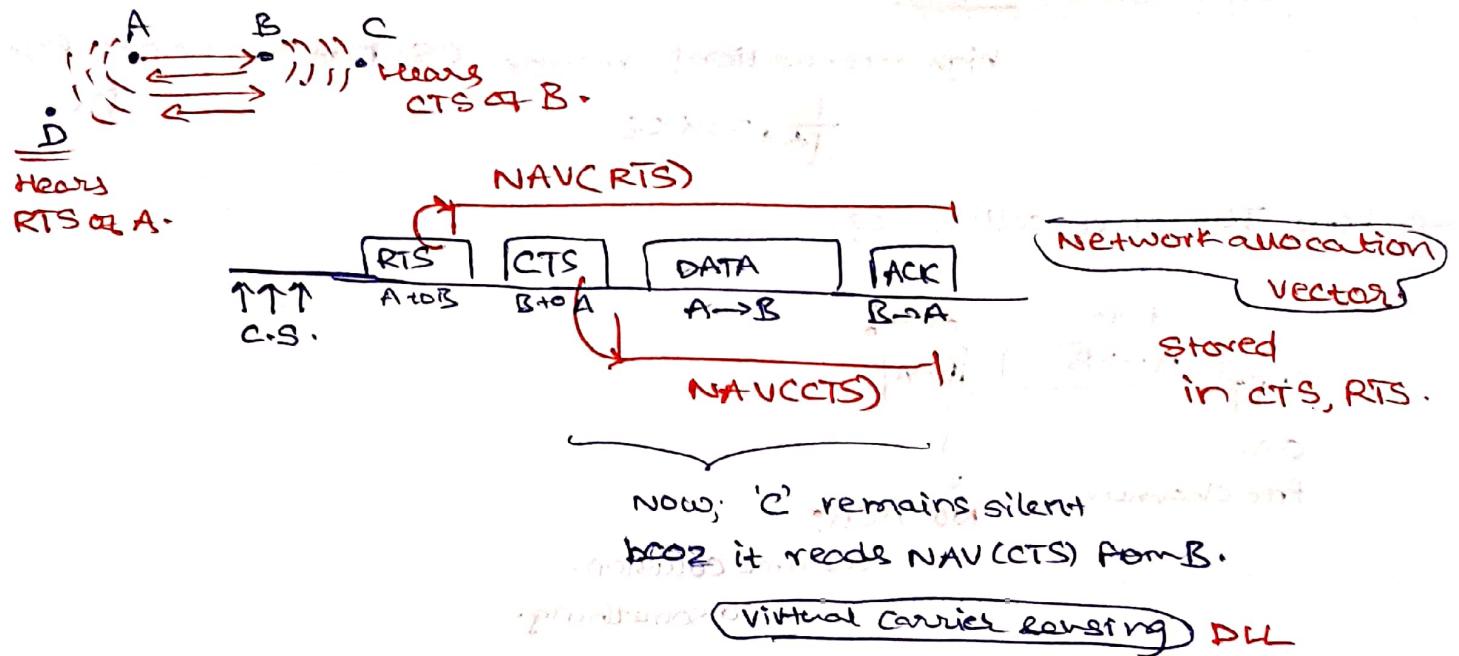
(B) will tell (C), that it is expecting (A) to transmit!

RTS: request to send. tells B, that it wants to send packet.

CTS: Clear to send.

Short frames

↳ B tells A, to continue!



Rule: Any node hearing an RTS/CTS should remain silent for the NAV in the msg.

To summarize:

option-1:-

1. Don't use RTS/CTS:-

when traffic is less.

Just carrier sensing

↓
Send data

↓
lookout for ACK

↓
retransmit, if no ACK.

But; why won't I send RTS, CTS....

they are just short frames ha...

Short frames → bits lol!

Kanis...

"we need RTS CTS for a wide range of coverage!".

• we use lower modulation techniques (BPSK); which will work with low SNR.

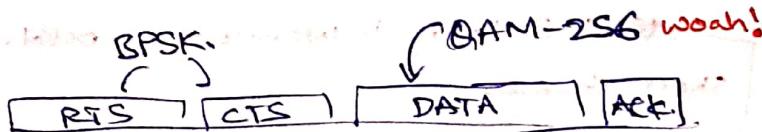
duration is long!

not trivial durations!

use; RTS, CTS, CarrierSensing and everything.

NO CTS → assume collision.

Time Scale:

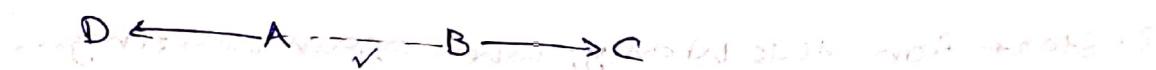


PROTOCOL:-

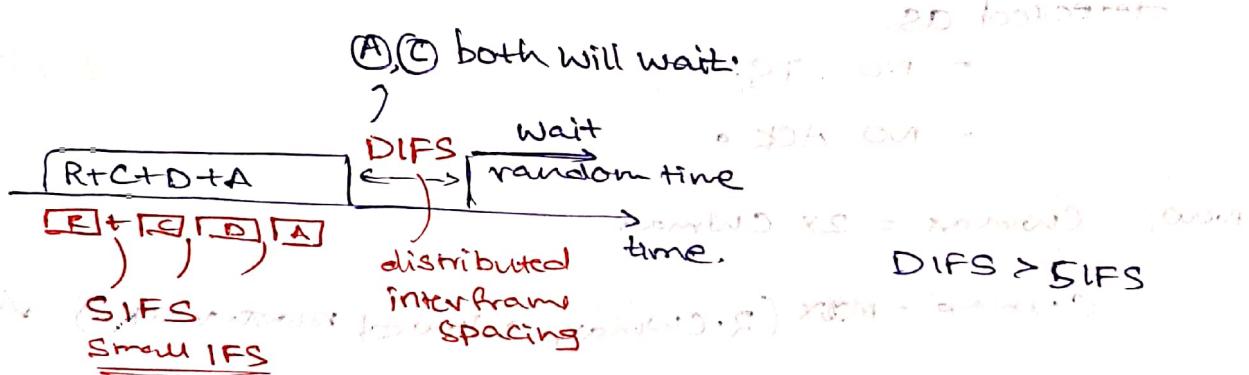
CSMA-CA - (wifi)

Collision Avoidance.

Exposed terminal Problem:-



But $A \rightarrow D$; B hears the RTS & stay quiet.
or due to carrier sensing.



SIFS is needed, bcoz;

When RTS: $A \rightarrow TX$, $B \rightarrow RX$

CTS: $A \rightarrow RX$, $B \rightarrow TX$

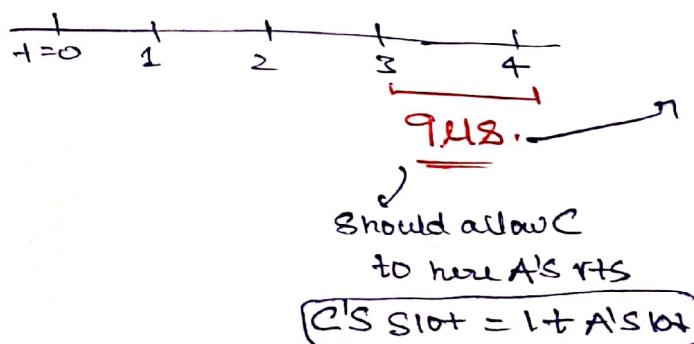


* random wait time!

contention window (CW)

Bad disagreement.

divide time into slots



propagation delay
+
possible offset
+
time to carrier sense.

→ How long to wait?

initialize $W \in (0, CW_{max})$

& wait for $W \times 9 \mu s$.

1. decrease W by 1; forever ideal slot.

! 2. freeze W ; if the channel is busy.

3. Start from this W only; when again waiting...

the node waiting from before, should preferably transmit.

Hence; Freeze W .

→ collision!

detected as

- NO CTS; either hits slot (J), A

- NO ACK either hits slot (J), B

Now; $CW_{max} = 2 \times CW_{min}$

$CW_{min} = \max(2 \cdot CW_{max}, \text{allowed max-value}) \quad \checkmark$

$W = \text{unif}(0, CW_{max})$.

repeat!

→ 1. wait for DIFS

2. decrement W every ideal slot

3. if channel busy; freeze W

Transmit, when $W=0$.

* Frame of CSMA-CA:-

IEEE 802.11 b,g,n,ac,ax



data → code → sent over channel.

rate $\frac{3}{4}$ → Sent N bits; $\frac{3N}{4}$ are information bits.

$S/6$ → $S/6$ are information bits.

			Bandwidth:-	Speeds:-
802.11g	64-QAM , rate $\frac{3}{4}$		20 MHz	54 Mbps
11n	64QAM	$S/6$	40 MHz	150 Mbps
11ac	256QAM , $\frac{3}{4}, \frac{5}{6}$		160 MHz	866 Mbps
11ax	<u>1024QAM</u> , $\frac{3}{4}, \frac{5}{6}$ ↑SNR		160 MHz	1.2 Gbps
				<u>SISO rates...</u>

* DIFS > SIFS

$$DIFS = SIFS + 2 \times \text{slot time} \quad (\mu\text{s}).$$

$$PIFS = SIFS + 1 \times \text{slot time}.$$

* Regarding QoS;

we can give preference to video call latency; by using smaller CWmax limits
 But Cross-layer interaction.

1920-21-10-10-10-10

1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10
1920-10-10-10-10-10
1920-10-10-10-10-10
1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10
1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

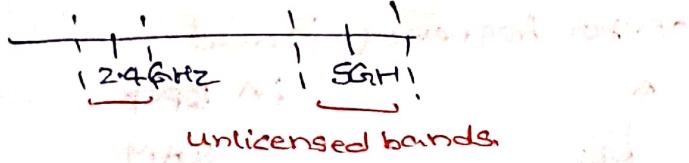
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10

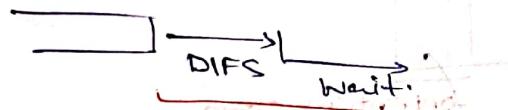
1920-10-10-10-10-10
1920-10-10-10-10-10

1920-10-10-10-10-10
1920-10-10-10-10-10



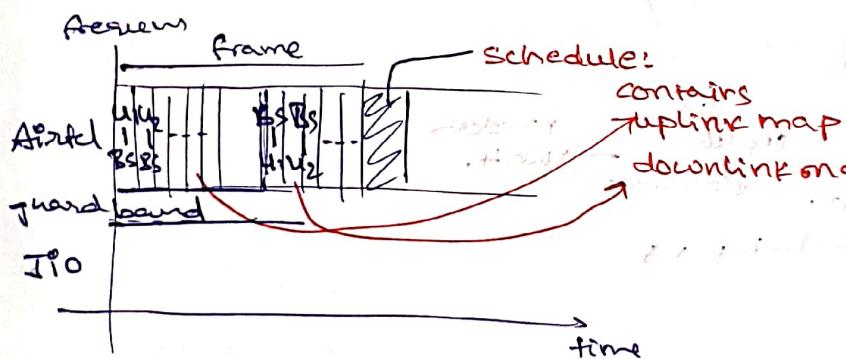
* saw I got a licensed bandwidth. Do I still have to do CSMA-CA?

Hii! I'm Airtel



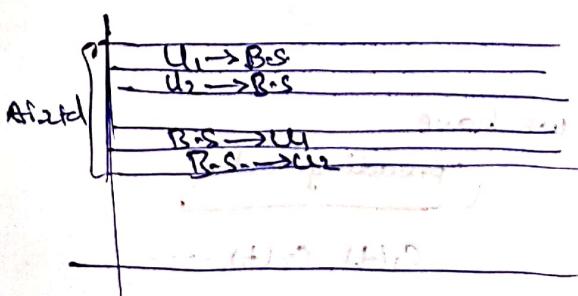
Can I save on this? I've got crores of people... :P.
phew! lets see...

→ TDMA: time division multiple access



(on)

→ FDMA: Frequency division multiple Access

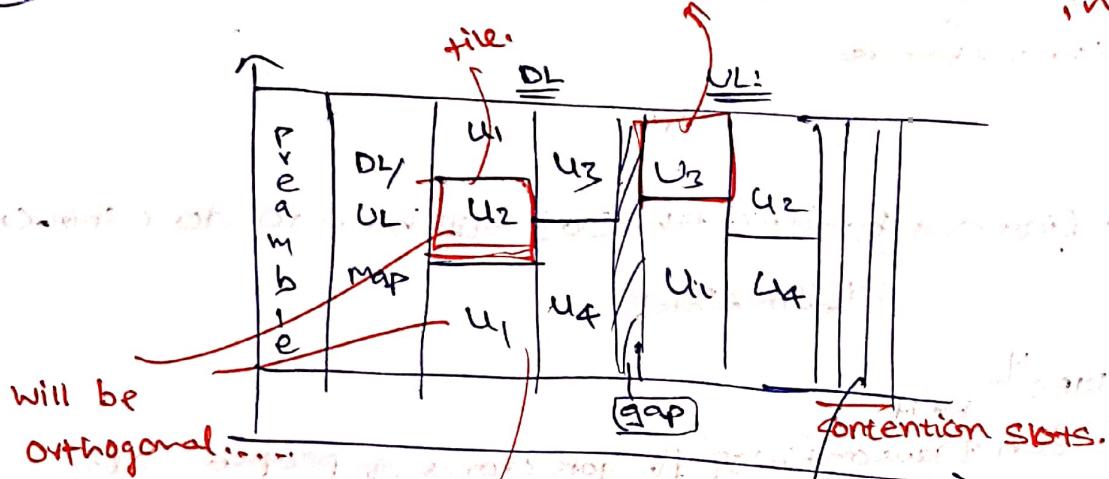


looks nice na? ya...

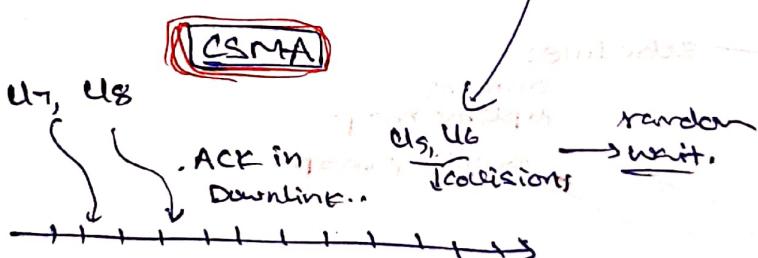
* 4G uses OFDMA: orthogonal Frequency division Multiple Access.

OFDMA

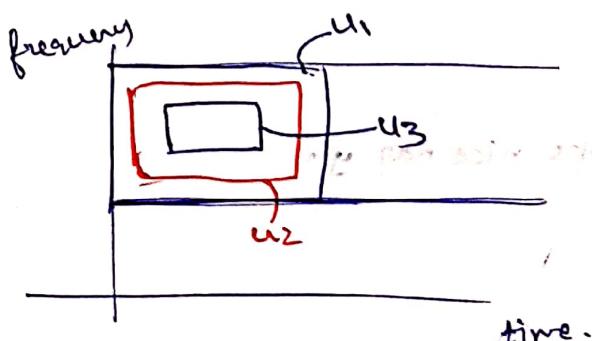
center frequency is used
in BPSK
QAM



this map is dynamic; caters to demand.



CDMA:- Code Division Multiple Access:-

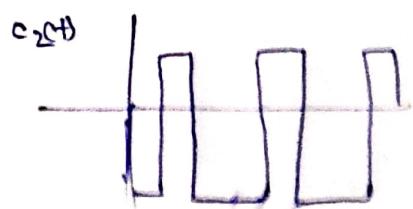


we have

Spreading codes

$$c_1(t), c_2(t), \dots$$

like



All spreading codes
are mutually
orthogonal

$$|G_1(H)|^2 = 1$$

$$|G_2(H)|^2 = 1$$

$$\int G_1(t) \cdot G_2(t) dt = 0$$

$$u_1 = c_1(t) \cdot A \cos 2\pi f_c t$$

$$u_2 = c_2(t) \cdot (A \cos 2\pi f_c t)$$

different constellation points.

somehow,
CDMA also
mitigates
multi-path issues.

* to recover from $r(t) = u_1 + u_2$,
dot product with $c_1(t) \cdot \cos(2\pi f_c t)$.

CDMA, OFDMA
robust to multipath

$$u_1 \cdot c_1(t) \cdot \cos(2\pi f_c t)$$

$$= \frac{A}{2} (1 + \cos 4\pi f_c t) \cdot 1.$$

remove by
LPF.

$$u_2 \cdot c_1(t) \cdot \cos(2\pi f_c t)$$

$$= -\frac{A}{2} c_1(t) c_2(t) (1 + \cos 4\pi f_c t)$$

remove by
LPF

∴ we have

$$\frac{A}{2} + -\frac{A}{2} c_1(t) c_2(t)$$

Now! average over time!

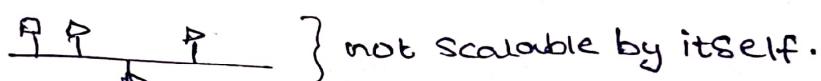
$$= \frac{A}{2}. \quad \therefore \text{we recovered } \underline{c_1(t)'s \text{ eqn.}}$$

Switching:-

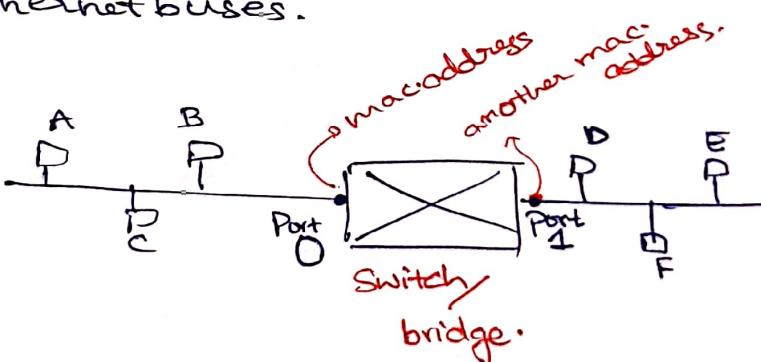
at Layer 2 (Ethernet).
also at Layer 3 \rightarrow routers.

- * use MAC addresses to switch, in layer-2.
 (Amplify only
necessary packets)
- use IP address; for layer-3.
- * the LAN bus; can have only so many devices connected...
 since more devices \rightarrow Much more collisions.

∴



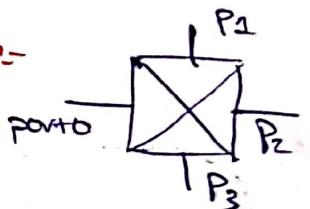
- * we use bridges/repeaters/L2 switches to connect different Ethernet buses.



every ethernet module manufactured, must have different MAC addresses

2^{48} possibilities.

- * The switch is supposed to send messages from one side source to other side destination, & ignore src-dest. on same bus.
 can have more than 2 ports!-



- for the switch to behave in such an intelligent manner, it needs FORWARDING TABLE.

This needs to be dynamic; bcoz, what if device A is removed from one bus, & connected to another bus?

FORWARDING TABLE:-

Dest.	Port#

1. Initially empty.

Dest.	Port#
A	0

2. Say a frame has come-in;

SRC	DEST
A	B

Now; it registers the sender, with the appropriate port.
append (A, 0)

3. Since B isn't in our forwarding table, we forward it via all other remaining ports.

4. Like this, switch builds up its forwarding table.

5. Set an expiryTimer on each entry. reset timer ^{to full} whenever we receive a frame, from the same MAC.

This allows the BRIDGE to delete unexisting MAC; & update port#

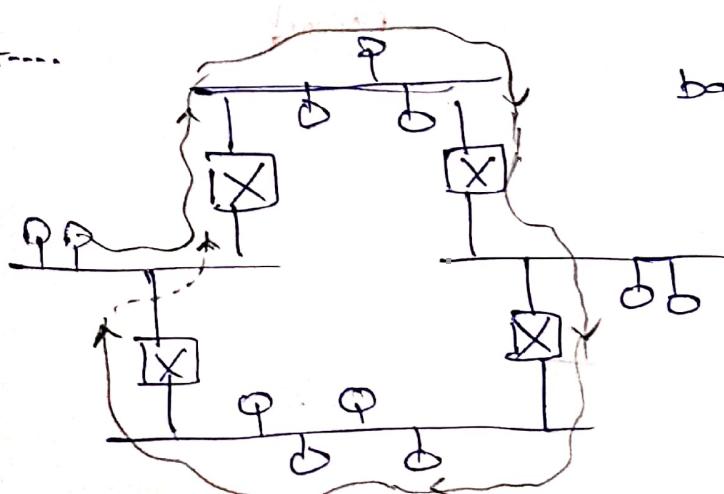
in case a device is connected to another.

What about collision detection?

* A huge issue- loops

what if a loop exists? Some frames will be repeatedly

Sent....



bandwidth will be exhausted like this.

→ SPANNING TREE PROTOCOL by radia perlman.

* every bridge

↳ has root port.

* Every LAN

↳ has designated port.

NO confusion in deciding the route of frames via BRIDGES.

→ any message; is sent into all LANs initially; & no loops!

eventually; the SWITCHES become better at deciding & then; we'll have efficient repetition of frame.

1. elect a root bridge. (one with lowest Bridge ID).

2. Every other bridge finds which Port is closest to root & make it rootport!

bridges have many ports

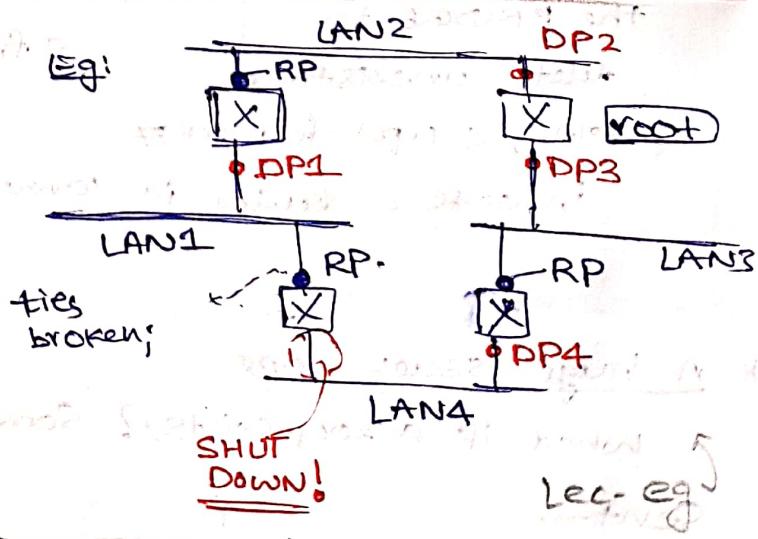
(we have lower neighbour switch ID, for tie-breaking)

3. All bridges connected to a LAN, ethernet bus.

elect one among them to forward frames of this LAN. (designated port) the one closest to root.

4. Any port; that is not a root port, nor a DP; is shut down.

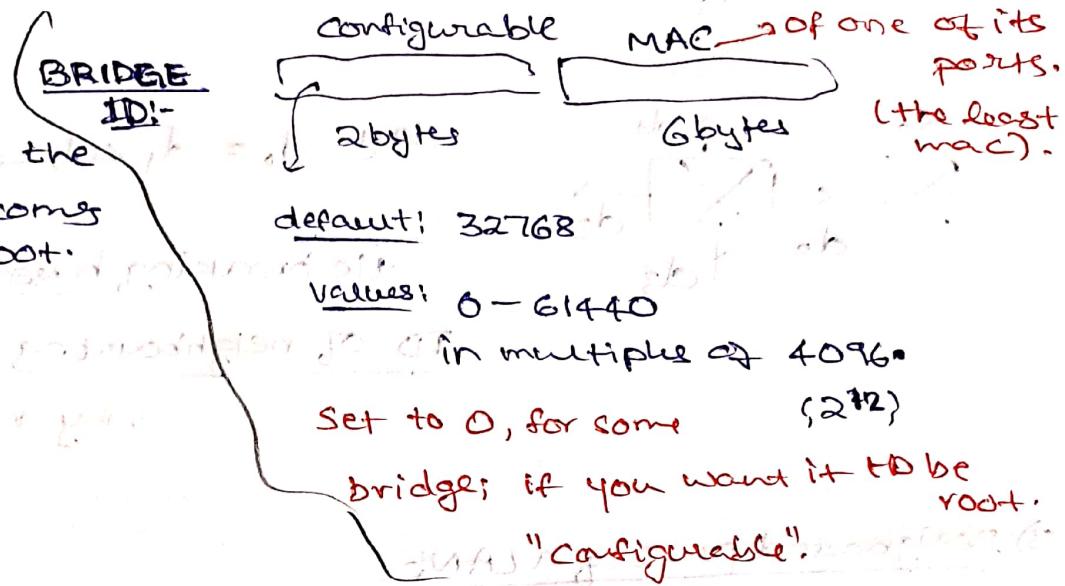
In this manner; we avoid loops.



Details:-

1) Elect root Bridge

each bridge has Bridge ID:-



- the bridge with the lowest ID becomes root.

Bridge ID is 16 bits, so it can have values from 0 to 65535. The value is stored in multiple of 4096.

Set to 0, for some bridge; if you want it to be root.

→ How do bridges talk:-

during the election algorithm:-

- each bridge broadcasts

(Y, d, X)
 minimum and d to the smallest ID till now
 distance from the smallest ID

∴ initial:-

SW1: (1, 0, 1)

SW2: (2, 0, 2)

!

SW4: (4, 0, 4)

SW1: (1, 0, 1)

SW2: (1, 1, 2)

SW3: (1, 1, 3)

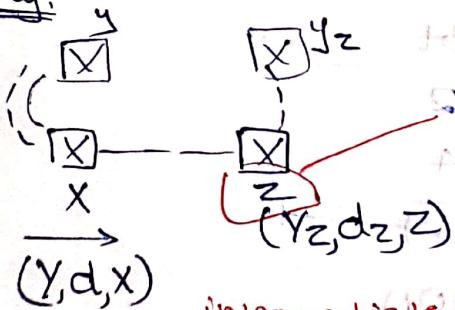
SW4: (2, 1, 4)

converge:-

1, 0, 1
 1, 1, 2
 1, 1, 3
 1, 2, 4

distances;
 also depend
 on wire
 speed.

say:



if $Y < Y_z$:

then $(Y_z, d_z, z) \leftarrow (Y, d + \text{dist}(X, Z), Z)$

if $Y = Y_z \wedge d + \text{dist}(X, Z) < d_z$:

varies
 based
 on
 throughput
 rate of
 ethernet.

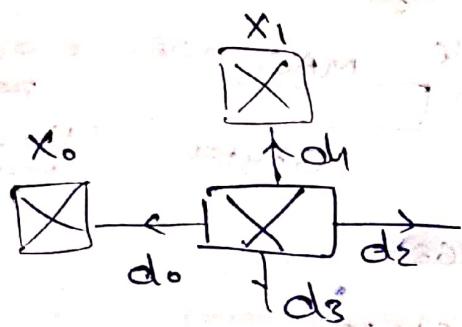
intermediate

step:-

$$d_z = d + \text{dist}(X, Z)$$

2) root port of a bridge

Smallest distance to root \rightarrow root port of bridge!



$$\text{say } d_0 = d_1 < d_2 < d_3$$

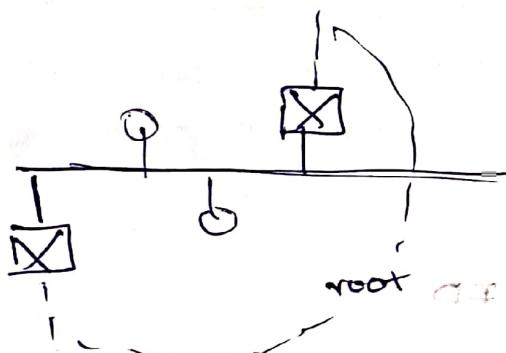
tie breaking, based on

ID of neighbouring switches of ports.

Why not MAC of ports themselves?

maybe some profit in Spanning tree..

3) Designated port of a LAN:



- See the distance of port from root.

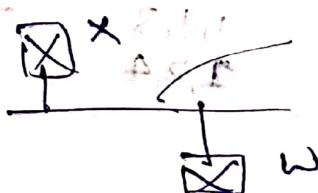
tie breaker:-

see the minimum SWITCHID of corresponding ports.

of bridges

Finally; PORTS which aren't "root port" or "designated port" are (with a switch) (with a LAN).

* in practicality;



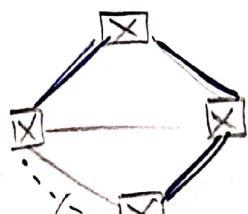
dist(X,W) \neq 1

depends on speed of wire, etc.

Speed	Cost
10 Gbps	2
1 Gbps	4
100 Mbps	19
10 Mbps	100

• okay! But is this it? Will layer-2 switching help me interconnect globally?

→ Shortcomings of Ethernet switching:



direct port,
not used!

1) path choosed is not optimal

2) Some links are not used at all.

3) Forwarding tables are $O(n)$

bcoz FLAT ADDRESSING.

no meaning in
MAC address.

Layer-3 (IP): Hierarchical Addressing

4) Stability is an issue for large SPT.

detected when
'Hello' messages
stop.

If some switch/link fails; whole SPT is reconstructed.

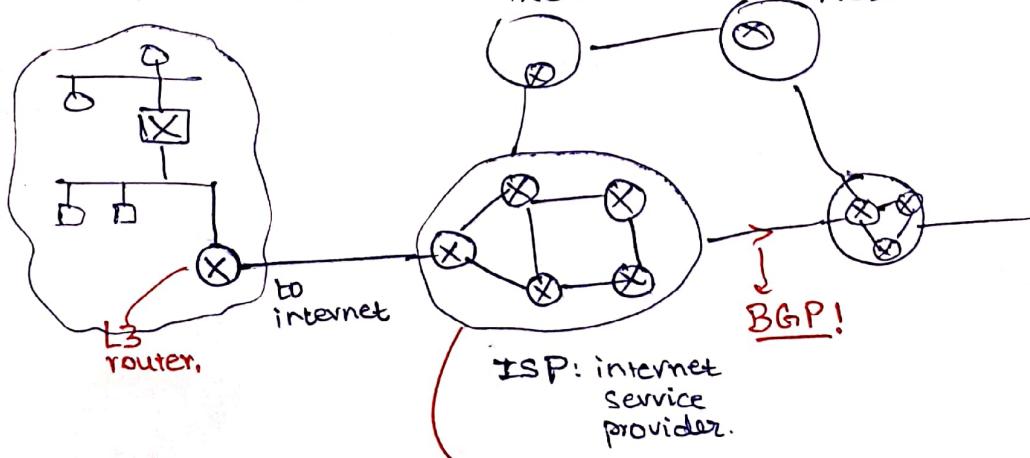
large Network → more reconstruction.

Some more issues:

- No common addressing Scheme. Across globe

- No common communication protocol.

→ The bigger picture:



Autonomous systems AS.

- can choose any internal routing protocol.

- Eg: least hops, least latency...

* Intra-domain routing : within AS.

will enable us
to go less than
 $O(n)$ tables.

Inter-domain routing: B/w AS.

BGP: border gateway protocol.

Hence Global Scalability.

→ How do routers exchange information between them? What is the job of routers?

→ Forwarding packets to destination

→ INTRA-DOMAIN ROUTING:-

→ Principles

1. Distance Vector based : RIP

transmit forwarding tables routing Information protocol
using Bellman-Ford algorithm

2. Link State Routing: OSPF

transmit neighbouring open shortest path first
topology information IS-IS

→ Intermediate System to Intermediate System.

* All these will

- use shortest path

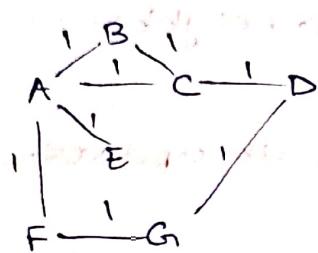
- Avoid loops.



→ Distance vector based protocol

→ Forwarded to destination IP

Distance-Vector routing: We just store the next node to be used, to reach dest.



* A sends out its 'table' to neighbours.

not BROADCAST.

A:

Dest.	Next hop	distance to dest
A	-	0
B	B	1
C	C	1
E	E	1
F	F	1

initially will be there.

hears from B:-

B	-	0
A	A	1
C	C	1

from C:-	
C	-
A	A
B	B
D	D

* Now; after hearing from neighbours, A updates its table:-

A	-	0
B	B	1
C	C	1
E	E	1
F	F	1
G	F	2
D	C	2

Bellman-Ford Algo.

Like this; converges for every node.

∴ we have the forwarding tables for every node.

(↳ these will be shortest paths)

→ if G-Flink is down!

G,F both send

for G:

Dest	Next hop	dist
A	F	∞

→ node D will hear this;

for F:

Dest

G

dist
∞

& update accordingly.

As a chain reaction, all nodes' tables get updated.

- Trigger updates: Event and link may need updating, rather than periodic.
- Some Event triggers a routing update to neighbours
- Eg: ~~F-G~~ link failure TRIGGERS an update.

- Periodic updates:

periodically keep sharing the routing tables.

0		A
1		B
2		C
3		D
4		E
5		F
6		G
7		H

→ Count to infinity Problem:

* X-A-B after convergence:

A:	dest next cost	B:
X	2	X A 2
B	1	A A 1

Now X-A is broken:

A: $\xrightarrow{\text{dest next cost}}$ sends trigger update
B: hasn't updated yet

X - ∞ X A 2

sends B to B + 1 \Rightarrow A-A = 1 (infinity)

sends B to A $\xrightarrow{\text{dest next cost}}$ sends next cost.
periodic update.

A:	dest next cost	B:
X B 3	$\xrightarrow{\text{dest next cost}}$ sends	X - ∞
B B 1		A A 1

A: $\xrightarrow{\text{dest next cost}}$ sends	B: $\xrightarrow{\text{dest next cost}}$ sends
B: $\xrightarrow{\text{dest next cost}}$ sends	A: $\xrightarrow{\text{dest next cost}}$ sends
A: $\xrightarrow{\text{dest next cost}}$ sends	B: $\xrightarrow{\text{dest next cost}}$ sends
B: $\xrightarrow{\text{dest next cost}}$ sends	A: $\xrightarrow{\text{dest next cost}}$ sends

3	$\xrightarrow{\text{dest next cost}}$ sends	4
5	$\xrightarrow{\text{dest next cost}}$ sends	6
7	$\xrightarrow{\text{dest next cost}}$ sends	8
8	$\xrightarrow{\text{dest next cost}}$ sends	9
9	$\xrightarrow{\text{dest next cost}}$ sends	10

Count to ∞ .

RIP protocol sets 16 as infinity.

i.e. can't reach destination

Propositions:-

split-horizon:-

Don't advertise information about destination to the neighbour, if it is the next hop for that destination.

So; B doesn't tell A about X.

problem solved! is it? what about 3-member loop?

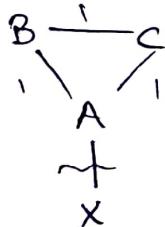


split-horizon with poison reverse:-

Node tell its next hop to the destination, that distance to destination is ∞ .

So; B tells (X, ∞) to A.

*



now; 3-membered loop is coming, even with split-horizons.

→ Routing Information Protocol RIP:-

- Distance vector based.
- cost of all links = 1.
- max distance = 16 → infinite.

Advantages:- Simple to implement. (not in my opinion)

Disadvantages:-

- count to infinity & routing loops
- convergence of routing tables is slow.
- since 16 is max, can't have large networks.

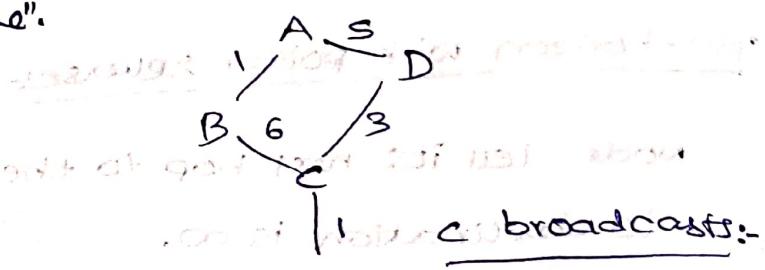
→ Link-State Routing

Each node broadcasts information about cost to its neighbours.

∴ each node has a picture of network topology completely.

* Apply beautiful dijkstra's for min paths, from that node to graph.

& decide "next node".

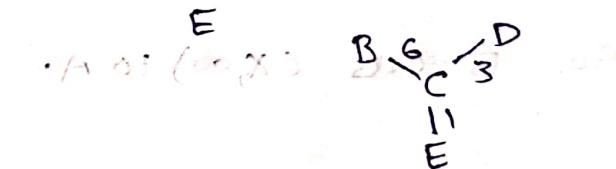


→ on link failure

* $A \cancel{\rightarrow} F$

A, F broadcast that their link has failed.

→ All nodes re-run Dijkstra's.



1) Advantages:

- no cost-to-infinity, routing loops.

- convergence of routing tables is quick.

2) Disadvantages:

- A bit complex than DV. (I personally don't think so)

Finally, what should be the cost on the links?

Speed?

Latency?

Queuing?

Ethernet bus;

joining two points.

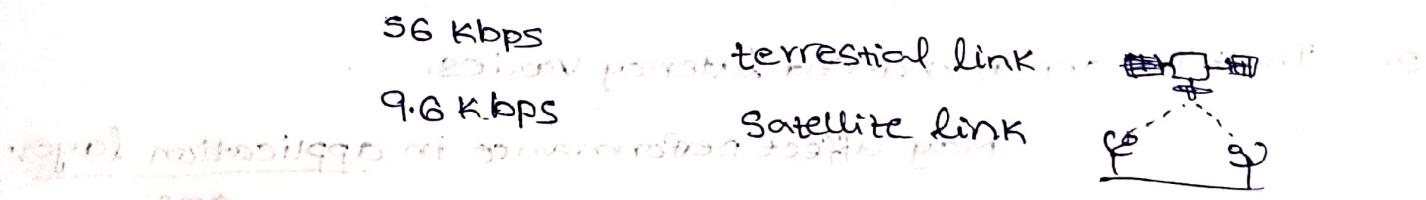
→ choosing link weights:

* DV; LSR choose shortest paths.

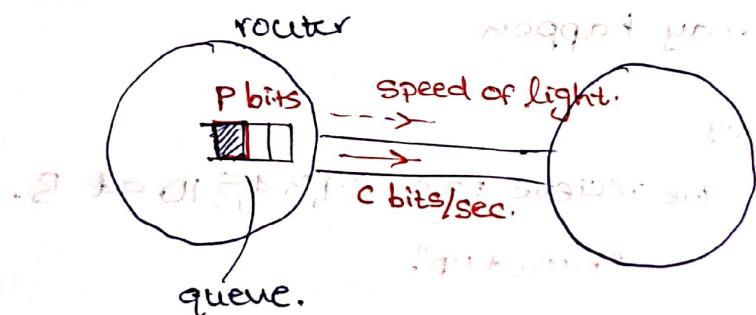
Hence, every Network Engineer has to decide... what are the weights...

lets ask god aka ARPANET (lol)

had two basic link speeds



* A link:-



$$\text{latency} = \text{queuing delay} + \text{transmission delay} + \text{speed of light delay.}$$

latency = queuing delay + transmission delay + speed of light delay.
taken time for last bit of message to reach demand for the link.

- well, looks perfect analogy for cost on links.

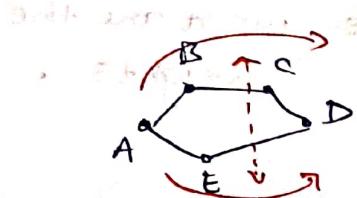
but latency is sooo dynamic.

→ take time window average, latency avg. of all pkts

is the link weight.

* Many problems arised, by this analogy...

1. Under heavy load, routing paths oscillate



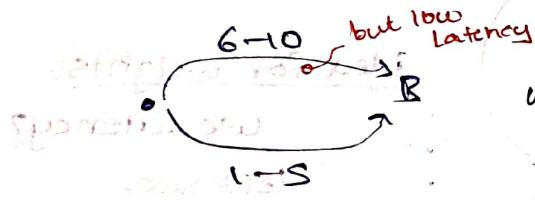
Queuing goes up \rightarrow increase link weights



2. In this manner, end-end latency varies.

May affect performance in application layer.

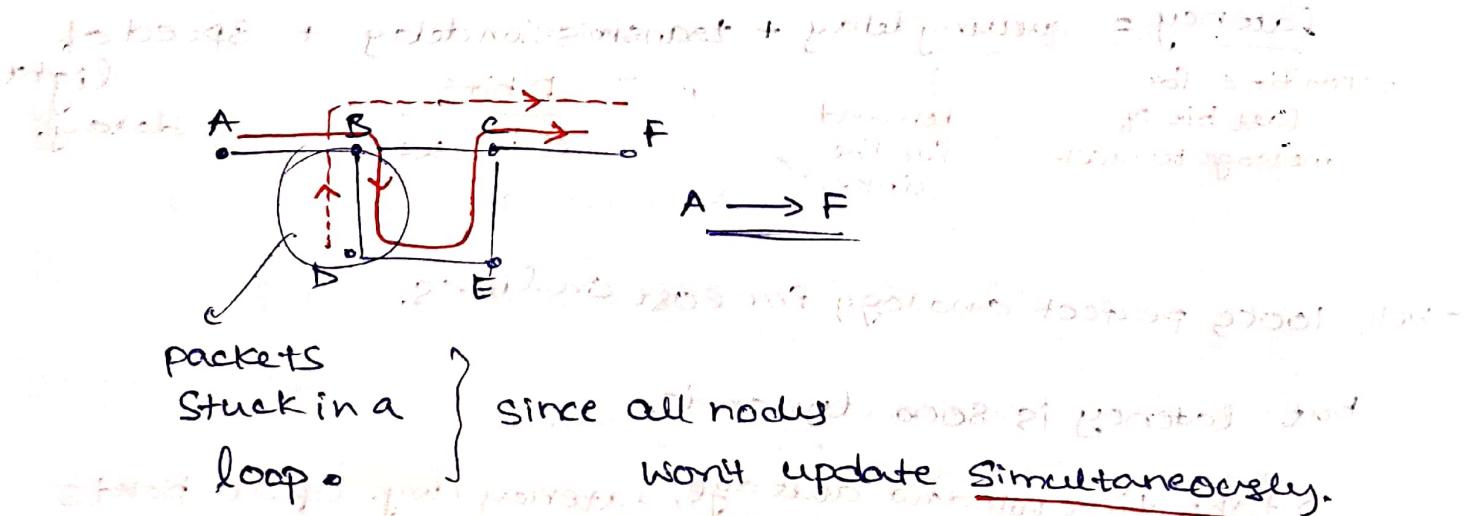
3. Reordering of packets may happen



we receive 1, 2, 6, 7, 8, 9, 3, 4, 5, 10 at B.

"missed up".

4. Routing loops are possible....

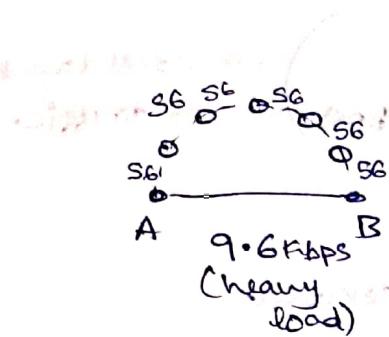


After convergence, we won't have these loops.

Approaches for finding consistent paths

5. Range of link weights were too wide.

Due to which some links were penalized too much.



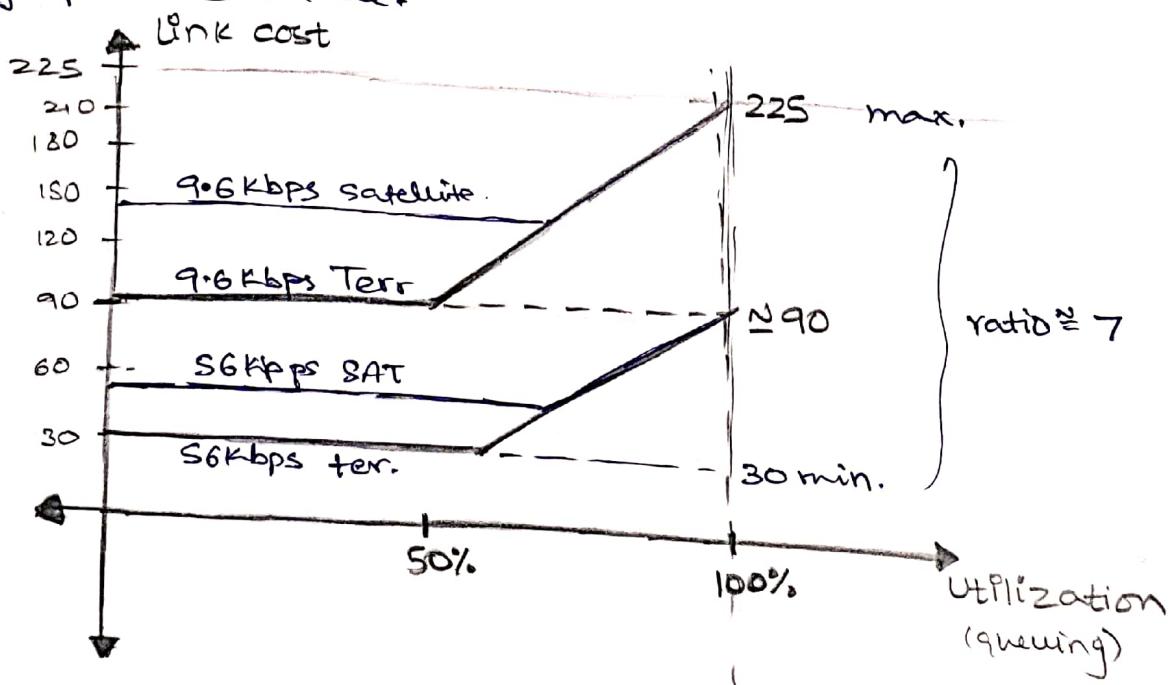
126 links of 56 kbps had same weight as one 9.6 kbps.

why use 126 entities, when we have a direct edge.

6. Satellites were penalized too much. (low utilization)



Hence, graph was made.



& weights changed infrequently.

Today;

OSPF protocol:-

based on link state routing

$$\text{wt. of link} = \max(1, \frac{10^8}{\text{line speed(bps)}})$$

no consideration
of traffic.

then how?

Setup a NOC - network operations center

change weights due to traffic manually!

Creating static and dynamic path selection based on traffic conditions

dynamic path selection by protocols

Protocol based on OSPF, BGP, IS-IS, RIPv2

Bandwidth

link weight = $\frac{10^8}{\text{line speed}}$

Latency

link weight = $\frac{10^8}{\text{line speed}} \times \text{latency}$

Link weight = $\frac{10^8}{\text{line speed}} \times \text{latency} + \text{latency}^2$

Link weight = $\frac{10^8}{\text{line speed}} \times \text{latency} + \text{latency}^2 + \text{latency}^3$

Link weight = $\frac{10^8}{\text{line speed}} \times \text{latency} + \text{latency}^2 + \text{latency}^3 + \text{latency}^4$

97/128

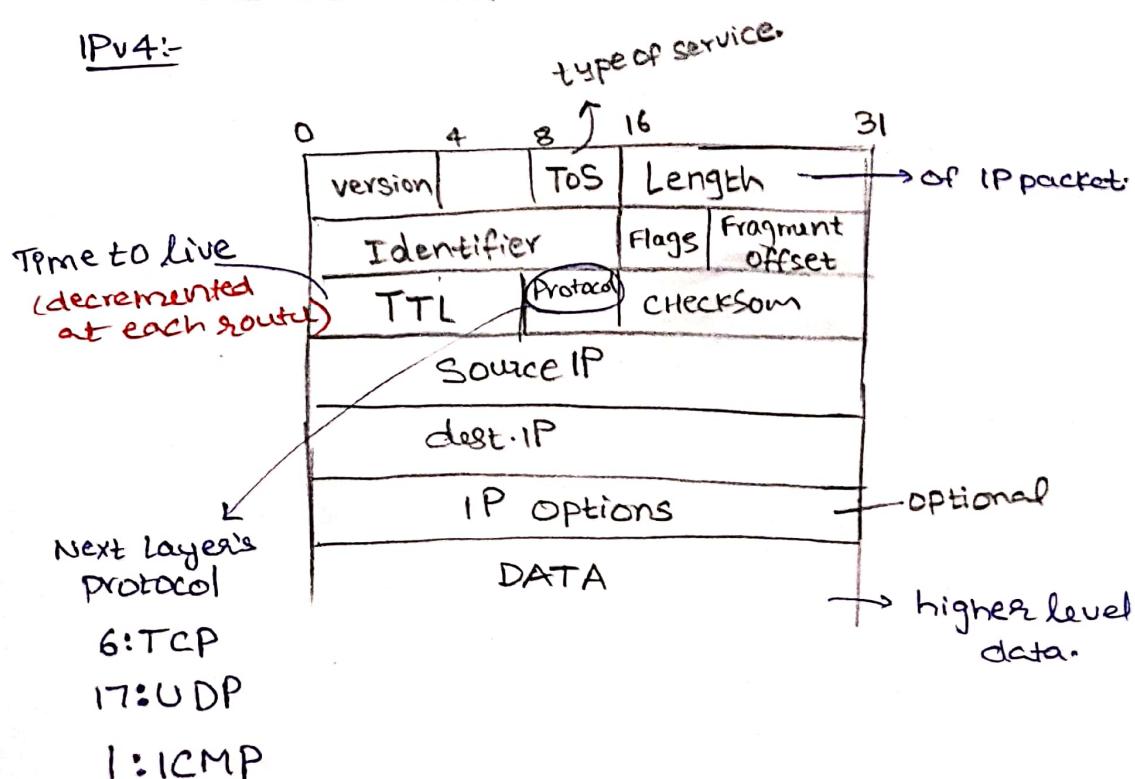
4. 1951-1952

Q18:- IP Addressing:- Layer-3: Network layer

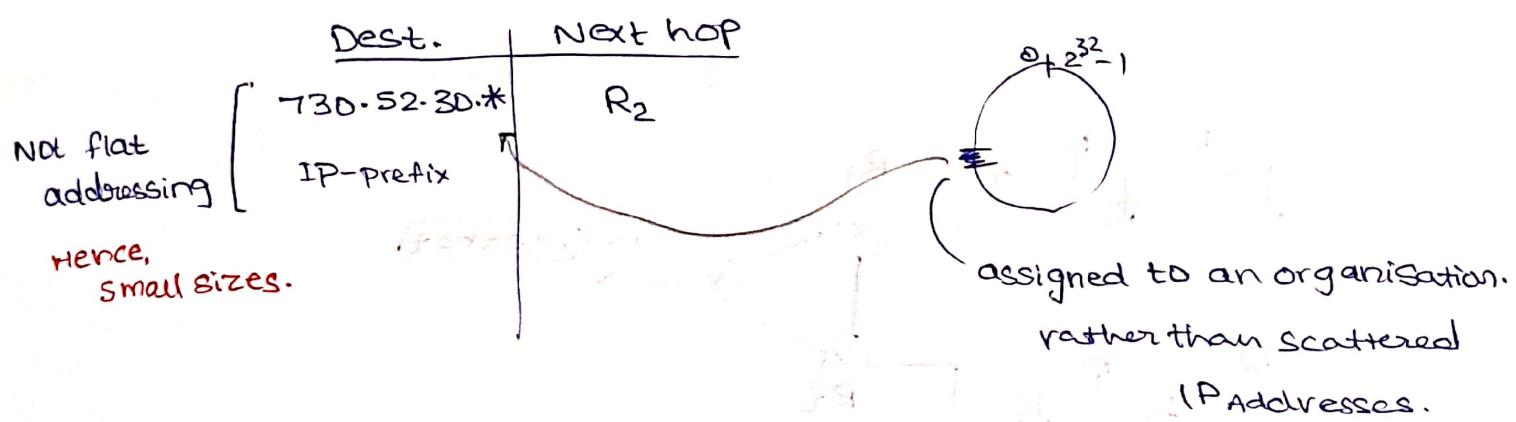
- * MAC address (48 bits) for layer 2.] Built into NIC.
- * IP address (IPV4: 32; IPV6: 128 bits) for layer 3.] assigned via software.
- * 2^{32} not at all sufficient. → Well use something called NAT; (Network Address translation)
- * IP address as
 $(8\text{bits}) \cdot (8\text{bits}) \cdot (.) \cdot (.)$
- * Ex: 255.255.255.255 } all 1s; reserved for broadcast.
 ➤ blocked by routers.
 ➤ works for a network.
- * 10.*.*.* }
 192.*.*.* } private IP address
- * public IP address } Should be unique, only one host must use this IP addr.

→ Layer 3 headers:-

IPv4:-



* Routing tables of routers be like,



Like machine's IPs be like

730.52.30.1

730.52.30.2

Network part | host part

Big networks have lots of hosts

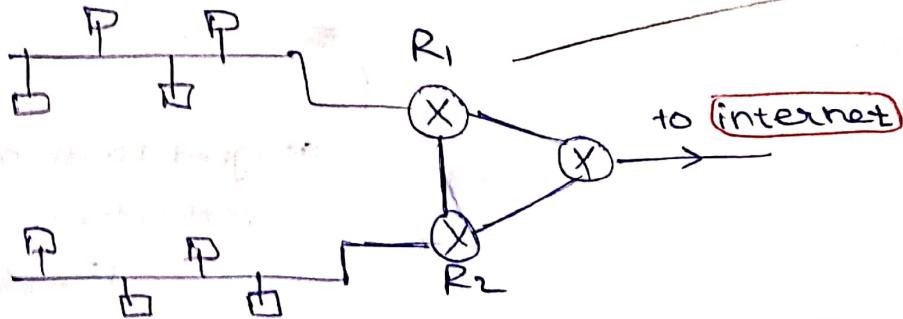
* class A: big MNCS 8 bits 24 bits 2^8 devices!

class B: 16 bits | 16 bits

class C: 24 bits | 8 bits, 2^8 devices here...

→ Subnetting:

- * we are given a 'slice' of IP addresses, how to divide among LANs



maybe this is Class C:-

24 bits
Network (fixed!) {
1 for R1
2 for R2 } we are subnetting!

* Subnet mask:

says which bits in IP address to use,

to decide which LAN to route to.

Eg: to consider first 25 bits:

MASK = 11111111.11111111.11111111.10000000 (binary)

255.255.255.128 (human readable notation)

* Subnet address - Is this unique?

S₁ for LAN1 (M1 is mask)

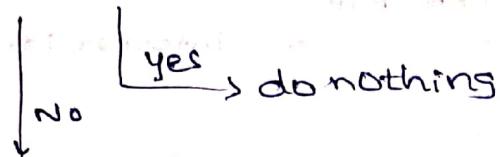
S₂ for LAN2

* Routers be like:-

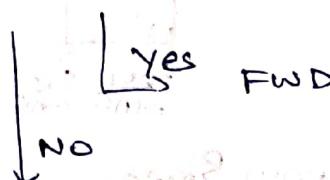
Router R1:-

dest. IP address is 'D'.

IS $(D \text{ AND } M_1) == S_1$?

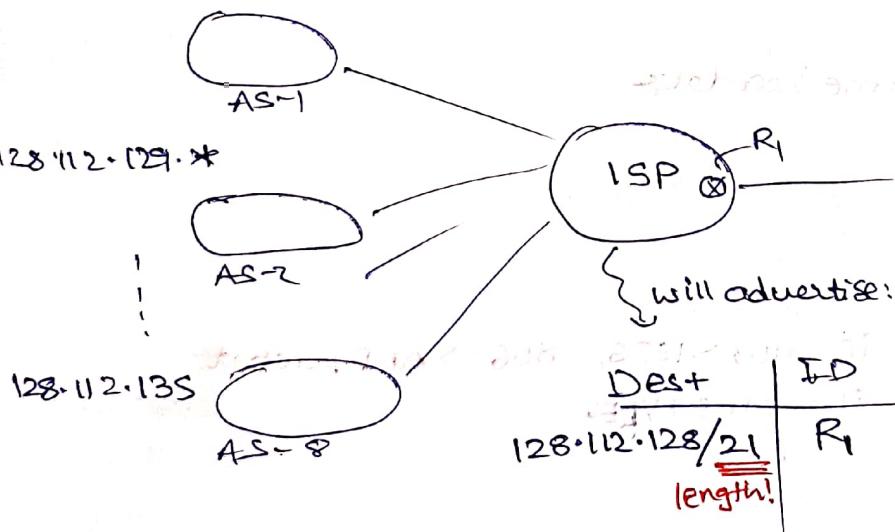


IS $(D \text{ AND } M_2) == S_2$?



→ Super netting :- helps to reduce size of forwarding tables.

128.112.128.*



128: 1 0000000
129: 1 00000001
130: 1 00000010
135: 1 00000111
supernet!

* IP prefixes

a.b.c.d/N

first N bits for prefix.

* CIDR:- Classless Inter Domain Routing

no class A B C ...

- use any prefix length (N)

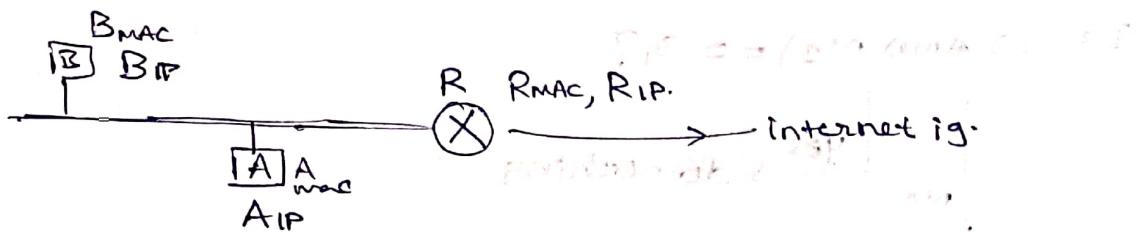
L-19 ARP & DHCP :-

Assign IP to host

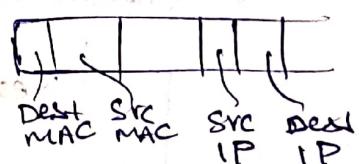
(find MAC of dest.)

→ Address Resolution Protocol :- Remember, IP addresses are software assigned.

whereas MAC is built in.



frames will be:-



* A wants to send to B

A Knows BIP; but doesn't know BMAC.

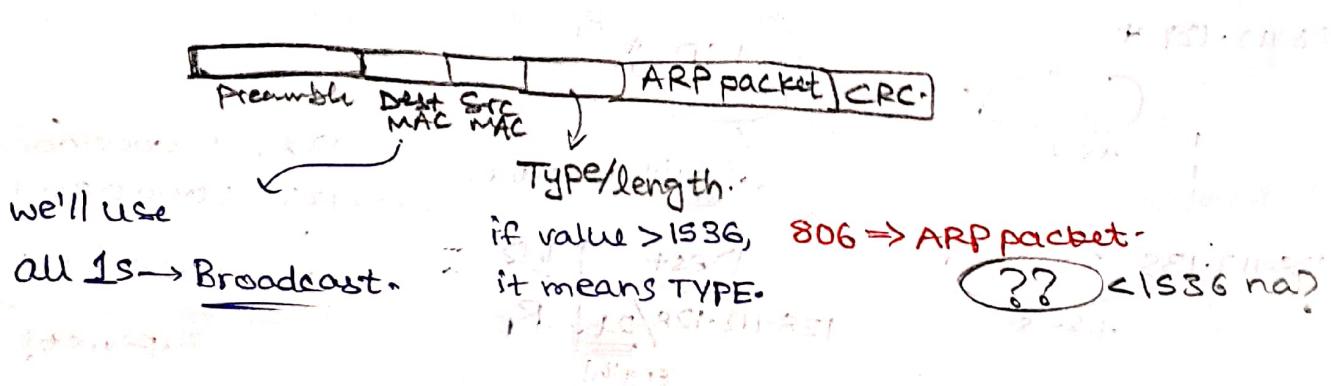
How? Server?

→ ARP to the rescuer:-

Where's ARP in the stack?



* Recall the ETHERNET frame header:-



* ARP request (by A):-

(A → all)

SENDER MAC; SENDER IP;

A mac

A IP

Target MAC Target IP

00000000--

BIP

* ARP reply
(from B to A)

unicast here.

Sender MAC , Sender IP

Bmac

Bip

Target MAC , Target IP

Amac

AIP

Host A

* Information: B_{MAC} ; stored in ARP cache of A, with a timeout.

Q1) How does A know, if dest. IP belongs to same network?

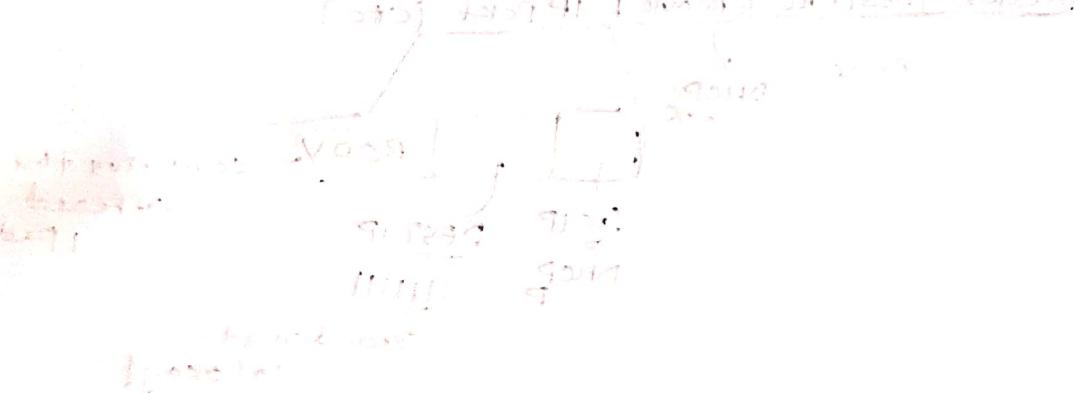
A: using subnet mask.

if $(\text{DEST IP}) \text{AND} (\text{MASK}) = (\text{AIP}) \text{AND} (\text{MASK})$
then same network.

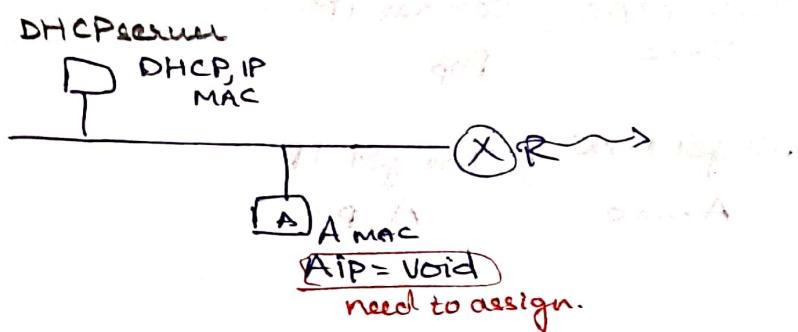
Q2) If DEST IP is different; we'll need routers RIP, RMAC.

A: RMAC; by using ARP.

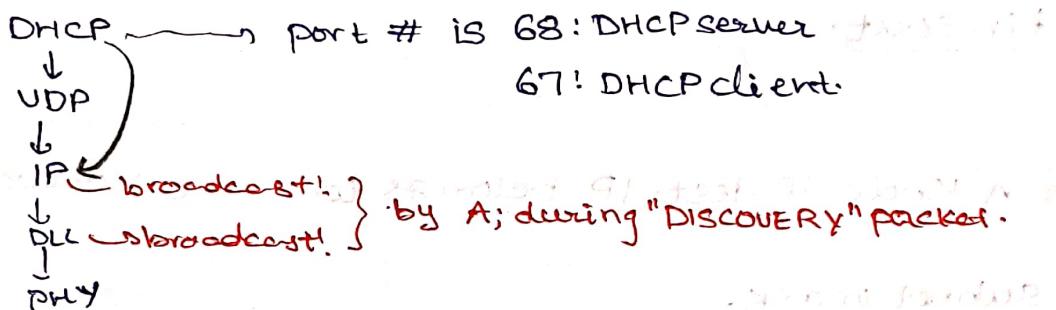
RIP, RMAC.
we'll get it
from DHCP
server



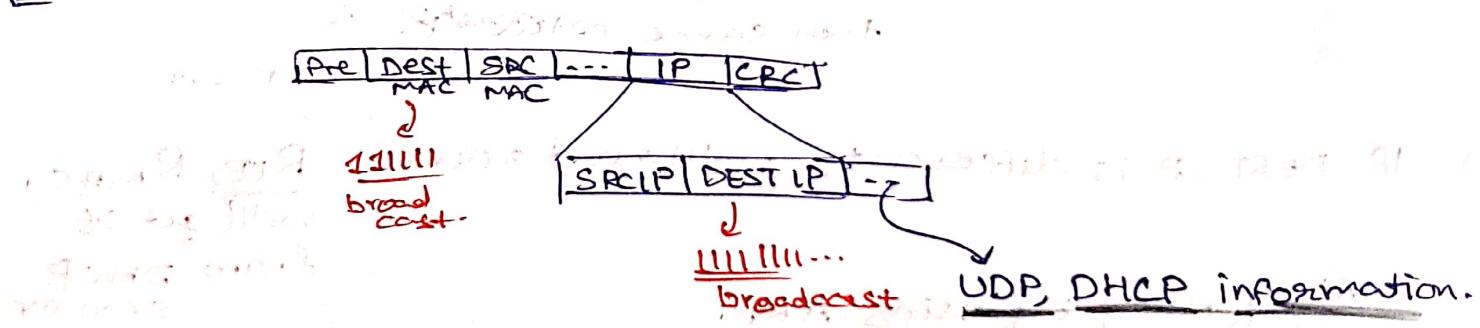
→ Dynamic host configuration protocol:-



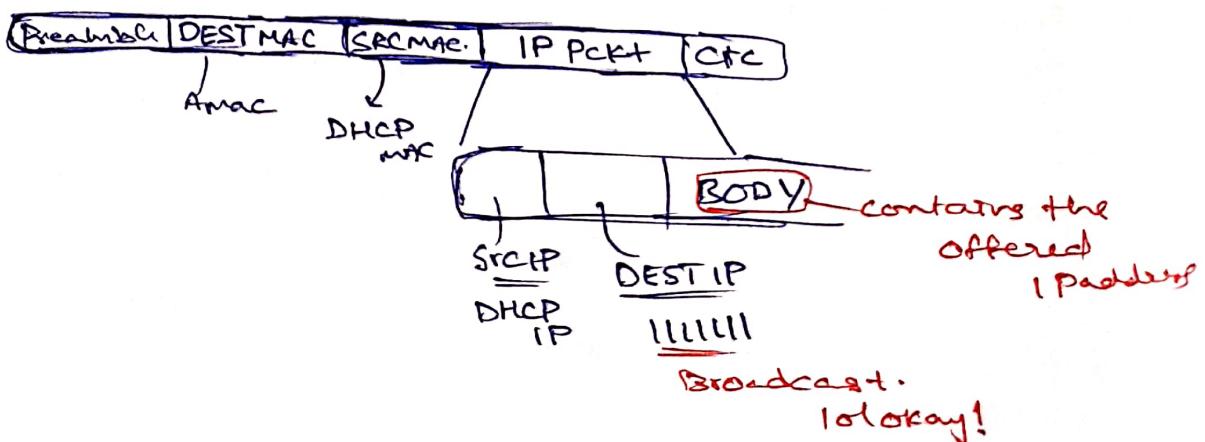
* in protocol stack:-



1. A sends out a 'DISCOVERY' packet:-



2. DHCP server replies with an Offer-

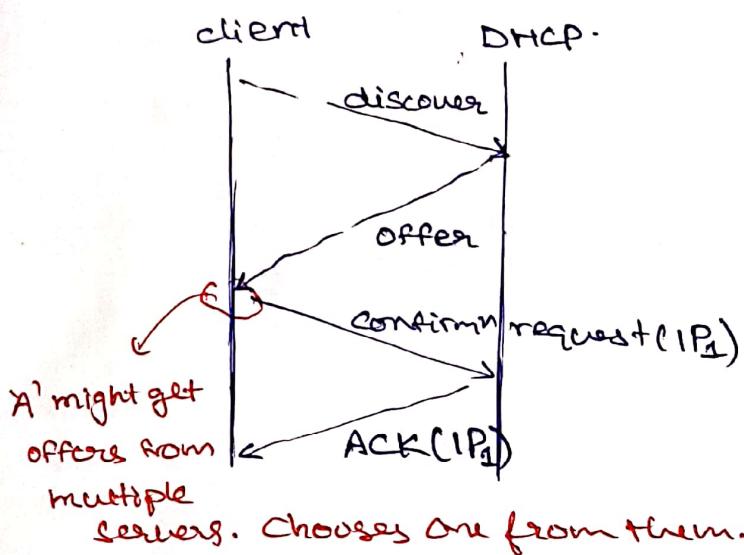


30

Host 'A' checks whether offered IP is taken... by ARP.

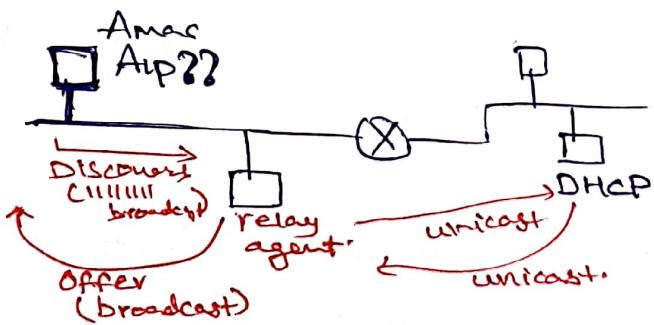
& then sends

confirmation
request.



* In absence of DHCP server; we'll have a relay agent; bcoz L3 routers, block broadcast signals!!
 (L3 broadcast)

Otherwise,
 RIP internet.



1986-08-20 10:00 AM - 10:30 AM

Wet day outside

Cloudy overcast

100% RH

90% RH

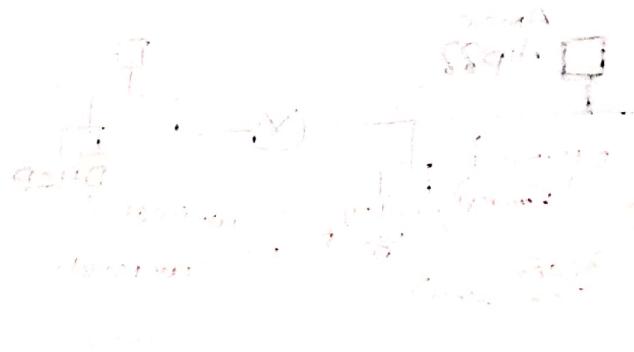
100% RH

80%

70% RH

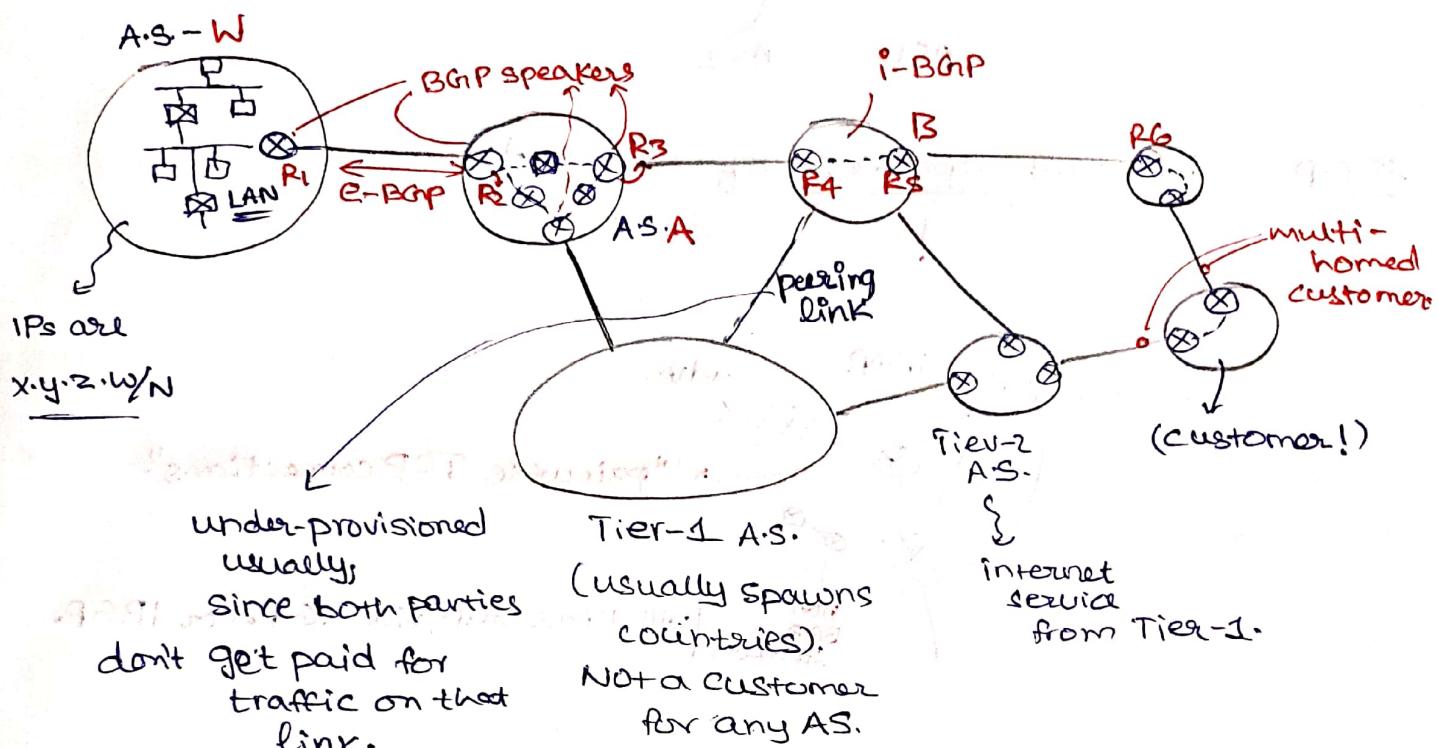
Small orange patches and yellowish-green. The leaves are still green.

Leaves are back to normal, greenish-yellow
Orange spots



Lec 20:- Border Gateway Protocol :- also HW3.

> Interdomain Routing: Between ASes. ^{for} the glue that brings the world together.



BGP advertisements:-

eBGP $R_1 \rightarrow R_2$: $x.y.z.w/N$

Prefix A.S. Path other

• Best attributes

$P_0 \rightarrow P_1$: number of nodes in P_1

$$R_3 \rightarrow R_4 : \underbrace{x.y.z.w/N}_{\text{---}} \quad \underline{A-N} \quad \underline{\text{---}}$$

prefix suffix other attr.

$$R_5 \rightarrow R_6: \quad x \cdot y \cdot z \cdot w / N \quad B - A - w \quad ,$$

Smith

path.

Uroplatus fimbriatus (Gmelin) (partim)

meaning: if RG gives a pkts to

الآن، أنا أعلم أنني أنت

for sure this path'll be followed.

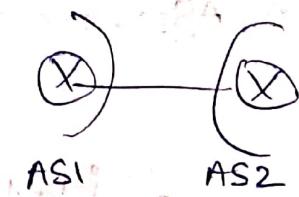
Le point de l'heure

* Tier-1 ASes; don't need to use peering links for

Flexibility is there.../

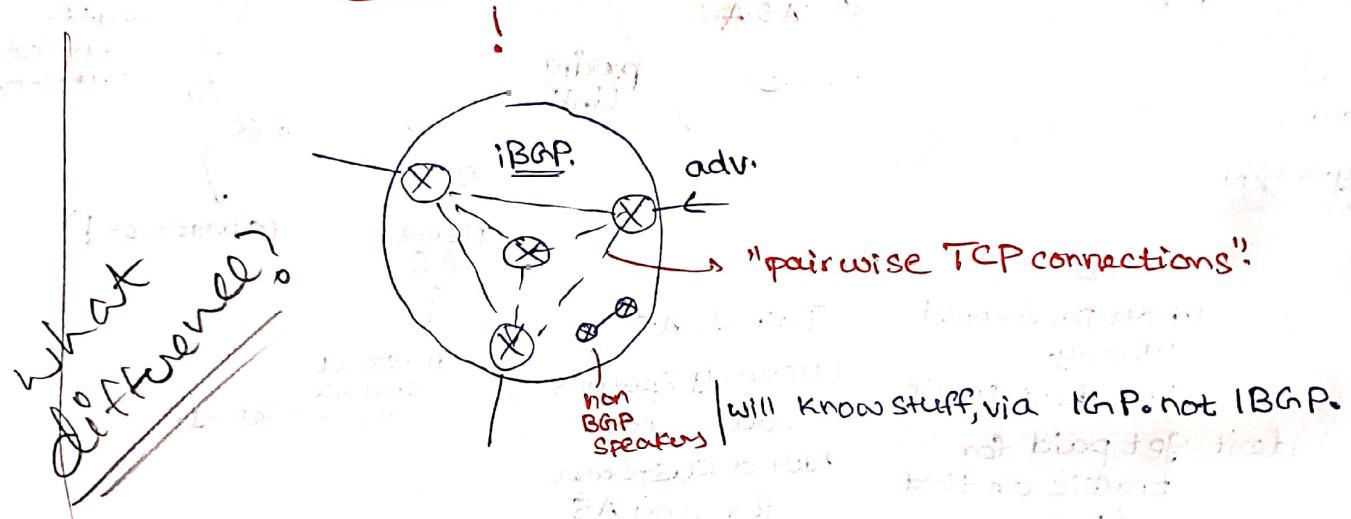
1974-1975
Yearly average rainfall

eBGP :- Protocol between BGP speakers in neighbouring A.S.



TCP # 179.

iBGP :- b/w BGP speakers of same A.S.



IGP :- Interior Gateway protocol } intra-domain routing!

Between every router of an AS. based on whether speaker or not.

DU, LSR

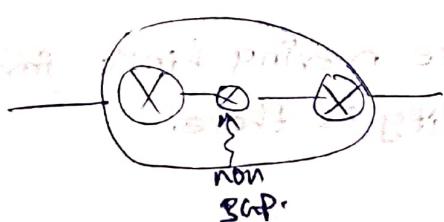
these are concepts.

NOT protocols.
hummm...

→ Steps to determine BGP paths:- Protocol:-

- 1) e-BGP speakers learn AS-paths from neighbours.
- 2) e-BGP nodes share learned information via iBGP with other nodes inside A.S.
- 3) BGP speakers choose routes to various IP prefixes.
- 4) The speakers; insert chosen routes into IGP. Since all routers; whether speakers or not; should be able to send pkt to every IP destination.

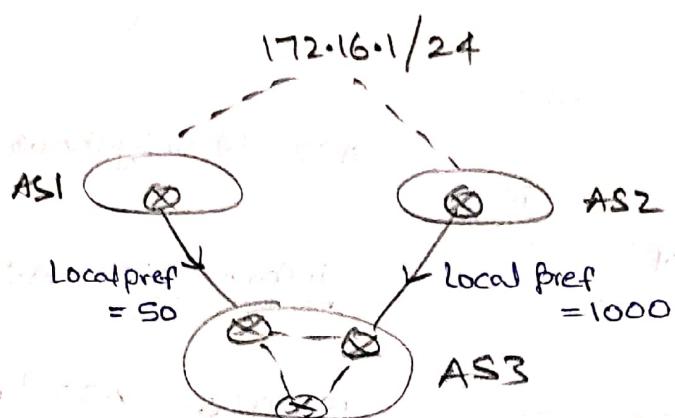
5) eBGP speakers can advertise



newly inserted routes to neighbouring AS.

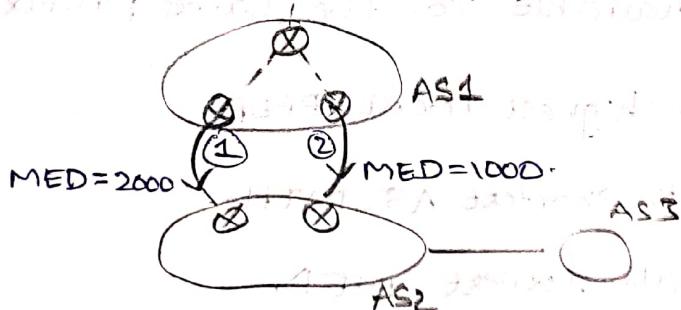
→ BGP attributes:-

1) LOCAL PREF:-



- > Admin of AS3 adds LOCAL-PREF himself. Usually left default.
- > Higher pref → more preferred path.

2) MULTI-EXIT DISCRIMINATOR:- MED

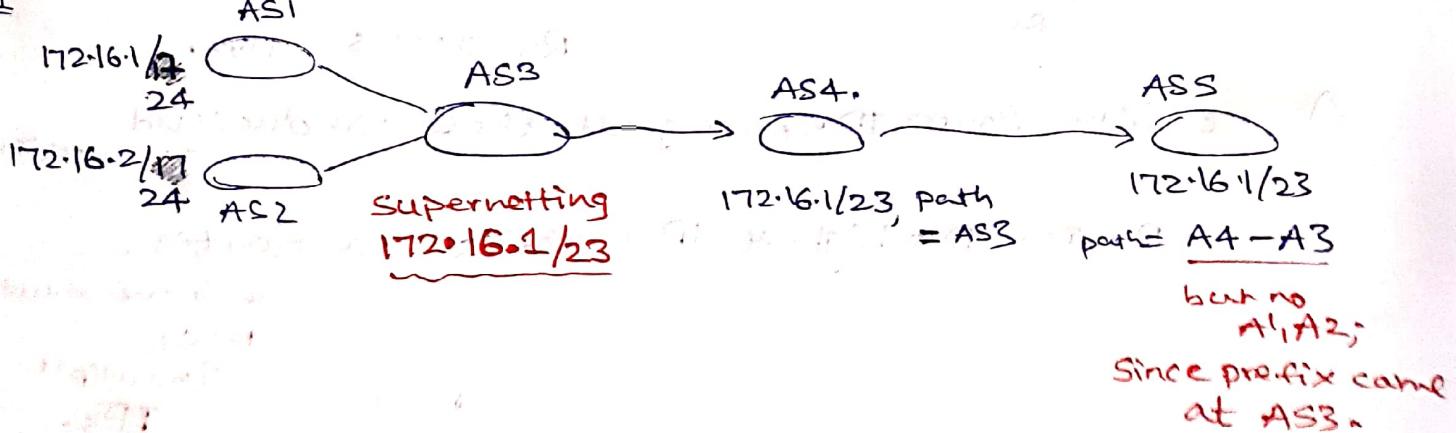


- > AS1 tells AS2 that use ② over ①. Lower MED → more preferred path.
- > lower MED, more preferred path.

3) AS-PATH:-

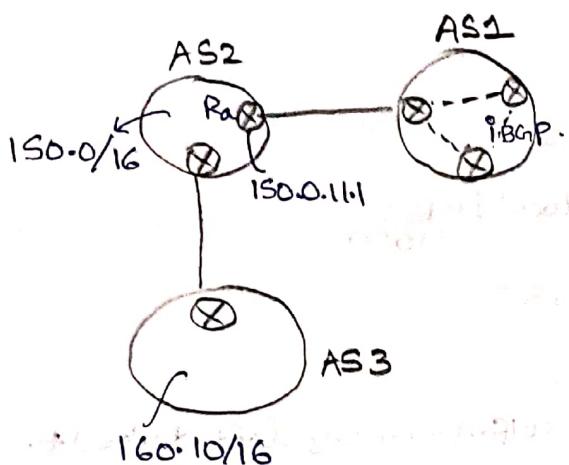
List of AS; all the way to the destination, who advertised that prefix first.

Eg:-



4) NEXT_HOP:-

the ip address of the BGP speaker; in the neighbouring AS.



AS2 is neighbour of AS1:-

information AS1 has:- (circulated by iBGP)
(forwarding tables)

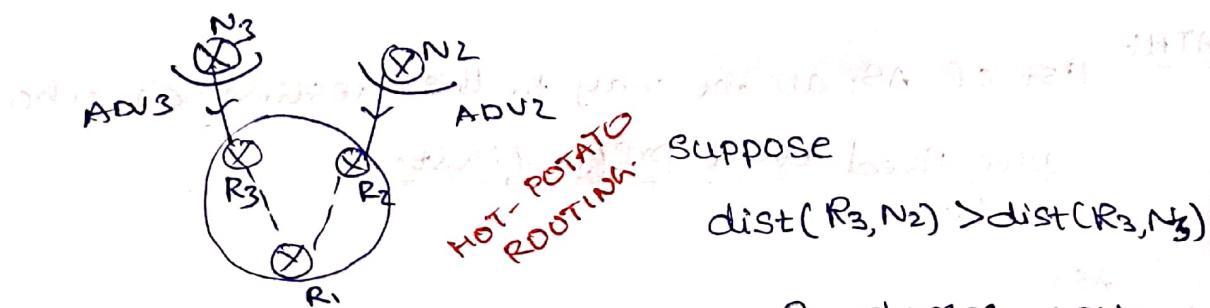
150.0/16 | AS2 | NEXT_HOP = 150.0.11.1

160.10/16 | AS2=AS3 | NEXT_HOP = 150.0.11.1

→ Rules to choose routes, in BGP:-

- * each BGP Speaker decides which AS-route to use, among many available for the same prefix.

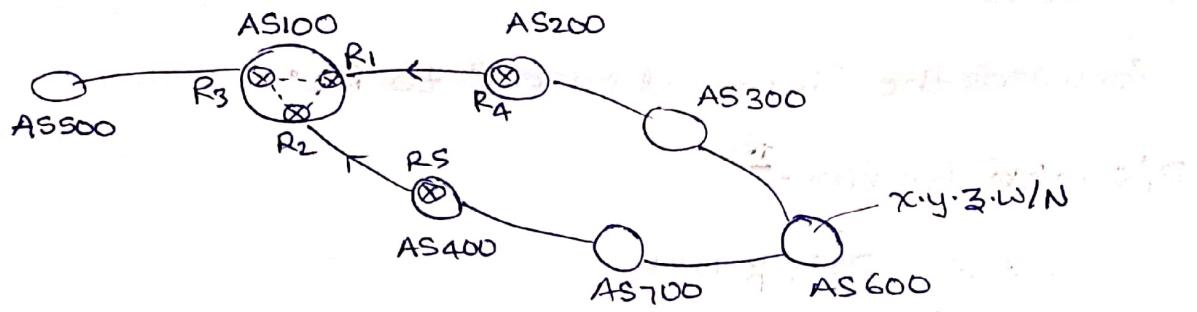
- a) use route with highest LOCAL-PREF.
- b) choose path with shortest AS-PATH.
- c) choose path with lowest MED.
- d) choose path learned via EBGP over one learned via iBGP.
- e) choose path with lowest IGP metric to NEXT-HOP.



- f) use least Router ID. (Among all, those who've sent advertisement)

RouterID = highest IP address on a router,

can have multiple NICs
multiple IPs.



* Look inside AS100!

Say both paths have same LOCAL-PREF and MED

Path 1: AS200 - AS300 - AS600

Path2: AS400 - AS700 - AS600

R_1 chooses path 1: eBGP over iBGP as the best route

R₂: chooses Path 2: $\text{B} \rightarrow \text{C} \rightarrow \text{D}$ at $t = 0.25$ s and $t = 0.5$ s.

R3: ? : both paths learned over iBGP.

if $\text{dist}(R_3, R_i) < \text{dist}(R_3, R_2)$ then R_3 goes for path 1.

Hot potato routing.

AS100 interior:-

Can Ra talk
with internet?
lets work
on it.

Solutions:-

Also; How R₃ will send to R₁?
Uation:- Same as below
(16P)

1) Encapsulation:-

BGP interaction with
IGP?

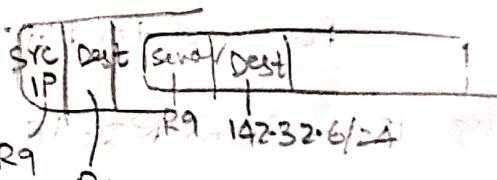
Speaker
Speaker

Say R9 has a destIP as 142.32.6/24 for a packet. But it has no BGP information. But, every router gotta send packets to world, not only speakers.

Ex: Rg's FWD table:-

Now Rg encapsulates this packet by a new IP layer :-)

DST	Next HOP
R1	R11
R2	R12
R3	R3
R11	R11
R12	R12



present in FWD table!

(PTO)

right
be
wrong.
Re is
internal.
Re will
do this.

- R_i receives this packet, strips off outer layer, forwards the "internal packet" to R₄.

R_i's table be like:-

DST	Next
R ₂	
R ₃	
:	
142.32.6/24	R ₄

∴ Encapsulation means: Non-speaker routers have no clue about outside IP prefixes, but they do know that, they need to ask BGP-speaker router for help, & thus forwards their packet to speaker router.

2) Pervasive BGP: All this work to allow internal routers to talk to internet. All routers are treated as BGP speakers.

* Say there's a unique exit for 142.32.6/24.

E also for interaction of BGP and IGP. (Speaker talk).

R₉s:-

BGP table:-

Prefix	Gateway/Exit
142.32.6/24	R ₁
:	

IGP FWD table:-

Dest	Next
R ₂	
R ₃	
R ₁	
R ₄	
R ₁₂	

forward to R₁₁.

- On receiving IPpkt from R₃; if Dest.IP = 142.32.6/24; then R₉ will do above thing. (Recursive lookup)

I understood wrong at first.

these solutions are for BGP-IGP interaction.

I thought Internal Router - BGP
interaction...

3) Tagged IGP:-

- * Internal routers may not be BGP speakers.

But IGP allows addition of tags.

- * R₁ can insert, into its own IGP. Not a new BGP table.

142.32.6/24
Prefix

R₁ → gateway router
TAG

insert into IGP.

propagated to all routers, by IGP (or LSR).

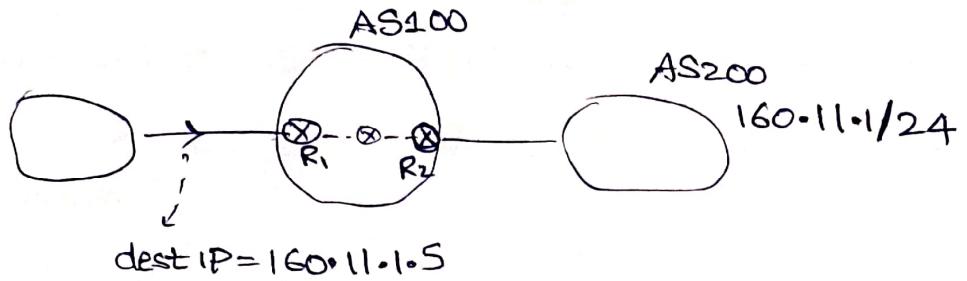
Ex: R₂'s IGP table be like

DST	Next	Tag	Cost
R ₁	R ₁₁	-	S
R ₂	R ₁₂	-	T
142.32.6/24	-	(R ₁)	1
142.32.6/24	-	(R ₂)	2

Now; R₂ will send a pkt: DestIP = 142.32.6/24

then R₂ will look up; find matches & forward to closest "tag".

Hot potato.



NOW; R₁ receives this packet.

- R₁ looks up BGP table & decides to Fwd to R₂.

But How! To calculate next hop?

How this interaction of BGP & IGP?

- Encapsulation: R₁ encapsulates & sends to R₂

- Pervasive BGP:- R₁ sends as it is, intermediate routers look up 2 tables; every time they receive this pkt; to determine NEXT.

- Tagged IGP:- R₁ sends as it is. only one. Internal routers look up IGP table; decide via hot potato & forward to NXT.

(Visit HW3.)

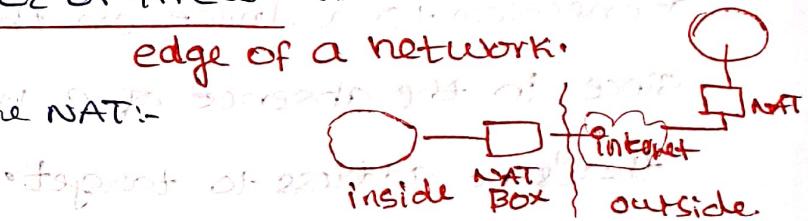
→ Network address Translation (NAT) :-

- * To access the internet; we need an public IP address. But we can use private IP in our private networks.
- * NAT's idea is to allow multiple devices to access the internet; through a single public address.
Diagram of NAT: A tree with root 192.*.*.* and branches 10.*.*.* and 192.*.*.*
- * NAT is a process in which one or more local address is translated into global address. Also port number is translated
 - The corresponding entries of IP address and port number is made, in the NAT table.
 - NAT operates on router or firewall at edge of a network.

* There are 3 ways to configure NAT:-

1) static NAT:-

- a private IP mapped to a legally registered IP.
i.e. public IP.
- Generally for web hosting.
- not for organisations, as many devices need internet.
costly to afford those many public IP.



2) Dynamic NAT:-

- a private IP is assigned a public IP, from a pool.
- if the pool of public IP is not free; then further IPPKts will be dropped.

Ex:-

we got 2 public IPs. Hence we can't afford 3 hosts.

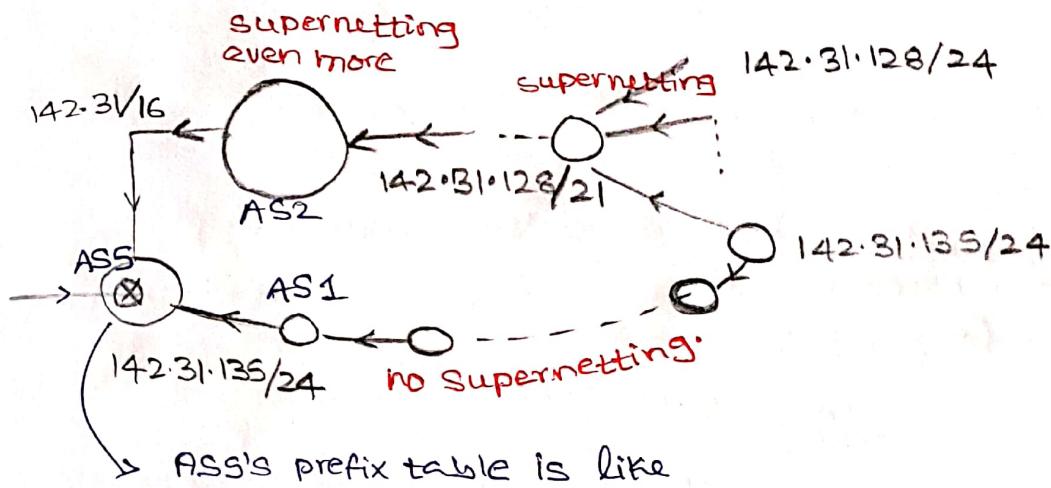
3) Port Address Translation (PAT):-

- Known as NAT overload.
- many private IPs to single public IP; port numbers are used to differentiate traffic.
- most often used, cost effective. Thousands can use single public IP.

→ NAT security:-

- * NAT provides privacy.
 - NAT allows you to display a public IP address while on a local network, helping to keep data & user history private.
 - NAT transfers packets of data from public to private addresses, it also prevents anything else from accessing the private device.
- * NATs make it difficult for an outside device to initiate a connection to a device on the private side of NATs since in the absence of a mapping in NAT table, there's no address to target.

→ longest prefix matching (BGP mei) :-



ASS's prefix table is like

Prefix	Next hop
142.31.16/16	AS2
✓ 142.31.135/24	AS1

fast lookup possible
due to hardware aids;
CAM: content
addressable
memory.

* NOW ASS gets a pkt with dest IP = 142.31.135.24.

"choose longest prefix matching to decide next hop!"

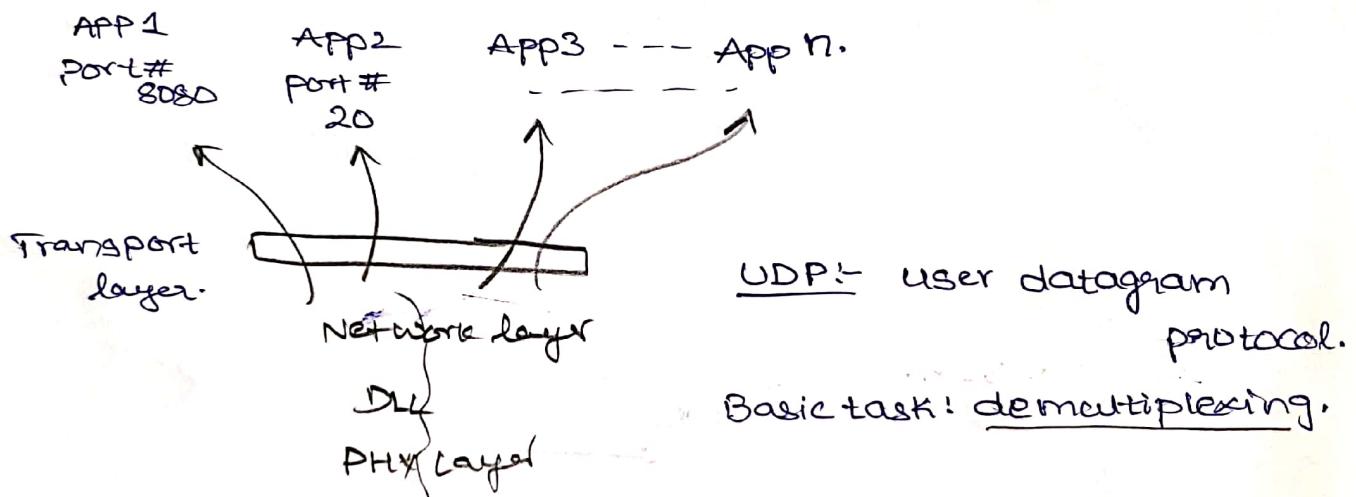
Hence AS2.

- This'll prevent routing loops. And also we'll have a deterministic path (instead of choosing randomly).
- longer prefix implies we are "closer" to destination.

Layer 4: Transport layer:-

- * Many applications are running simultaneously on our PC; bcoz we have **port numbers**.
Each application → Allotted a port number.
- * The transport layer acts as a 'DEMULTIPLEXER'; send pkts **(implemented in OS)** to their appropriate applications.

Hence the name 'transport' inside OS.

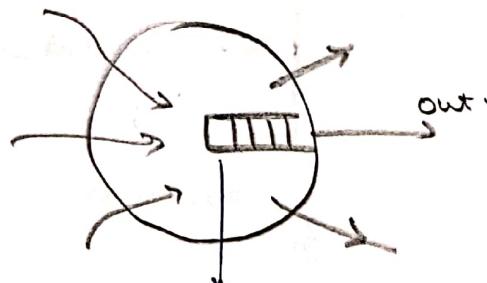


TCP:- transfer control protocol:-

- * Does a-lot- of stuff.
- * There's no client TCP & server TCP. Its just TCP. Both behaviours are part of this.
 - Reorders received segments based on seq#.
 - retransmits lost packets. (maybe lost due to queue full at routers).
 - 2-way data transfer.
 - Acks are present. for sent data.
 - Congestion control, after inferring congestion via different mechanisms.
!! routers don't have L4.
 - Flow control.

- * TCP is purely an end-to-end congestion control protocol.
 - This means, it does not require any special information from routers or explicit information about other TCP flows.

- * Recall the structure of a router:



Van Jacobson
 Sally Floyd
 1980s
 TCP.

this is usually a FIFO queue.

also "drop tail mechanism".

i.e. when it is full it'll drop further incoming packets.

in case of congestion,

Tcp will infer this

~~fix by packet loss due & reduce segment sending rate.~~

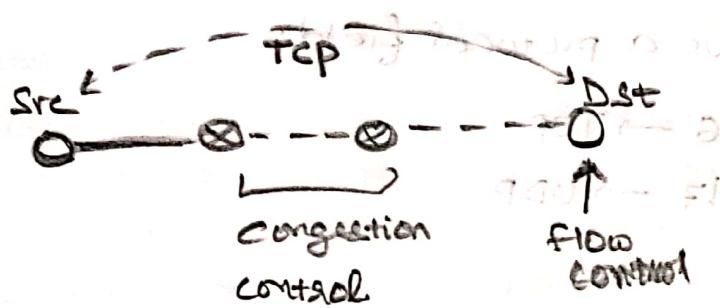
"packet loss due to congestion".

3 Dup.ACK method.

(decreases congestion window).

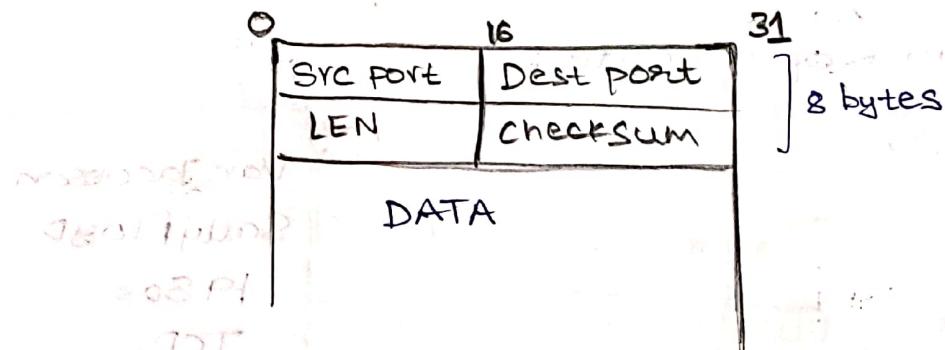
Flow control:

when there is congestion at the end host (rather than intermediate router) then the end host can inform the sender. No need of inferring this flow congestion.

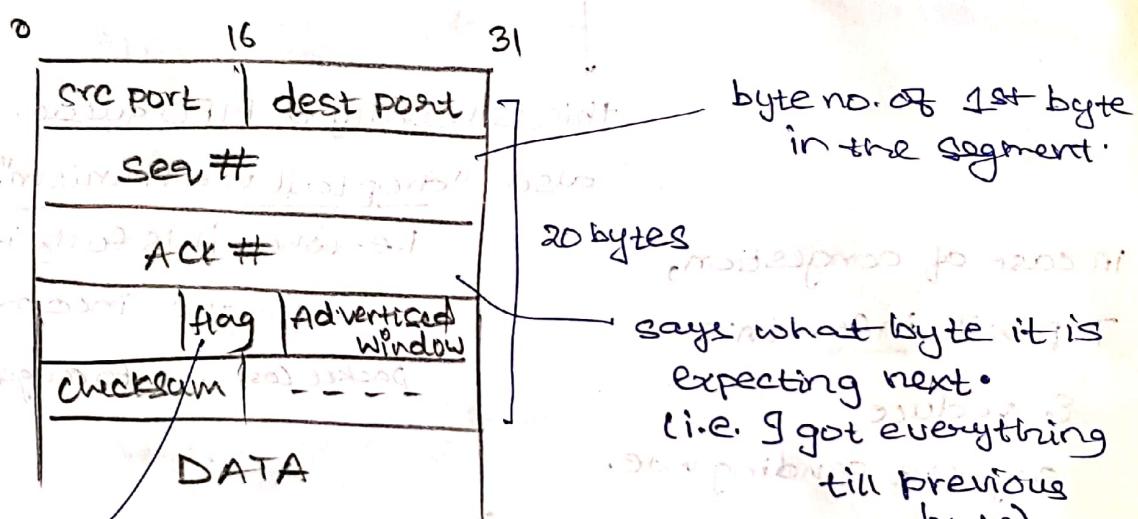


→ Headers of the transport layer:-

1) UDP:- this is a plain protocol.



2) TCP:-



will have various bits

set to 0/1.

SYN | FIN | RESET | PUSH | URG | ACK

Start connection | end connection | to say if ACK is there or no.

* IP header will have a protocol field:-

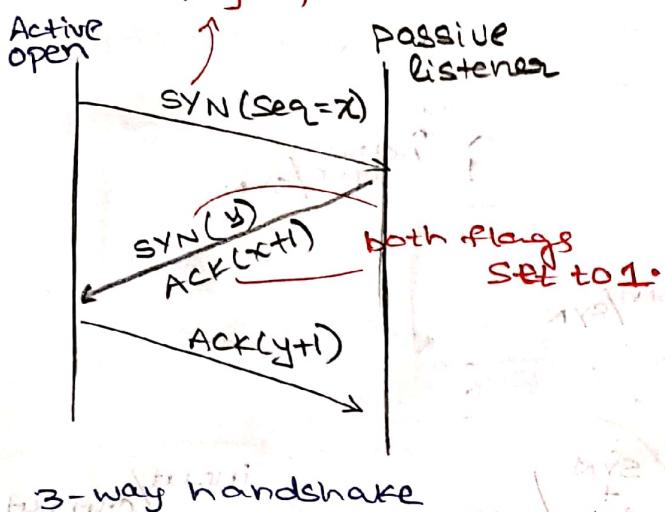
6 → TCP

17 → UDP.

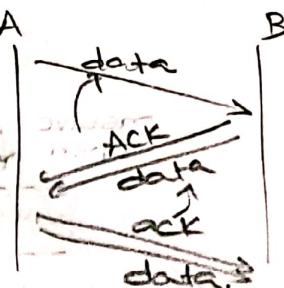
QUESTION:-

→ connection establishment:

flag = 1, in header.



• Typically:



* SYN → synchronization

* x,y → randomly selected not 0.

* The SYN, FIN segments, typically shouldn't have any data. But taken to have 1 byte.

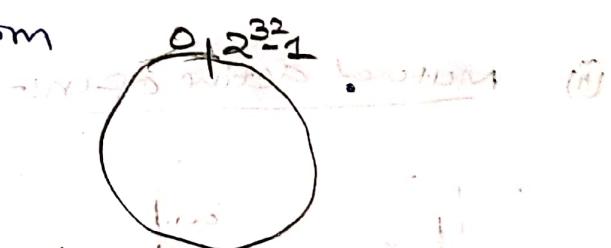
Hence ACK is x+1 bytes!

* if ACK = K; then it means "I have everything till K-1".

* x,y are randomly selected from

Bcoz say if we take in a
ordered fashion;

We'll maybe have overlaps...



↳ [NOT SURE HERE].

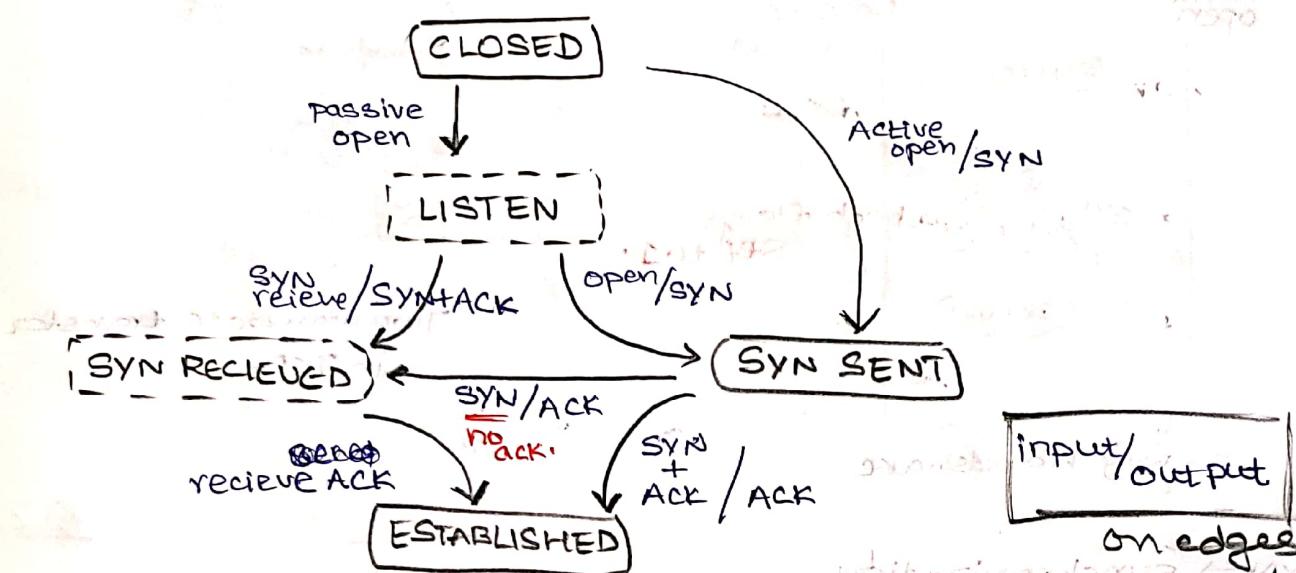
verbatim from book:-

TCP specification requires that each side of a connection select an initial starting sequence at random. The reason for this is to protect against two incarnations of the same connection reusing the same sequence numbers too soon - that is, while there is still a chance that a segment from an earlier incarnation of a connection might interfere with a later incarnation of the same connection.

→ State Diagrams-

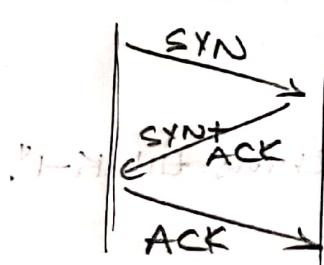
→ SYN

- opening a connection:- same diagram should be used for a client & a server.

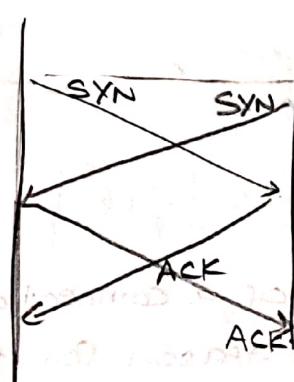


* this state diagram encompasses 2 things :-

- 3-way handshake: 1 round trip between 2 hosts



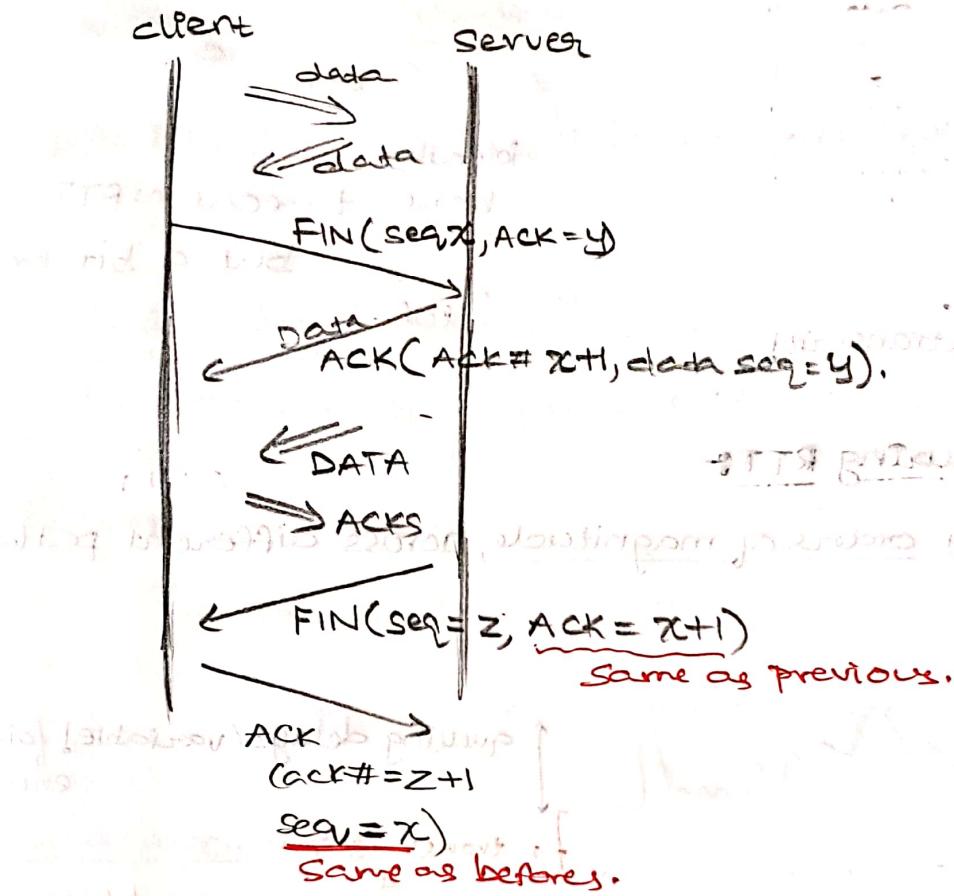
- mutual active open:-



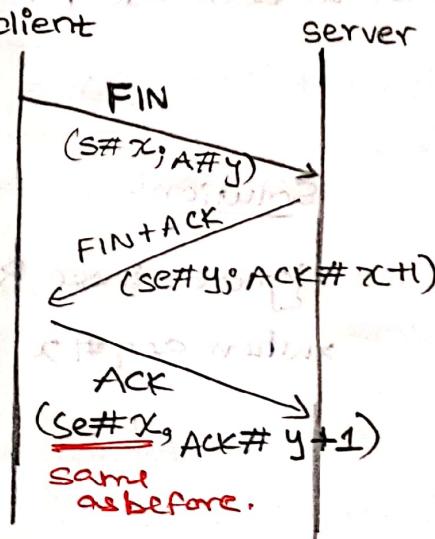
2) terminating a connection:- (use FIN flag).

State diagram becomes convoluted.

i) Half-close:- (only one end closes)



2) 3-way handshake termination:-

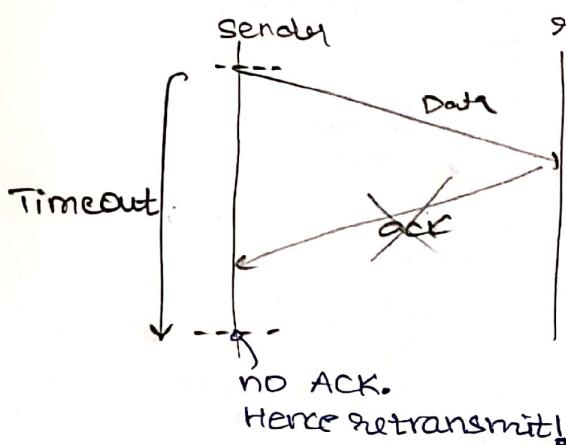


Both sides must close the connection.

If only one side closes, it means it won't send data, but it will receive.

→ TCP timeout

- * When does the sender retransmit?

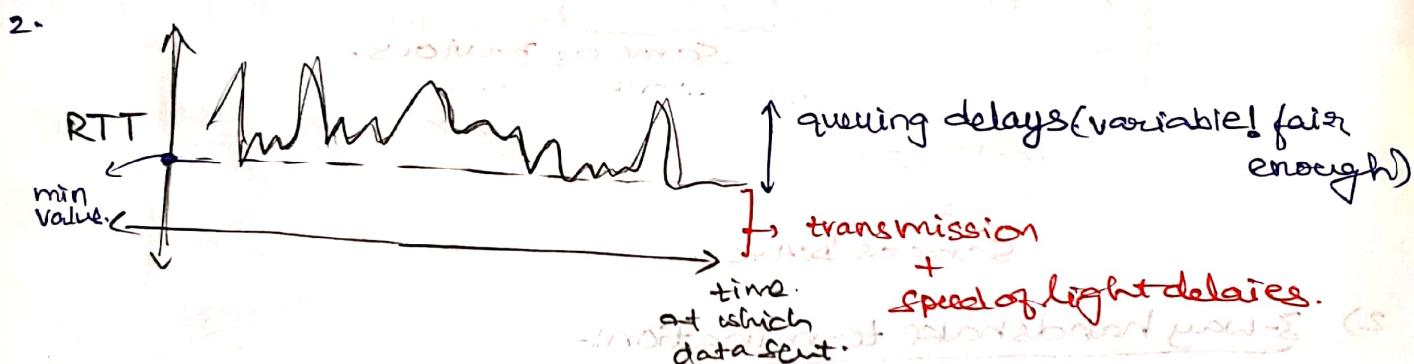


ideally,
have $\text{timeout} \approx \text{RTT}$
but a bit more.

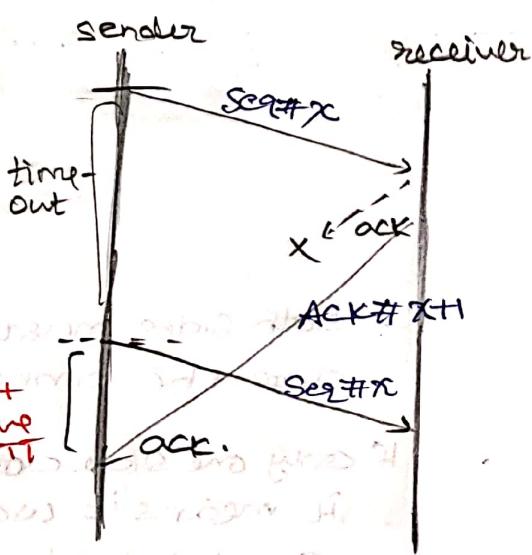
lol...

→ issues in measuring RTTs

1. RTT varies by orders of magnitude, across different paths.



3. Packet loss makes it difficult to estimate RTT.



Solution:-

ignore these RTT measurements
when seq#x is 'retransmitted'

* so, how DO we measure avg. RTT or something like this:-

1) old algorithm:-

sample RTT \leftarrow most recent

RTT measured

$$\text{Est. RTT} = \alpha \cdot \text{Est. RTT} + (1-\alpha) \text{Sample RTT} \quad \alpha \in (0,1)$$

"Smoothing".

$\alpha \approx 1 \rightarrow$ less importance to current RTT

$\alpha \approx 0 \rightarrow$ lot of importance to current RTT.

then,

$$\text{Timeout} = 2 \times \text{Est RTT}$$

mmm.... has some issues....

2) New Algorithm:-

assume gaussian distribution on RTT.

1. estimate Est. RTT as before.

$$\text{Est. RTT} = \alpha \cdot (\text{Est RTT}) + (1-\alpha) (\text{Sample RTT})$$

$\alpha \rightarrow \frac{7}{8}$

2. estimating deviation:-

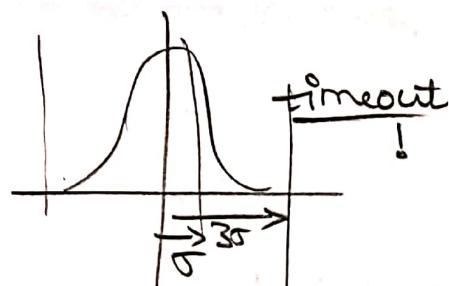
Standard deviation \rightarrow intricate
RMS value....

use mean deviation!

$$= \frac{1}{N} \sum |x_i - \bar{x}|$$

$$\text{Dev} = \text{Dev}(1-\beta) + \beta \cdot |\text{diff}| \quad \beta \rightarrow \frac{1}{4}$$

where diff = Sample RTT - Est. RTT



Set timeout
= mean +
 $N \times \sigma$
my midsem
approach!
Any

3) calculate Timeout, from RTT & Dev.

$$\text{timeout} = \mu \times \text{EstRTT} + \phi \times \text{Dev}$$

1

μ
mean RTT

4

(mean dev); not std dev

"estimated"

RTT measured at each = 298

RTT variance = 298 - 298

mean dev = 298

$$298.298 \times 298 = \text{timeout}$$



2

example 2

RTT measured = 298.298 ms

RTT variance = (298.298 - 298.298) × (298.298 - 298.298) = 0 ms²

minimum timeout = 298.298 ms

maximum timeout = 298.298 ms

random =

RTT

measured = 298.298 ms

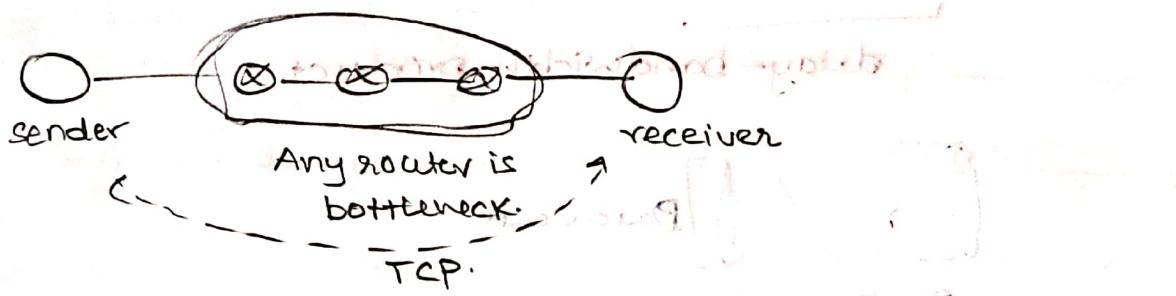
variance =

$$298.298 - 298.298 = 0$$

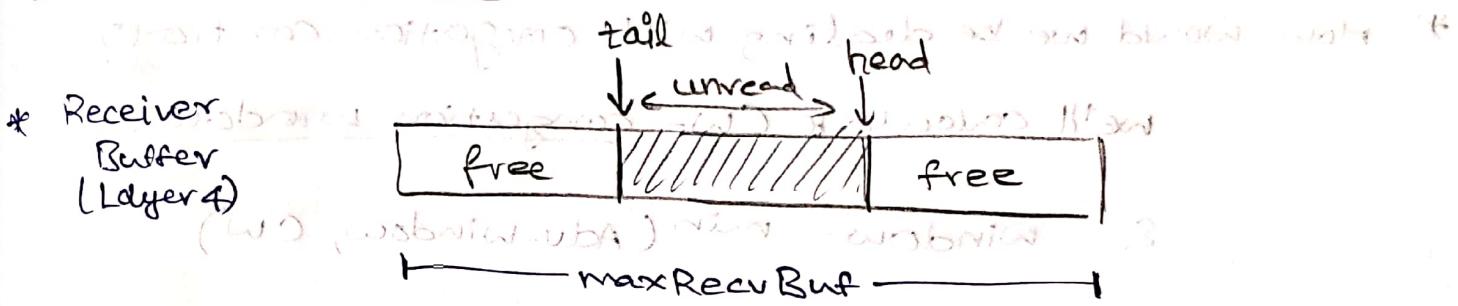
$$298.298 \times (298.298 - 298.298) = 0$$

RTT measured = 298.298 ms

→ congestion control: hosts need to infer packet loss. → greater RTTs.



lets see at receiver's end:



Flow control: congestion control at receiver.

- NO need to 'infer'!

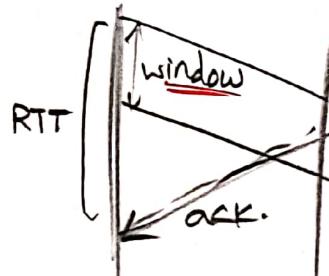
- * A field in L4 Header → Adv. window = space left in receiver buffer.
So, the sender knows!

Sender has a "window".

Window = Max amount of data (in bytes) which is outstanding, i.e. sent into the network but not acknowledged.

- * Say RTT is constant

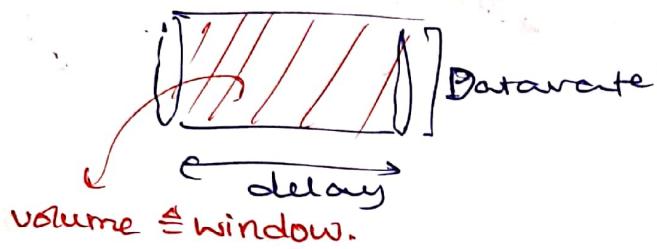
$$\text{Data rate} = \frac{\text{Window}}{\text{RTT}}$$



∴ we'll have window < Adv. window (for flow control).

* $\text{window} = \text{RTT} \times \text{datarates}$

delay-bandwidth product:



* How would we be dealing with congestion control?

we'll calculate CW- congestion window.

$$\Sigma \text{ window} = \min(\text{Adw}-\text>window, \text{CW})$$

→ window to be backed off by congestion window
instantly back off

at first stage = Congestion window → released of bits after
retransmission → if no retransmission → no release
of bits → if retransmission → release of bits
"window" = a set release

of bytes (set of bits) ready for transmission = window

and increment with each time slot, principle

new window size

get full connection

delayed acknowledgement

(Forward) ACK only contains bits > window next

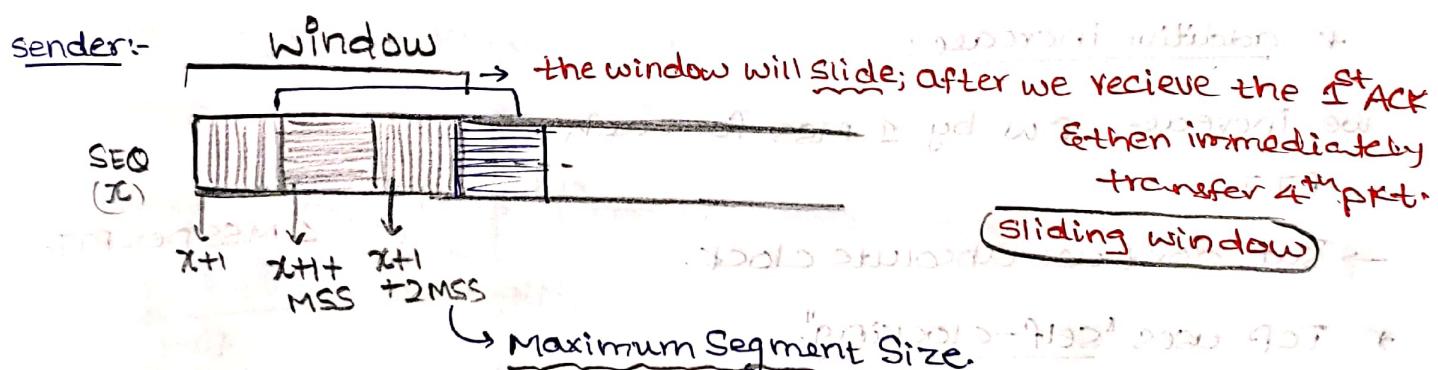
TCP Congestion Control

- * problems with congestion at routers
 - packets are dropped.
 - queuing delays; more RTT.
- * Router will not distinguish UDP, TCP packets. It is responsibility of TCP to do congestion control.
 - VOIP → UDP protocol.

* How do we infer Congestion?

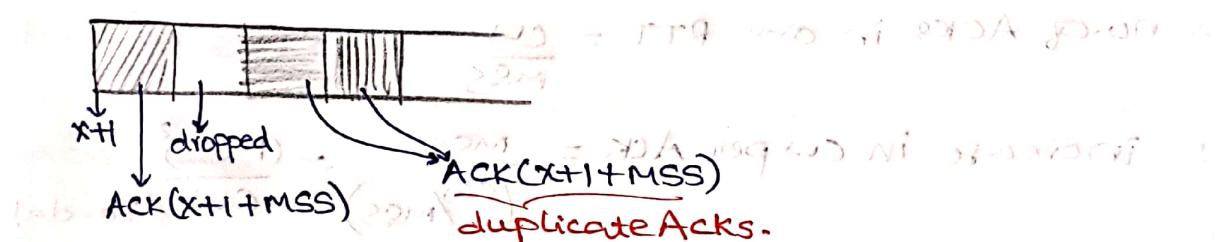
- large RTT values
- packet loss. How to infer packet loss?
 - Timer expiry
 - 3 Duplicate ACK approach.

→ Detecting Packet loss:-



If initial window = 3 MSS; we'll send out 3 segments & wait for 3 ACKs.

receiver :-



cumulative ACKs :-

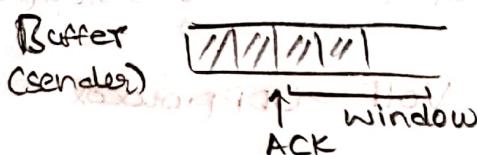
Ack#z (sent) ⇒ means; "I have every byte till (z-1)!"

(PTO)

Duplicate ACKs:-

- ACK# of some previous ACK is sent again.
- Each dup ACK says "one more segment received".

* where as; at the sender side; window starts at the latest



→ Principles of congestion control: For now, assume $window = CW$; & take AdvWindow is sufficient.

• If no congestion, increase CW conservatively.

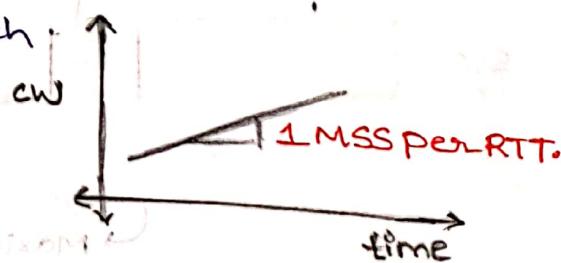
$$\text{DataRate} \approx \frac{CW}{RTT}$$

well have something called slow-start tool!

∴ as $CW \uparrow$; dataRate \uparrow . (yay!)

* additive increase:

we increase EW by 1 MSS for each RTT.



→ But, don't use Absolute clock.

* TCP uses "self-clocking".

i.e. increase CW on arrival of ACKs.

$$\text{No. of segments in one RTT} = \frac{CW}{MSS}$$

$$\therefore \text{No. of ACKs in one RTT} = \frac{CW}{MSS}$$

$$\therefore \text{Increase in CW per ACK} = \frac{MSS}{\frac{(CW)}{MSS}} = \frac{(MSS)^2}{CW} \quad \text{Ta-da!}$$

∴ on receiving an ACK;

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

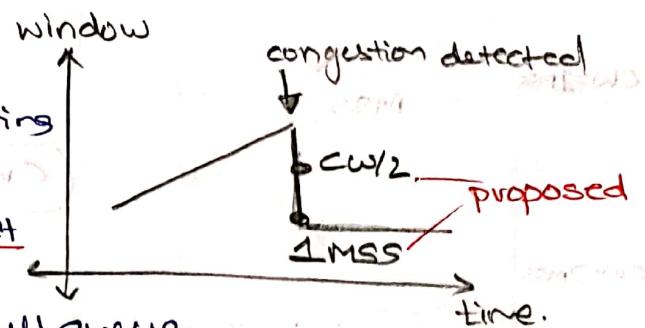
; Assuming no congestion.

2. If congestion is detected, decrease CW aggressively.

- aggressively, bcoz already data incoming rate exceeds outgoing at router.

Hence, small decrease won't help,

the already full queue.



TCP Tahoe \rightarrow CW = 1

TCP Reno \rightarrow CW = $\frac{CW}{2}$ MULTPLICATIVE DECREASE.

TCP Vegas \rightarrow & some other mechanisms.

TCP Africa, cubic \rightarrow we ain't studying.

AIMD:

additive

- multiplicative increase

multiplicative decrease.

3. Slow start - (Multiplicative increase)

* Initially, we donot know what the optimal datarate/CW is.

It might be 10 Gbps ($\approx 10^3$ MSS).

Starting from 1 MSS & additive increase \rightarrow very slow!

* Hence, be aggressive initially.

- double the CW per every RTT.

again! using self-clocking.

Congestion avoidance in each RTT!

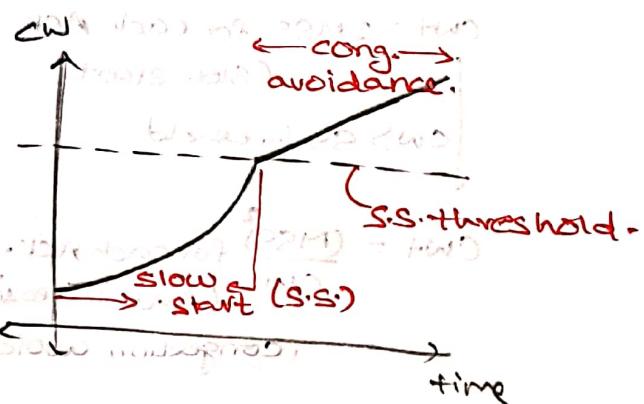
$$\text{ACKS} \# = \frac{CW}{MSS}$$

since, doubling!

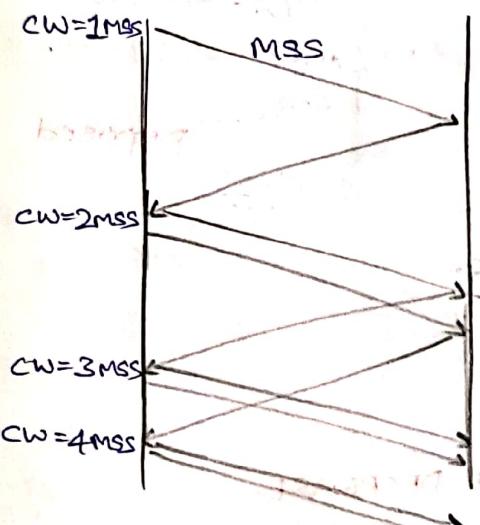
$$\begin{aligned} \text{increase per ack} &= \frac{CW}{(CW/MSS)} \\ &= 1 \text{ MSS} \end{aligned}$$

$$CW_t = 1 \text{ MSS}$$

Per Ack. } self-clocking!

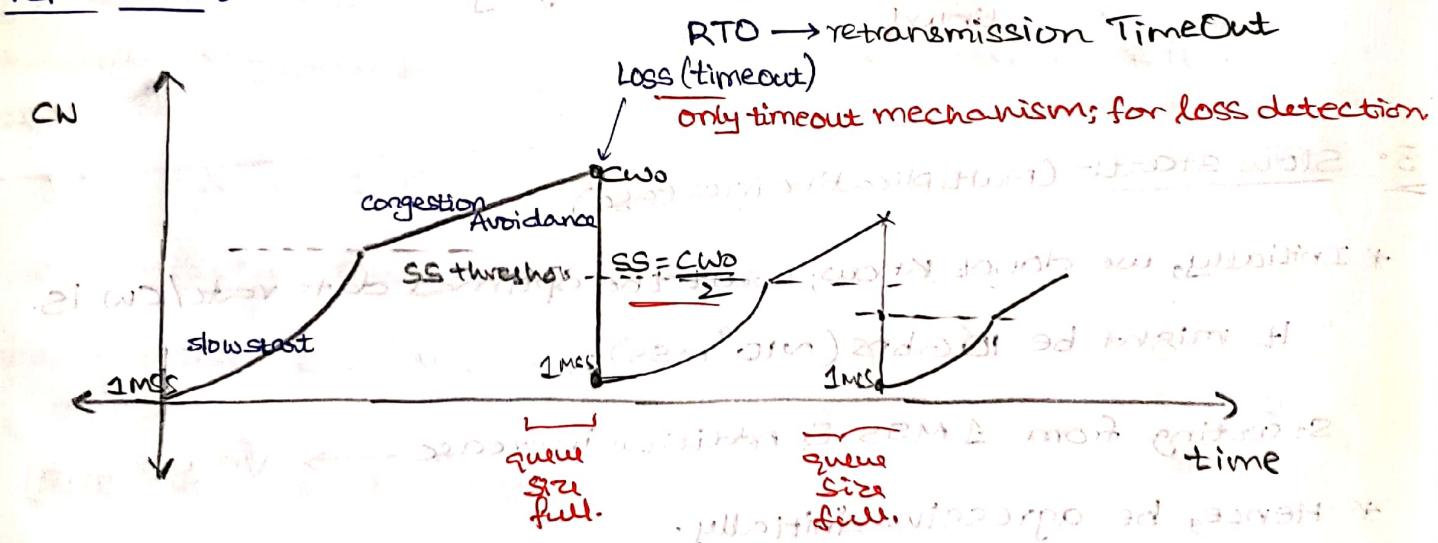


* in Slow Start phase:



$CW_t = 1MSS$ per ACK received.

→ Tcp Tahoe:



rules:

$CW = 1 MSS$ (Initialization)

$CW_t = 1 MSS$ for each ACK
(Slow start)

$CW > SS_threshold$

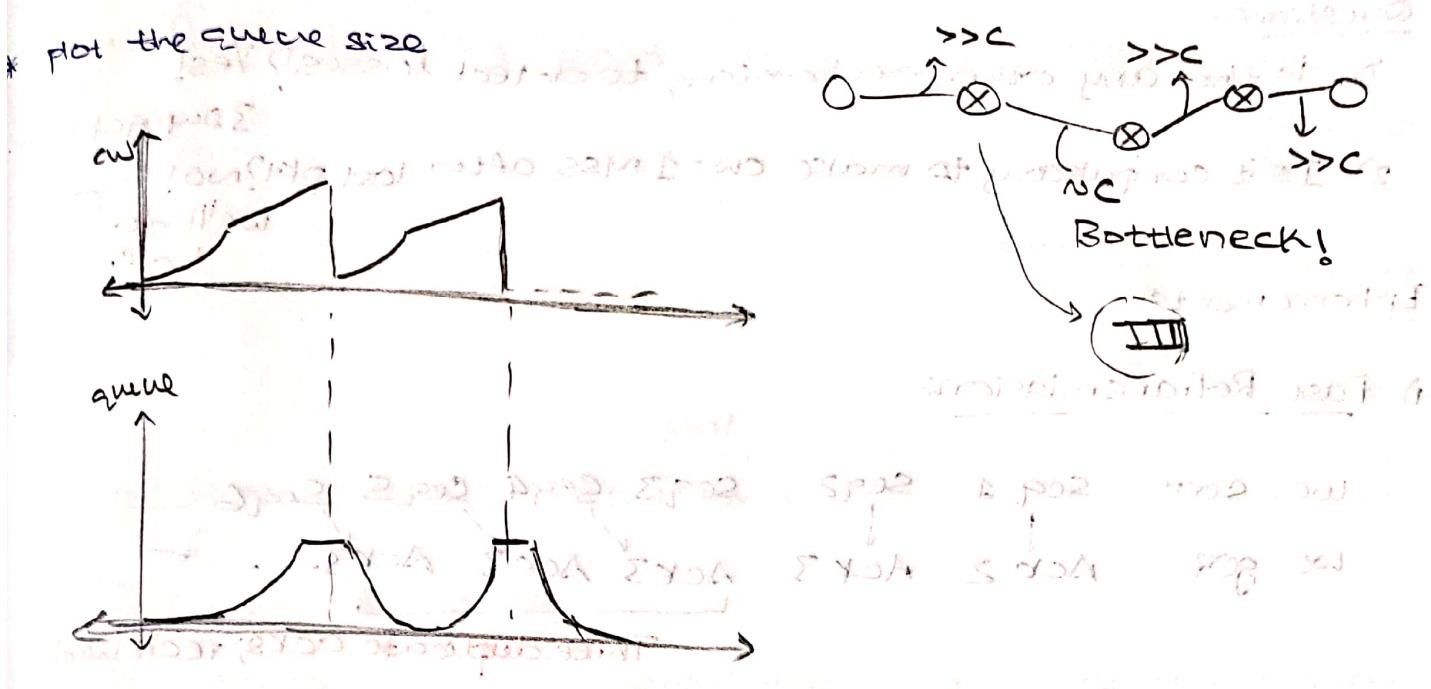
$CW_t = \frac{(MSS)^2}{CW}$ for each ACK
(Add. increase)

(congestion avoidance)

$SS_threshold = \frac{CW}{2}$

min. $SS_threshold = 2 MSS$

gave in a RFC
(Request for
comments)

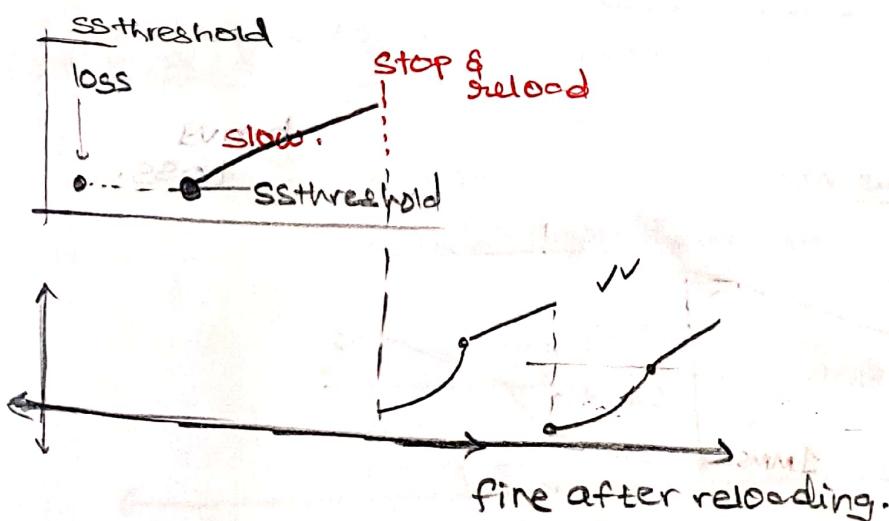


- * Two TCP Tahoe flows... okay! (See in slides; nothing special graphs).

→ Stop & Reload Phenomenon:-

- > We observe that sometimes, downloading is slow; so we reload the page, now it's fast!
- Because, due to bad luck, at the start itself, we might detect packet loss.

So; we'll be in additive increment mode later.



Questions:-

- 1) Is there any other mechanism to detect losses? Yes!
3 Dup ACK way.
- 2) Is it compulsory to make $CW = 1 \text{ MSS}$ after lost pkt? No!
we'll do half.

Enhancements:-

1) Fast Retransmission:-

We sent seq 1 seq 2
We got Ack 2 Ack 3

lost...
seq 3 seq 4 seq 5 seq 6 ...
Ack 3 Ack 3 Ack 3 ...

Three duplicate ACKs received
before RTO expiry.

Declare seq 3 lost.

reset RTO.

new way to
find lost seqs.

2) Fast Recovery:-

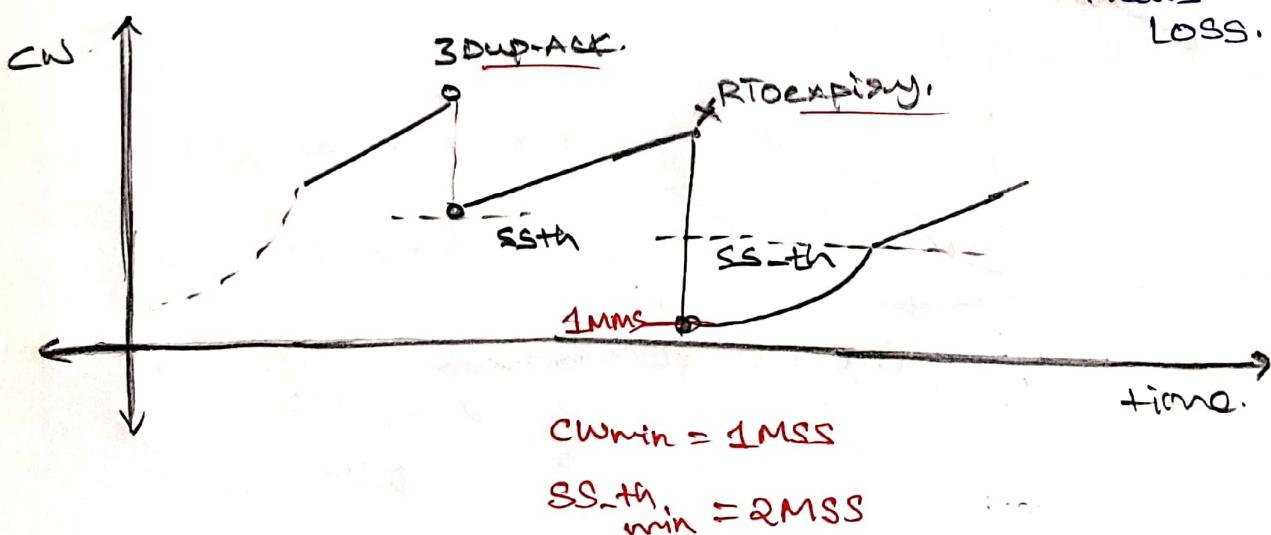
when loss is detected by 3 DUP ACKS; $CW = \frac{CWS}{2}$

Ex the usual rule:- when loss detected by $CW = 1 \text{ MSS}$

Also $SS_{th} = \frac{CWS}{2}$; so we'll have
Add. incr.

RTO expiry

meaning the
other ACKs are also
absent now!



- * TCP Reno; is TCP Tahoe + Fast retransmission + Fast Recovery.
- * CW = $\frac{C_W}{2}$, if 3DupAck
- * CW = 1; if RTO / expiry
retransmission TimeOut.
- * some miscellaneous points from RFCs;

① In Slow Start; SS

if the ACK acknowledges 'N' bytes of new data;
rather than 1MSS; then

$$C_W+ = \min(1MSS, N)$$

In Congestion Avoidance; CA

every ACK; do

$$C_W+ = \frac{(MSS)^2}{C_W}; \text{ but in one RTT; increase by } 1MSS.$$

that means,

take this caution.

② Initial RTO = 1 Sec or more

Loss by Timeout:-

$$ss_th = \max(2MSS, \underline{\underline{\frac{C_W}{2}}})$$

$$RTO = \min(2 * RTO, \max_{RTO})$$

CW is made 1MSS only.
You confused before!

this is RTO
for the same
Seq#

Hardcoded
Value

after sender
retransmits.

TCP Vegas: (1994-95)

These are notes - 2023-24 - Page 7

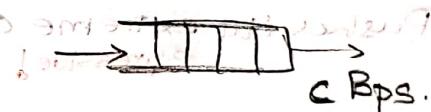
uses RTT values also; to infer congestion. **Nice!**

- * The previous two TCPs try to achieve the best possible rates. They would increase packet loss for themselves & also, maybe for other VoIP- UDP packets.
- * TCP-Vegas is RTT sensitive. It will try to use optimal band width & also allow transfer for other packets. Less pkt drops.

- * Detect loss & modify window, SS_th as in TCP Reno.

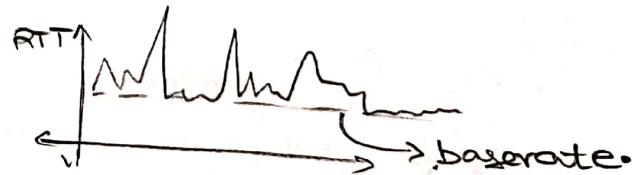
Slow start rules are the same.

- * Model Path as a queue



> queue is empty:

then $RTT = \text{baseRTT}$



$$\text{Expected Rate} = \frac{CW}{\text{baseRTT}}$$

$$\text{Actual Rate} = \frac{CW}{RTT} \leq \text{Expt. Rate}$$

Diff = expt. rate - Actual rate.
→ "smoothed out" remember!

∴ Hence, during CA :- congestion avoidance.

'Diff < α:-'

do A.I.

(additive increment)

+ 1 MSS per RTT

(or) 1 MSS per RTT

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

'Diff > β:-'

do additive decrease.

- 1 MSS per RTT

α < Diff < β:-

(W is not modified!)

data rate comparison

* We are optimal enough!

Diff = expt. rate - Actual rate

(TCP Vegas) approach 4.

$$= \frac{Cw}{\text{base RTT}} - \frac{Cw}{\text{RTT}}$$

$$= \frac{Cw}{(\text{base RTT})(\text{RTT})} (RTT - \text{base RTT})$$

Say queue length B

$$\approx \frac{B}{\text{Base RTT}}$$

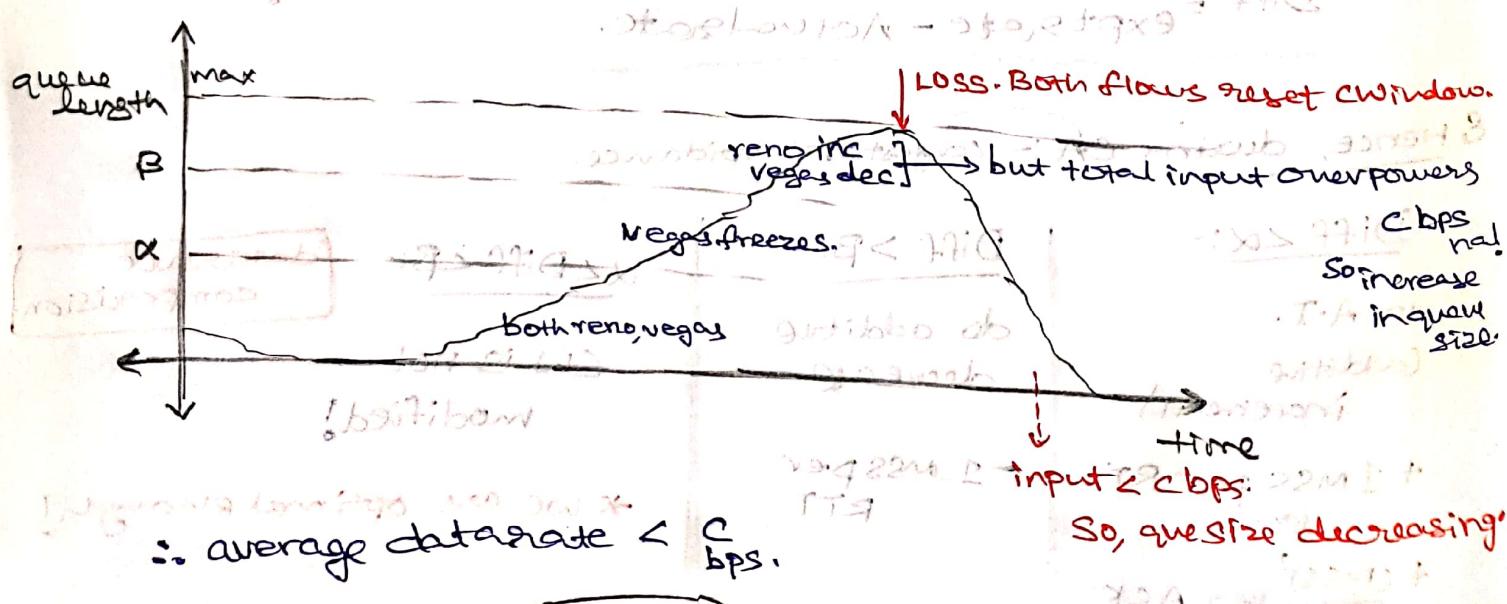
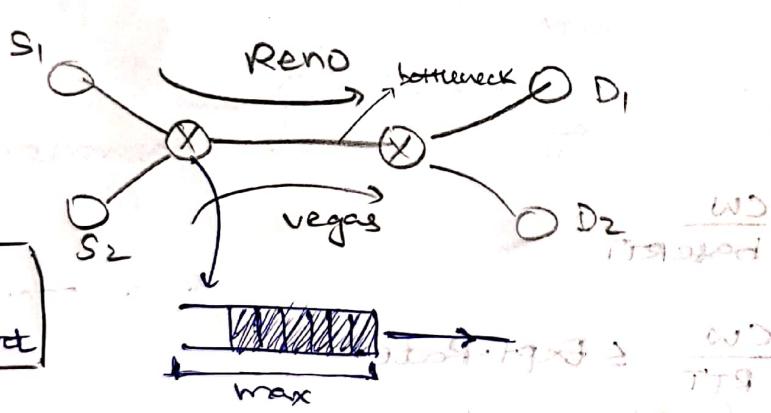
\therefore this difference is \propto to queue length.

Q2) What if TCP Vegas flows compete with TCP Reno flows?

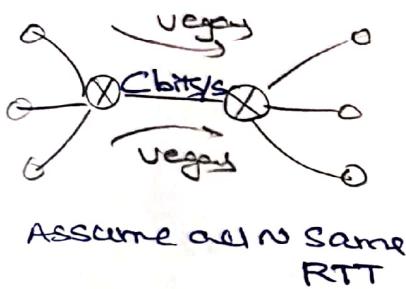
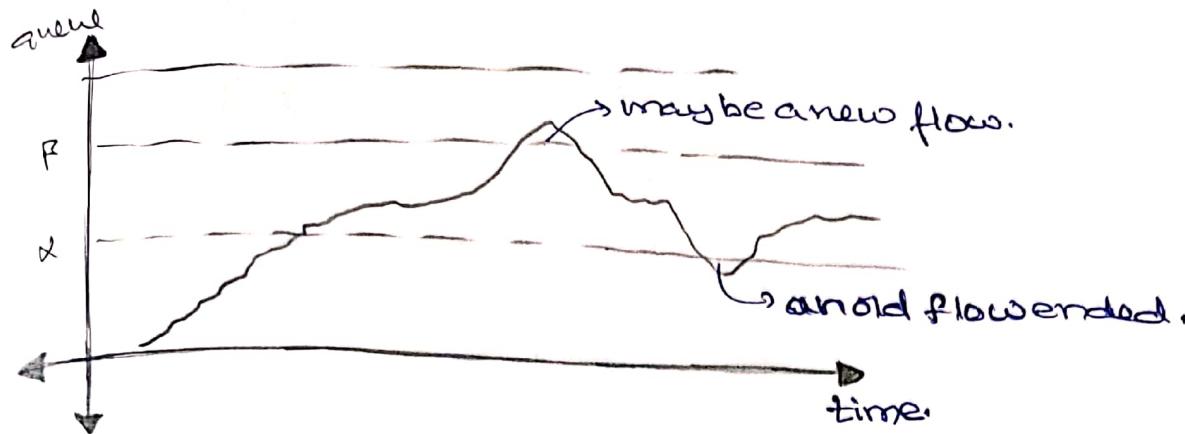
Vegas \rightarrow gets satisfied

Reno \rightarrow Pushes till extreme of packet loss situation: aggressive!

this results in pkt loss for all other flows of the common router.
It's evil Reno!



All vegan flows; then what?
lol 😊



A TCP flow is a s-tuple

$(\text{src}_{\text{IP}}, \text{dst}_{\text{IP}}, \text{src}_{\text{port}}, \text{dst}_{\text{port}}, \text{Protocol})$
field

- two devices
may have multiple
flows b/w themselves
Awesome!

benefits:

a) NO pkt losses. (TCP & VoIP & ... any flow)

b) queuing delays are predictable.

b/w $\frac{q_1}{C}(\alpha's)$ & $\frac{q_2}{C}(\beta's)$.] Some ~~appln~~ need this.

c) Throughput can be higher than all Reno case bcoz
the queues never empty here.

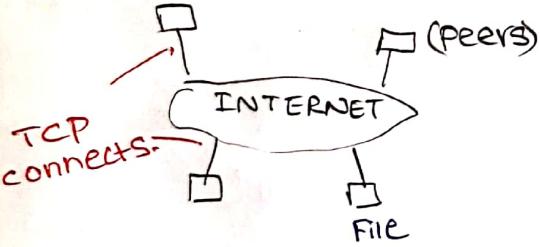
i.e. datarate N Cbits/sec.

max.
data
rate

Claim: All vegas is 50% better than all reno.

- * we have unlimited flexibility in application layer. we can define our own protocols here; if needed.
- * various protocols here include SMTP, HTTP, HTTPS, FTP, SSH, IMAP,
Simple
mail
transfer
protocol
LDAP, POP,
TLS/SSL.....

→ Peer to peer networks:-

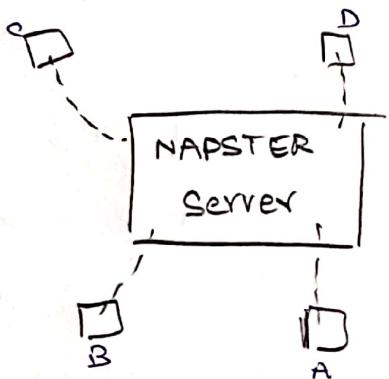
- * usual, starting internet was used on a client-server basis.
(peer-peer interaction was like mails; which needed server to store database).
- * A P2P should look like :-


- 1) Bootstrapping
 - How does one user 'get into' this web.
How does he start.
- 2) What sort of service?

1990's NAPSTER (digital content sharing platform).

(lot of music was pirated! so shutdown main server.)

*



Server:-

- stores list of files; available with each user.
- Eg: f1 | IP A
f3 | IP B
- whichever peer requests for a file, gets a reply from server containing IP of file's host.
- C directly connects to A over TCP.

* Bootstrapping:-

- A logs in
- A gives details of files it has.

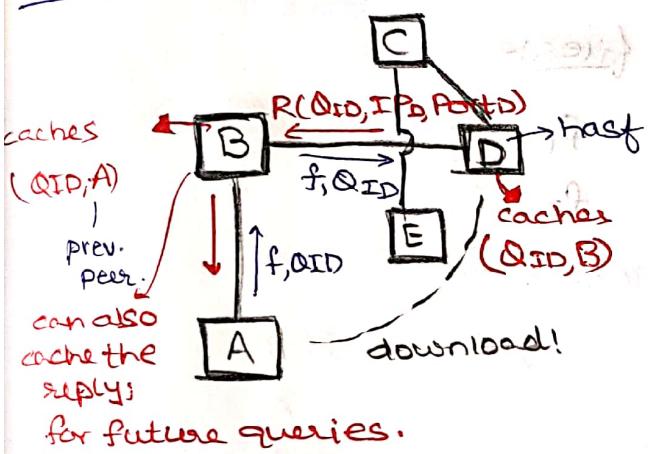
* Drawbacks:-

centralised → single point of failure.

↳ not fault tolerant

↳ legally, easy to takedown.

Gnutella:- (next gen P2P idea)



* Bootstrapping:-

(1) Application might have IPs of some other peers.

(2) Lookup IPs from one or more websites.

& connect to them.
(TCP!)

* suppose A wants file 'f'.

can't search complete internet: Gotta do limited search.

LIMITED BROADCAST ($TTL=n$)

(this is application TTL.)

maybe suppose $n=2$.

First method:-

Query from A has "f", QID (unique ID for query).

Reverse reply; over the path which query travelled.

2nd Method:-

Query has IP_A as well. D will reply directly to A.

V0.4 → $TTL=7$

V0.6 → $TTL=4$ } since network became dense.

Next Idea:-

- * If N nodes in the network, can I query only $\log(N)$ nodes
to find where data is stored?



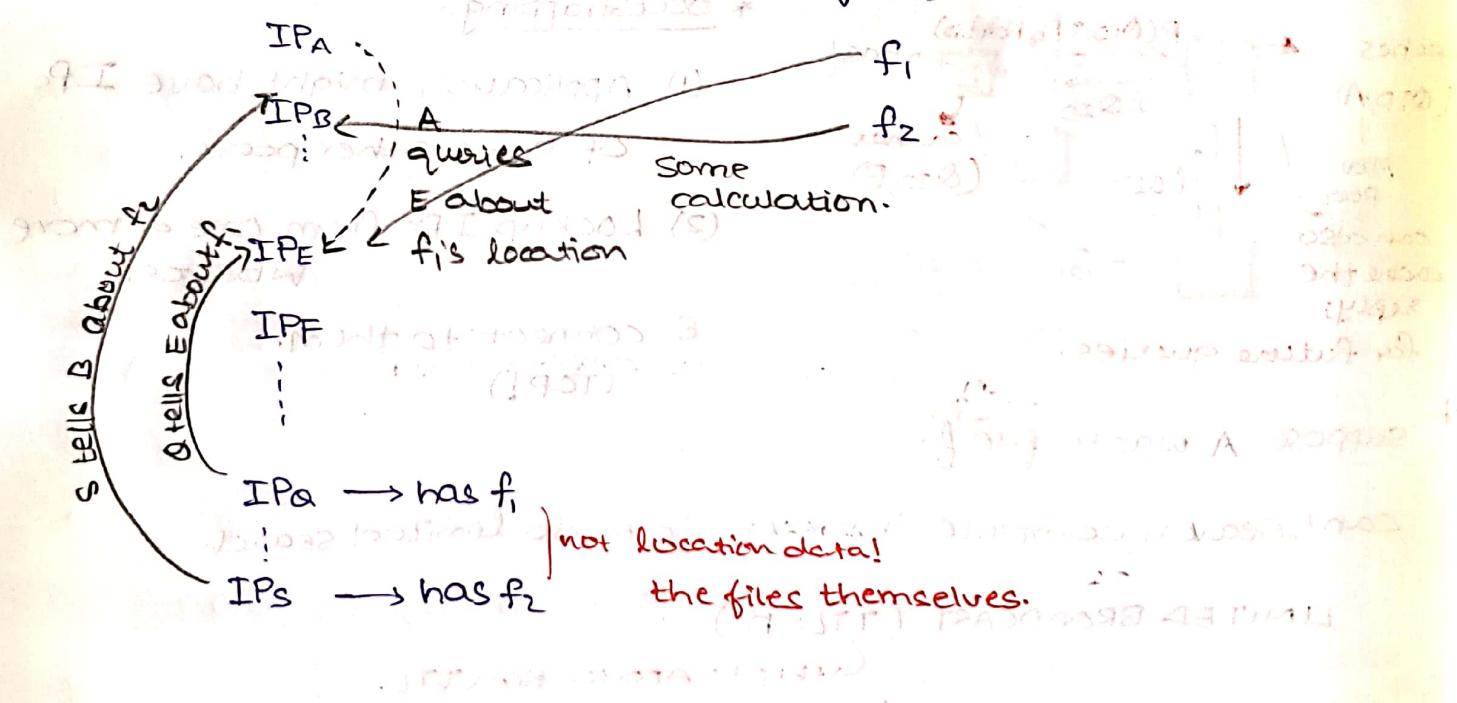
N nodes
log N levels

there's some
hope...

- * lets store (location data) of files, at some peer whose identity can be calculated based on file name.
then, every peer can calculate this specific peer & query him.

Eg: IP addresses :-

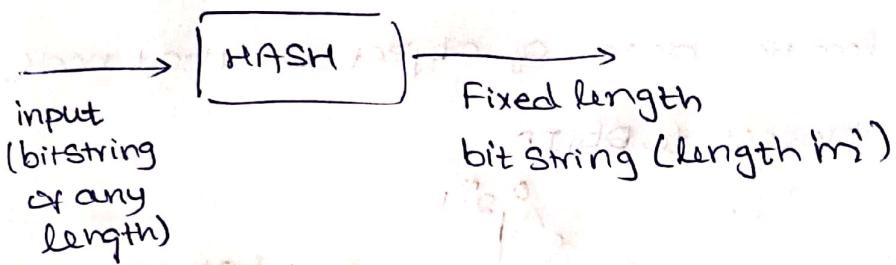
files:-



- * Here note that every peer should have reasonable equal load as others. Our mapping of $f \rightarrow \text{node}$ should be as such!. Even though in IPwise bunch of nodes are close or far apart, they gotta have equal chances for being responsible for ' f '.

- * Better use hash functions -

→ hash function:



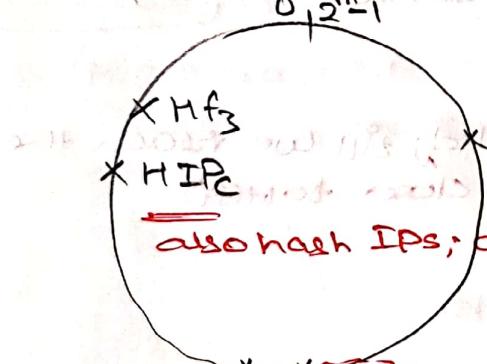
* properties of hash:

- uniformity: map inputs uniformly to 2^m space.
- Avalanche effect: small change in input, gives huge change in output.
- collision: very less probability for hash collision.

* Now; $f_1 \rightarrow H \rightarrow h(f_1)$

$$f_2 \rightarrow H \rightarrow h(f_2)$$

$m = \text{hash output length}$.



Hf_1 HIP closest!

now; this means that

Peer-E will store all data regarding f_1 .

$Q \xrightarrow{\text{tells about } f_1} E$

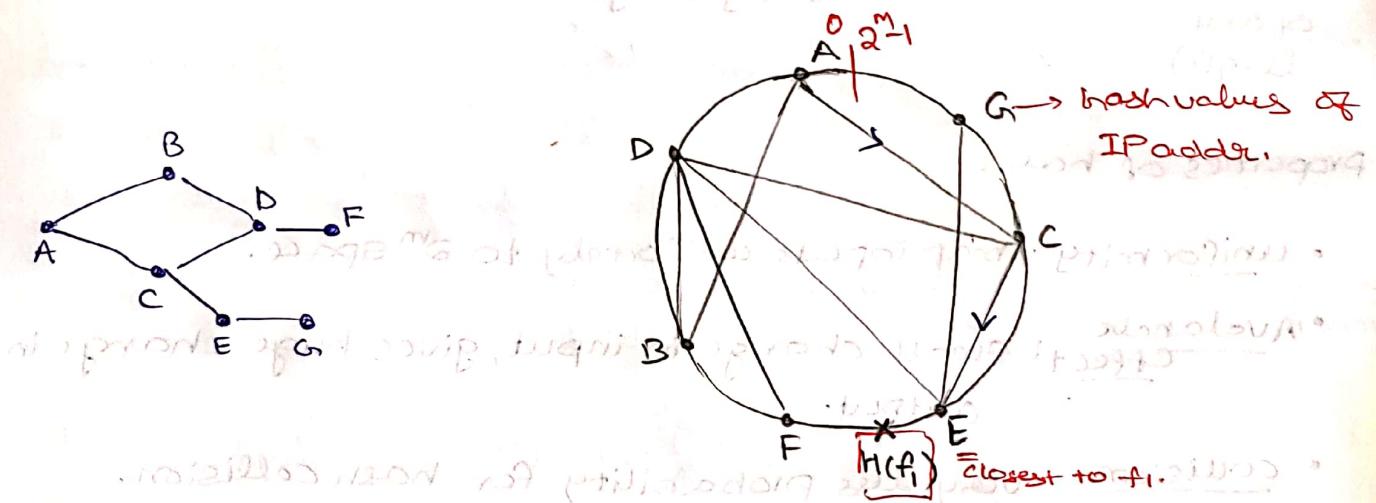
$A \xrightarrow{\text{queries about } f_1} E$

L29:

* nodeID = Hash(IP)

objID = Hash(f)

- we want to store location data of object on that node, whose nodeID is closest to ObjID.



* Our idea: route query message closer to H(f) in this virtual space till we reach E.

- Suppose A wants to find out where file f_i is present.
- A sends query to the node in its leafset that is closest to $H(A)$.
- This query is recursively forwarded; till we reach the "node closest to $H(f_i)$ ".

leafset:
of A; peers connected to A

But this need not always work.

We'll mandate a special structure among peer networks; for this algo to work.

- * 1st query forwarding from A to E:
then, E connects to A.
has IPA.

PASTRY:-

The leaf set of a node x includes $\frac{L}{2}$ closest nodes on either side of node x in the virtual space.
 2^m Hash space.

Eg: $L=4$:



for node A. Like this, for EVERY node.

!@!!@!

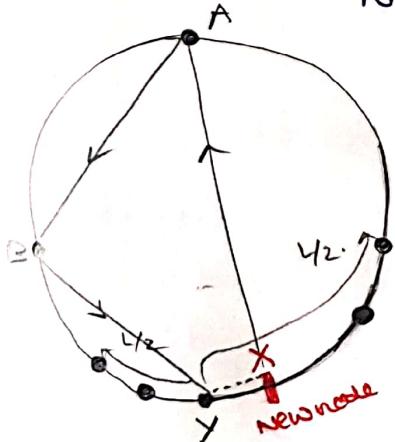
- * This PASTRY ensures that our "closer to H(f)" routing protocol will reach E - the closest node to H(f).
- * PASTRY ensures that leaf set of Any node is $L+O(\log N)$.
- * Now, how am I going to construct PASTRY?

• Recursively

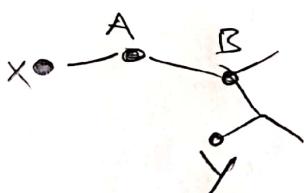
its also
better to
have some
more nodes.

Suppose P2P network already has the property PASTRY.

Node X is added.



Say; X; via bootstrap,
is connected to A:



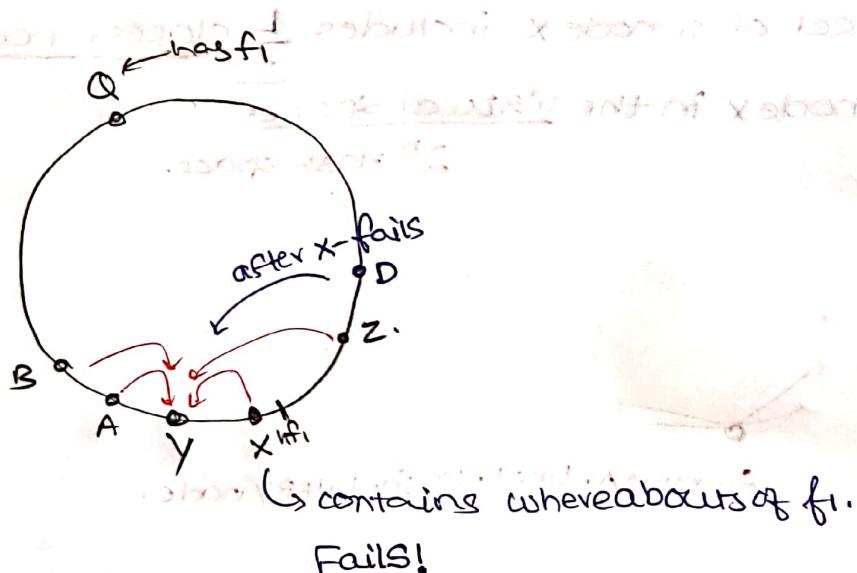
- 1) node X routes a message to find node, whose nodeID is closer to $\text{Hash}(\text{IP}_x) = Y$.
- 2) From Y, X determines its L_2 neighbours on either side.

This way of storing data 'everywhere' is

DHT.

Distributed Hash Table.

→ Fault tolerance, what if a node fails...



- * Let's be a bit redundant!. Every node stores the file info of 2 other nodes, on its either side!
 - Eg: Y stores not only of itself's; but also of B,A,X,Z.
 - * All peers send [keepalive] messages to peers of their leafset. if any node within 2 distance fails; (X here) then Y finds about a further node - W & stores its info.

L30:

Domain Naming System:-

DNS.

* www.google.com 192.168.1.13
userfriendly! bath! we need!

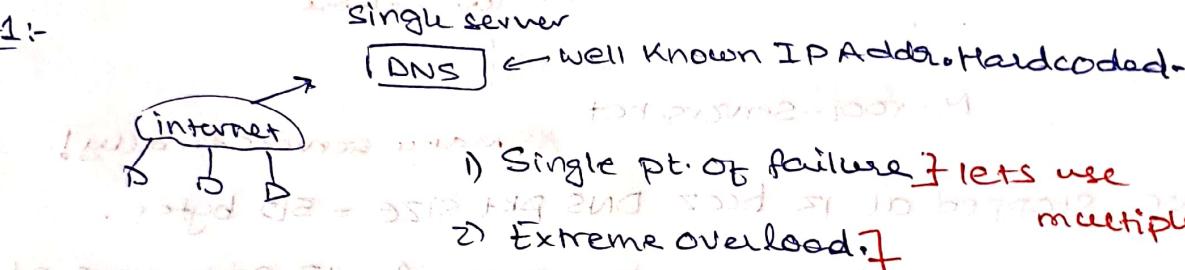


DNS

gotta be
robust

Importance is obs.

* idea-1:-

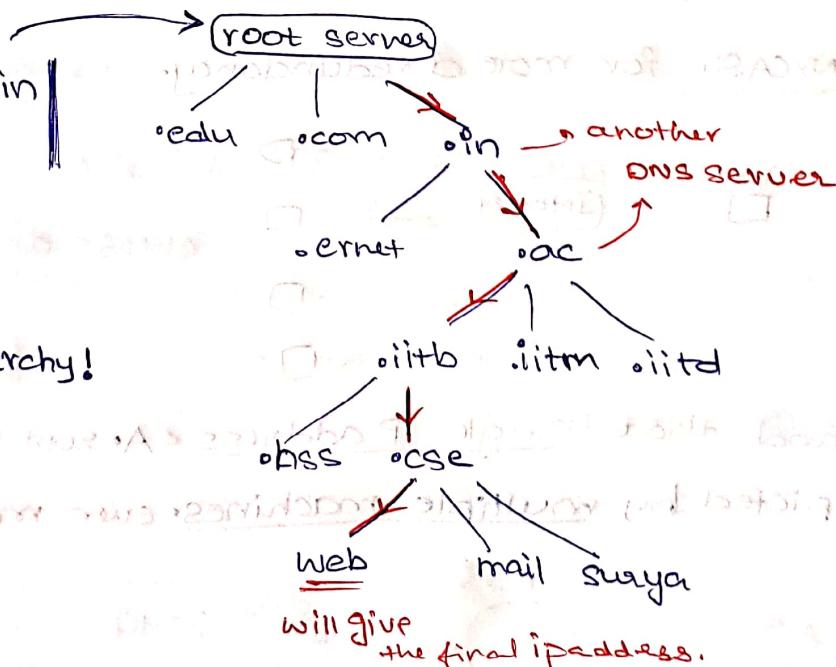


* idea-2:- (what runs now-a-days)

Let's maintain
Hierarchy.

URL:-

www.cse.iitb.ac.in



→ Root server :-

1. well provisioned - bandwidth, CPU... cooling!

*to prevent DoS (denial of service) attacks. To prevent overload.

- 2. Redundancy - We've got multiple root servers.

* 13 root servers globally.

A. root-servers.net — 10 machines

M.root-servers.net

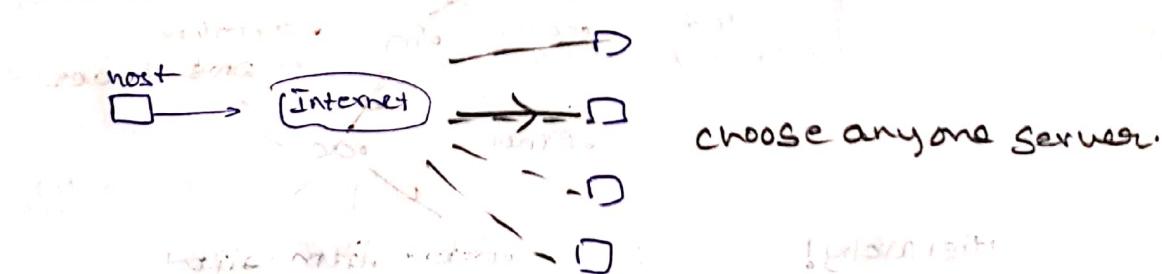
many servers finally!.

* we stopped at 13, bcoz DNS pkt size = 512 bytes.

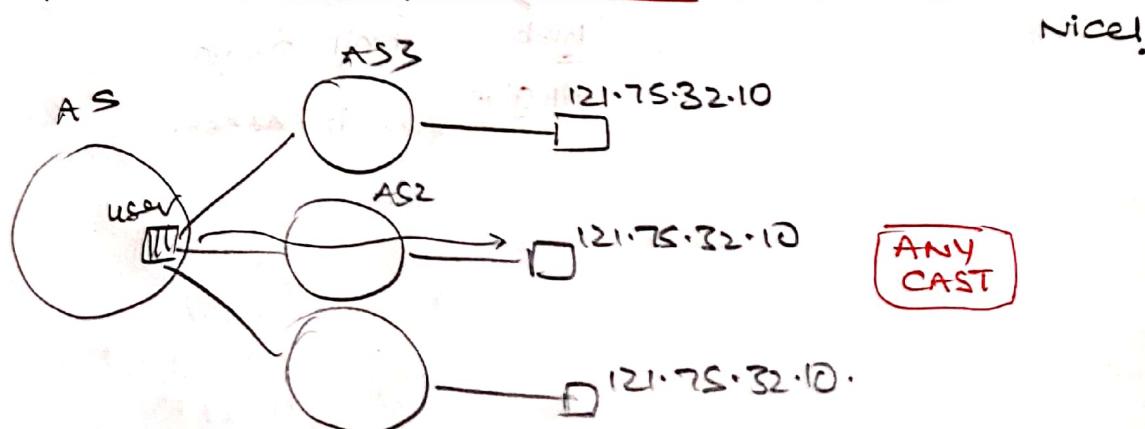
has space for 13 addresses at most.

not convinced!

- * use ANYCAST for more redundancy.



- understand that single IP address <A root server. net> is depicted by multiple machines. even more redundancy!



(allowed only for DNS servers)
(not part of a AS).

→ Resource Record: Inside any DNS server:-

< Name, Value, Type, < Class, TTL > >

URL value says how 'Value' is to be interpreted

 (meaning, internet) IN

 (meaning, internet) in host systems.

* Type :- < name, meaning of value ->

A < A, IP address, interpretation >

NS Name Server (the host, running the DNS Service in domain corresponding to 'Name' field)

CNAME Canonical Name (alias of host, specified in domain in 'Name' field)

MX (name of host, running mail server in domain specified in 'Name')

Note: Linux command - \$ dig www.google.com a few times. What does it do? used to gather DNS info.

Domain Information Groper.

DNS runs on UDP (port 53 for server).



Eg: root server has Resource Records

$\langle \text{edu, as.nstid.com, NS} \rangle$

name

points to name of one of DNS servers of edu domain.

$\cdot\text{edu}$

$\langle \text{as.nstid.com, 192.5.6.32, A} \rangle$

Server as.nstid.com has

$\langle \text{princeton.edu, dns.princeton.edu, NS} \rangle$

$\langle \text{dns.princeton.edu, 128.112.129.5, A} \rangle$

next level server has

$\langle \text{cs.princeton.edu, ---, A} \rangle$

$\text{dns.cs.princeton.edu}$ has

$\langle \text{penguin.cs.princeton.edu, ---, A} \rangle$

$\langle \text{www.cs.princeton.edu, corewals.cs.--, CNAME} \rangle$

$\langle \text{cs.princeton.edu, mail.cs.princeton.edu, MX} \rangle$

$\langle \text{mail.cs.princeton.edu, 128.12.1.5, A} \rangle$

* Look up example:

DNS root server IP, manually config or DHCP.

Local DNS Server

