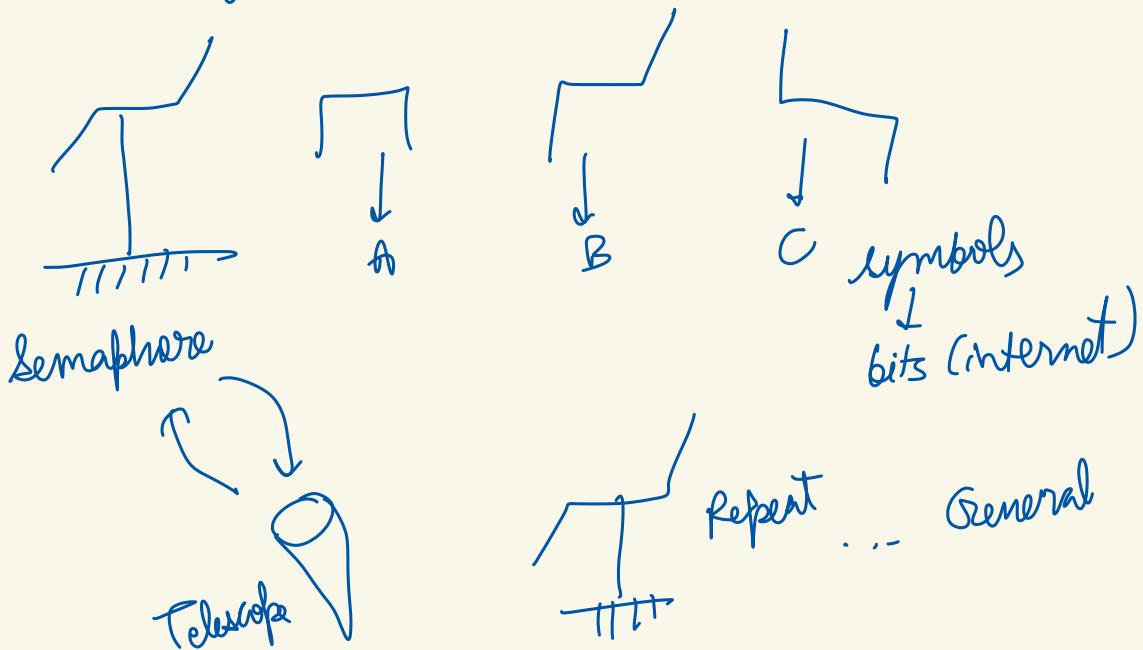
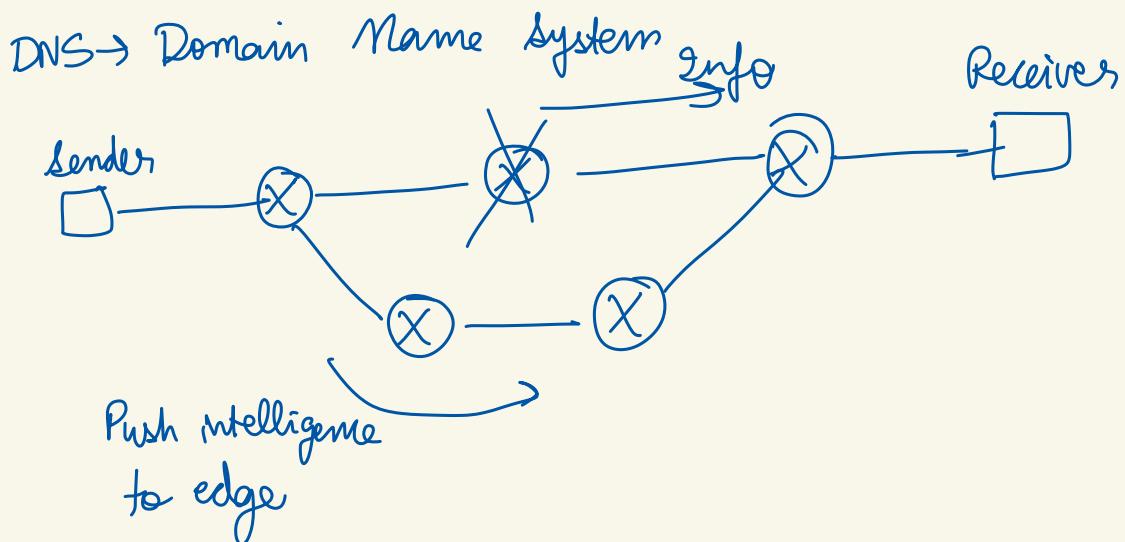




**CS 348**

# US- ARPANET - DEFENCE

Advanced Research Projects Agency NET  
(forerunner to the contemporary internet)



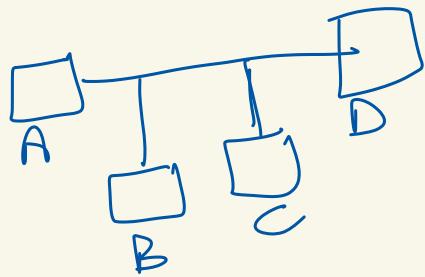
Napoleonic's semaphore was the world's first telegraph network.

(le système Chappe)

Physical Layer: Communicating bits from one node to next using signals.

Protocol: Rules for formatting and transmission of data

Bus topology



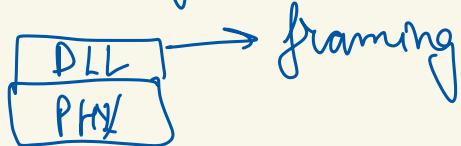
Collisions because of shared medium

Medium Access Control (MAC)

CSMA → CS → Collision Detection  
CA → Collision Avoidance

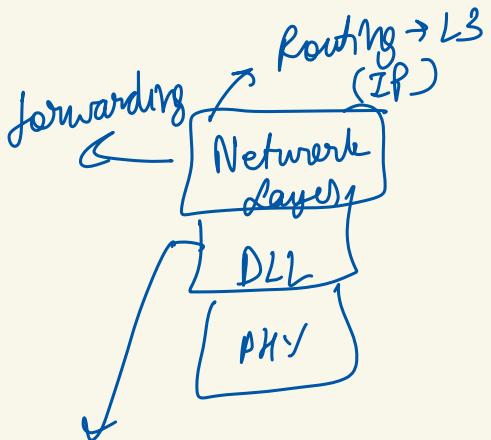
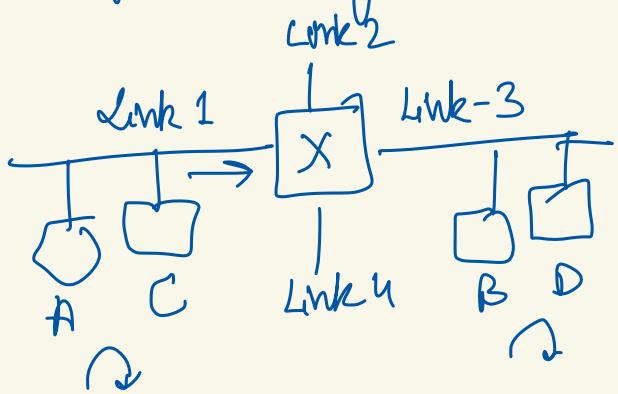
Carrier Sense  
multiple Access

MAC → part of data link layer

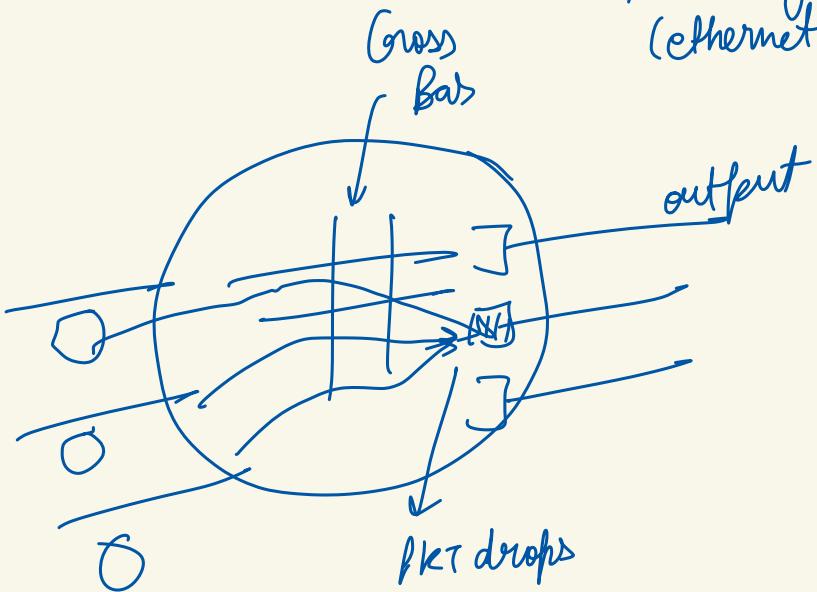


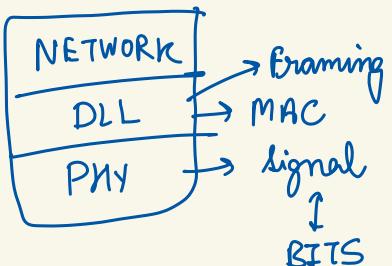
framing

Topology = graph of network

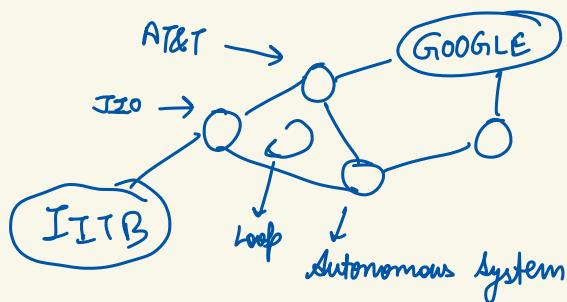
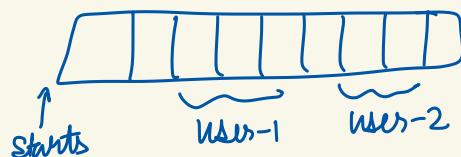


switching  
(Ethernet)





BIT STRING



**Simplex** : One-way communication

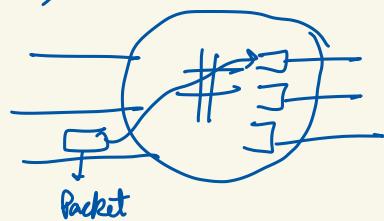
**Full-duplex** : simultaneous communication in both directions

**Half-duplex** : Communication in both directions but not simultaneous

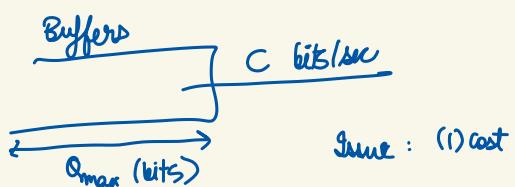
**BGP**: Border Gateway Protocol

Message :

- L2  $\rightarrow$  frame (1 unit of data)
- L1  $\rightarrow$  symbol
- L3  $\rightarrow$  packet
- L4  $\rightarrow$  (TCP) (UDP)
- Segment
- Datagram



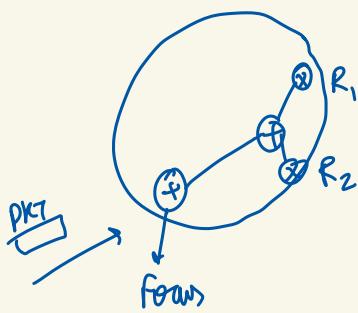
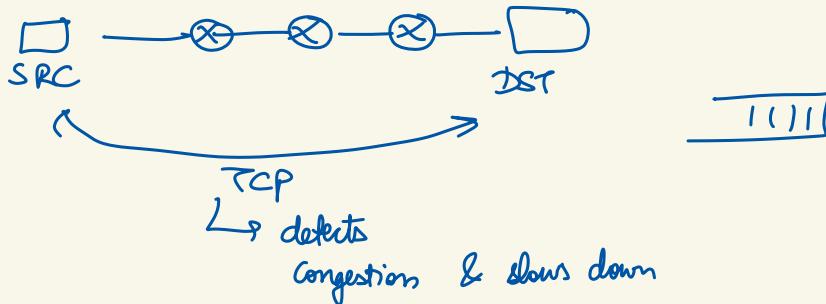
Latency  $\rightarrow$  Delay



$$\frac{Q_{\max}}{C} \rightarrow i^{\text{th}} \text{ router}$$

$$\text{End-to-end delay} = \sum_{i=1}^N \frac{Q_i}{C(i)} + \sum_i d_i$$

speed of light delay  
+ transmission delay



#### L4: Transport Layers

TCP: Transmission Control Protocol

Congestion control  
↓  
video call

Reliability  
↓  
file transfers

UDP =

#### Layer 5: Application Layer

Web (HTTP)      Email (SMTP)      VoIP      Text Messaging      P2P

L4 :

TCP

UDP

IP (Internet Protocol)

L3 :

WIFI, WLAN, Ethernet ...

L2 :

L1 :

## Design Protocols in Modules

Each sub-problem handled by some protocol.

OSI-5

Advantages of Layering / Modularity

APP2

(1) Ease of Development → Only certain problems handled by particular layers

TRANS

(2) Debugging

Netw

(3) Many applications and many physical layer technologies.

DLL

(4) Ease of modifying → only change 1 layer to address a problem

Phy

Protocol  
layering :

APPL

Email, whatsapp, YouTube, FTP

TRANS

TCP, UDP

NET

IP

DLL

WiFi-DLL, Bluetooth, 9G, Ethernet

PHY

RFC  
Request  
for Comments

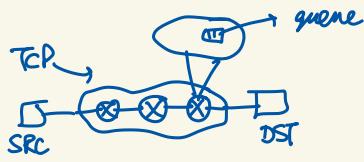
### Advantages:

- Development & debugging (Ease)
- Easy to modify one layer without breaking entire system
- Can have different choices for each layer (Compatibility)

### Disadvantages:

- Opaqueness about other layers

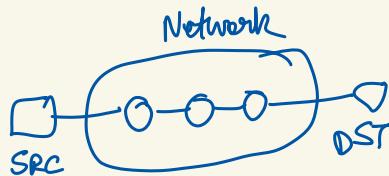
May not be a drop (guess)



- (2) Redundancy of Tasks  
 (3) Suboptimality

Ex:

$\dots \rightarrow P$



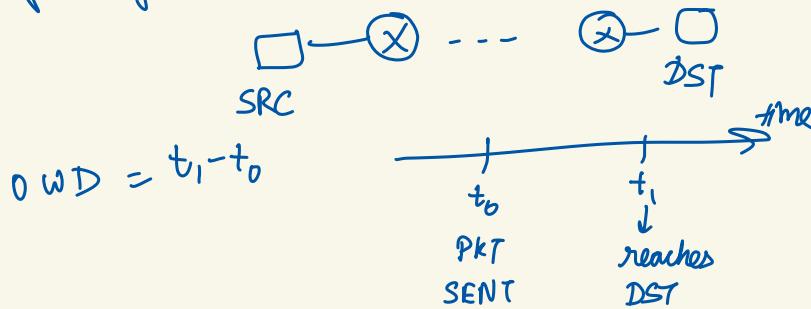
Best-effort: No guarantees on QoS

Quality of Service : amt of packet drops

Latency (delays)

Latency Delay:

(i) One-way delay



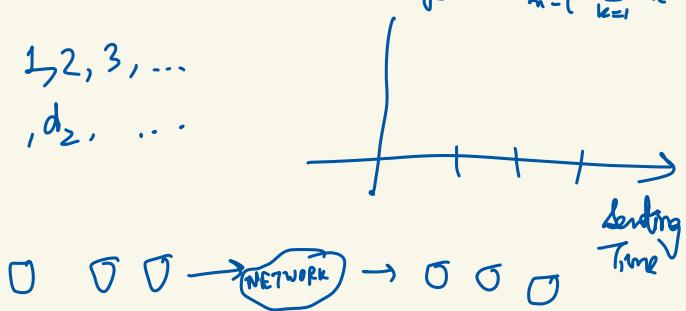
Jitter: Variability in end latencies

PKT S : 1, 2, 3, ...

OWD d<sub>1</sub>, d<sub>2</sub>, ...

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

$$\text{Average Jitter} = \frac{1}{n-1} \sum_{k=1}^n e_k$$

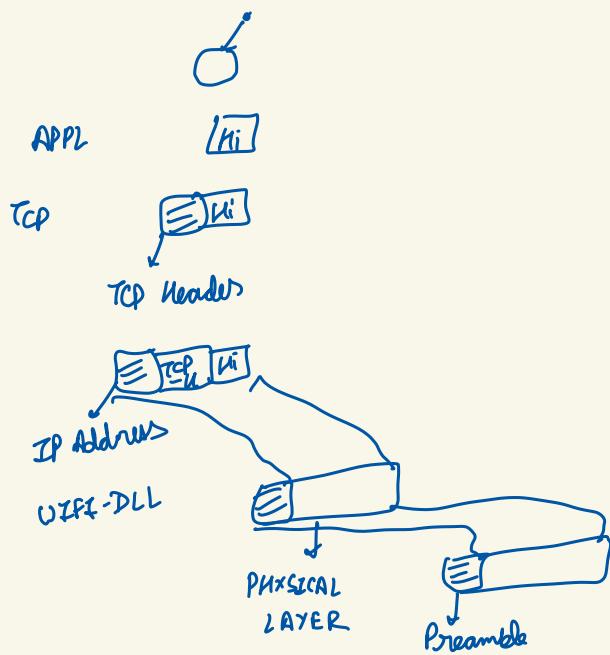


# Telephone Network

First one (by Graham Bell) → voice

Low jitter  
RTT, QoS

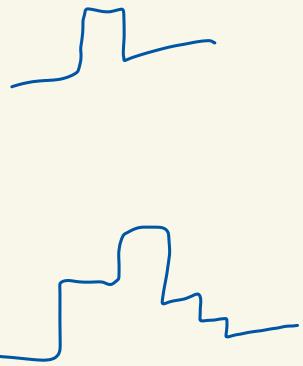
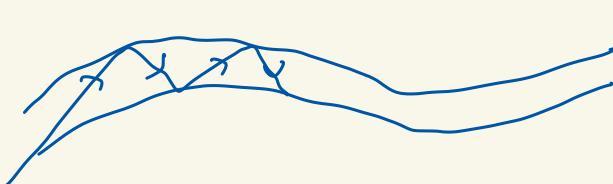
No data loss



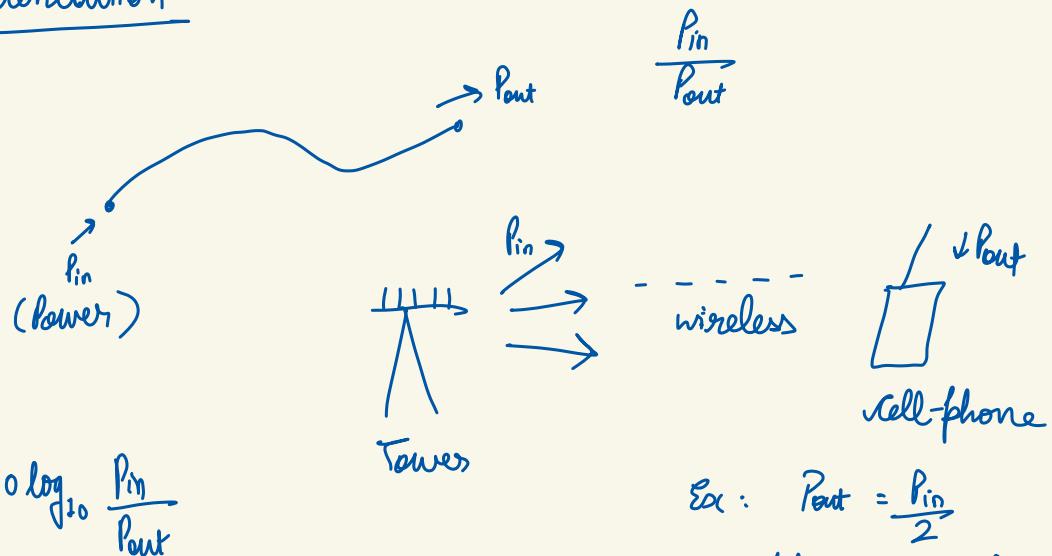
## PHY LAYER

Optic Fibre  $\rightarrow$  submarine cables

Single-mode better than multi-mode



## Attenuation



$$10 \log_{10} \frac{P_{in}}{P_{out}}$$

$\rightarrow$  Attenuation in decibels (dB)

$$P_{out} = \frac{P_{in}}{10} \quad \text{Atten} = 10 \text{ dB}$$

$$\text{Ex.: } P_{out} = \frac{P_{in}}{2}$$

$$\begin{aligned} \text{Atten} &= 10 \log_{10} \frac{P_{in}}{P_{out}} \\ &= 10 \log_{10} 2 \\ &\approx 3 \text{ dB} \end{aligned}$$



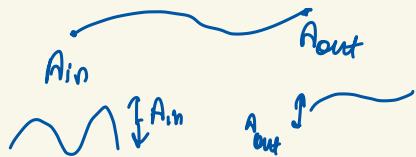
$$10 \log_{10} \frac{P_{in}}{P_f} = 3$$

$$10 \log_{10} \frac{P_f}{P_{out}} = 3$$

$$6 = 10 \log_{10} \left( \frac{P_{in}}{P_f} \times \frac{P_f}{P_{out}} \right) = 10 \log_{10} \frac{P_{in}}{P_{out}}$$

$\frac{P_{in}}{P_{out}}$	dB
2	3
4	6
10	10
100	20
1000	30
10000	40

Power  $\propto (\text{Ampl.})^2$   
Attenuation



$$= 10 \log_{10} \left( \frac{A_{in}}{A_{out}} \right)^2 = 20 \log_{10} \frac{A_{in}}{A_{out}}$$

Absolute power in decibel scale

1 mW as a reference

Power P (watts) in  $\text{dBm}$

Ex: (1)  $P = 1 \text{ mW} = 10^{-3} \text{ W}$

$$10 \log_{10} \frac{P}{10^{-3}}$$

$$10 \log_{10} \frac{10^{-3}}{10^{-3}} = 0 \text{ dBm}$$

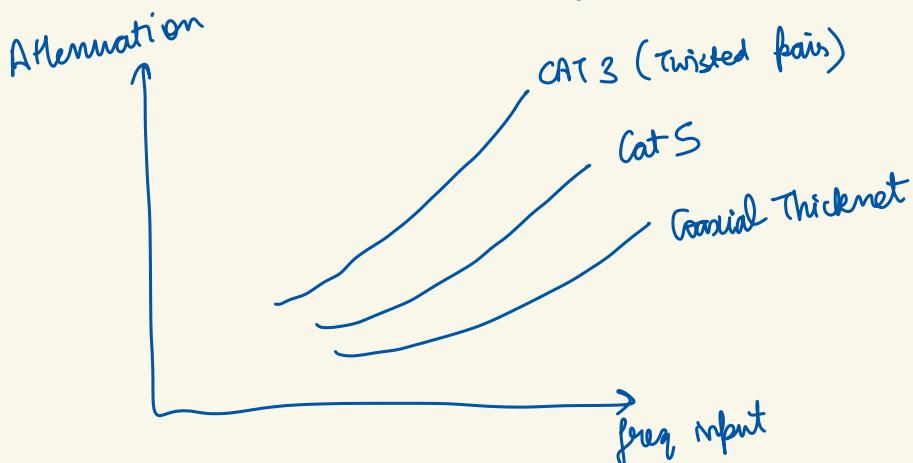
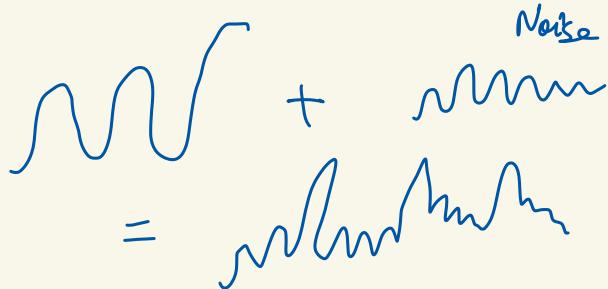
(2)  $P = 2 \text{ mW}$   
(3)  $P = 10 \text{ mW}$

$$\begin{aligned} 10 \log_{10} 2 &= 3 \text{ dBm} \\ 10 \log_{10} 10 &= 10 \text{ dBm} \end{aligned}$$

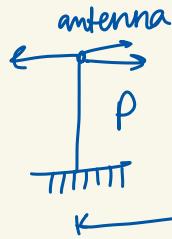
$$(4) P = 1 \mu\text{W} = 10^{-6} \text{ W}$$

$$10 \log_{10} \frac{10^{-6}}{10^{-3}} = -30 \text{ dBm}$$

Received Power  
\_\_\_\_\_  
Noise power



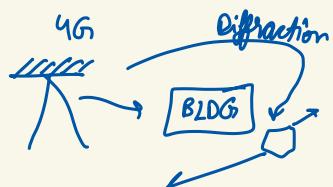
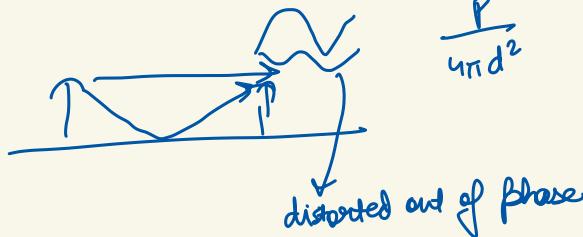
### Wireless

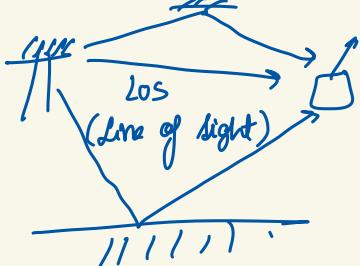


$$P_{rec} = \frac{PA}{4\pi d^2}$$

$$\text{Reality} \rightarrow P_{rec} \propto \frac{P}{d^\alpha}$$

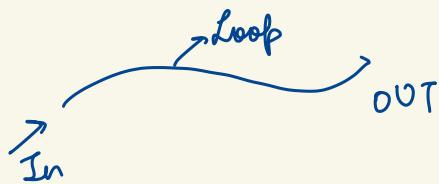
$$5 > \alpha > 2$$



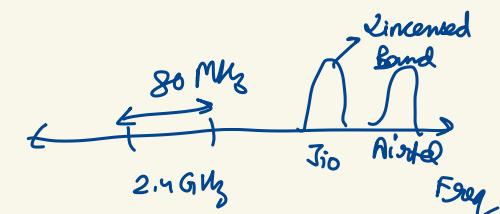
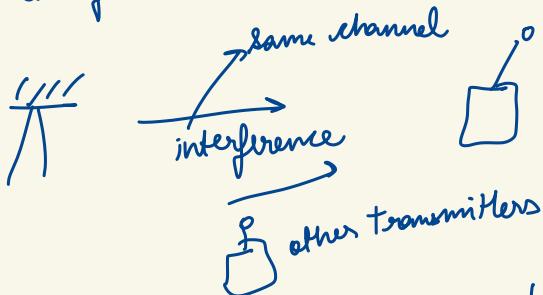


Shannon → Entropy  
→ Capacity  
↓  
Max data rate  
for communication over  
a channel

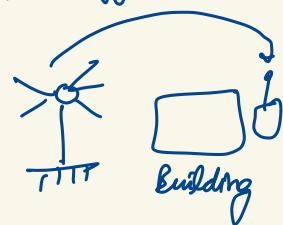
### Wireless Channels



- (1) Attenuation much higher than wired
- (2) Interference



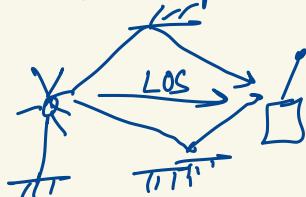
### (3) Diffraction



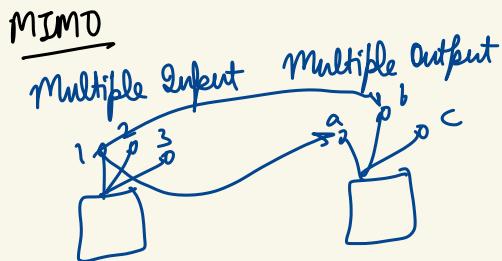
wifi  
unlicensed  
band

LOS = line  
of sight

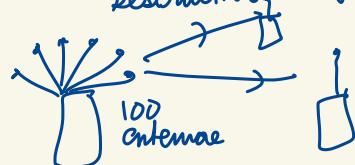
### (4) Multi-path



Constructive and  
Destructive Interference



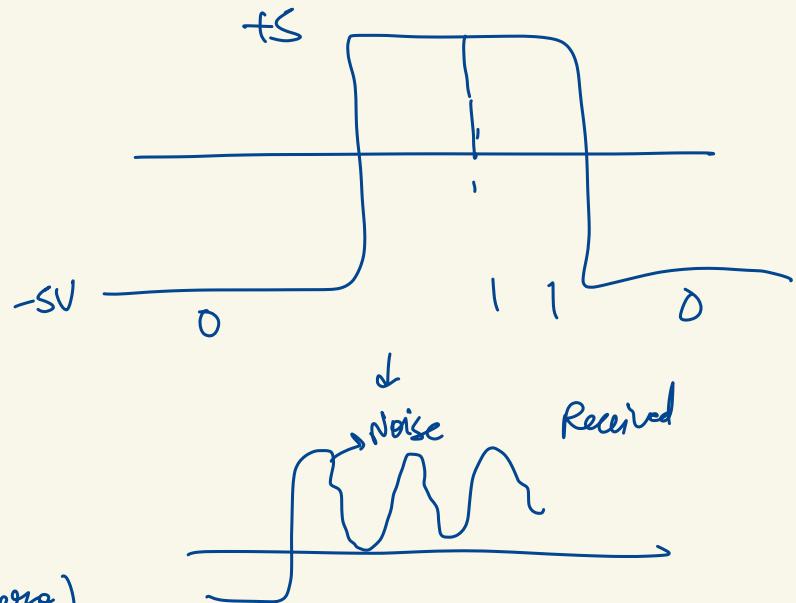
Massive  
MIMO



## Signalling

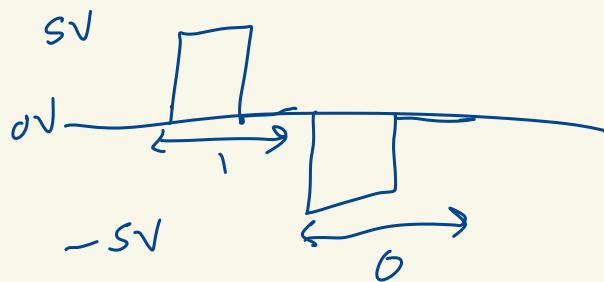
### Wired

Non-return to zero (NRZ)

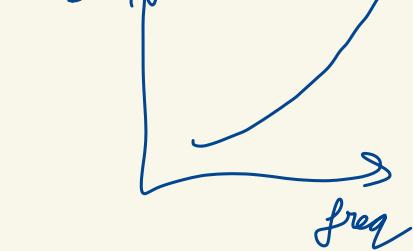


RZ (return to zero)

transmit the signal, then return to zero in each pulse

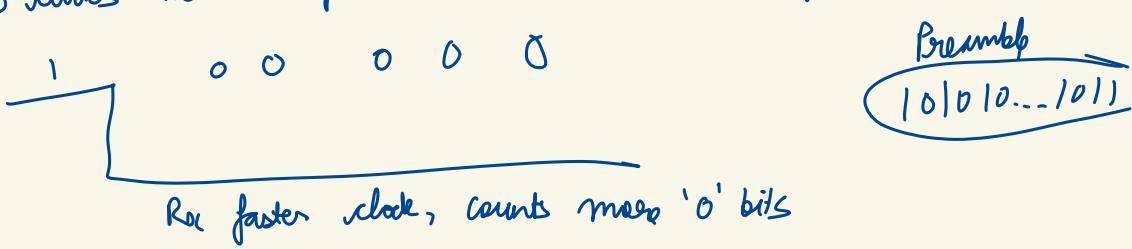


dB (Attenuation)

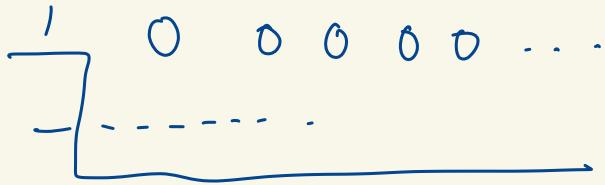


## Issues with NRZ

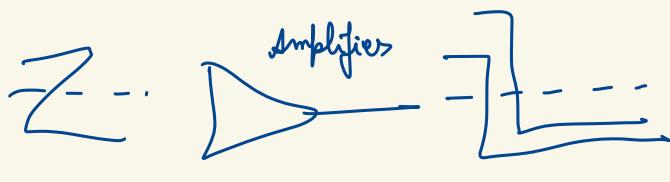
- (1) Clocks not in sync (frequencies can be different) at Tx/Rx



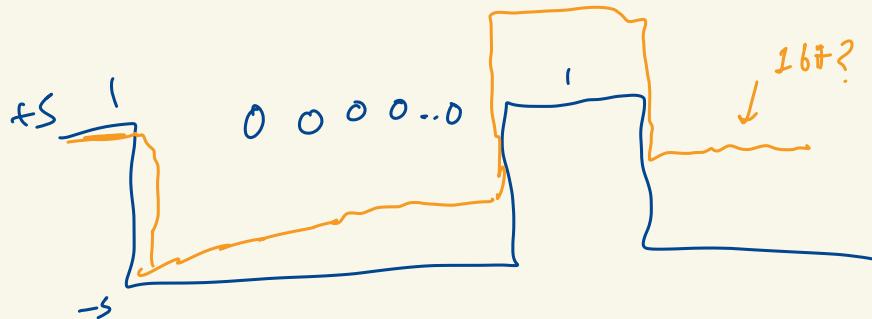
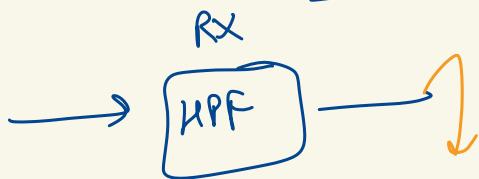
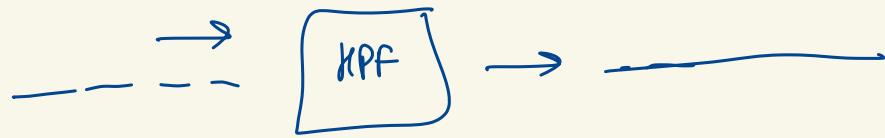
## (2) Baseline wander



Average signal power at the receiver is taken as baseline -  
if it is high (because of more 1s) then noise affects the further 1s



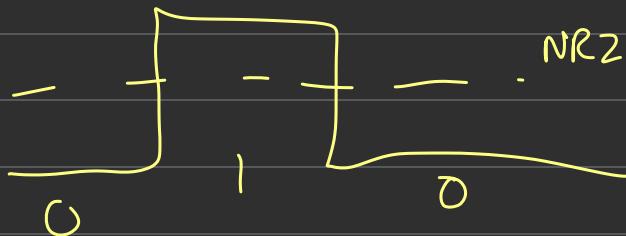
High-pass filter → Removes DC (offset) and low freq. signals



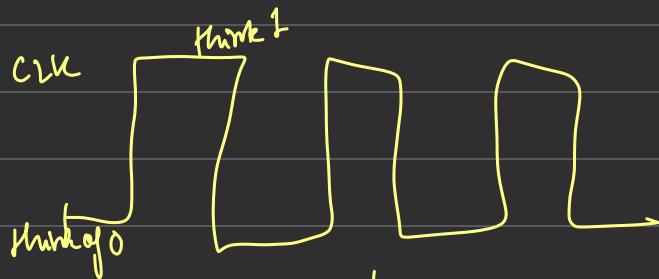
A high pass filter is an equalisation tool that removes all frequencies below a set point

Bandpass filter allows signals within a selected range of frequencies to be heard/decoded.

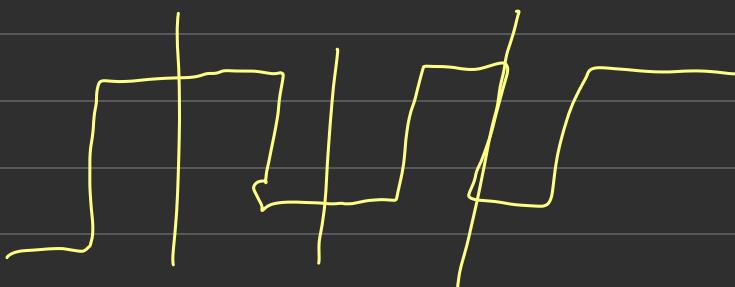
### Manchester Coding



$$\begin{array}{r|l} \text{XOR} & \\ 0 & 0 \\ 0 & 1 \end{array}$$



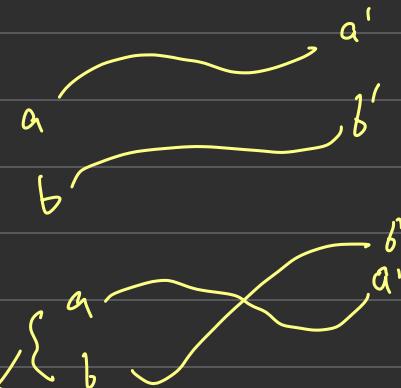
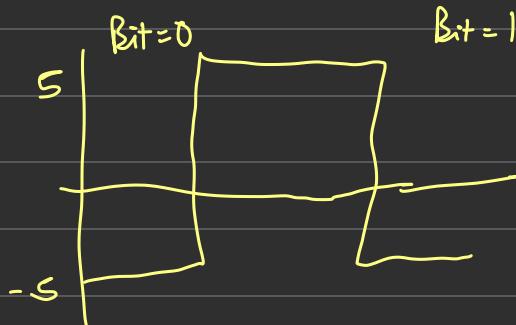
- (a) signal transition every bit period  
(b) Avg signal per bit period is 0



Manchester  
coded signal

Line Coding → Wired media Bits → Signals

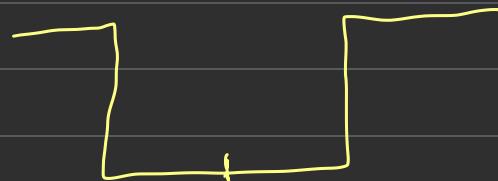
### Manchester Coding



If flipping happens, then receiver thinks:

Need to know polarity

Ideas: Idea → Default first bit at 0.



T<sub>x</sub>

R<sub>x</sub>



End of preamble  
can be used to  
convey about polarity

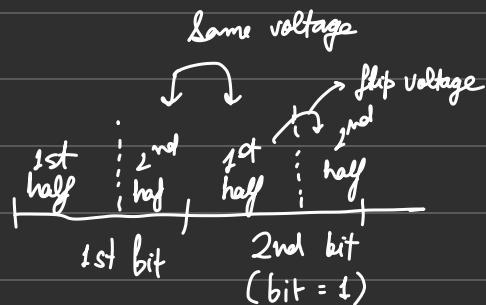
Ethernet       $101010\dots11$        $\underbrace{\hspace{1cm}}$       Consecutive 1s

← Receiver might  
start receiving  
later

## Differential Manchester Coding

Rule:  $\text{bit} = 1 \rightarrow$

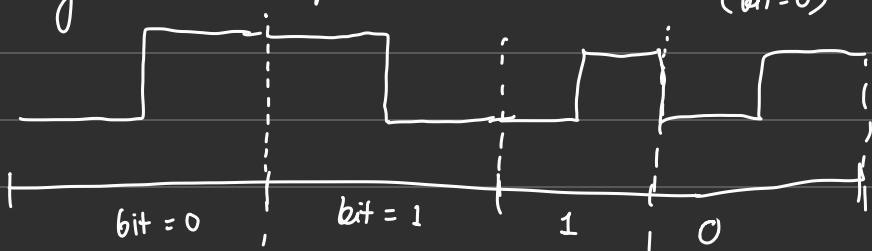
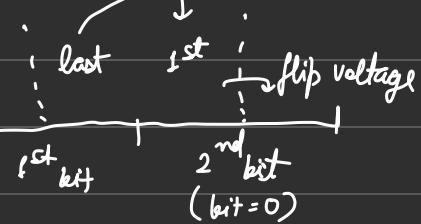
Voltage in 1<sup>st</sup> half of bit period  
is same as last half of  
prev. bit.



If  $\text{bit} = 0 \rightarrow$

Voltage in 1<sup>st</sup> half of bit is opp.  
to that of prev. bit period

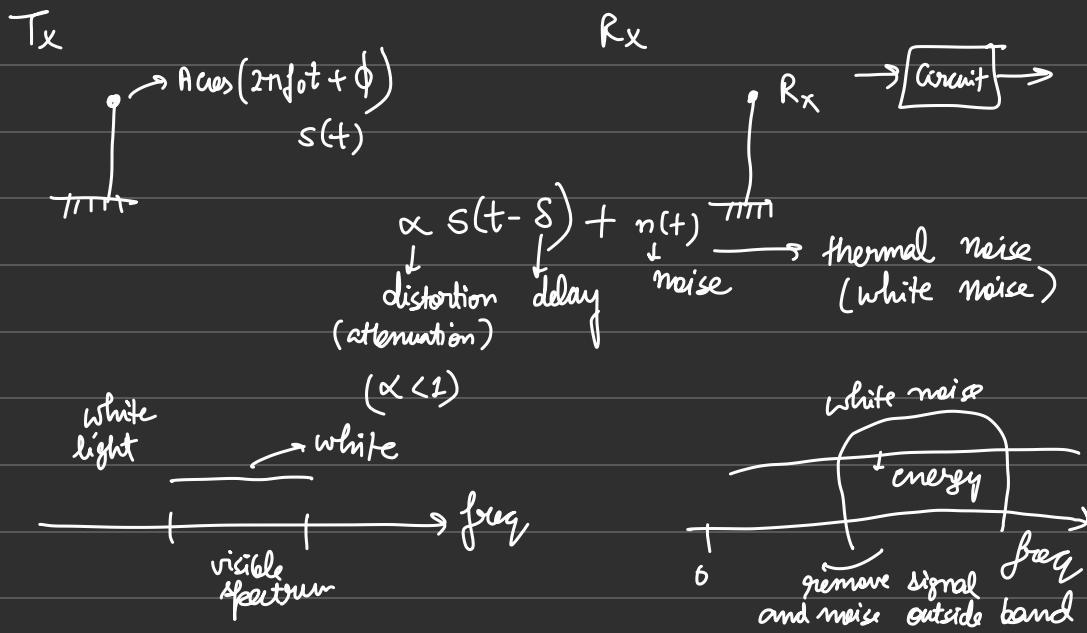
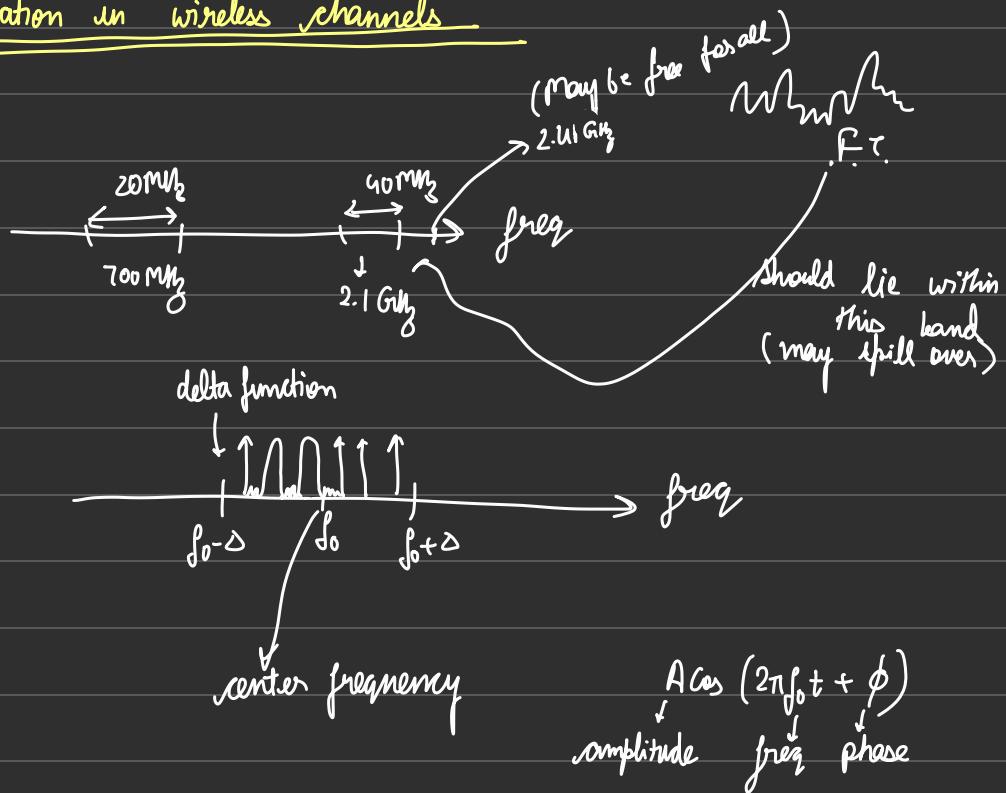
opposite



Even if this is flipped over, the message transmitted remains the same.

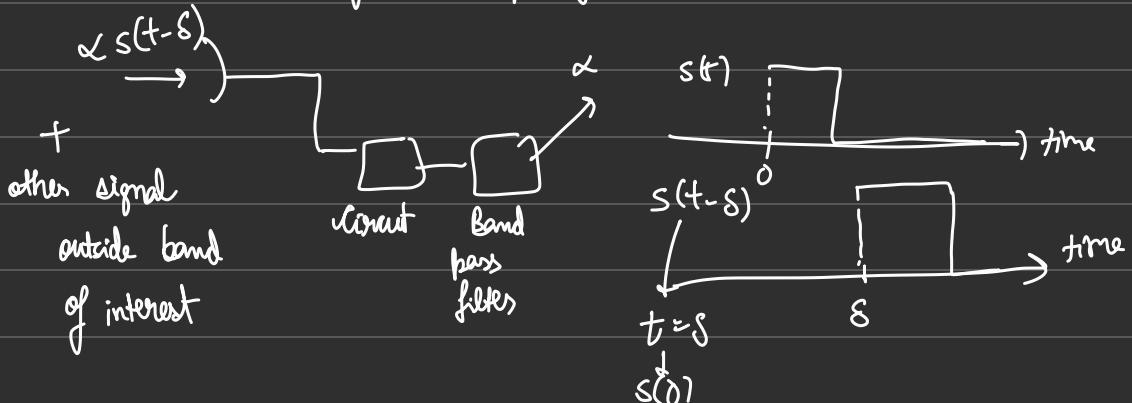
Default signal for the first bit.

## Modulation in wireless channels

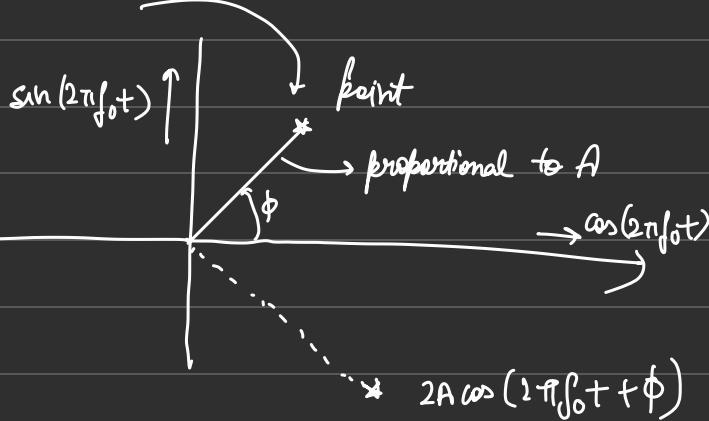


$$\alpha s(t-\delta) + n(t)$$

noise after band pass filter



$A \cos(2\pi f_0 t - \phi)$  → vector space



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

In this vector space

$$a(t) = \dots \quad 0 \leq t \leq T$$

$$b(t) = \dots$$

Inner product

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

3d space

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

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

Dot product

$$\langle a, b \rangle = a_x b_x + a_y b_y + a_z b_z$$

↳ Information about the angle

## Unit Vectors

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

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

$$T = \frac{1}{f_0}$$

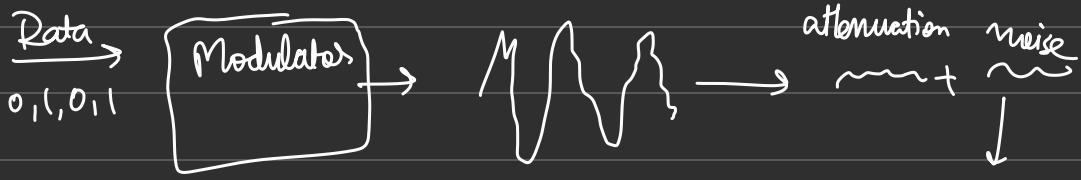
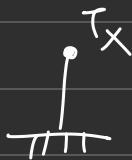
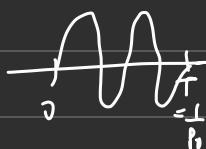
Q.: Inner product of these?

$$\langle s_1(t), s_2(t) \rangle = \frac{2}{T} \int_0^T \underbrace{\sin 2\pi f_0 t}_{\frac{1}{2} \sin 4\pi f_0 t} \underbrace{\cos 2\pi f_0 t}_{2 \text{ time periods in } 0 \text{ to } T} dt$$

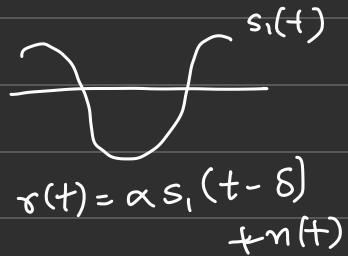
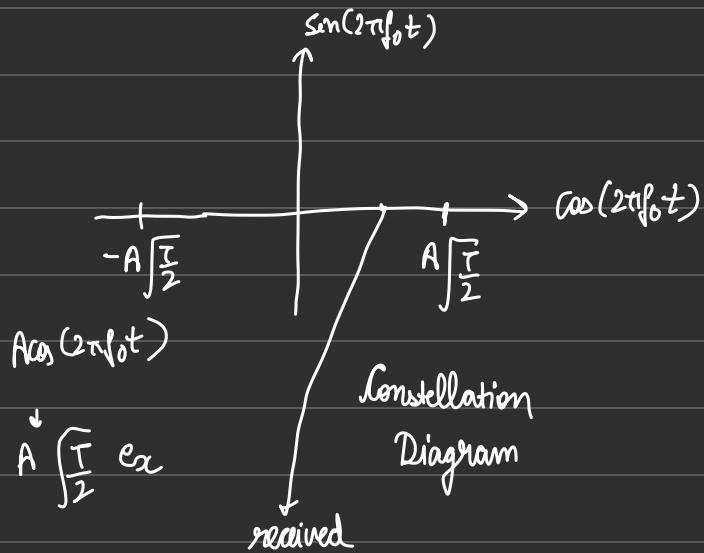
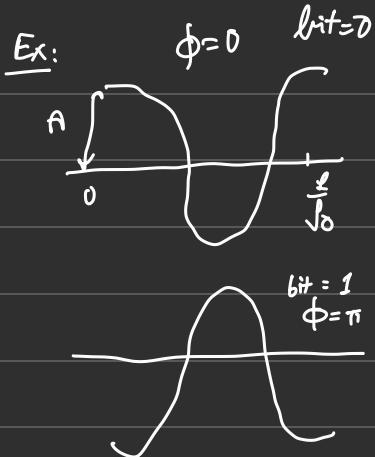
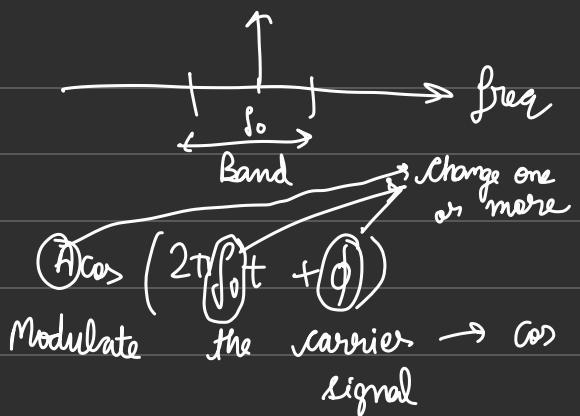
$$\langle s_1(t), s_1(t) \rangle = \frac{2}{T} \int_0^T \cos^2(2\pi f_0 t) dt$$

$$= \frac{1}{T} \int_0^T 1 + \cos(4\pi f_0 t) dt$$

$$= 1 \quad \text{Integration} = 0$$



Demodulator  
of Data



Noise : white, Gaussian, Additive



$$s_{i,x} = \langle s_i(t), e_x \rangle$$

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

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

$$s_1(t) = A \cos(2\pi f_0 t)$$

$$s_{1,x} = \langle s_1(t), e_x \rangle = \int_0^T A \cos(2\pi f_0 t) dt$$

$$= A \sqrt{\frac{T}{2}} \underbrace{\left( \sqrt{2} \cos(2\pi f_0 t) \right)}_{\sin(4\pi f_0 t)}$$

$$s_{1,y} = \langle s_1(t), e_y \rangle = \int_0^T \sqrt{\frac{2}{T}} \underbrace{\sin(2\pi f_0 t) dt}_{\text{integral over } T \text{ is 0}} \underbrace{(A \cos(2\pi f_0 t))}_{\sin(4\pi f_0 t)}$$

$r(t) \rightarrow r_x, r_y ?$

$$r_x = \langle r(t), e_x \rangle \xrightarrow{s_m} \cos(2\pi f_0 t)$$

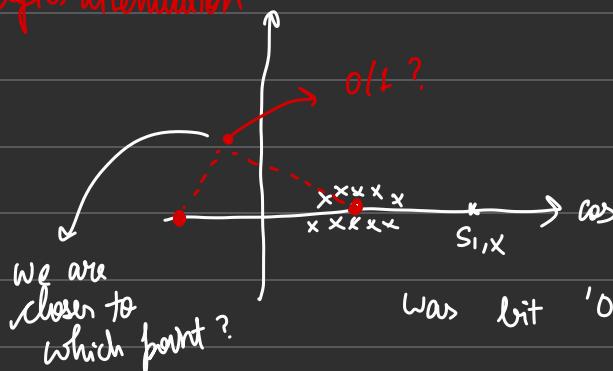
$$r_y = \langle r(t), e_y \rangle \xrightarrow{s_m} \sin(2\pi f_0 t)$$

$(s_{0,x}, s_{0,y})$

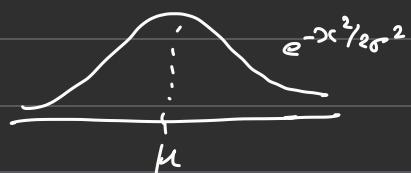
$$\text{Suppose } r(t) = \alpha s_1(t) + n(t)$$

$n_x = \langle n(t), e_x \rangle$   
 $n_y = \langle n(t), e_y \rangle$   
 Each is identically distributed, independent gaussian random variable

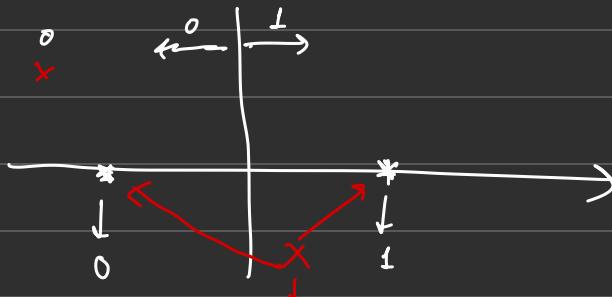
• After attenuation



Was bit '0' sent or '1'?



Suppose we know constellation after attenuation and without noise

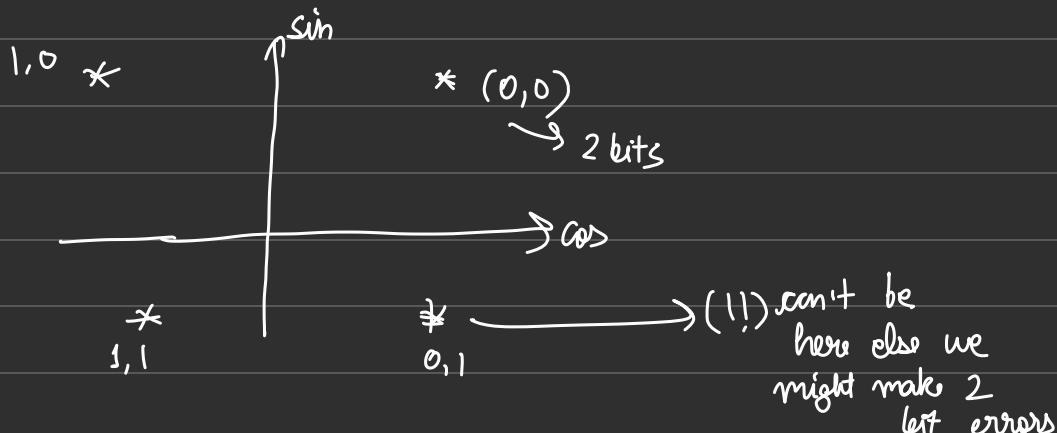


received  $\rightarrow$  closer to 1  $\rightarrow$  So 1 was received

Anything on left side  $\rightarrow$  0  
right side  $\rightarrow$  1

QPSK

Quadrature PSK

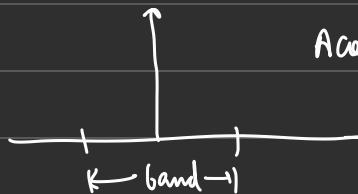


PHY  
↓  
DLL

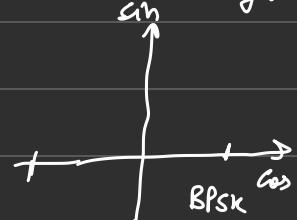
$s(t)$  →  
 $\tau_x$

$\tau(t)$   
↓  
signal noise

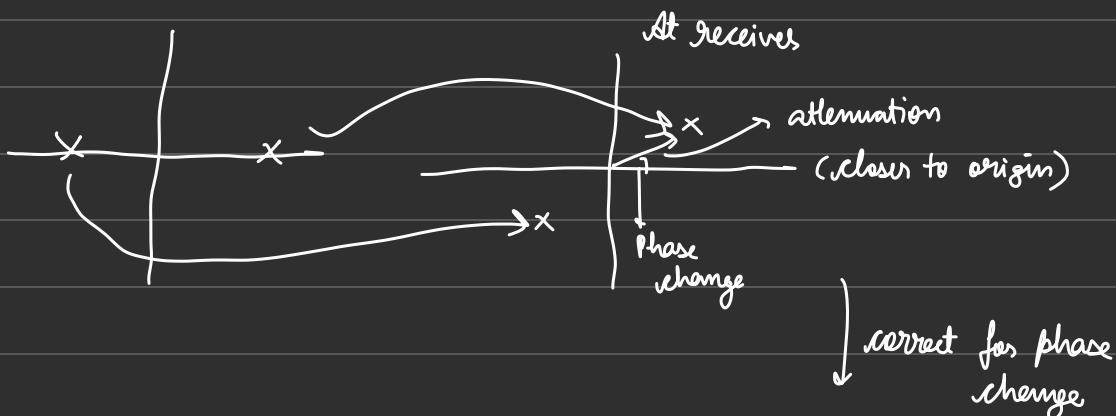
Constellation  
Diagram



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



signal → channel →  $R_x$   
 $\tau_x$



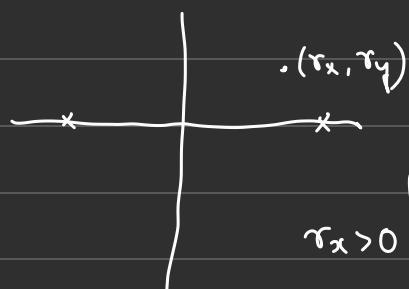
$\tau(t)$

$$\tau_x = \langle \tau(t), e_x(t) \rangle$$

$$\tau_y = \langle \tau(t), e_y(t) \rangle$$



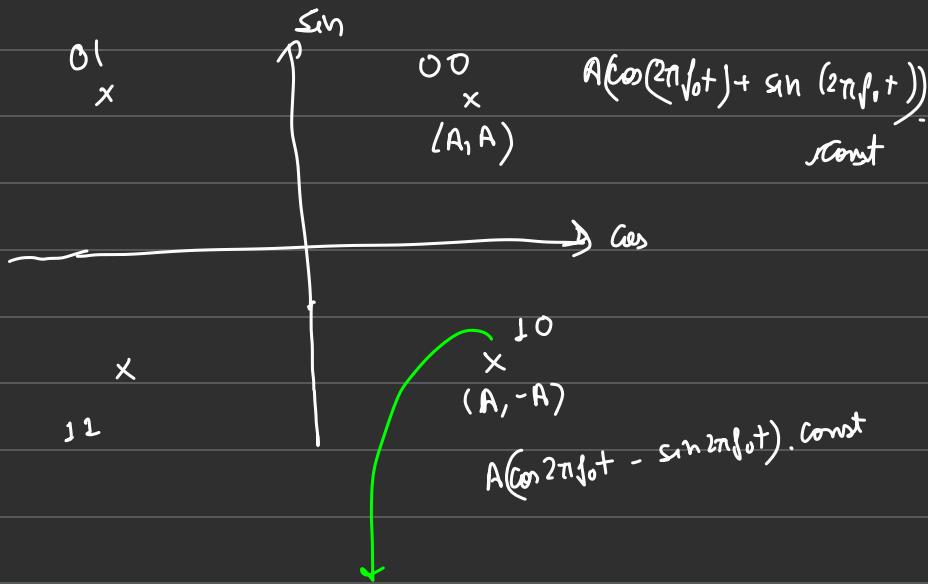
(rotate back  
by  $\phi$ )



AWGN → Additive white Gaussian noise

BPSK → Binary Phase Shift Keying  
 $\tau_x > 0 \quad \tau_y \leq 0$

QPSK



so? Not a good idea

because  $00$  &  $11$  are near

hence, we might have 2 bit errors instead of one

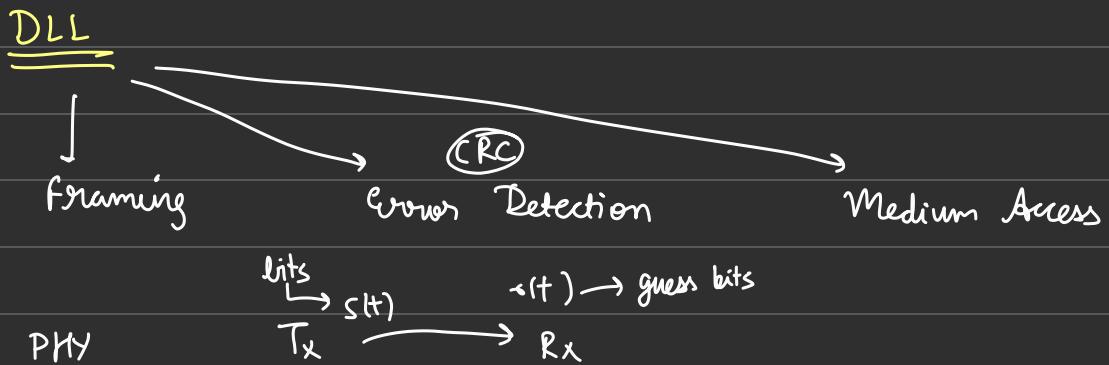
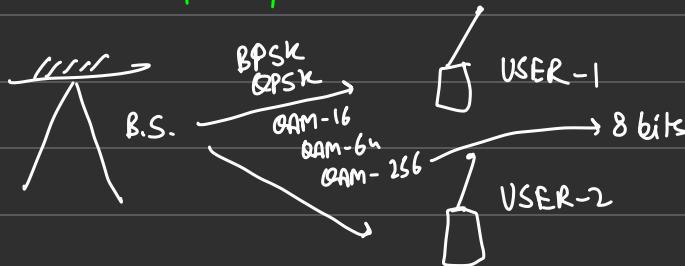
QAM-16





Given  $P_s$  (bit errors) <  $\epsilon$

Attenuation?



data 1011 01011 011010

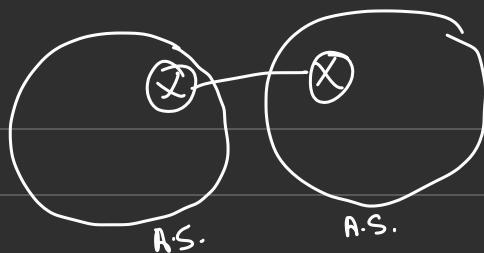
IP adder?



HDLC

High level Data link control

WAN  $\rightarrow$  Wide Area Network



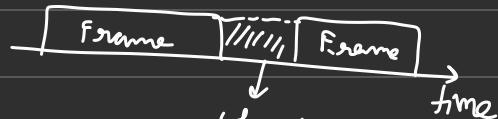
## Synchronous Mode (HDLC)



Always → transmit

HDLC  $\rightarrow$  Wired (NRZ, Manchester...)

Default Seq: 0111110



what to  
send here?

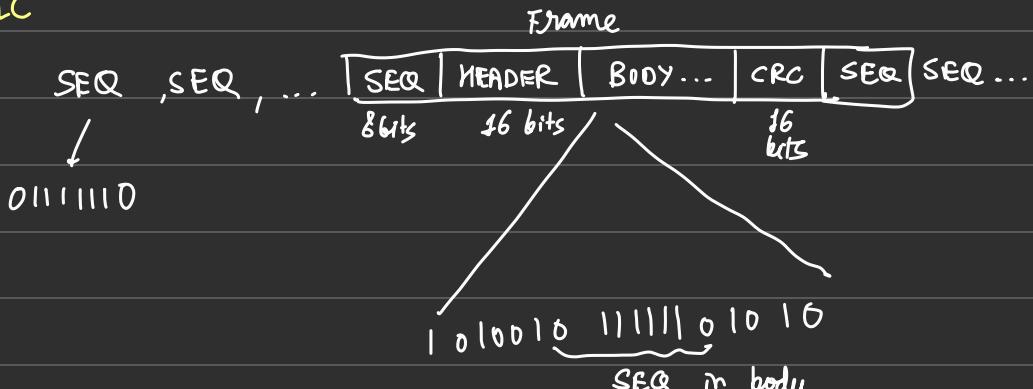
The diagram illustrates a frame structure with the following fields:

- seq**: 8 bits.
- Frame**: 16 bits.
- Header**: 8 bits.
- Body**: The main data payload.
- CRC**: 8 bits.
- seq**: 8 bits.
- check**: A cyclic redundancy check field.

Bit stuffing is indicated at the bottom left, and cyclic redundancy check is indicated at the bottom right.

DL

## HdLC



Sol<sup>n</sup> - Put zeros at every alternate position  $\xrightarrow{\text{So } n \text{ bits}}$  lot of redundant bit.

*sol<sup>m</sup>:* Every time after 5 ones  $\rightarrow$  add a 0.

At transmitter :

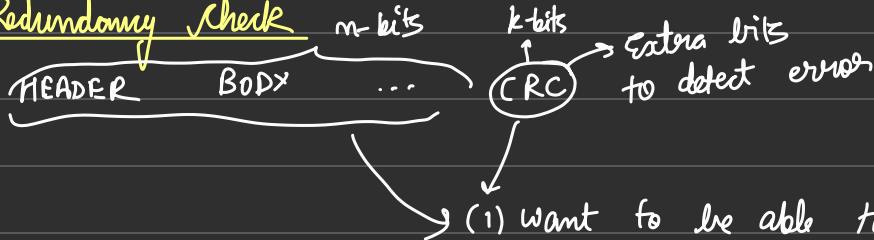
00 111 11 0 1010 111 11 0 1010 111 111 010

At receiver:

If seq 0 11111

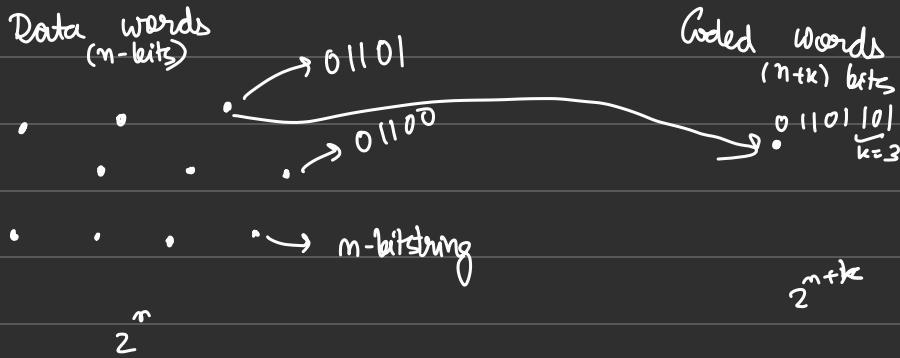
$\left\{ \begin{array}{l} 0 \Rightarrow \text{assume bit stuffing and remove} \\ 10 \Rightarrow \text{seq at end of frame} \\ 11 \Rightarrow \text{Error, so discard frame} \end{array} \right.$

### Cyclic Redundancy Check



(2) Creation of CRC and Verification of CRC should be computationally efficient

(3) for given 'k', CRC should be computable for any 'n'



Hamming Distance : Given two codewords  $a_1, a_2, \dots, a_m$   $b_1, b_2, \dots, b_m$

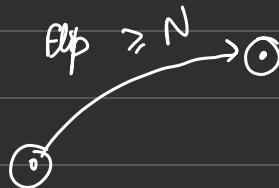
# bits in which they differ

Min. Hamming Distance of a code.

pairs of codewords

Error Detection: If min. HD is  $N$  (of a coding scheme) then we can always detect  $(N-1)$  or less bit errors.

We detect that there were errors but may not know position of bit errors.

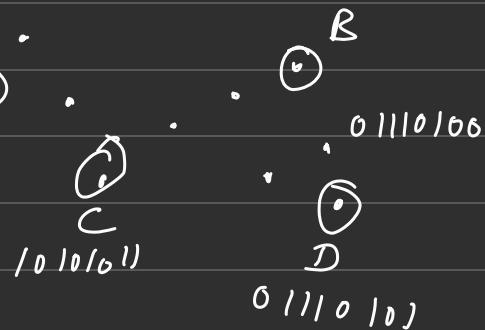


### Error Correction

Suppose min. HD is  $(2t+1)$

and  $t$  bit errors is ' $t$ ' or

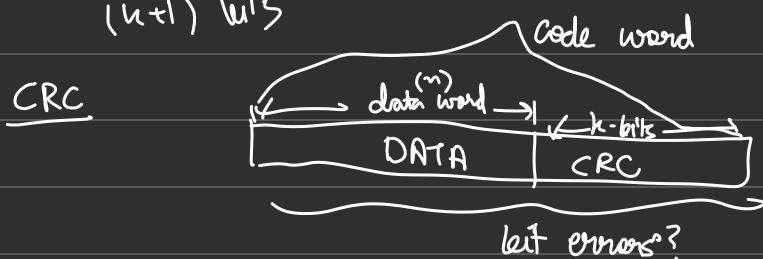
less. Then by mapping received code word to the nearest valid code word corrects all errors.



Graal's fields

1101

Division  
 $(n+1)$  bits



<sup>2</sup>  
dataword  
space

$$\text{Codeword space} \rightarrow 2^{m+n}$$

## Galois Fields

		0, 1		XOR
x	0		+	0
0	0	0	0	0
1	0	1	1	0

Data: 110110 K=3

Divisor / Generators : 1101 (k+1 bits long)

depending on the first bit of  
the current dividend

Codeword =

$110110 \underbrace{111}_{\text{CRC}}$

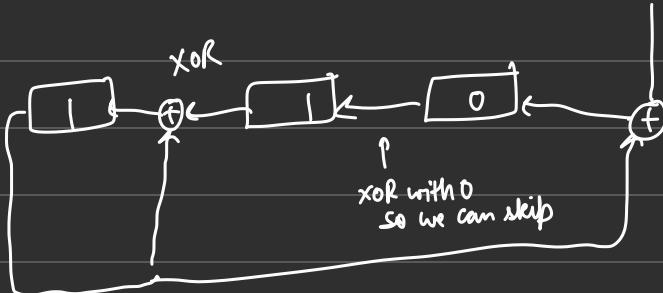
Remainder

Initially with all  
zeros

1101

Data

110010 | 000



At receiver:

Kodeword

Data | CRC

## Concatenation

## Polynomial representation of Bitstrings

$$\begin{array}{l} \text{11 O} \\ \downarrow \quad \downarrow \quad \downarrow \\ x^3 \quad x^2 \quad x^1 \quad x^0 \end{array} \rightarrow 1 \cdot x^3 + 1 \cdot x^2 + 0 \cdot x + 1 \cdot x^0$$

$x^3 + x^2 + 1 \rightarrow C(x)$

(Divisor / Gen. Poly)

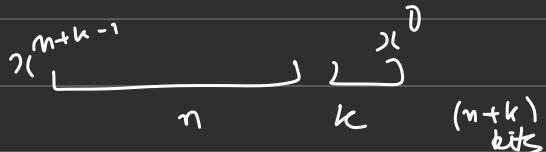
$$\begin{array}{ccccccc}
 \text{vars?} & & & & & & \rightarrow p(x) \\
 \hline
 110 & 110 & 111 & & & & \\
 000 & 001 & \downarrow & 00 & \downarrow x^0 & \rightarrow E(x) & \text{Errors bitstring} \\
 & & x^3 & & & &
 \end{array}$$

$$\underline{Q} \quad \frac{P(x) + E(x)}{C(x)} = 0?$$

If not zero then detected errors

### Types of Errors

(1) Single bit errors



$$E(x) = x^i ; \text{ for some } i \in \{0, 1, \dots, n-k+1\}$$

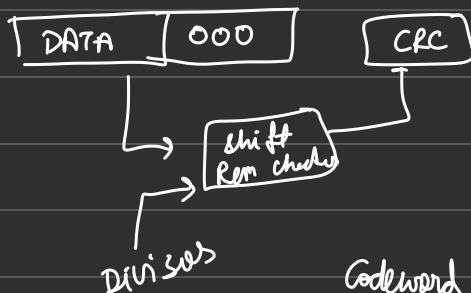
$$\frac{P(x) + E(x)}{C(x)} = \frac{P(x)}{C(x)} + \frac{E(x)}{C(x)}$$

~~$= 0$~~

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

$$C(x) = [x^n + \dots + x^2]$$

$x^{n+m}$  highest       $x^2$  lowest



TRANSMIT



$$\text{Min HD} = d$$

↓  
fewer than  $d$  errors  
always detected  
\* May detect errors  
of  $n, d$  bits  
No guarantee

↓  
 $\text{Rem} = 0 ? ?$   
 $\text{Rem} \neq 0 \Rightarrow \text{Error}$

## Polynomial Arithmetic

$$1101 \rightarrow 1 \cdot x^3 + 1 \cdot x^1 + 0 \cdot x + 1 = x^3 + x^1 + 1$$

GIF(2) : Addition / Substr (xor)

Data: 110110  $\rightarrow x^5 + x^4 + x^2 + x$

$$\begin{array}{r} x^3 \\ \hline 1101 \sqrt{110110 | 000} \\ \hline x^3 \\ \hline x^3 + x^2 + 1 \\ \hline \text{Remainder} \\ = x^2 + x + 1 \\ = (111) \end{array}$$

If A(x) is divisible by B(x)

$$A(x) = B(x) \cdot D(x)$$

We will send

$$110110 \mid 111 \rightarrow x^2 (x^5 + x^4 + x^2 + x) + x^2 + x + 1$$

$P(x) = \text{codeword}$   
 $C(x) \rightarrow \text{divisor (generator)}$   
 $E(x) \rightarrow \text{error polygon}$

$P(x) + E(x)$  is received  
 $E(x) = x^6 + x^3$



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

$P(x)$  chosen so that  $C(x)$  divides it

Q1 Does  $C(x)$  divide  $E(x)$ ?

## Types of errors



(1) Single bit errors

$$E(x) \in x^i \text{ for some } i \in \{0, 1, \dots, n-k-1\}$$

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

then  $C(x)$  does not divide  $E(x)$

$$\frac{x^i}{x^k + \dots + 1} \quad \left| \begin{array}{l} C(x) \mid D(x) \stackrel{?}{=} x^i \\ I \\ C(x)(x^m + \dots + x^k) \\ m \geq q \\ x^{k+1} + \dots + x^q \end{array} \right.$$

$$k+m > q$$

(2) Two bit errors

$$E(x) = x^j + x^i \quad (j > i) \\ = x^i (x^{j-i} + 1)$$

Suppose  $C(x)$  is of form  $x^k + \dots + 1$

$$\frac{E(x)}{C(x)} = \frac{x^i (x^{j-i} + 1)}{x^k + \dots + 1}$$

know  $x^i$  terms  
not cancelled out

Ques Does  $C(x)$  divide  $x^{j-i} + 1$

Def : (order of a Polynomial):

The smallest  $r$  such that  $C(x)$  divides  $x^r + 1$  is called its order.

Known  $C(x) = x^k + \dots + 1$   
s.t. order is  $2^k - 1$

Ex.  $k=16$ ;  $C(x) = x^{16} + \dots + 1$  can find  $C(k)$  s.t.  
it will not divide  
any  $x^p + 1$  for  $p < 2^{16}-1$

(3) Odd number of errors

$$E(x) = x^i + x^i + \dots$$

\ odd numbers of terms

If  $C(x) = (1+x)(\cdot)(\dots)$  then all odd errors detected

If  $C(x)$  has even # terms

$$\text{NDLC CRC} = \underbrace{x^{16} + x^{15}}_{x^{15}(x+1)} + \underbrace{x^2 + 1}_{(x+1)(x+1)} \\ = x^{15}(x^2 + x + 1)$$

want to show

$$\begin{aligned} E(x) \\ = C(x)D(x) \\ \downarrow \\ \text{not possible} \end{aligned}$$

$$x^i + x^i + \dots +$$

\ odd # terms

$$F(1) = 1 + \dots - 1 \quad |$$

$$\text{Case (1)} \quad C(x) = \\ (1+x)(\dots)$$

$$C(1) = (1+1) \quad (\because \overline{0})$$

RHS is 0

Case (1)  $C(x)$  has even # terms

$$\begin{aligned} C(1) &= 1 + \dots + 1 \\ &\quad \text{even} \\ &= 0 \end{aligned}$$

(4) Burst of errors :

Data / CRC : 101 ...  $\underbrace{1101}_{\text{interference}} \dots 11$   $\rightarrow$  consecutive errors

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

$$\frac{E(x)}{c(x)} = \underbrace{x^i [x^{l-1} + x^{l-2} + \dots + 1]}_{x^k + \dots + 1} \quad \begin{matrix} x^i \text{ does} \\ \text{not} \end{matrix}$$

If  $l-1 < k$  then all factors of  $c(x)$  cannot be canceled.

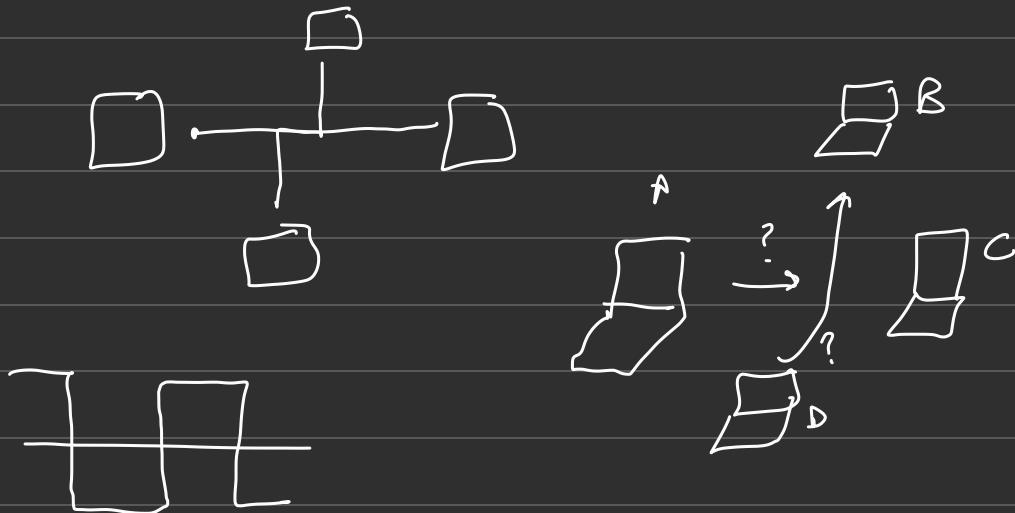
cancel  
any factor in  
the denominator

So  $c(x)$  of form  $x^k + \dots + 1$

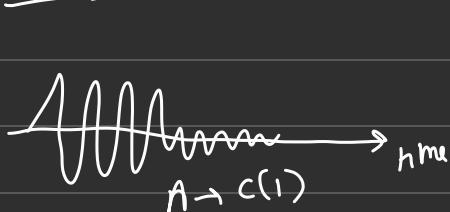
detects all bits of error of length ' $l$ ' such that  $l-1 < k$

Ethernet:  $x^{32} + x^{26} + x^{23} + x^{22} + x^{16} + x^{12} + x^{11} + x^{10} + x^8 + x^7 + x^5 + x^4 + x^2 + x + 1$   
CRC-32:  $x^{32} + x^{26} + x^{23} + x^{22} + x^{16} + x^{12} + x^{11} + x^{10} + x^8 + x^7 + x^5 + x^4 + x^2 + x + 1$

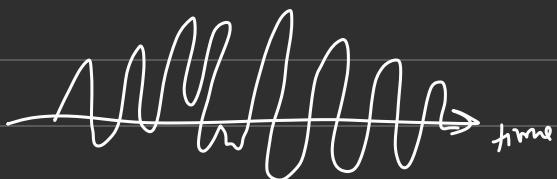
## Medium Access



A + C

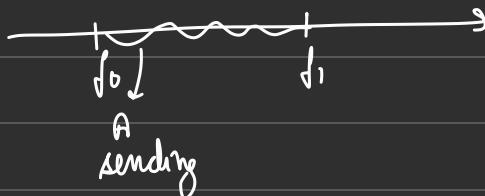


D → B



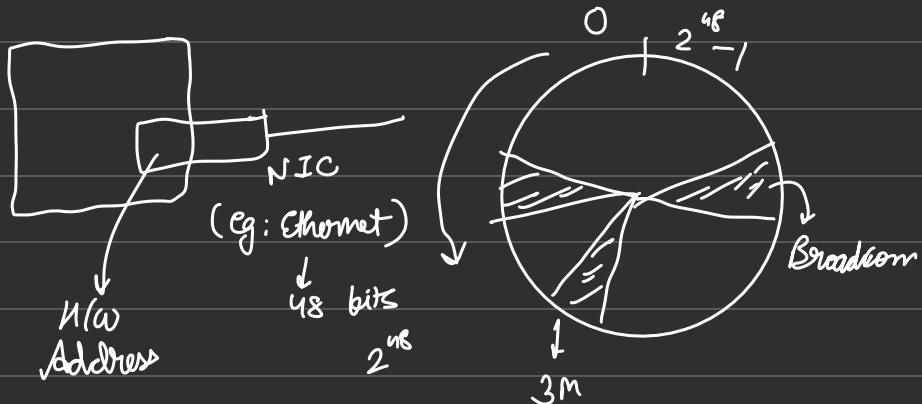
Q) what is being sent, by whom, for whom?

Identifiers

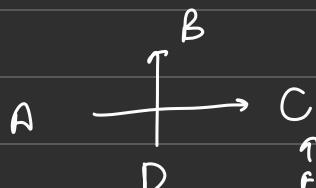


Drawbacks : (1) If  $f_0, f_1$  have to jam  
(2) Throughput of A is  $\frac{1}{2}$  compared to if A uses the full band.

MAC Layer (DLL) → Use hardware address



Interference?

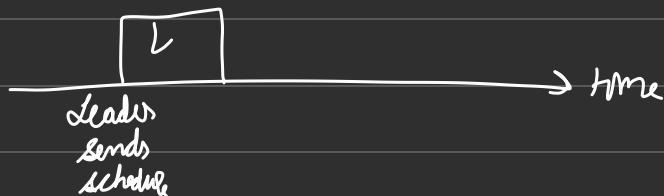


Dif. Freq Bands?

Who uses what bands?  
(hidden terminals A and E are hidden)



Idea: One node is leader



### Token based

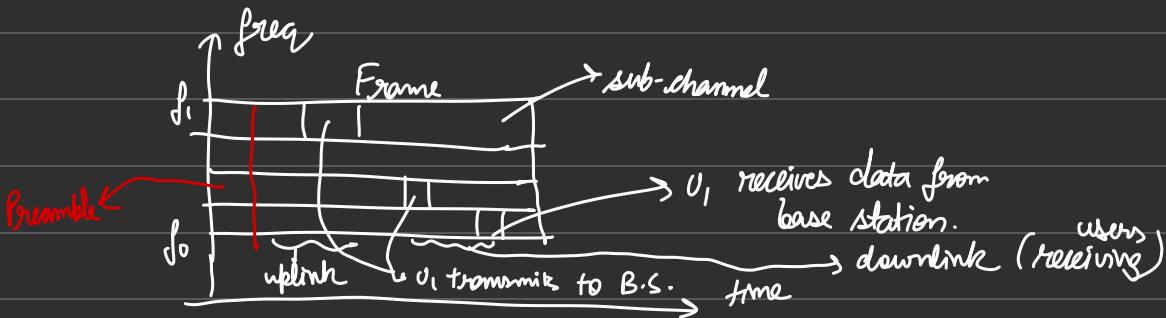


Token release message  $\rightarrow$  say who giving to

### Medium Access



options:  
 (1) Central coordinator  
 (2) Decentralised



Ex.  $BS \rightarrow U_1$  (less attenuation) }  $f_1$ , reverse for  $f_2$   
 $BS \rightarrow U_2$  (more attenuation)

Some column divided into  
different slots

### Issues with central coordinator

- single pt. of failure
- may not be possible in certain situations

wireless

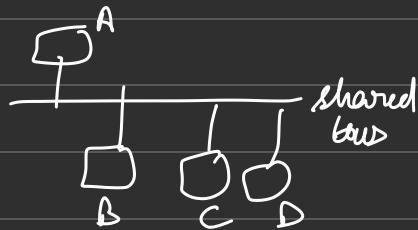
- unlicensed bands

### No central coordinator

want

- (1) Plug & play
  - Connect and disconnect

- (2) No central coordinator



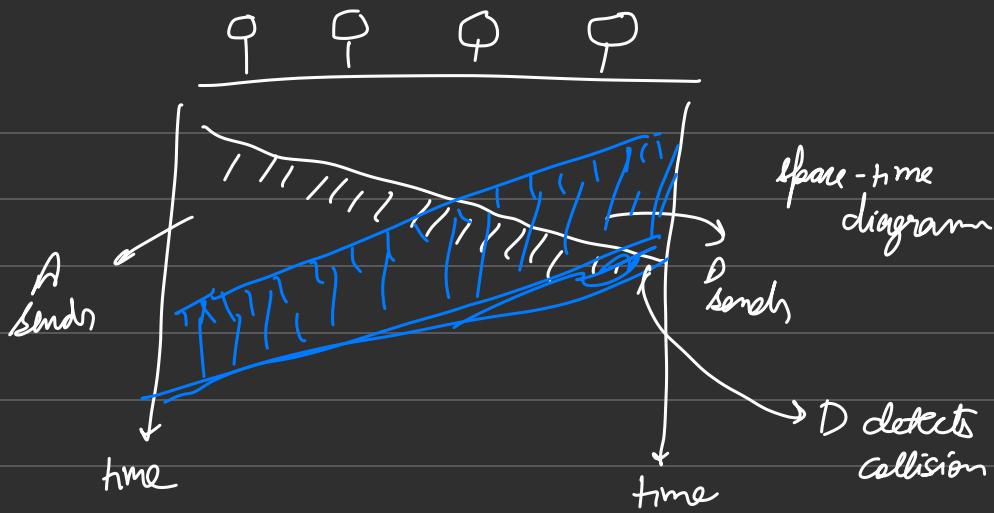
Idea: suppose most of time 1 node transmits, only rarely we have collisions

Collisions: multiple nodes transmit simultaneously

→ signals add up at receiver.

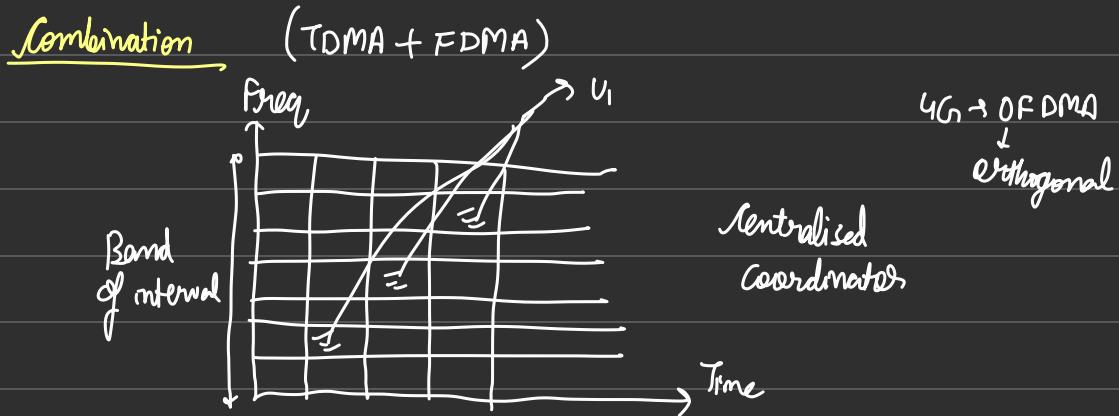
↳ Neither can be deciphered -

If collision occurs → stop transmitting.



Ethernet allows large frame size so that

(A) sender can detect collisions.



Token Passing :

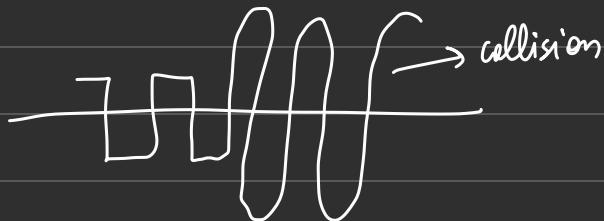
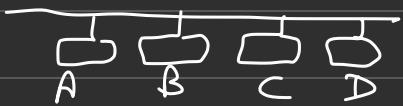


- Only token holder can transmit
- Pass on token every so often

Ethernet:



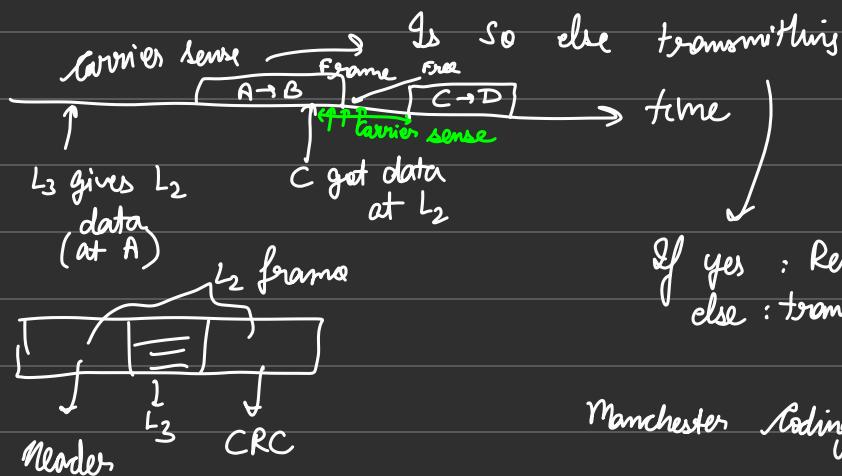
CD: Energy level



IEEE  
802.3

Ethernet

WiFi : 802.11



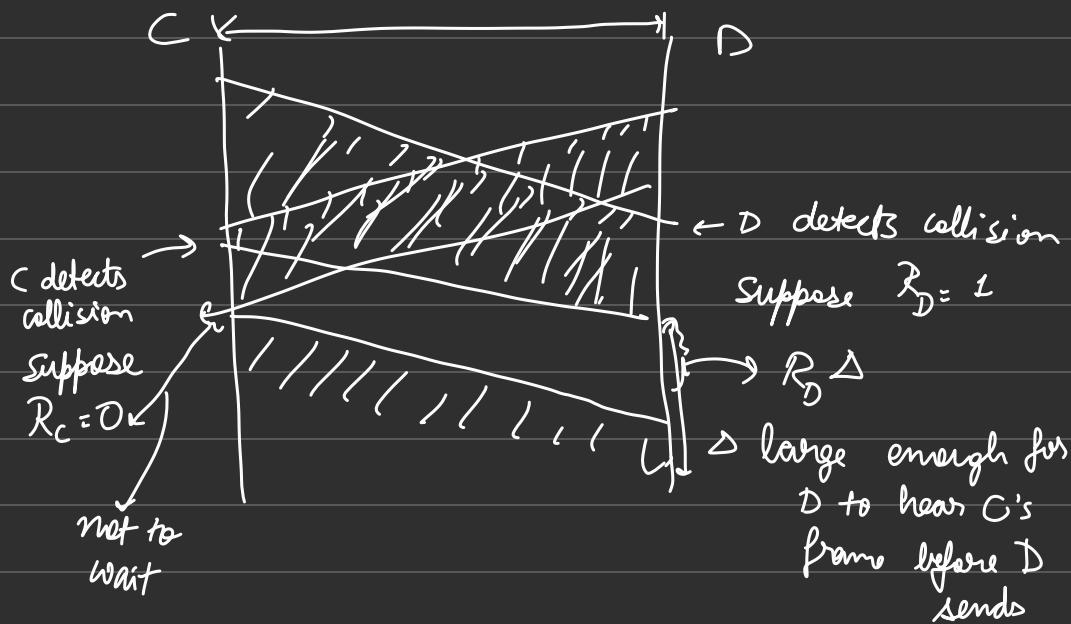
Jammering signal  $\rightarrow$  long enough for all to hear collision.

Random wait time

Each node throws random numbers

$R$ , unif. from some range

wait for  $R_D$  before trying to send again.



Now show if  $D > RTT$ , then D will not transmit before C.

$$RTT = 2 \times \text{time for signal to go from one end of network to the other})$$

Max cable length = 2500m

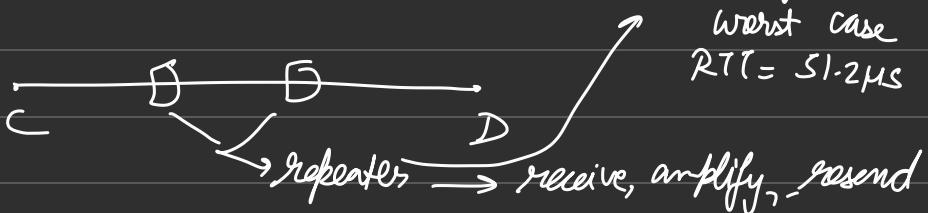
Speed of signal =  $2 \times 10^8$  m/s

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

$$RTT = 2 \times OWD = 25 \mu s$$

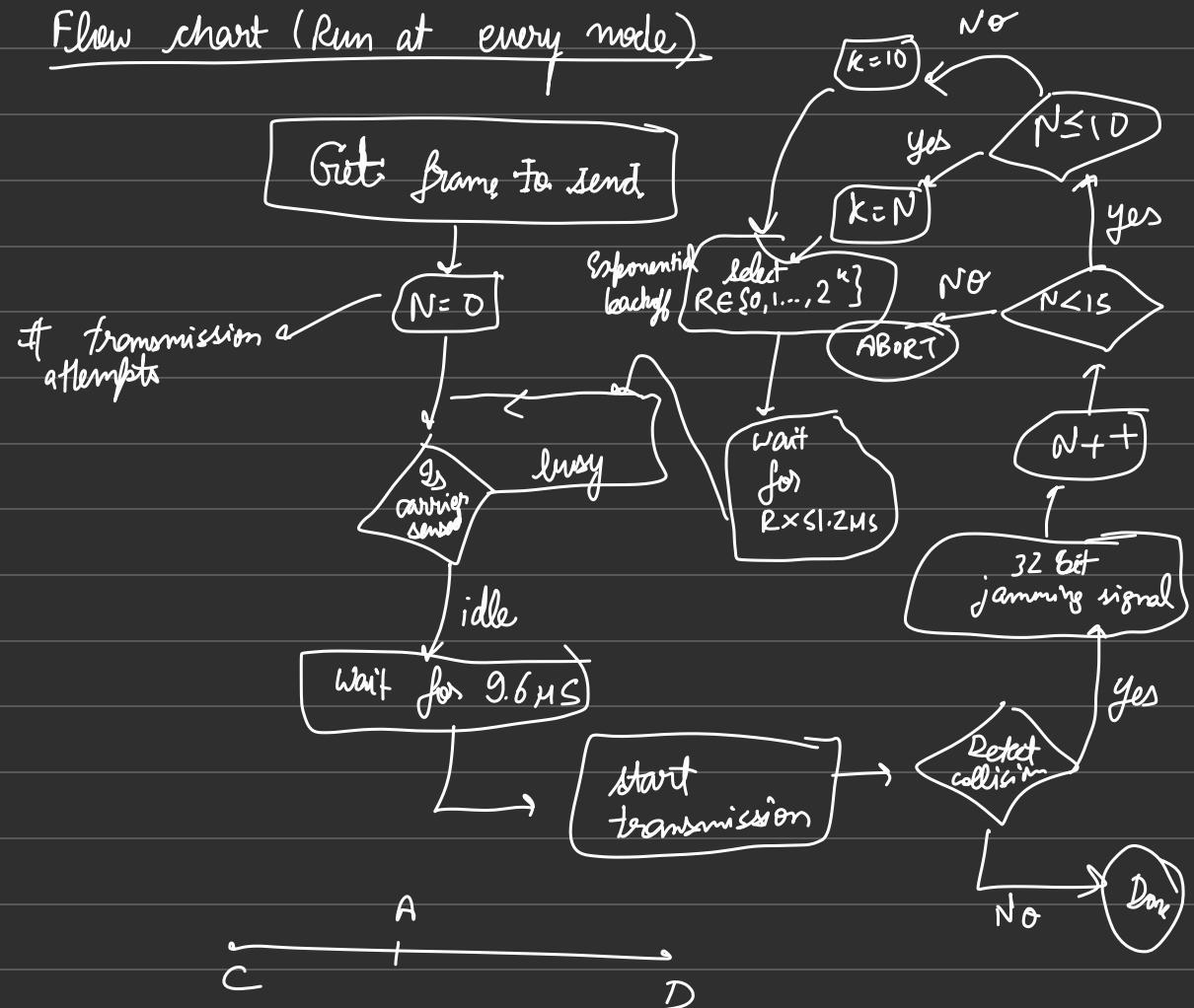
Including these worst case

$$RTT = 51.2 \mu s$$



$\Delta$  chosen by ethernet = 51.2 μs.

Flow chart (Run at every node)

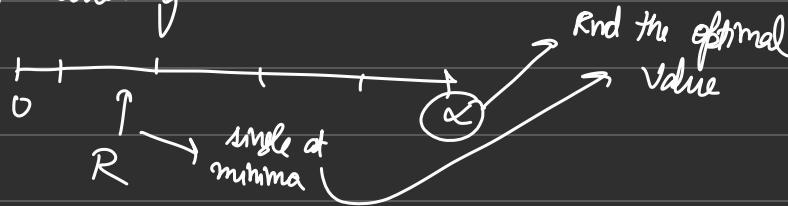


(Interference gap, for receiver to finish processing prev. frame)  
Ex: CRC check

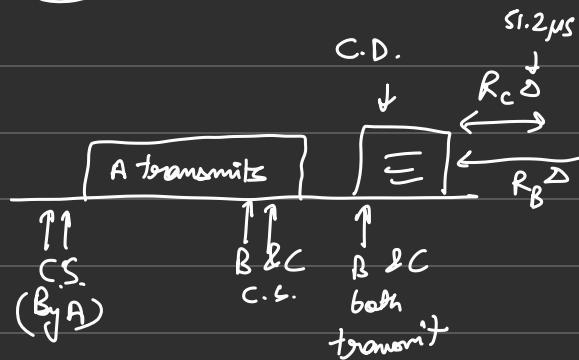
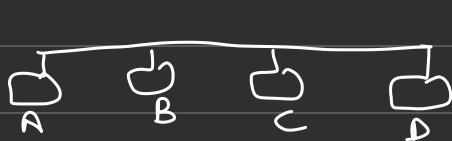
## Exponential Backoff

Don't know how many colliding

Suppose M colliding



## CSMA-CD Ethernet



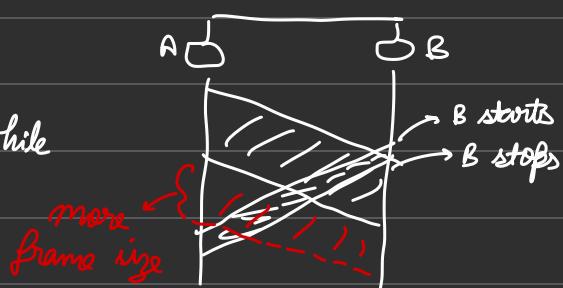
$$R \in [0, 1, 2, \dots, 2^k]$$

$k$  increments up to 10 for every collision (for some frame)

## Limits on frame size (Ethernet 802.3)

(1) Min. frame size

Want A to detect collision while transmitting



$\Delta$  = Worst case R.T.T

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

original standard 10Mbps

Frame size 64 bytes

$$= 512 \text{ bits}$$

Time to transmit @ 10Mbps

$$= \frac{512}{10} = 51.2 \mu\text{s}$$

## (2) Max frame size

### (a) bit errors



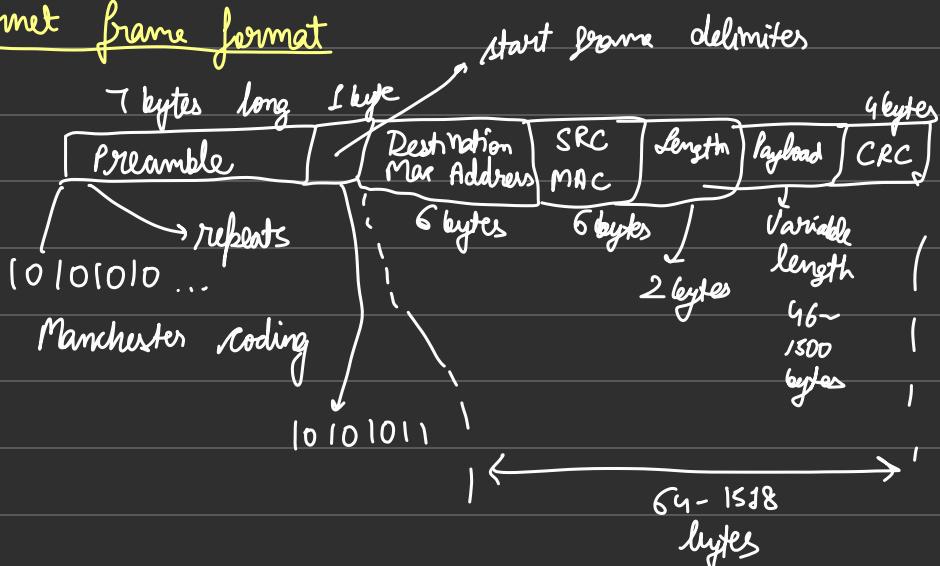
$\sim 10^3$  bits

High prob. of at least one bit error

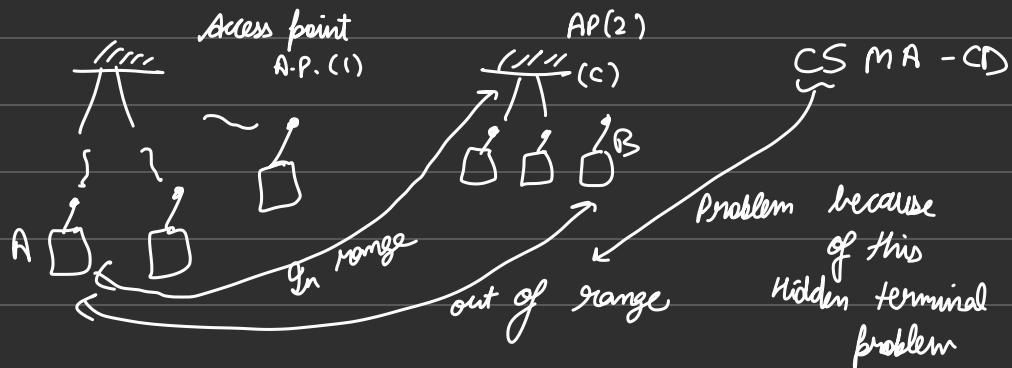
(b) Don't want single user to monopolise the channel for very long

Max frame size = 1518 bytes.

## Ethernet frame format



WiFi (IEEE 802.11) 802.11 b/g/n/ac



Everyone can hear the other in ethernet

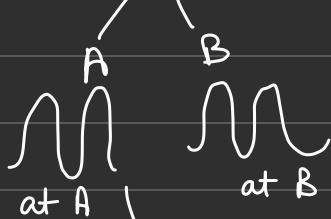


A cannot sense B's signal

(Hidden from each other)

A's signal at C      B's signal at C

Collision at C



Signal decay with distance  $d$

$$\sim \frac{1}{d^2}$$

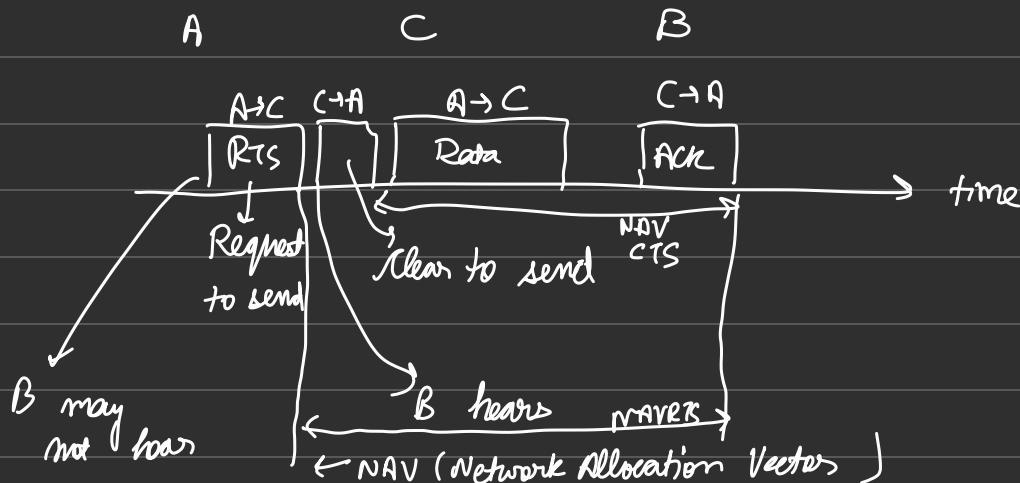
While transmitting A cannot hear anyone else

Collision detection does not work.

sol<sup>m</sup>: Receiver sends ACK



How to handle Hidden Terminal?



Symmetry is assumed (if B can hear C then C can hear B too)

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

Virtual carrier sensing

hearing CTS → B remains silent.

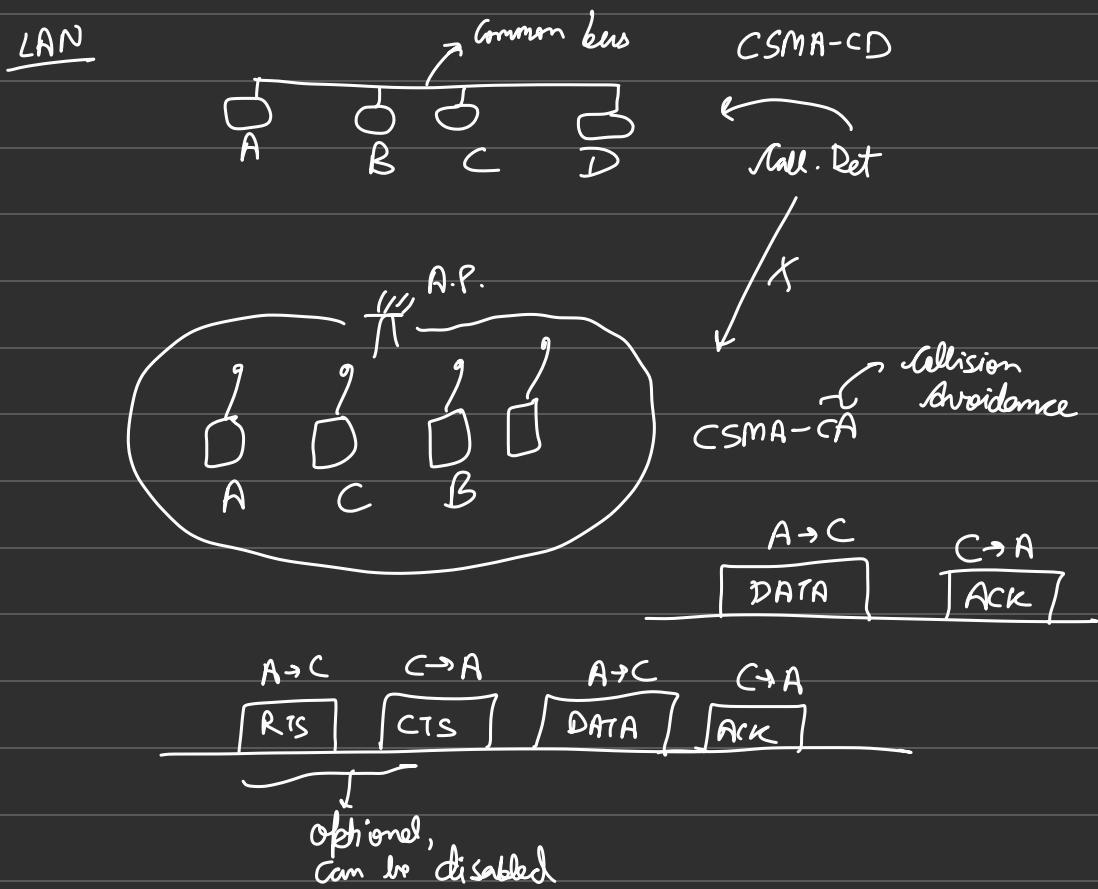
\* What if A sends RTS but does not get CTS ?

→ A assumes collision, tries to retransmit

\* What if RTS, CTS, DATA are sent but ACK didn't reach A ?

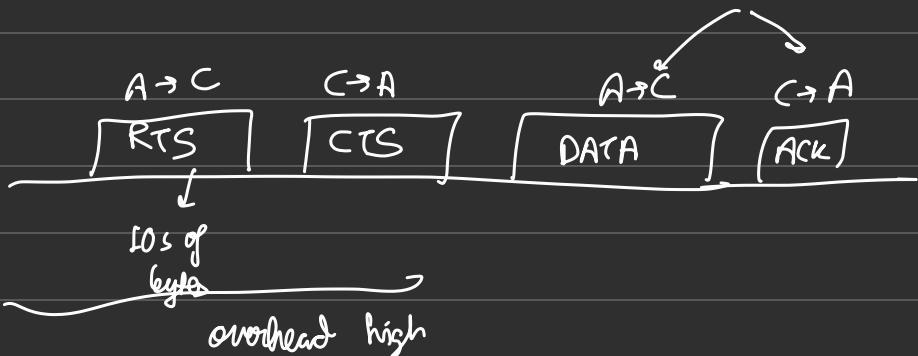
→ A assumes collision, tries to retransmit (starting with Data)

## WiFi MAC (IEEE 802.11)



QAM - 1024





### Exposed Terminal Problem :



A can hear only B.

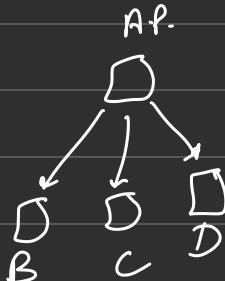
B hears A, C

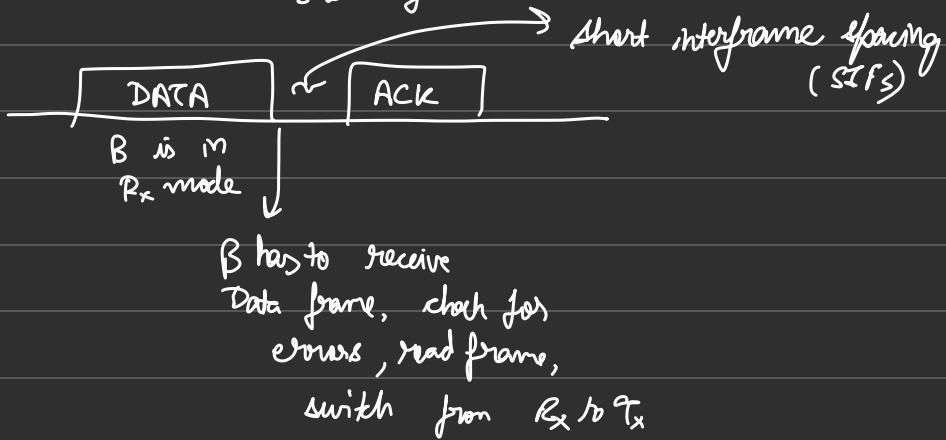
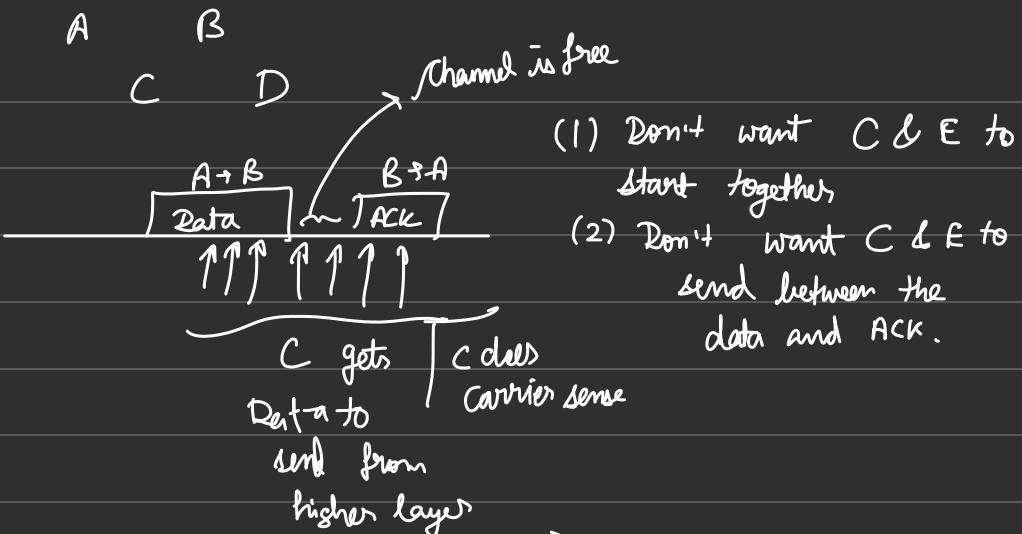
D hears only C

Case (i) RTS/CTS enabled  
If B starts first, C remains silent  
(due to RTS)

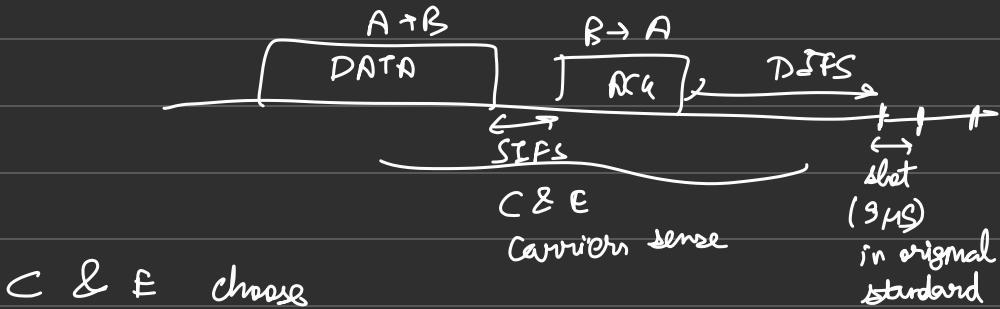
Case - (ii) RTS / CTS disabled  
B starts,  
C does carrier sense  
so C remains silent.

Usually :





Rule: C has to wait for atleast DIFS (SIFS + 2xt )



Contention window random variable

$$r_c \in \{0, 1, \dots, CW_c\}$$

$$r_e \in \{0, 1, \dots, CW_E\}$$

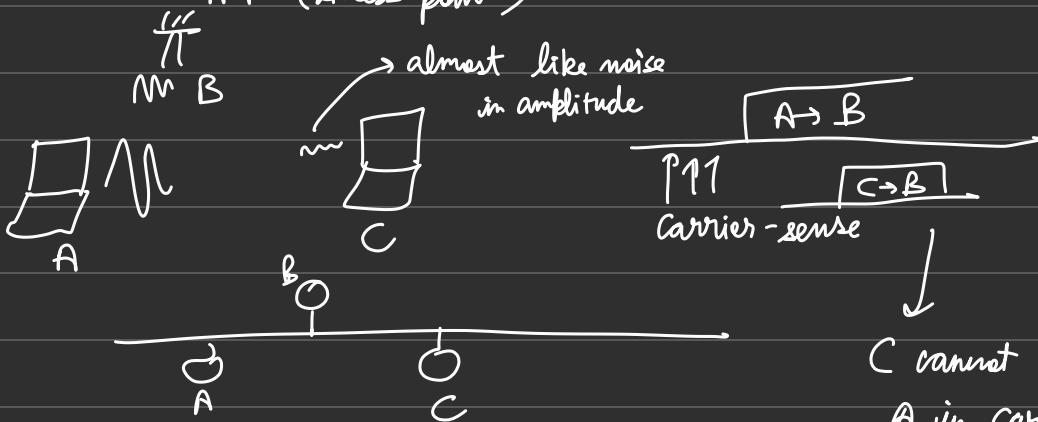
Initially  $CW_{max} = 15$  after limit  
= 1023

For collision  $CW_{max} = 2 \times CW_{max} + 1$

## Wireless Laws ( WiFi )

A.P. (Access point)

CSMA-CD will it work?



$$\text{decay in amplitude} \propto \frac{1}{d^2}$$

theoretically  $d^2$

practically decay is faster  
than  $d^2$

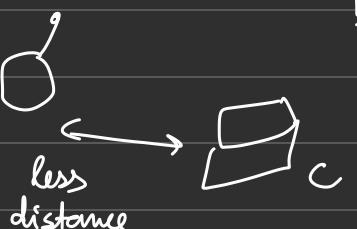
Collision detection

will also not work

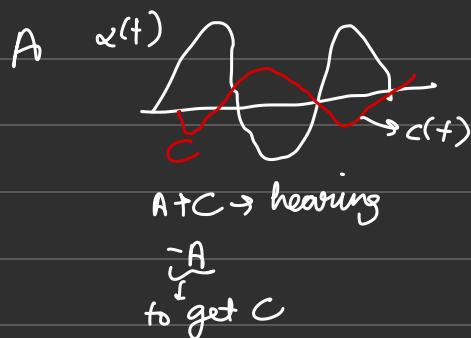
because difference in amplitudes is of orders of magnitude

A.P.

B



→ So C.S. work



$$A \text{ receives } \beta \cdot \alpha(t) + \theta \cdot c(t)$$

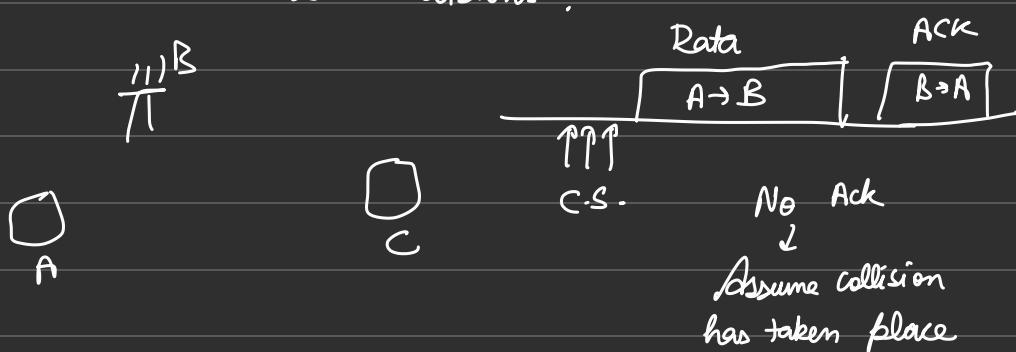
$\beta+1$        $\ll 1$

$$\beta \alpha(t) + \theta c(t) - \alpha(t)$$

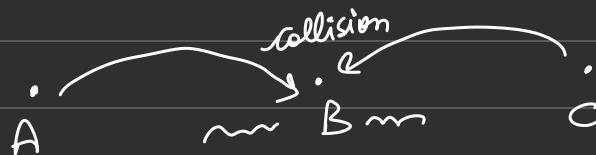
$$= (\underbrace{\beta-1}_{\text{still large}}) \alpha(t) + \theta c(t)$$

than  $\theta c(t)$

C.S. ✓ How to detect collisions?

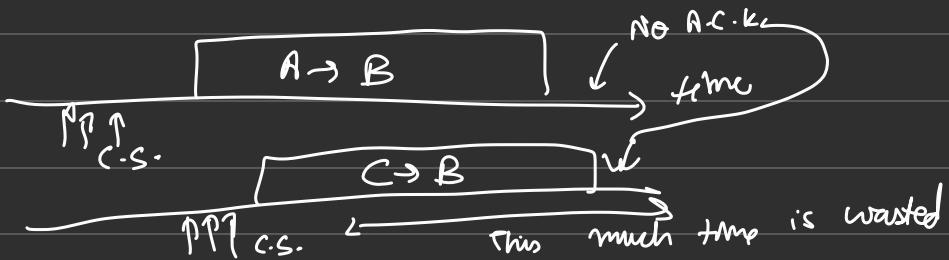


Hidden Terminal problem:



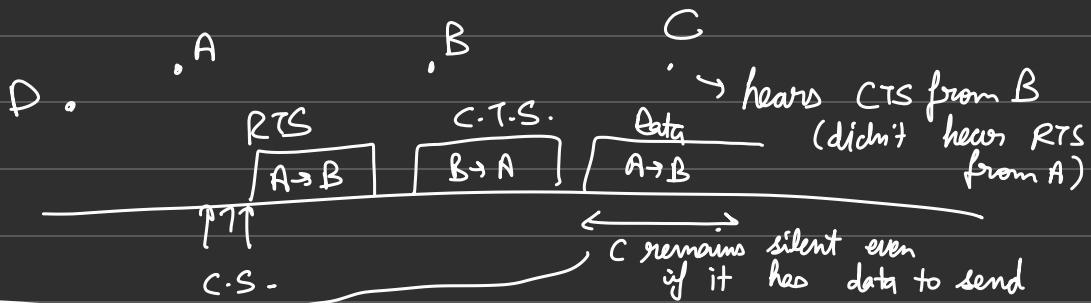
B can hear A and C

A and C cannot hear each other (hidden)



Solution to save the wasted time  $\rightarrow$  Virtual Carrier Sensing

RTS: Req. to send (short frame to tell the receiver it intends to send a data frame)



CTS  $\rightarrow$  Clear to send (short frame, saying that sender can transmit DATA)

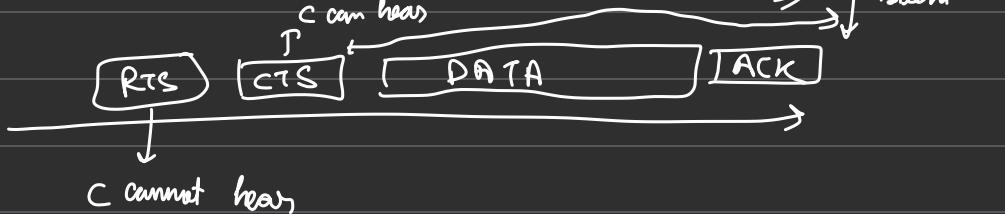
One solution: wait till  $B \rightarrow A$  acknowledge

Another solution: CTS has the length of waiting time

NAV : Network Alloc Vector

NAV.CTS (by B)

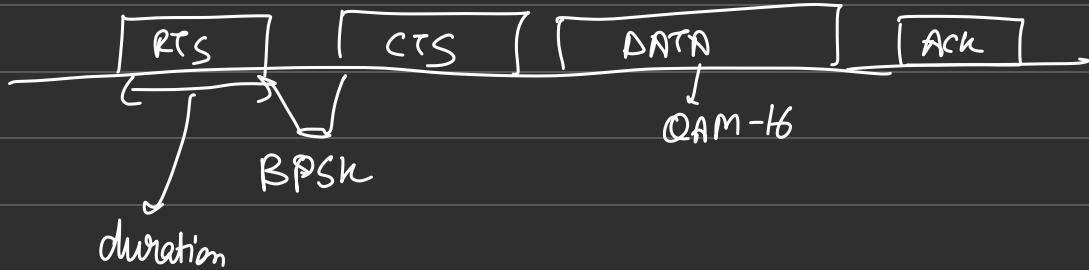
$\downarrow$   
based on NAV.RTS (by A)



Rule: Anyone hearing an RTS or CTS should remain silent for NAV mentioned in it.

option: (i) Don't use RTS / CTS : C.S. , DATA, Get ACK  
(ii) Use RTS / CTS : C.S. , RTS, CTS, DATA, ACK

↓  
try to retransmit if no ACK  
Usually sent using lower modulations (e.g. BPSK)  
and hence their duration is large.



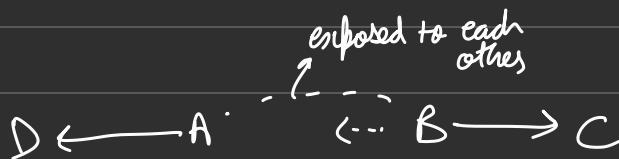
No, CTS, assume collision, try to retransmit later

## Wi-Fi

### CSMA - CA

Collision avoidance

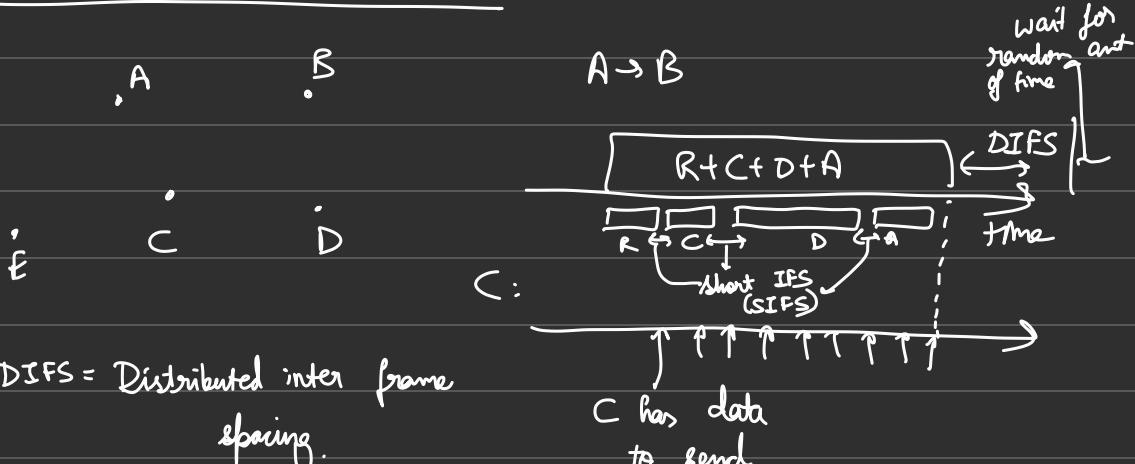
Exclusive Terminal Problem:



D can receive from A and C can receive from B simultaneously (in theory)

However, if  $A \rightarrow D$  then B remains silent due to C.S.

### Contention Window (CW)

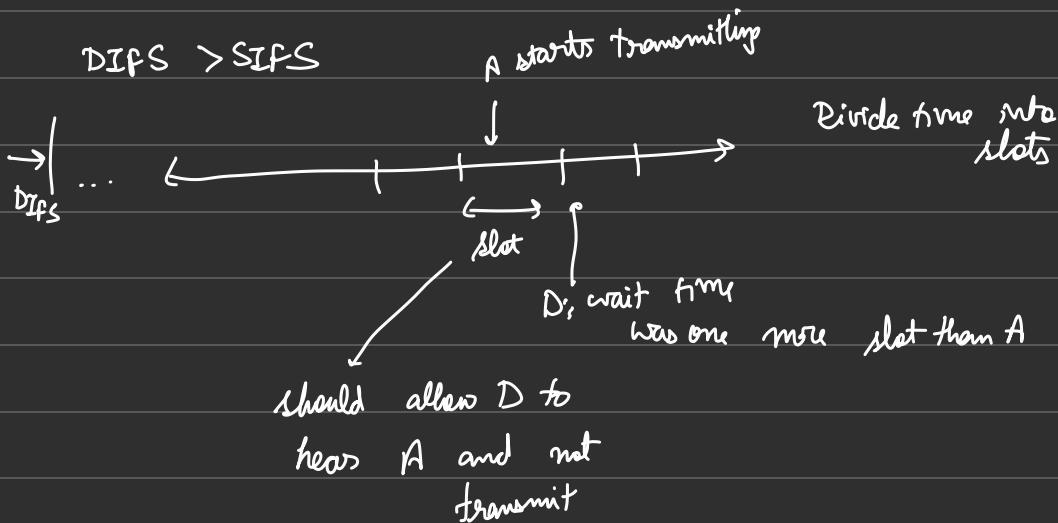


DIFS = Distributed inter frame spacing.

SIFS? : B hears RTS,  $R_x \rightarrow T_x$  circuit shift

A sent RTS,  $T_x \rightarrow R_x$  starts listening

↓  
decodes RTS  
(PHY/MAC processing)



prep. delay + possible offset (starting time for measuring data)  
+ time to CS.

C must hear A, carries sense if  $\rightarrow$  if A was transmitting  
prep-delay      c.s. time

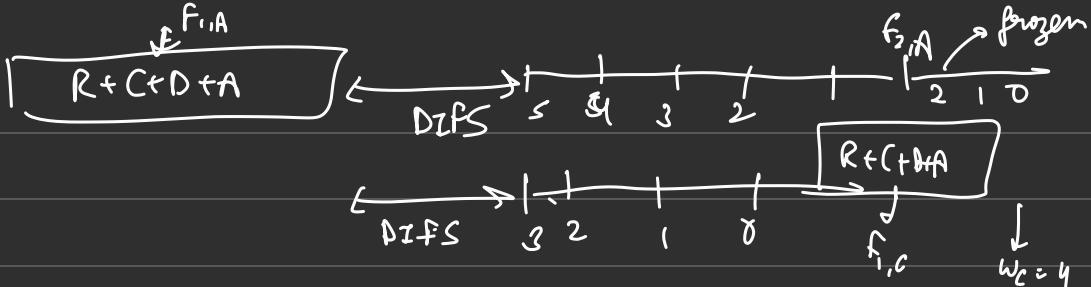
$+ RX \rightarrow TX \rightarrow$  No one else was transmitting

How long to wait?

Remaining waiting time  $w$

Initialize  $w \in \text{Unif}(0, w_{\max})$

For every idle slot, decrement  $w$  by 1.



Collision: sent RTS, no CTS

sent DATA, no ACK

$$CW_{max} = CW_{min} \times 2 \rightarrow \text{on collision}$$

$$CW_{min} = \max(CW_{max}, \text{max allowed value})$$

repeat procedure :

- $w \in \text{Unif}(0, CW_{max})$
- wait for DIFS free period
- decrement  $w$  every free slot
- if channel get busy
- freeze  $w$

Transmit when  $w=0$

IEEE 802.11

b, g, n, ac, ac  
latest

- MAC address of sender and receiver

Framing  
↓ check sum (CRC)

Frame :





rate  $\frac{3}{4}$  : send  $N$  bits,  $\frac{3}{4} N$  information bits  
 rate  $\frac{5}{6}$  : "  $\frac{5}{6} N$  information bits

802.11 g → 64 QAM, rate  $(\frac{3}{4})$  (54 Mbps) Ch width 20 MHz  
 ↓ bit rate

$\left\{ \begin{array}{l} 11n \rightarrow 64 \text{ QAM, rate } (\frac{5}{6}) \quad (150 \text{ Mbps}) \quad 40 \text{ MHz} \\ 11ax \rightarrow 256 \text{ QAM, rate } (\frac{3}{4}, \frac{5}{6}) \quad (866 \text{ Mbps}) \quad 160 \text{ MHz} \end{array} \right.$

$11ax \rightarrow 1024 - \text{QAM, rate } \dots \frac{5}{6} (1.2 \text{ Gbps}) \quad 160 \text{ MHz}$

MIMO  
 multiple input  
 multiple output

↓  
 SISO  
 single input  
 single output

DIFS > SIFS : To give priority to someone already transmitting

$$\text{DIFS} = \text{SIFS} + 2 \times \text{slot time}$$



$$PIFS = SJIFS + \text{slot\_time}$$

↓  
point coordination function. (PCF)

D → distributed CF



Higher priority possible?

$$w \in \text{Clif}(0, C_{\max})$$

802.11e	<u>C<sub>Max</sub></u>	<u>Initial</u>	<u>Max Value</u>
Voice	3	7	7
Video	7	15	15
Others	15	1023	1023

Need

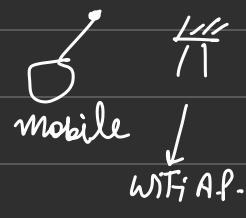
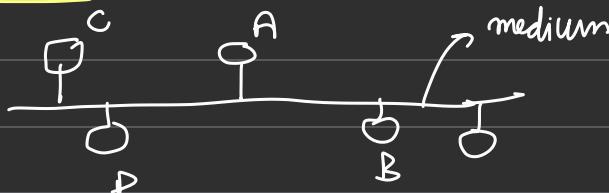
Cross-layer interaction

$$SJIFS = 10 \mu s \quad 11g \quad (\ln$$

higher for others

Max. . gets tries if collisions and then give up.

Medium Access (DCL)



## Ethernet LANs

CSMA : Carrier sense multiple access.

TDMA : Time division multiple access

Random access protocols : schedule for transmissions is not decided in advance.

Idea: If a node has data to send, create a frame and sends it.

Broadcast : all nodes on network are destination of frame  
unicast  $\rightarrow$  single node is destination of frame  
multicast  $\rightarrow$  a subset of nodes is destination of a frame (packet)

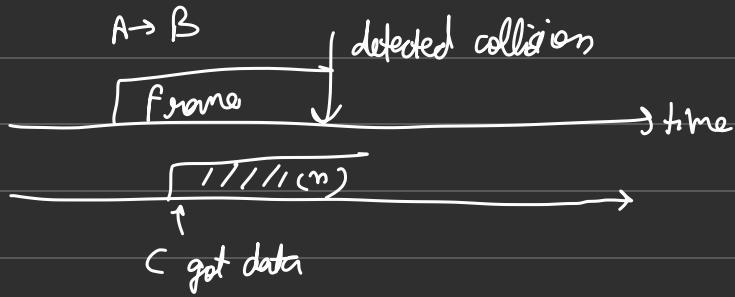
what if there is a collision?

A  $\rightarrow$  B at same time  
C  $\rightarrow$  D

(1) Collision Detection  $\rightarrow$  so that A and C know that a collision occurred.



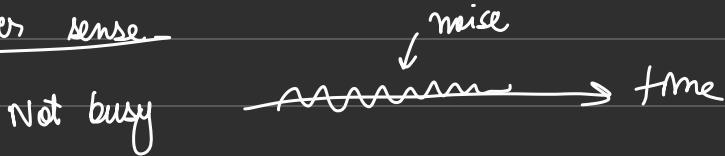
(2) Back-off for random time



(3) Carrier-sense : If a node has data to send, first sense the medium to see if a transmission is ongoing. Don't transmit if medium is busy.

CSMA-CD  
↓  
collision detection

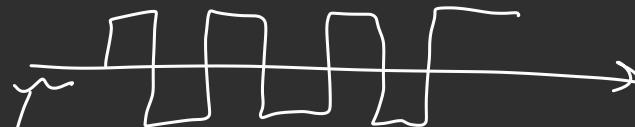
Carrier sense



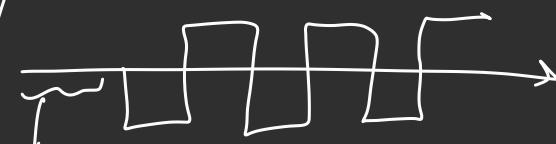
If energy of signal > threshold → conclude that channel is busy  
↓  
larger than noisy energy

## Collision Detection

At A :

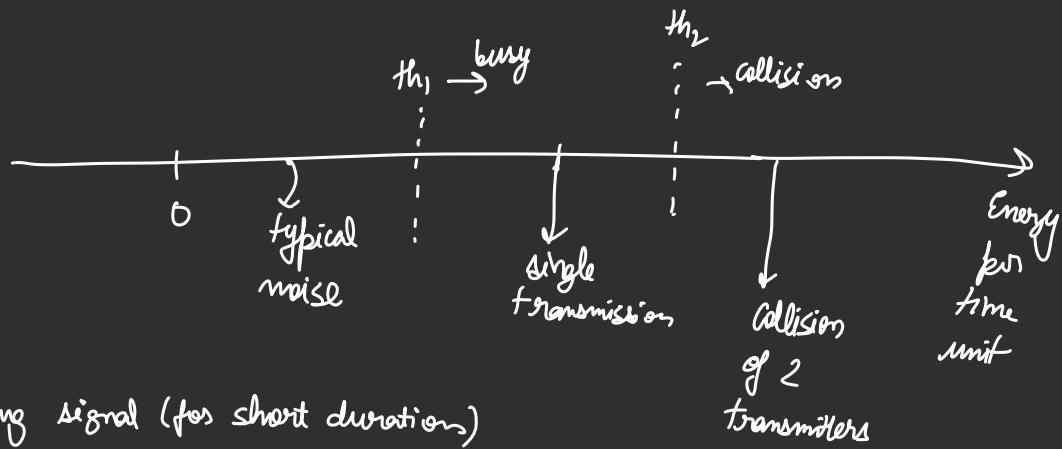


At C :



Carrier sense

Idea : Another threshold for collision detection ( $th_2$ )  
 $\text{energy} > \text{threshold}_2$  then collision



Jammering signal (for short duration)

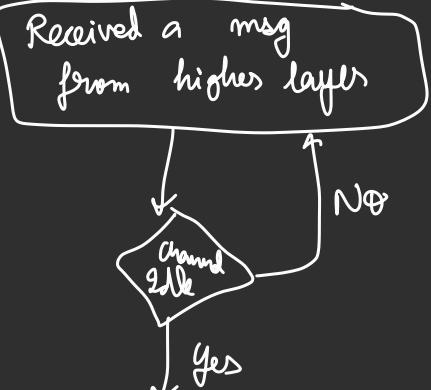
A sends



for C to detect

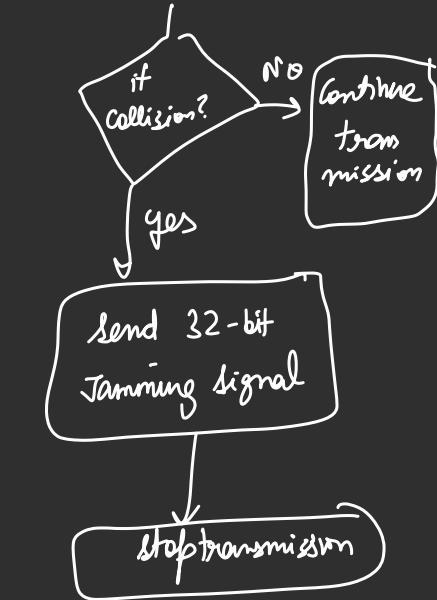
C also sends a jamming signal  
after collision detection

## Flow Chart



Transmit after a short gap

set to be 3.6 μS (if the receiver was the same  
for below-frame,  
he should be given some time for processing)



Frame Details in 802.3 → to signify that preamble has ended  
Start frame delimiter

Preamble	$S$ $F$ $D$	Dest. MAC ADDR	SRC MAC ADDR	$L$ $E$ $N$	Payload	CRC
7 bytes	1 byte	6 bytes	6 bytes	2 bytes		4 bytes

$2^{48} \rightarrow$  Mac addresses

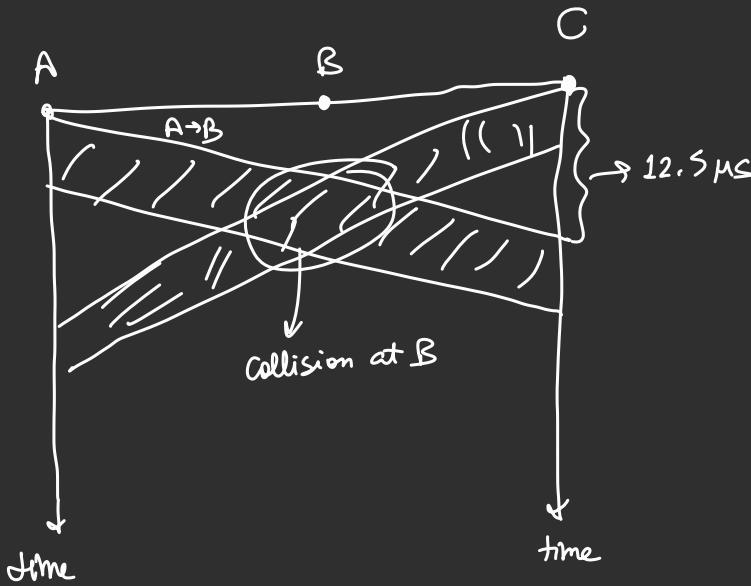
$2^{32} \rightarrow$  IP v4

64 - 1518 bytes

why min frame size?

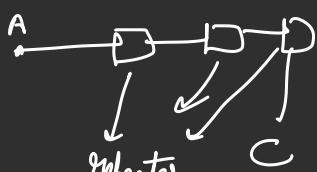
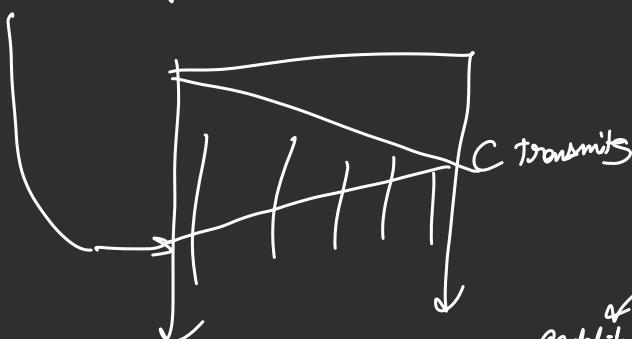
2500m (approx)

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



A must detect collisions before the frame ends, if a collision occurs at the receiver.

want A's frame to continue till the end



(replicates signal on one link onto other link)  
amplify and  
reproduce

≈ 50 μs RTT in the worst case.

64 bytes, 10 Mbps  $\rightarrow$  51.2  $\mu$ s  
 transmission rate

Max frame size

very large

large frames

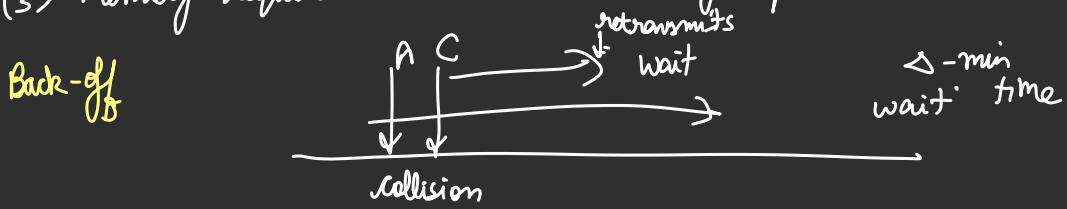
1 1 1 1 1 1 1 1

- (1) Higher prob. of getting a bit in error in frame  
 $'p'$   $\rightarrow$  suppose prob. of bit errors

$N \rightarrow$  frame length  
 Prob. of no errors  $(1-p)^N \rightarrow$  assuming bit errors are independent

entire frame will have to be discarded.

- (2) others need to wait longer to get a chance to transmit.
- (3) Memory requirements at NIC card go up.



Wait  $\in \Delta \cdot U$   
 $\downarrow$   
 $\{0, 1, \dots, 2^k\}$   $\rightarrow$  longer if more are colliding

Exponential backoff

1<sup>st</sup> collision :  $k=1$  ;

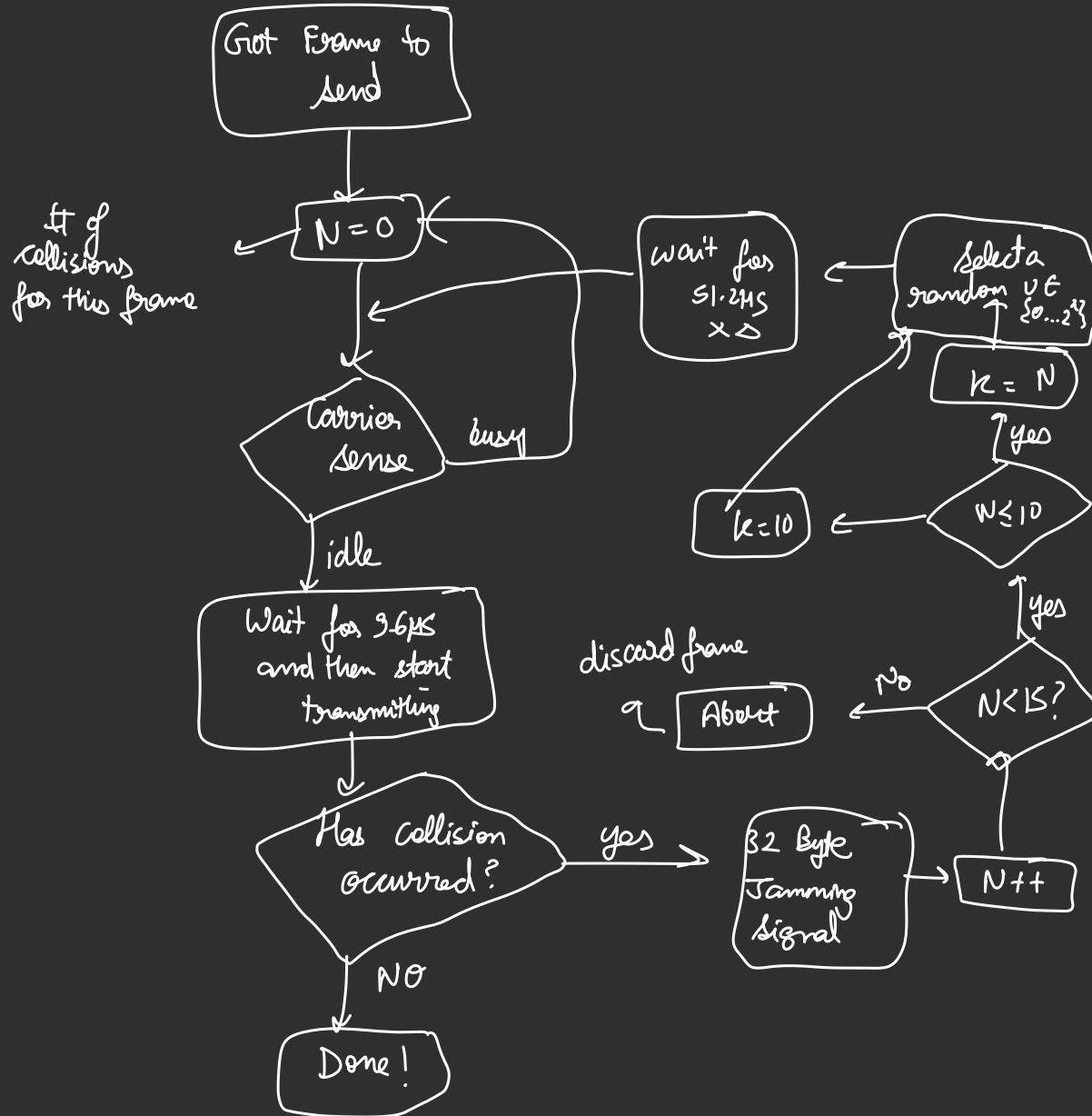
2<sup>nd</sup> collision :  $k=2$  ;

:

10<sup>th</sup> collision  $k=10$  ;

11<sup>th</sup>  $\rightarrow k=10 \dots 15^{th}$

# Flow Chart



MAC

CSMA-CD / CA → Collision  
↓      WiFi      Avoidance

Diagram illustrating the frequency bands for 80 MHz and 2.4 GHz.

## Collision Detection

## Decentralised



WiFi A.P

A. P.

9

六

No central  
Coordinator for MAC

Uplink : User  $\rightarrow$  B.S.

Downlink: BS  $\rightarrow$  Users

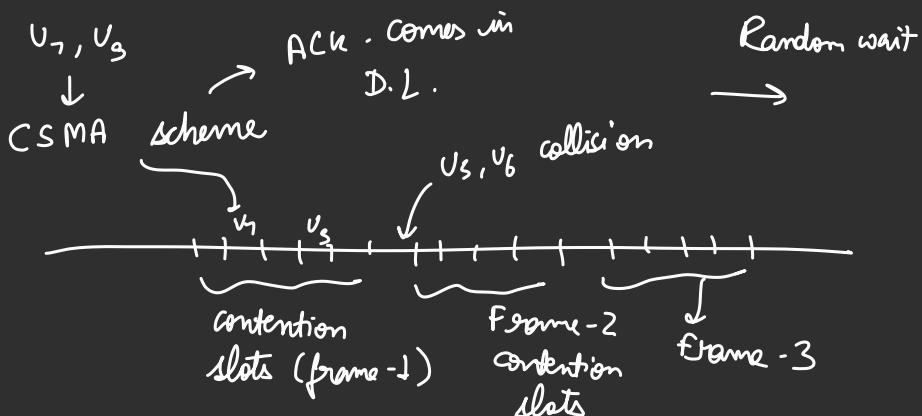
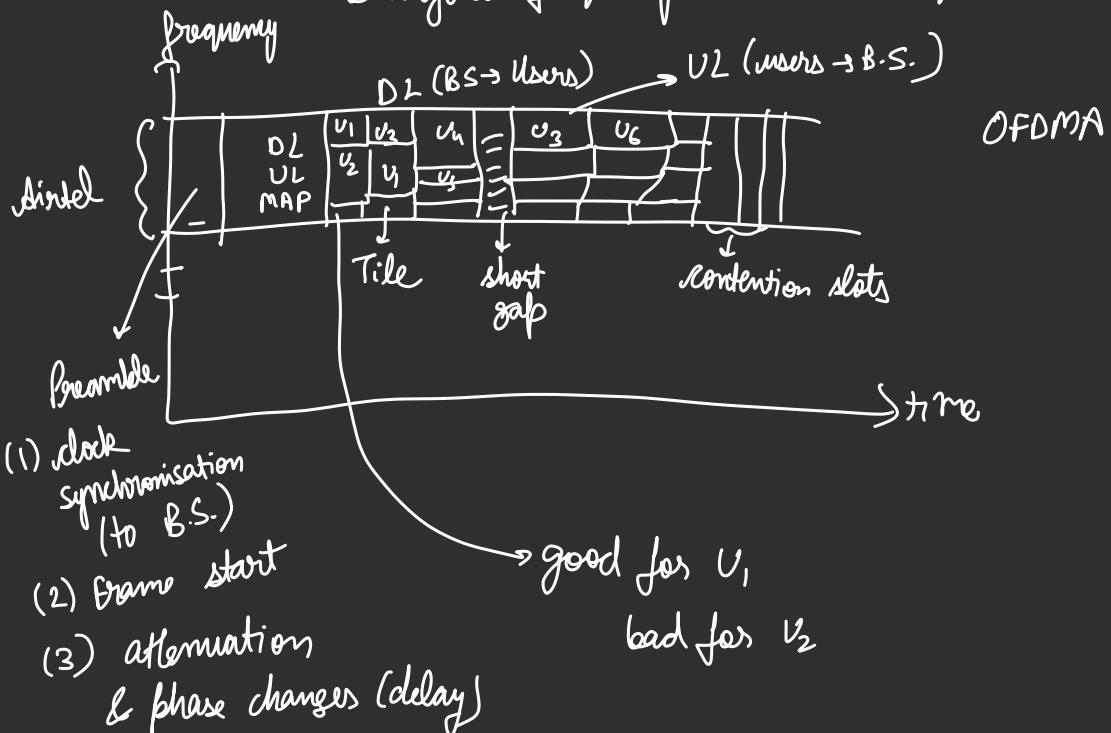
→ schedule  
UL-MAP  
DL-MAP

Time-division multiple access

## FDMA

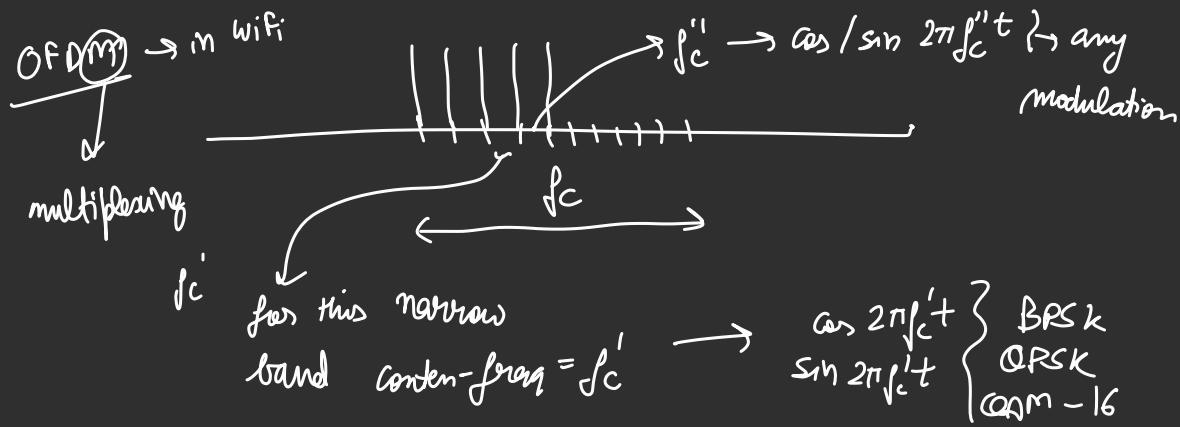
Frequency division  
multiple access

4G-LTE uses OFDMA  
 ↓  
 orthogonal frequency division multiple access

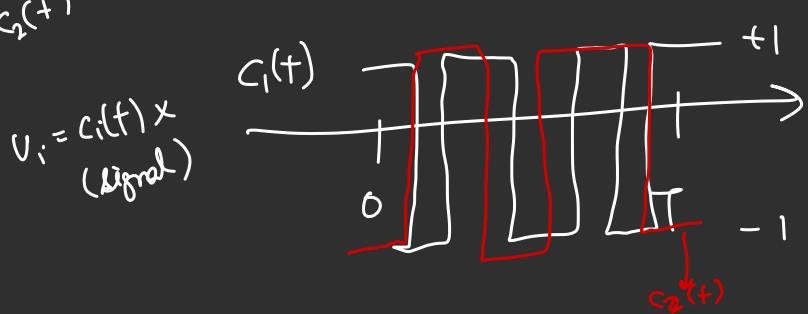
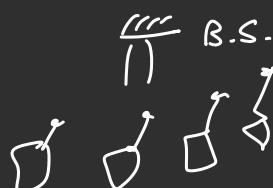
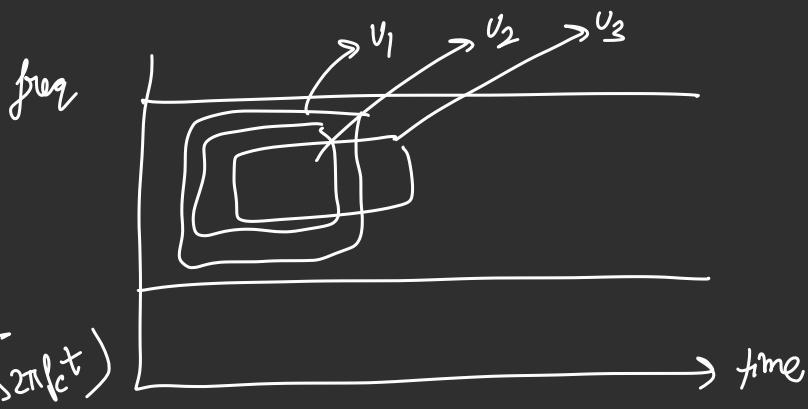


OFDMA

$$\xrightarrow{f_c} \mu_3$$



CDMA : Code Division Multiple Access



$$\begin{aligned} c_1^2(+1) &= 1 \\ c_1(+1) \perp c_2(+1) \\ \int_0^T c_1(+1) c_2(+1) dt &= 0 \end{aligned}$$

All codes are mutually orthogonal

B.S

$$\alpha(t) = \sum_{i=1}^n s_i(t)$$

# users

$$\begin{aligned}s_2(t) c_1(t) \cos 2\pi f_c t \\= -c_2(t) c_1(t) A \cos^2 2\pi f_c t \\= -\frac{A}{2} c_1(t) c_2(t) \underbrace{\left[ 1 + \cos 4\pi f_c t \right]}_{\substack{\text{low pass} \\ \text{filter remove}}}\end{aligned}$$

After low pass filters

$$\int_0^T \frac{A}{2} + \left( -\frac{A}{2} c_1(t) c_2(t) \right) dt = \frac{AT}{2}$$

orthogonal

CFDM, OFDM robust to multi-path

↓  
3G

↓  
Wifi, 4G

$$\begin{aligned}\alpha(t) \times c_1(t) \\ \cos 2\pi f_c t \\ s_1(t) \times c_1(t) \\ \cos 2\pi f_c t\end{aligned}$$

$$\begin{aligned}&= (c_1(t) A \cos 2\pi f_c t) \\&\quad (c_1(t) \cos 2\pi f_c t) \\&= A \cos^2 2\pi f_c t \\&= \frac{A}{2} \underbrace{\left[ 1 + \cos 4\pi f_c t \right]}_{\substack{\text{removed by} \\ \text{integration}}}\end{aligned}$$

(also by low pass filter)  
since  $2f_c$  is a high frequency

Switching → 12-switches (Ethernet)

↓  
MAC Address used to switch

13-switches (Routers)

↳ use IP address to switch

Intelligent  
isolation  
(for scalability)  
Switch  
Isolating different parts of LAN

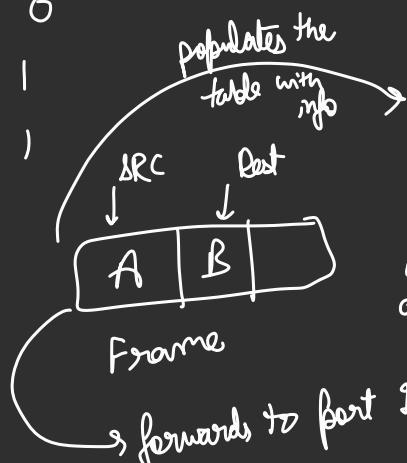
### FORWARDING Table

Dest	Port #
A	0
B	0
C	0
P	1
E	1
F	1

ports also have MAC addresses

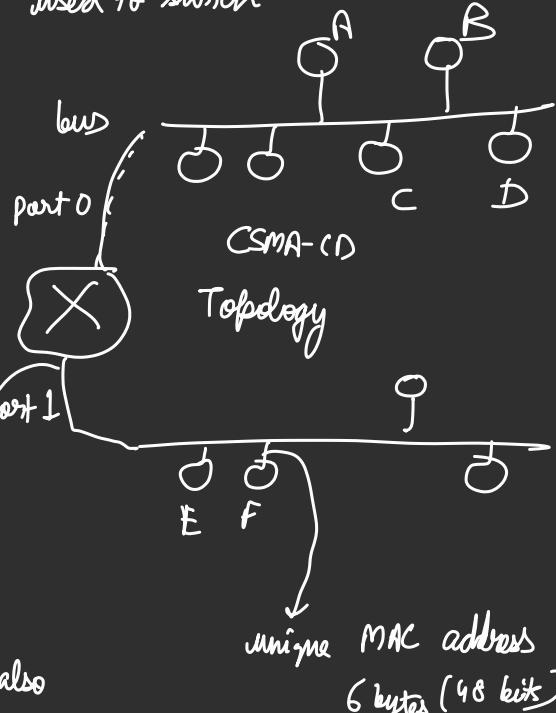
populates the table with info

Ethernet switches = bridges



[B | A]

→ received on port 0  
Don't forward.



### Init Table

Dest	Port #
A	0
B	0
F	1

Expiry time  
delete after expiry time

B | C

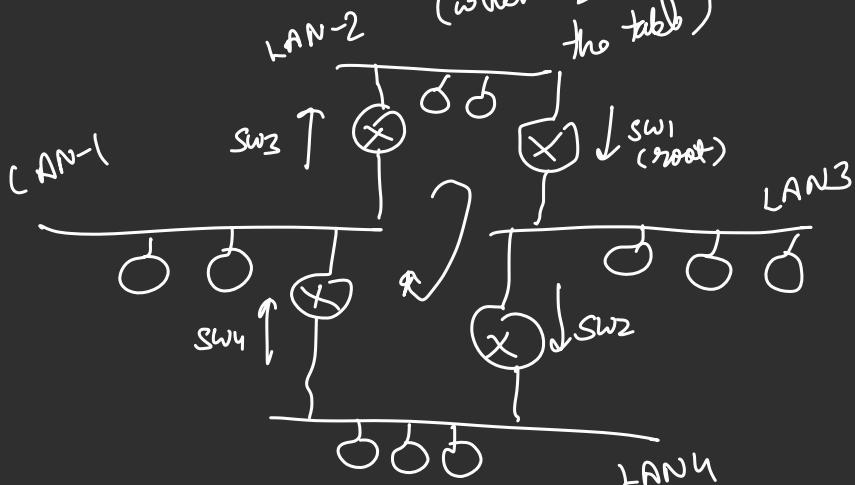
→ forwards to port 1



If dest is not in the table, then forward on all ports

A | B

↓  
Send it to all ports  
(when B is not present in the table)



Frame is going around forever (using resources)  
We need to avoid loops.

## Spanning Tree Protocol (Ratcliffe Perlman)

### SPT Protocol

- (1) elect root bridge
- (2) Each bridge finds which port is closest to root, assigns this port as root port. (tie breaking rule) → Designated port.
- (3) Any port which is not a root port or a designated port is disabled

## Details:

(1) elect root

Bridge ID

lowest becomes root

Each bridge tells its  
neighbours

$(y, d, x)$

the smallest  
ID heard  
till now

Configurable  
part (2 bytes)

6 bytes

MAC

Address

(smallest MAC  
of all ports)

default  
value : 32768

0 - 61440

(multiples of 4096)

after hearing  
 $sw_2, sw_3$

$sw_1 : (1, 0, 1)$

$sw_4 : (4, 0, 4)$

$sw_2 : (1, 1, 2)$

$sw_n : (2, 1, 4)$

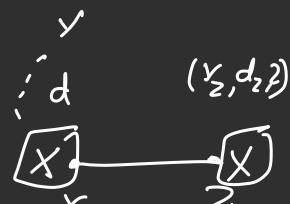
: Converge

$sw_1 : (1, 0, 1)$

$sw_2 : (1, 1, 2)$

$sw_3 : (1, 1, 3)$

$sw_4 : (1, 2, 4)$



$(y, d, x)$

if  $y < y_2$

then  $y_2 = y$

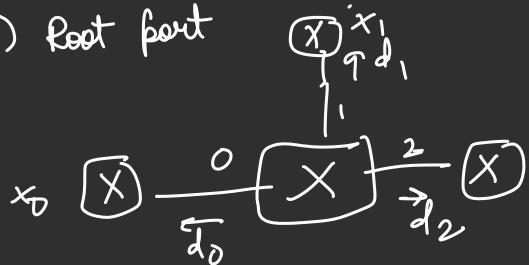
$d_2 = d + \text{dist}(x, z)$

if  $y = y_2$  but

$d + \text{dist}(x, z) < d_2$

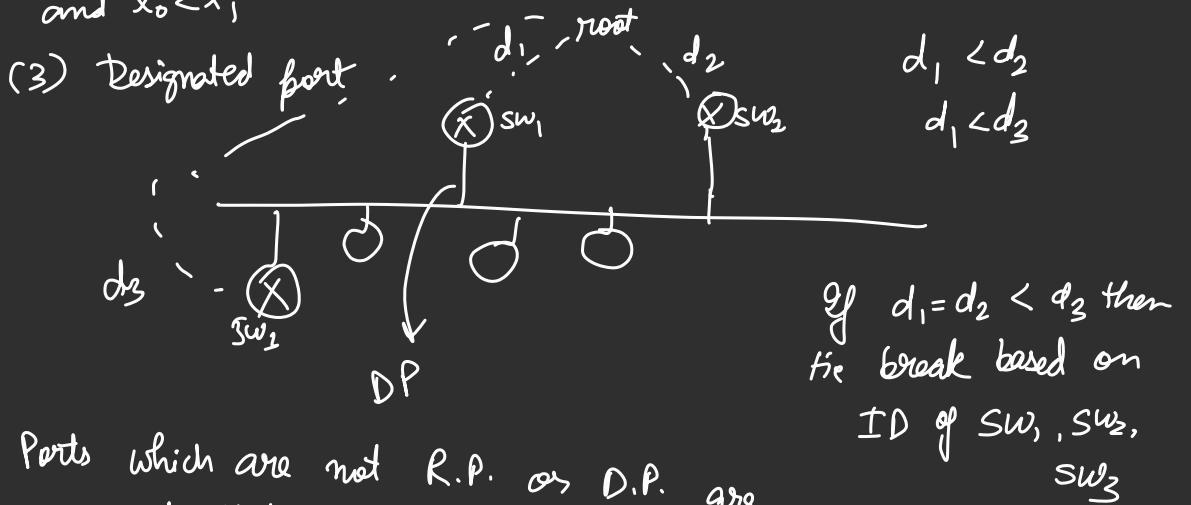
then  $d_2 = d + \text{dist}(x, z)$

(2) Root port

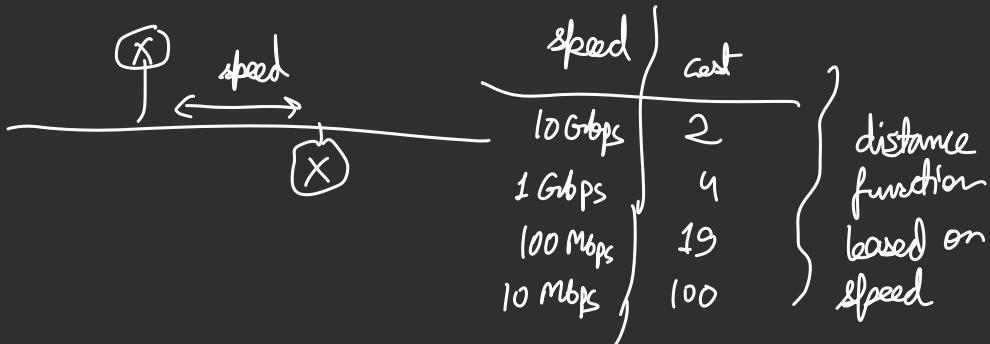


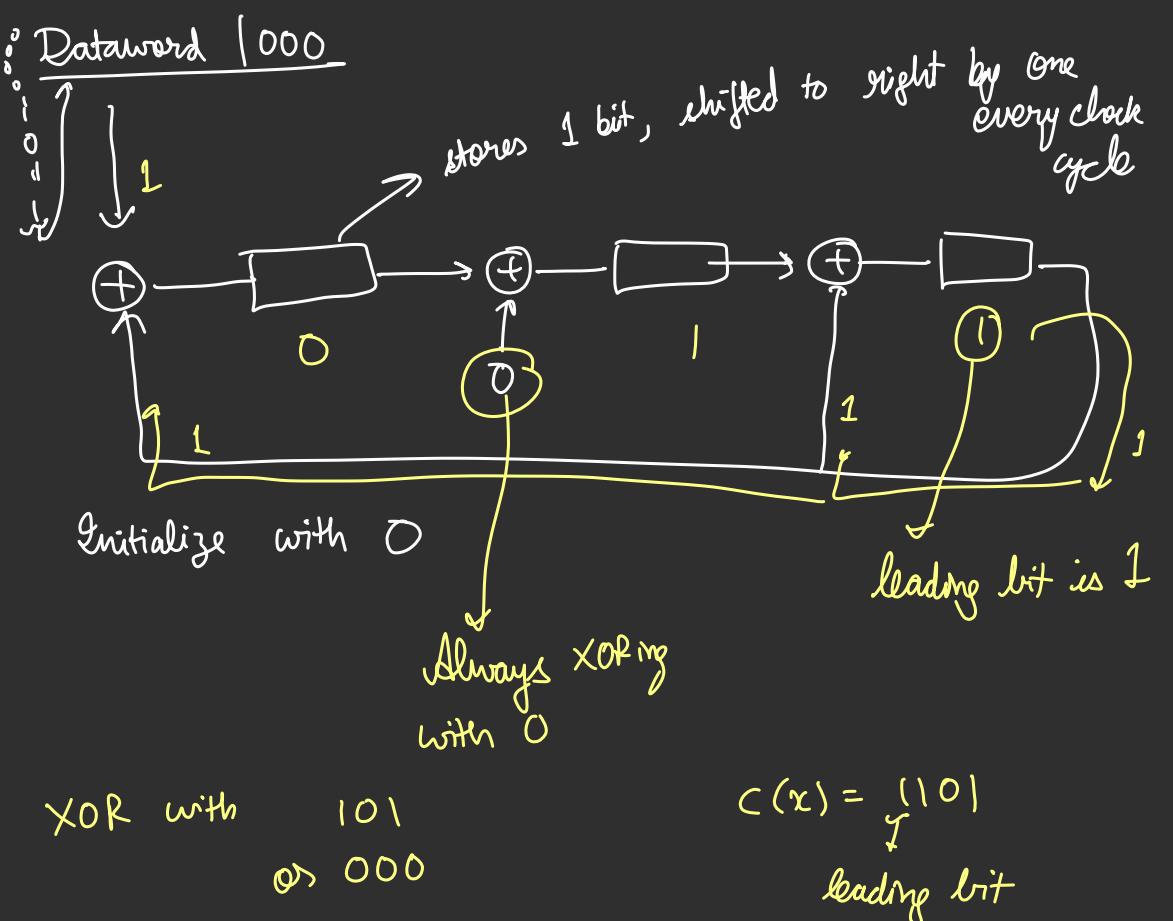
if more than one port  
has the <sup>some</sup> smallest distance  
to root, then tie-break  
based on ID (smallest)

If  $d_0 = d_1 < d_2$   
and  $x_0 < x_1 \rightarrow 0 = \text{root port}$



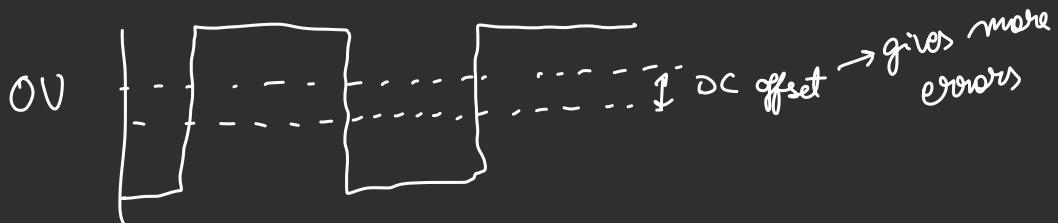
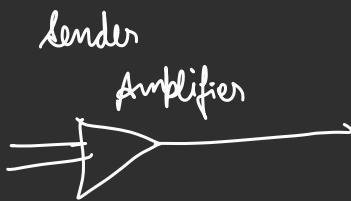
Ports which are not R.P. or D.P. are disabled.





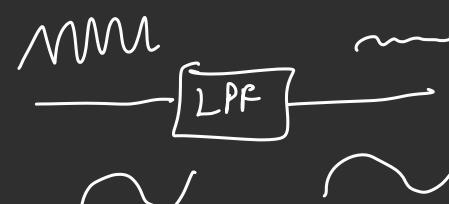
Last bits in the circuit will be CRC

Baseline Wander (Problem with NRZ)



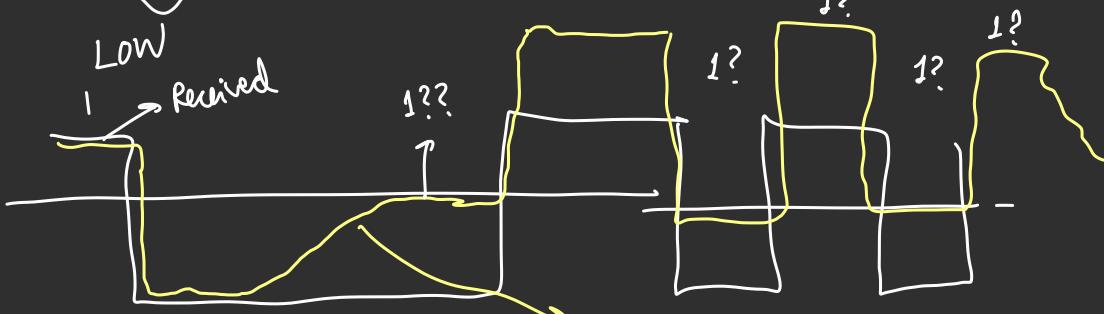
at receiver → 1?

High-freq



Low  
Received

→ Attenuated  
1?



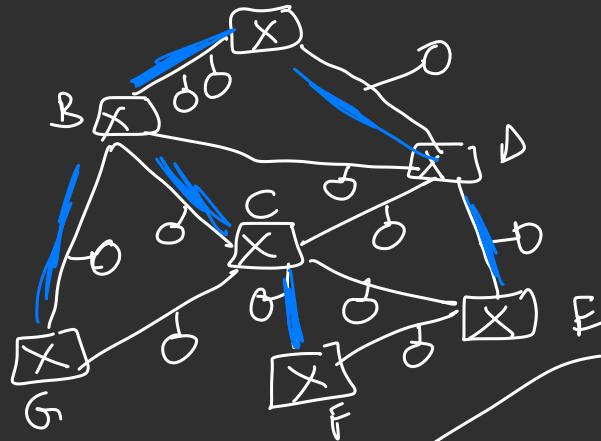
Because of high-pass filter

## Layer-3 switching

### Shortcomings of Ethernet switching:

$C \rightarrow E$ :  
 $C - B - A - D - E$

- path is not optimal (shortest path)
- not using some links - resource utilization poor
- $O(N)$  table size
- Stability



Flat addressing  
(Layer-2)

Layer-3 (IP): Hierarchical addressing

If root fails: SPT is reconstructed

→ root sends a frequent msg that other nodes should hear

→ Hello messages from root don't come

or other switch in spanning tree fails

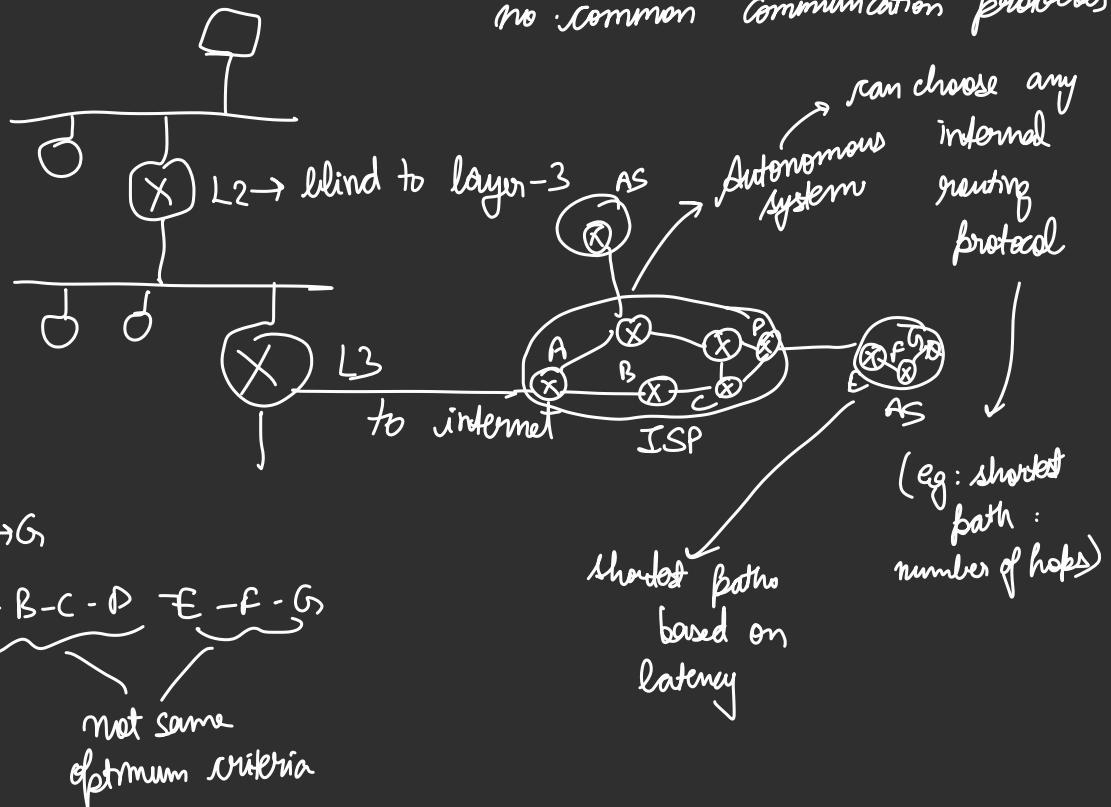
## II Spanning tree

### Fwd table

Dest	Port
host-1	0
host-2	1
⋮	⋮
all hosts	⋮
no hosts	⋮
	$O(N)$

Larger LANs  $\rightarrow$  more freq. SPT construction

Other issues : No common addressing scheme across globe .  
↓  
no common communication protocols



Intra-domain routing : Routing protocols within AS

Inter-domain routing : Routing protocols between AS

↳ Border Gateway Protocol (BGP)

Intra domain routing



RIP  
routing  
information  
protocol



link-state routing

↓  
OSPF : open shortest path first

IS-IS : Intermediate system to IS

Use shortest path  
avoid loops.

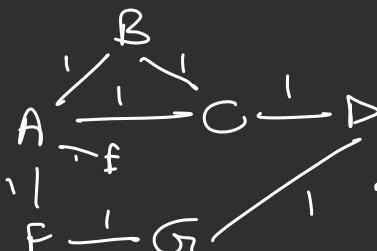
Distance Vector

A sends out:

(A, 0)

↓  
distance  
to myself is 0

A hears (B, 0) (F, 0) (C, 0)



Forwarding Table at A

Dest	Next hop	Cost
A	-	0
B	B	1
C	C	1
F	F	1
E	E	1

A sends out to own neighbours :

(A, 0), (B, 1) (C, 1) (F, 1)

hears from (B) (C, 1) (A, 1)

from (C) (D, 1) (B, 1) (A, 1)

from (F) (A, 1) (G, 1)

Rest	Neat	Cost	$G \xrightarrow{f} F \otimes G$
A	B	0	f fails
B	C	1	
{F}	F	1	f tells A that its distance to G = $\infty$
D	C	2	
E	E	1	
<del>FG</del>	F	2	x Invalid
C	-	$\infty$	A knows this to neighbours

replace  $\times$  by  $\oplus$ : -  $\oplus$   $\infty$  }  $\infty$

C tells A that G have distance 2 to G.

- Triggered update : Event triggers a routing update to neighbours  
Eg:  $F \xrightarrow{f} G$  fails which triggers an update  $(G, \infty)$  to A
- Periodic updates : Tell neighbours information in routing table about  $(\text{dest}, \text{dist})$  to various destinations

## Count to Infinity Problem

$x \xrightarrow{1} A \xrightarrow{1} B$

→  $x$ -knows      →  $A$  knows

Rest	Next	Cost	
A	A	1	
B	A	2	

Rest	Next	Cost	
X	X	1	
B	B	2	

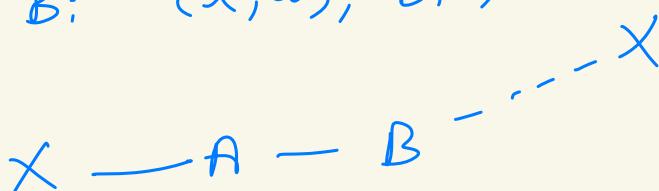
→  $B$  knows

Rest	Next	Cost	
A	A	1	
X	A	2	

If  $x$ -A link fails, A needs to send to B  
 B sends  $(x, \infty)$       These are sent at the same time  
 $(x, 2)$

B tells A:  $(A, 1), (x, 2)$

A tells B:  $(x, \infty), (B, 1)$



A will replace  $(A, -, \infty)$  by  $(x, B, 3)$   
 when B knows this if replace  $(x, A, 2)$  by  $(x, -, \infty)$

A tells B :  $(x, 3), (B, 1) \rightarrow$  B receives  
 and updates  
 B tells A :  $(x, y), (A, 1)$  Replaces  $(x, -, \infty)$   
 by  $(x, A, u)$   
 In A's table  $(x, B, 3)$  replaced by  $(x, B, 5)$

A tells B :  $(x, \leq), (B, 1)$

B tells A :  $(x, \leq), (A, 1)$

\* This is the routing loop and goes to infinity.

\* Routing Information Protocol (RIP) : max distance of 16  
 $\text{Cost} = 16 \Rightarrow$  Cannot reach destination.

### Split-Horizon

Do not advertise information about a destination to a neighbour if that neighbour is the next hop to the destination.

$x \xrightarrow{1} A \xrightarrow{1} B$

Dist	Next	Cost
B	B	1
X	-	$\infty$

Dist	Next	Cost
A	A	1
X	A	2

Not told to A under split horizon.

A tells B:  $(x, \infty)$

B tells A: nothing about dist to x.

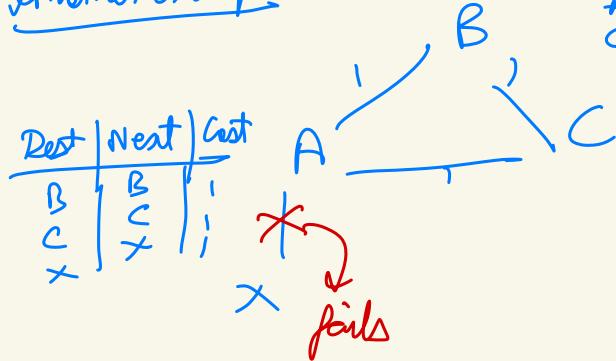
B updates  $(x, A, 2)$  to  $(x, -\infty)$

split horizon with poison reverse:

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

B sends advertisements  $(x, \infty)$  to A.

Another example:



Rest	Next	Cost
A	A	1
C	C	1
X	A	2

A sends  $(x, \infty)$  to B & C

Rest	Next	Cost
A	A	1
B	B	1
X	A	2

Suppose the message to C gets lost.

C tells B,  $(x, 2)$

C uses split horizon,  
so it does not give  
distance about X to A.

B put  $(x, C, 3)$  in its table

Again count-to-infinity problem!

B sends  $(x, \infty)$  to C

B tells A:  $(x, 3)$

A will update  $(x, B, 1)$

A will tell

C:  $(x, 4)$

RIP: Routing Implementation Protocol - D.V. based

Cost of all links are 1

distance vector

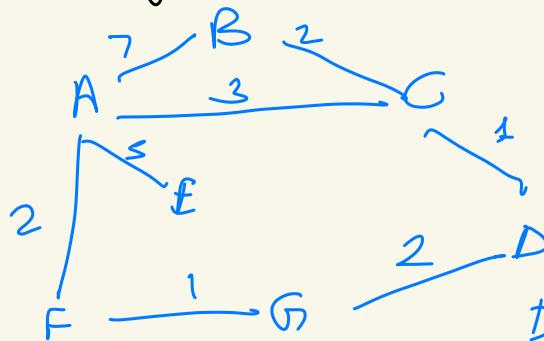
max. cost allowed to a destination is 16  
( $16 = \infty$ )

### Advantages of distance Vector:

Simple and easy to implement.

- Disadvantages
- cost-to-infinity & creating loops
  - convergence of routing tables takes some time.

### Link State Routing:

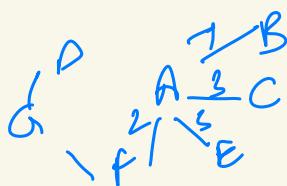


Each node sends to all others (broadcast)  
information about cost to immediate neighbours

D serves C to G to all.

### Dijkstra Algorithm:

Each node finds shortest path tree to all other nodes in network.



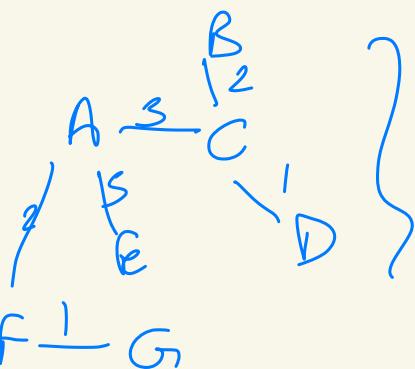
$$T = \{A\}$$

choose closest neighbours

$$T = \{A, F\}$$

$$T = \{A, F, G\}$$

$$T = \{A, F, G, C\}$$



On link failure

$$A \cancel{\xrightarrow[2]} F$$

A and F broadcast to all, that their link is failed.

All are going to run Dijkstra's algo.

**Advantages** - No routing loops, count-to-infinity problems, convergence to routing tables is fast.

**Disadvantages** - Algo is more complex than D.V.

$$T = \{A, F, G, C, D\}$$

$$T = \{\dots, E\}$$

$$T = \{\dots, B\}$$

our tree

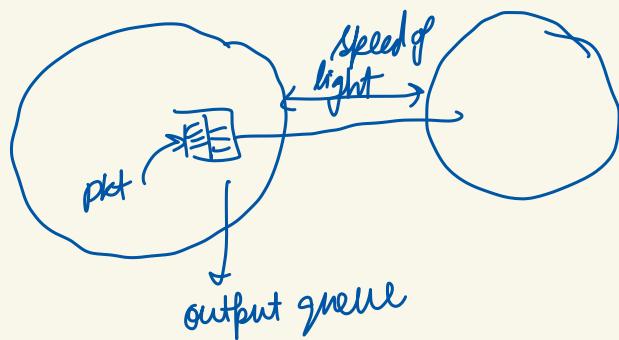
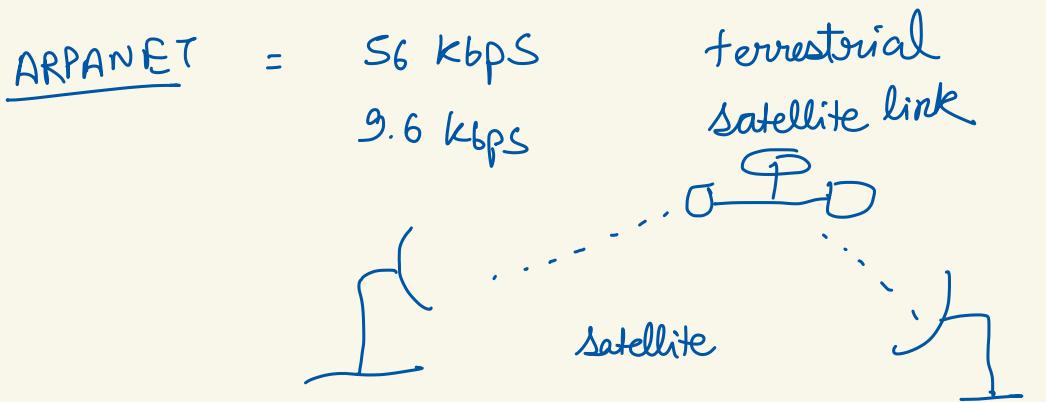
Finally

A's routing table:

Dest	Next	Cost
B	C	5
C	C	3
D	C	4
E	E	5
F	F	2
G	F	3

Q) what to use as cost?

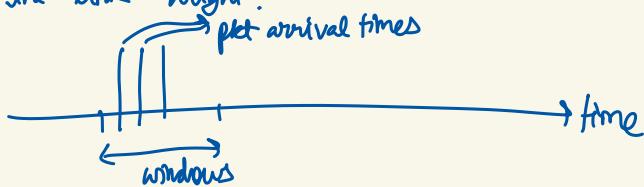
How to choose link weights?



Idea - I  
use latency

latency of pkt on this link = queuing delay  
+ speed of light delay  
+ transmission delay

Take time window, avg latency of all pkts in window  
is the link weight.



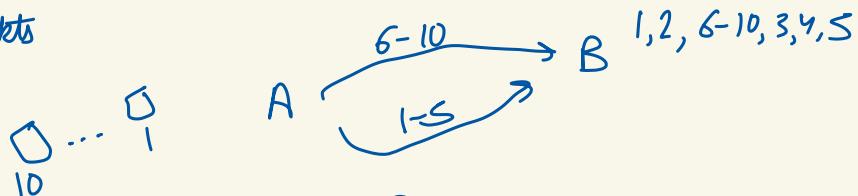
## Issues :

- Under heavy load, routing path oscillates

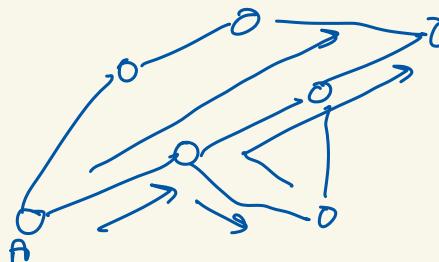
Queuing delays go up  $\rightarrow$  increased link weights  $\rightarrow$  new shortest paths  
 ↗ can paths used heavily

- end-to-end latency (eg : A to B) keeps varying  
 ↳ may affect application layer performance.

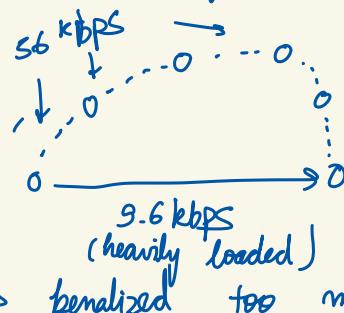
- reordering of pkts



- routing loops possible

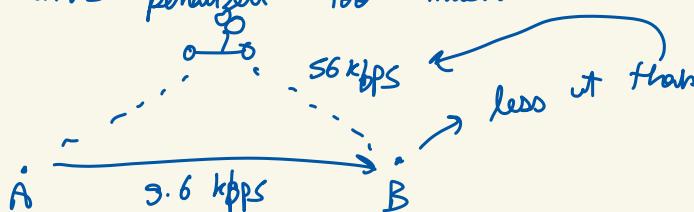


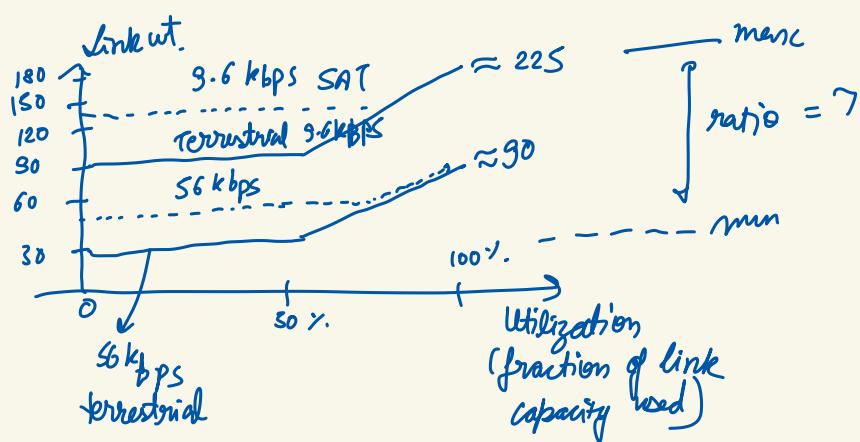
- Range of link weights was large due to which some links were penalized (due to very high weights) too much.



Found that 126 links of 56 kbps had some wt as single 9.6 kbps link

- Satellite links penalized too much





MTS change infrequently

Today : OSPF : ut of link max.  $\left( \frac{10^3}{\text{link speed (bps)}} \cdot 1 \right)$

NOC : network operations centers  
(AT&T)

## IP Addressing

MAC address (layer 2)

IP address (layer 3)

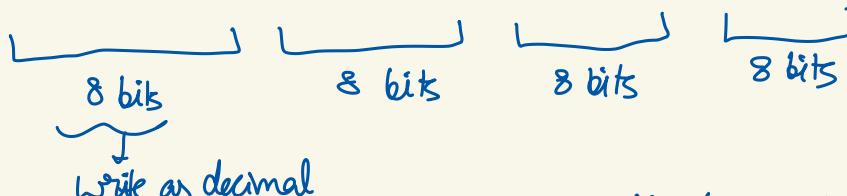
IPv4 : 32 bits  $\rightarrow 2^{32}$

IPv6  $\rightarrow 128$  bits

NAT Network Addr

Translation

(to reuse IP addresses)



Ex : 255. 255. 255. 255  $\rightarrow$  all 1's reserved for broadcast

{ 10. \* . \* . \*  $\rightarrow$  anything

private  
IP

public IP → should be unique in the internet

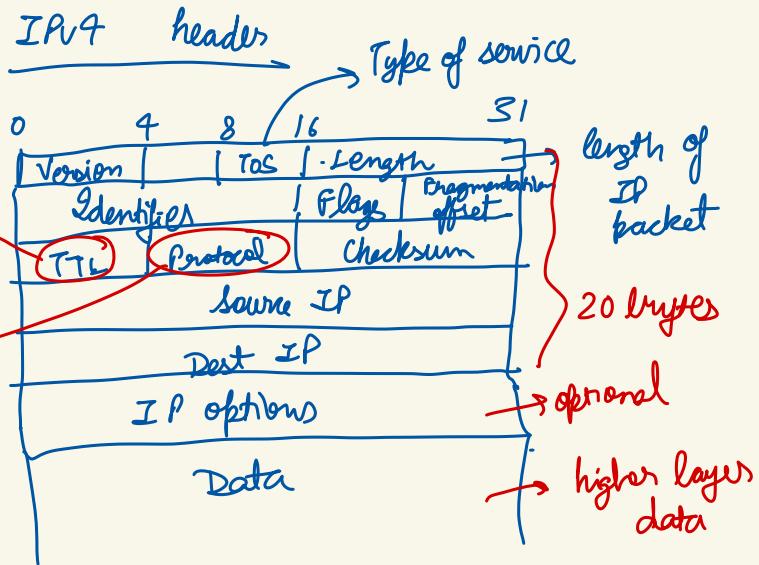
only one host must use that IP address

### IP header

Time to live  
(decremented  
at each "router")  
Go drop packet

Next layer  
protocol

- 6: TCP
- 17: UDP
- 1: ICMP



### Routing Table

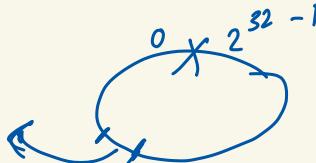
$$IPv4 \approx 2^{32}$$

small  
in numbers

Destination	next hop
730.52.30. <sup>**</sup> IP prefix	R <sub>2</sub>

JIT  
Bombay

slice



a. b. c. d  
octet

X

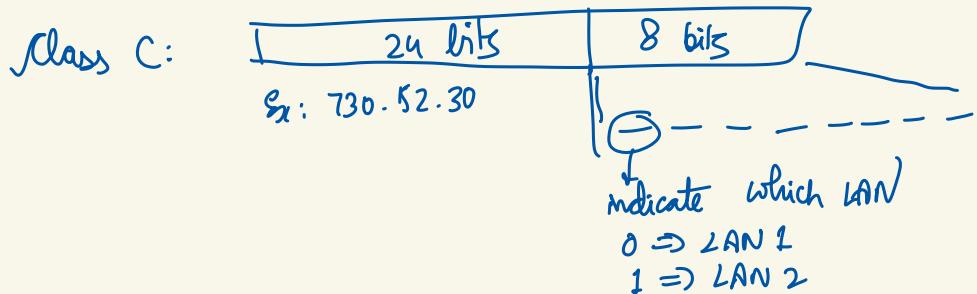
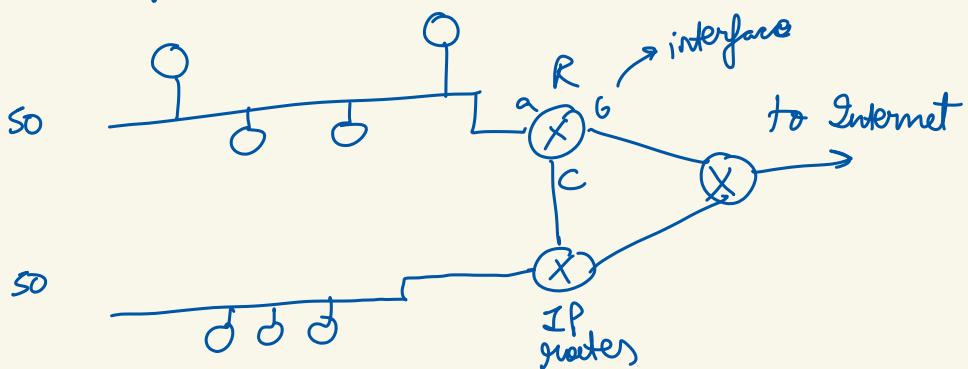
730.52.30.\*

assign to IIT

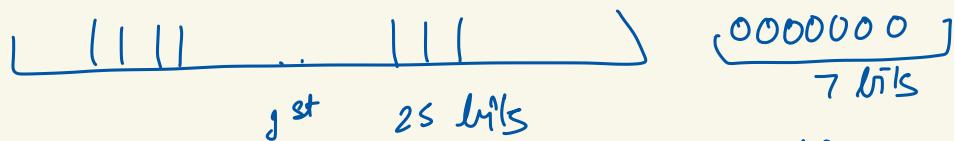
prefix of all machines

Class A	8 bits	$\underbrace{27 \text{ bits}}_{\text{host}}$	$2^u \text{ hosts}$
	network		
		host	
B:	16 bits	16 bits	
	network	host	
C:	24 bits	8 bits	
	network	host	

Subnetting: Given a slice of IP addresses; how to divide among LANs, setup config, interval greater.



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



Subnet address:  $S_1$  for LAN-1 ( $M_1$  is Mask)  
 $S_2$  for LAN-2 ( $M_2$  is mask for LAN<sub>2</sub>)

Router-R: suppose dest. IP is 'D'  
 $G_S(D \text{ and } M_1) == S_1 ?$

$$S_1 = 730.52.30. \underbrace{0}_{0}$$

Yes  $\rightarrow$  Do nothing.

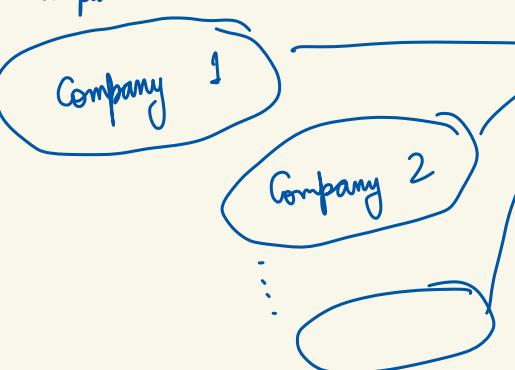
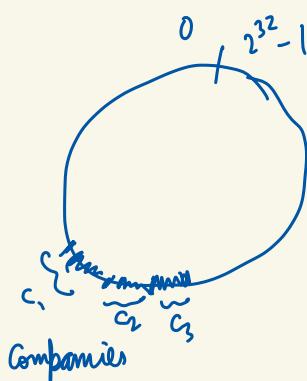
No  $\rightarrow G_S(D \text{ and } M_2) == S_2 ?$

$$(S_2 = 730.52.30.128)$$

1 is the leading bit  
~~leading bit~~

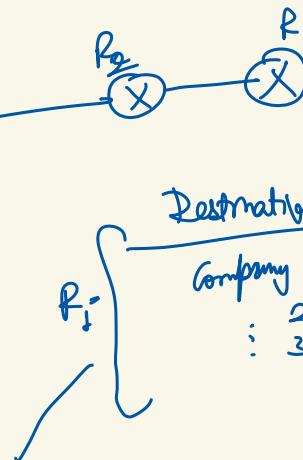
NO  $\downarrow$   
 Yes  $\rightarrow$  FWD to interface C

send to 'b'

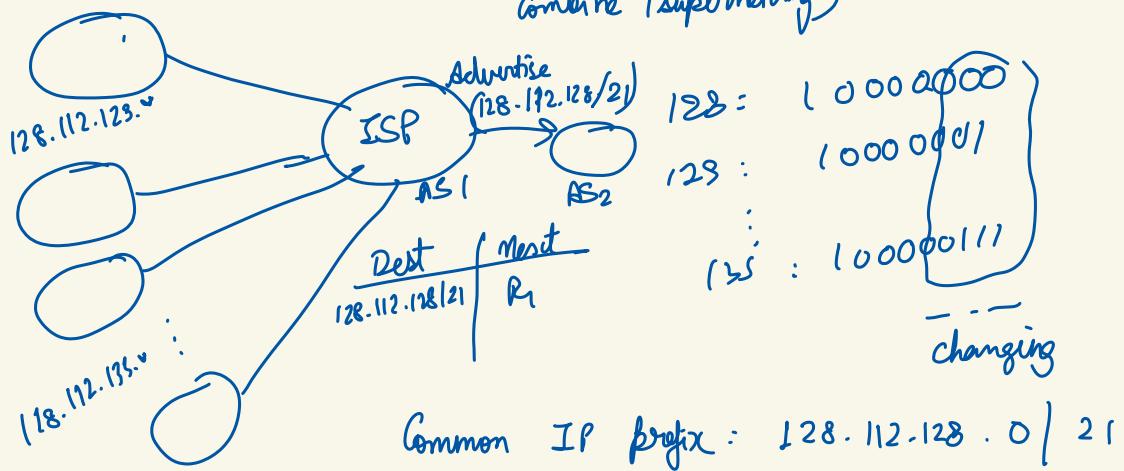


Destination      Host

Company 1	$R_2$
⋮	$R_2$
⋮	$R_2$



$128.112.128.x \rightarrow$  class C



If given dest. IP address D

If first N bits of D match

with the first N bits of a.b.c.d  
then D belongs to that prefix

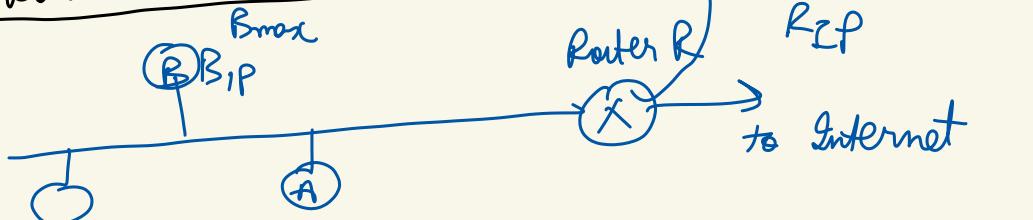
a.b.c.d / N

↓  
Consider N leading  
bits to get the  
IP Prefix.

CIDR: classless inter. domain routing

↳ specify any prefix length (N)

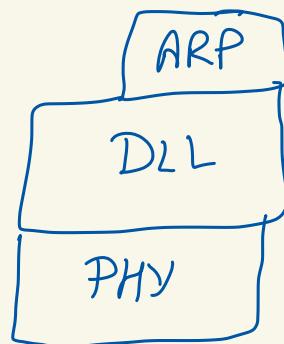
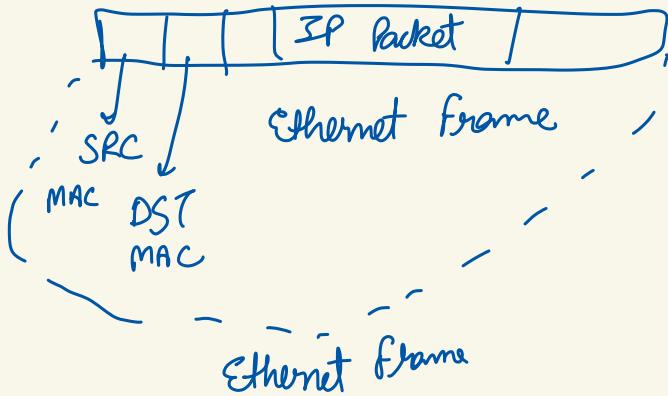
ARP: Address Resolution Protocol



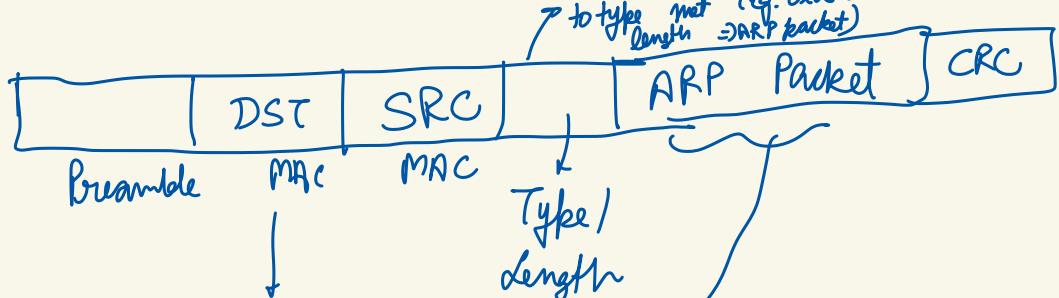
A wants to send an IP packet to B

A knows B's IP; A does not know B's MAC

ARP to the rescue



Ethernet Frame :



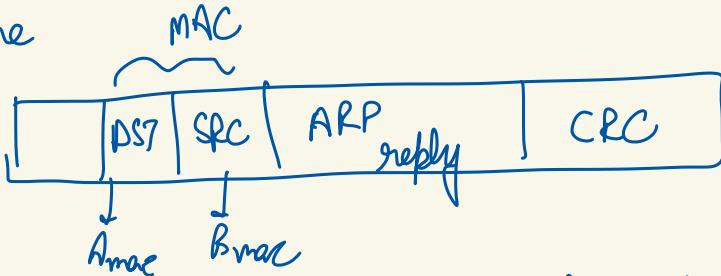
ARP Request : Sender MAC ; Sender IP  
(A<sub>MAC</sub>) (A<sub>IP</sub>)

Target MAC ; Target IP  
(All zeros) (B<sub>IP</sub>)

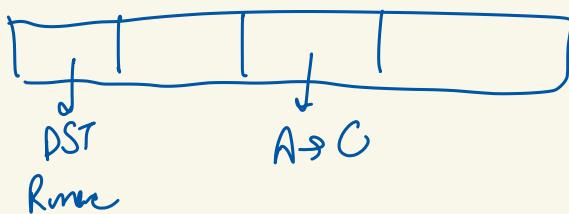
ARP Reply : Sender MAC ; Sender IP  
(From B→A) (B<sub>MAC</sub>) (B<sub>IP</sub>)

↓  
unicast  
frame

Target MAC ; Target IP  
(A<sub>MAC</sub>) (A<sub>IP</sub>)



Information (B<sub>mac</sub>) stored in ARP cache at A.  
with a timeout (orders of minutes)



(1) How does A know if DST IP belongs own network or not?

(2) If DST IP is not in own network , how to known R<sub>IP</sub>, R<sub>MAC</sub> .

$$(1) \quad A_{IP} = a_1 \cdot a_2 \cdot a_3 \cdot a_4$$

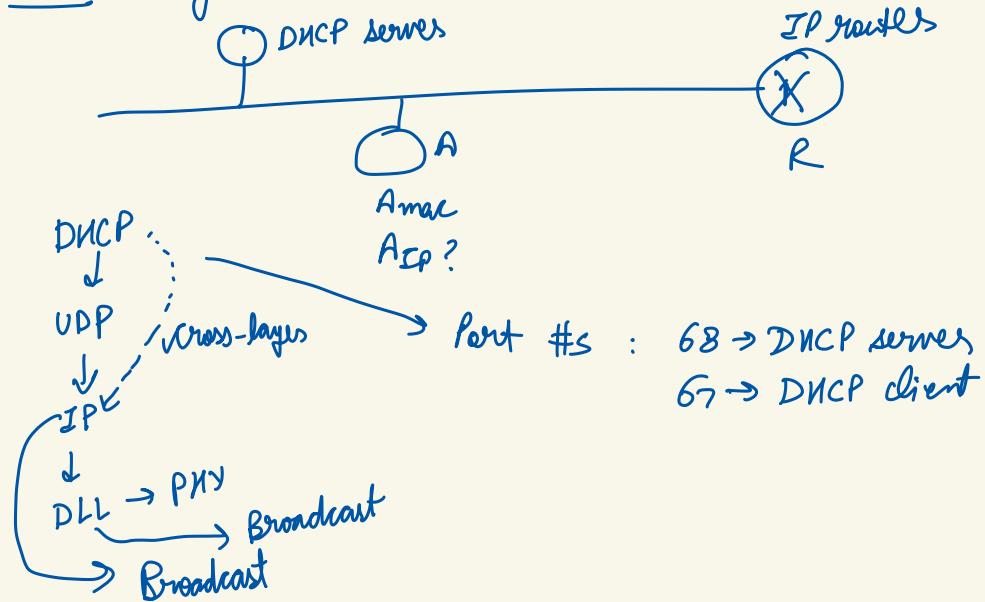
subnet MASK : Ex  $\underbrace{255 \ 255 \ 255 \ 0}_{\text{all } 1\text{s}}$

If  $\underbrace{\text{DST-IP}}_{\text{Network address}} \text{ AND MASK} = A_{IP} \text{ AND MASK}$

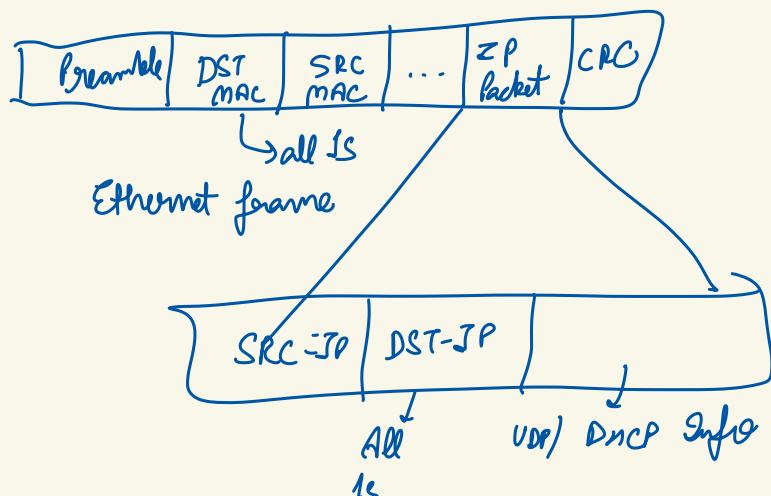
then DST-IP is in my network .

(2) Suppose know RIP, how to find RMAC  
Use ARP

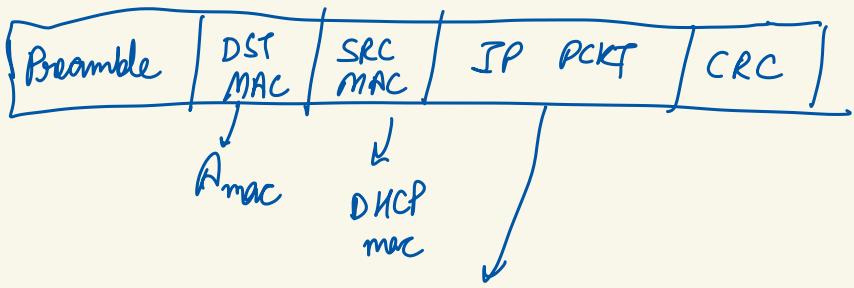
DHCP: Dynamic Host Configuration Protocol



'A' sends out DHCP 'Discoverer' packet

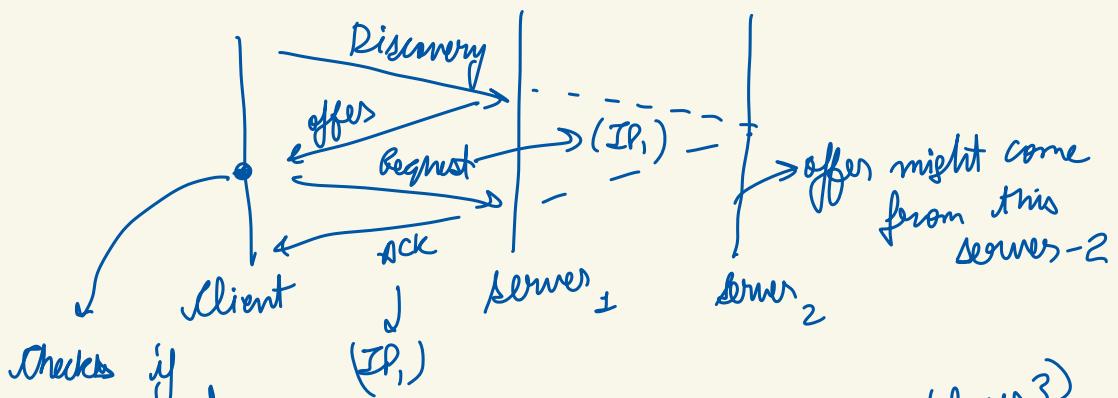


DHCP server replies (offer) → Contains a potential IP address for A.



we are yet  
 to configure Aip

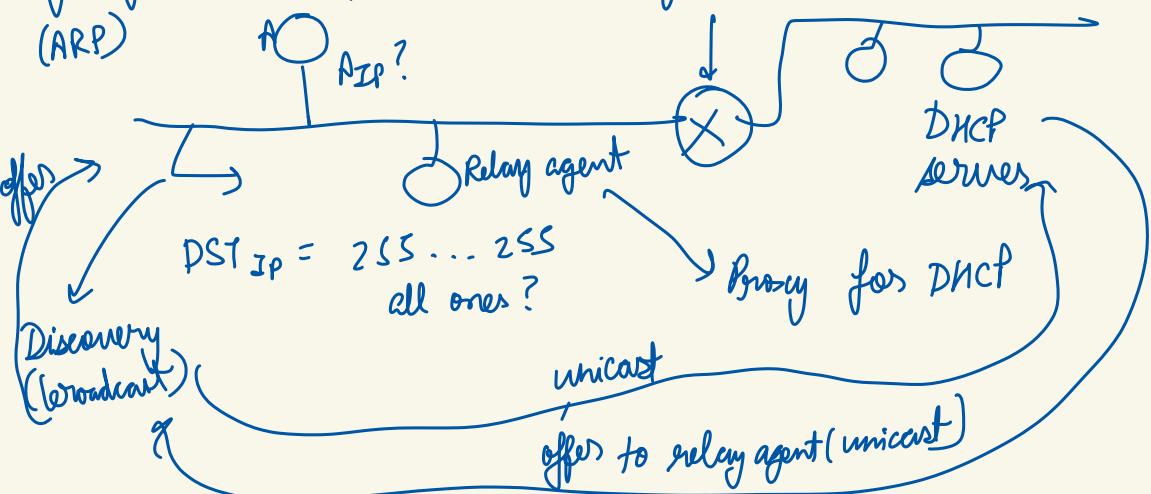
SRC IP: DHCP IP  
 ← DST IP: All ones (broadcast)



Checks if IP<sub>1</sub> is used by anyone else (ARP)

A mac  
A IP?

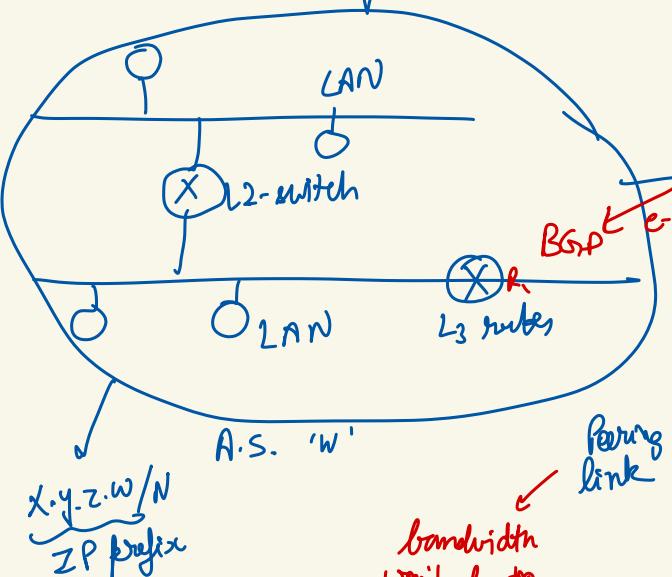
gateway ignores (layer 3) does not forward all 1s



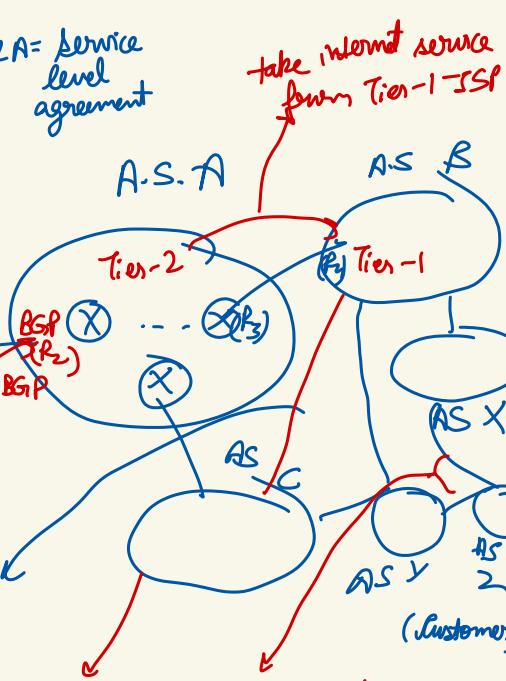
# BGP : Border Gateway Protocol

SLA = service level agreement

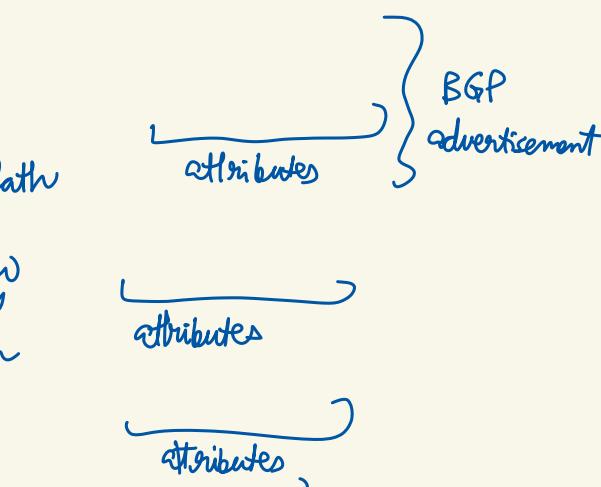
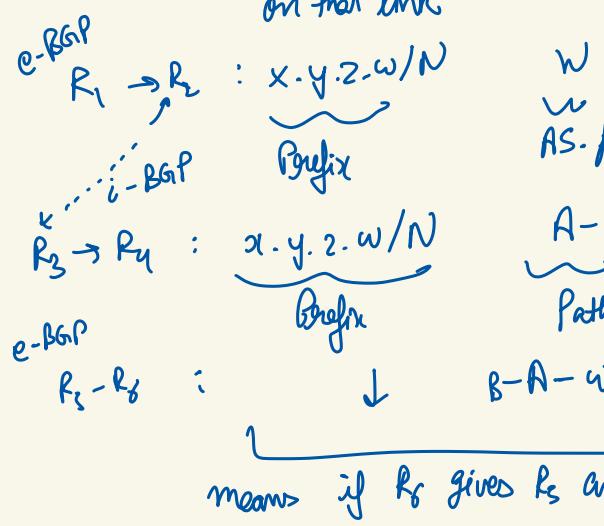
Inter domain routing: between ASes



bandwidth won't be too high  
usually under-provisioned because both don't get paid for delivering traffic on that link



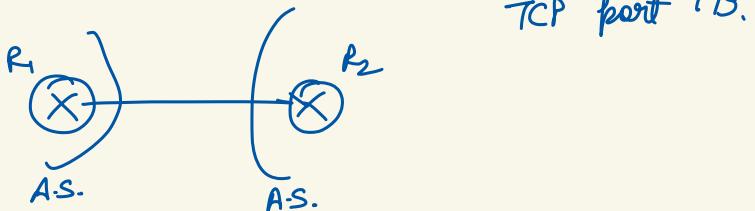
usually widely connected across country/globe, not customers for any other AS.



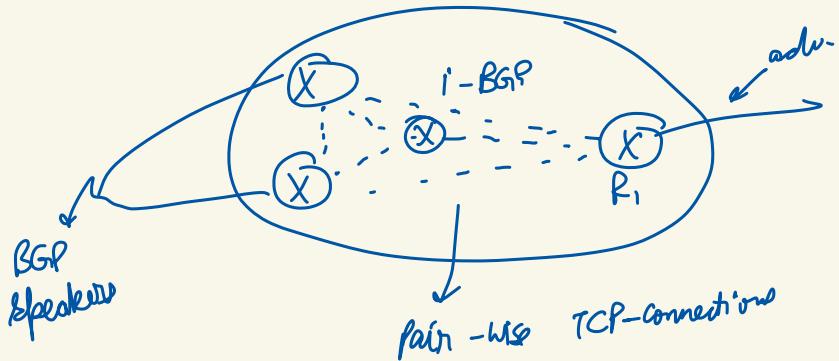
matching prefix, then it will definitely route it along AS path  
B-A-w

B does not have to tell C: that there is a path B-A or  
B-A-

e-BGP: Protocol between BGP speakers in different (neighboring) ASes



i-BGP : Protocol between BGP speakers within same AS-

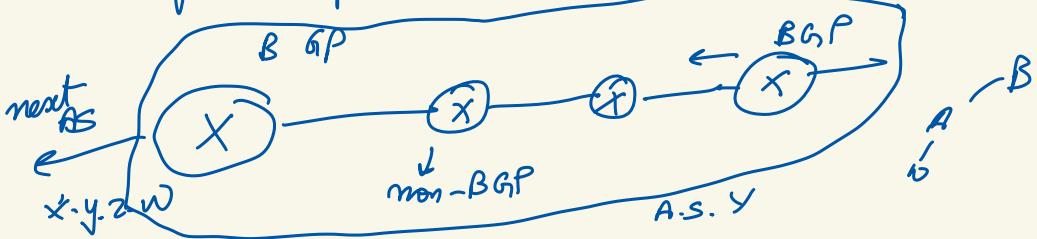


IGP : Interior Gateway Protocol  
(Intra-domain routing, LSR, DV)

Helps to decide BGP paths

- (i) e-BGP speakers learn AS-paths from neighboring routers in other ASes
- (ii) e-BGP nodes share learned information via i-BGP with other BGP speakers in own AS.
- (iii) BGP speakers select routes to various IP prefixes.
- (iv) Insert chosen routes into IGP (since all routers must be

able to forward pkts to a destination IP

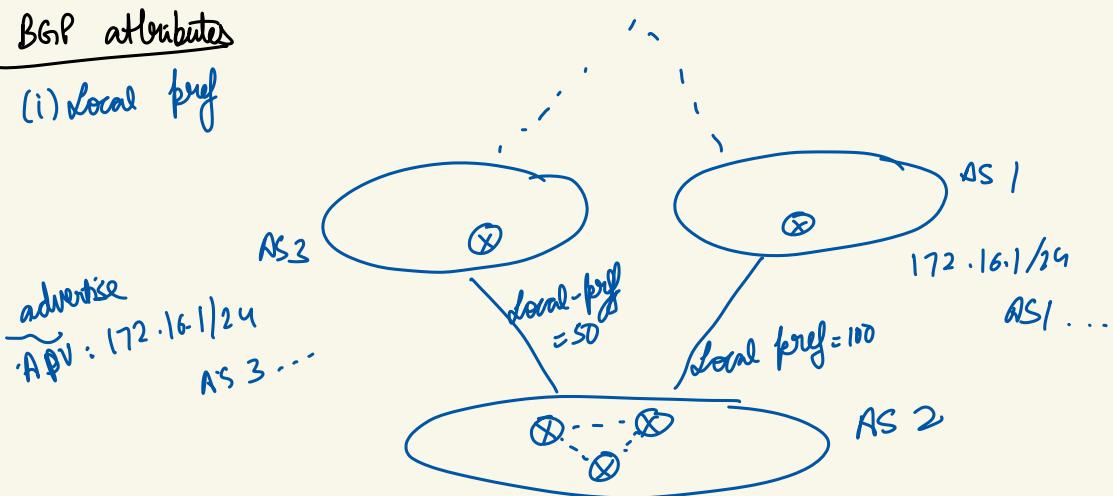


- (v) eBGP speakers can advertise newly created routes to neighbouring ASes

172.16.1/24

### BGP attributes

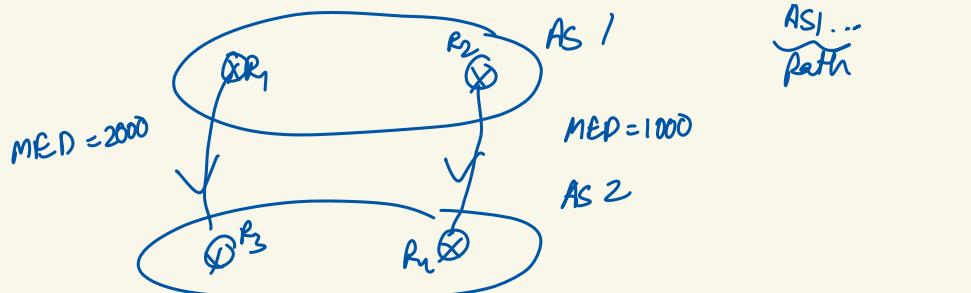
- (i) local pref



Admin of AS2 adds local pref himself

Higher local pref  $\Rightarrow$  more preferable path

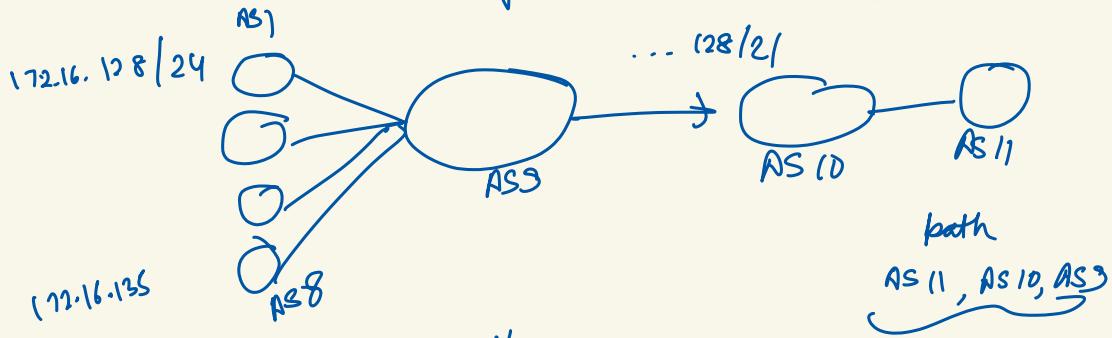
- (ii) multi-exit discriminators (med)



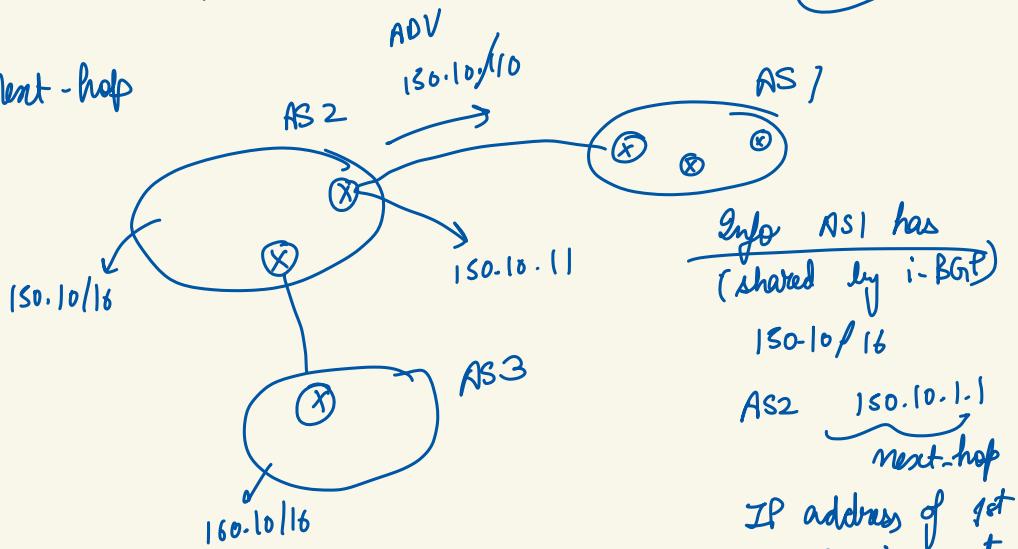
means: AS1 is telling AS2: I prefer that you send me packets for  
 172.16.128/24 on link R<sub>1</sub>-R<sub>2</sub> (lower MED)

(iii) AS-path : list of AS numbers all the way to the destination with that IP prefix

↳ whenever first adv that IP prefix



(iv) Next-hop



Rules to choose routes :

Each BGP speaker decides which AS route to use among many available for the same prefix

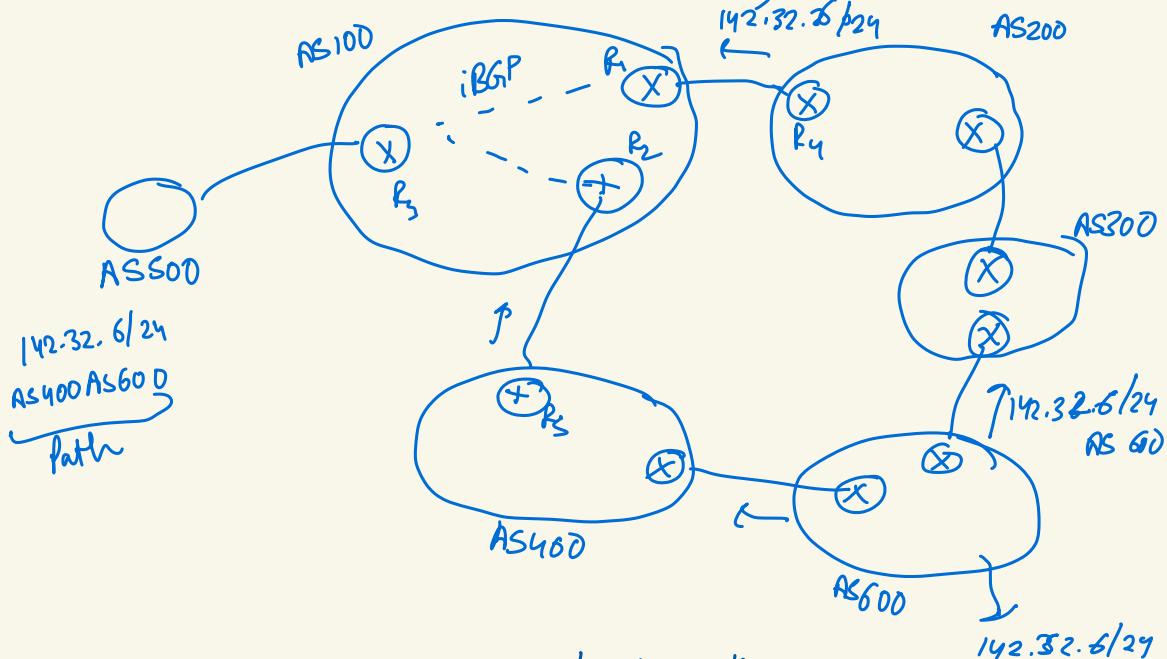
- (a) Use route with largest local pref
- (b) Choose path with shortest AS-path

160.10.10/16 AS2, AS3 150.10.1.1  
 path next hop

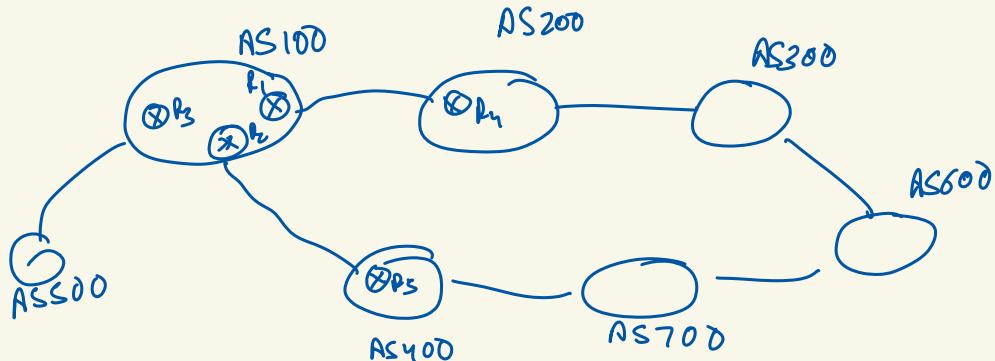
- (c) Choose path with lowest MED
- (d) Choose path learned over e-BGP over path learned over i-BGP
- 
- $x.y.z.w/N$  AS-PATH (ADV1)
- $x.y.z.w/N$  AS-PATH
- $x.y.z.w/z$  AS-PATH
- R<sub>2</sub> chooses to use ADV2 to send to x.y.z.w/z
- (e) Choose path with lowest IGP metric to next Hop
- $\text{dist}(R_3, R_2)$        $\text{dist}(R_3, R_1)$
- $R_3$  compares and chooses the smaller - /
- Suppose  $\text{dist}(R_3, R_2) > \text{dist}(R_3, R_1)$   
 $\Rightarrow R_3$  chooses ADV
- Not-potato routing**
- (f) Use router id (lowest router id among all BGP speakers who have sent these advertisements)
- Router-id = highest IP address on route

### BGP

1. Local Pref  $\rightarrow$  largest
  2. AS-Path  $\rightarrow$  shortest
  3. Med  $\rightarrow$  lowest
  4. eBGP learned route over iBGP learned one
  5. Not potato routing
  6. Router ID (lowest)
- AS200 AS300 AS600



- 1) If want all routers in AS100 to use the top route, admin of AS100 can set local-pref highest for it.
- 2) Suppose we want all BGP routers to use the lower route.  
AS 400 - AS 600
  - (i) Use local-pref
  - (ii) By default all will use the shortest AS-path (no need to set local-pref highest for this path).



Suppose Local pref and MED are same for both paths

Path 1: AS200 - 200-600

Path 2: AS400 - 700-600

$R_1$ : learned path 1 over eBGP path 2 over iBGP

$R_2$ : learned Path 2 over eBGP; path 1 over iBGP

$R_3$ : both over iBGP

$R_1$  uses Path 1

$R_2$  uses Path 2  $\rightarrow$  iGP metric

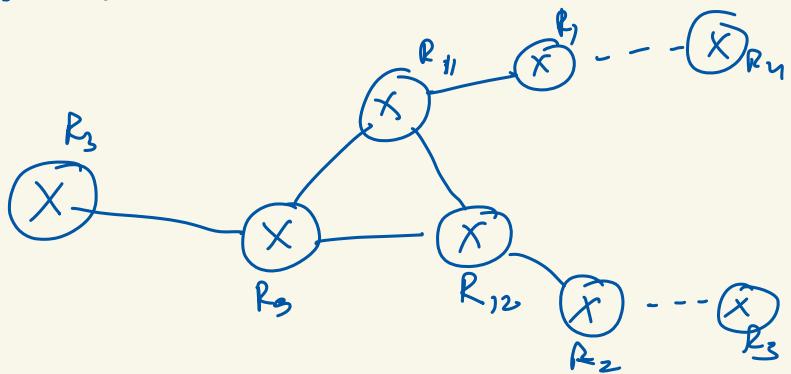
if  $\text{dist}(R_3, R_1) < \text{dist}(R_3, R_2)$  then

$R_3$  uses path 1 (hot potato)

if equal then lowest router id.

AS100

BGP interact with IGP?



Solutions:

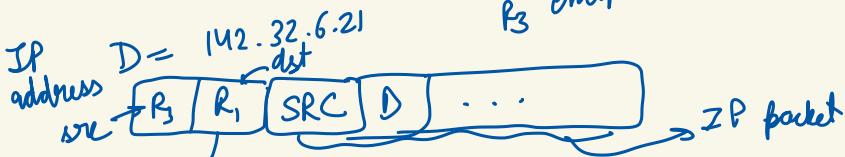
(1) Encapsulation

Suppose  $R_3, R_{11}, R_{12}$  have no BGP information

Ex: in  $R_3$ 's table

$192.32.6/24 \rightarrow ??$  ] external prefix not in the table

$R_3$  encapsulates



DST	next
$R_1$	$R_{11}$
$R_2$	$R_{12}$
$R_3$	$R_3$
$R_{11}$	$R_{11}$
$R_{12}$	$R_{12}$

another  
IP headers

R<sub>1</sub> receives this packet, strips off outer IP layers, finds the internal pkt to R<sub>4</sub>

R<sub>1</sub>'s table  
encapsulated

DST | Next  
R<sub>2</sub>  
R<sub>3</sub>  
⋮  
142.32.6/24 R<sub>4</sub>

(2) Pervasive BGP

All routers are BGP speakers in AS100

\* Suppose there is a unique exit for 142.32.6/24

BGP table

true for  
all

Prefix	Gateway/Exit
142.32.6/24	R <sub>1</sub>

On receiving

IGP

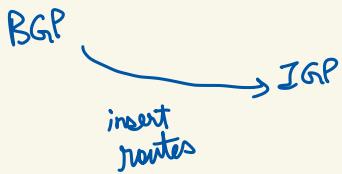
DST IP Addr	Next
R <sub>2</sub>	
R <sub>3</sub>	
R <sub>1</sub>	
R <sub>11</sub>	
R <sub>12</sub>	

SRC / D1 ...

from R<sub>3</sub>  
recursive lookup.  
forward to R<sub>11</sub>

### (3) Tagged IGP

Internal routers may not be BGP speakers



But IGP allows addition of tags  
R<sub>1</sub> can insert into its own IGP

142.32.6/24      R<sub>1</sub> → gateway router  
 prefec      tag  
 propagated to all routers  
 using TGP (say LSP)

At R<sub>n</sub> IGP Table

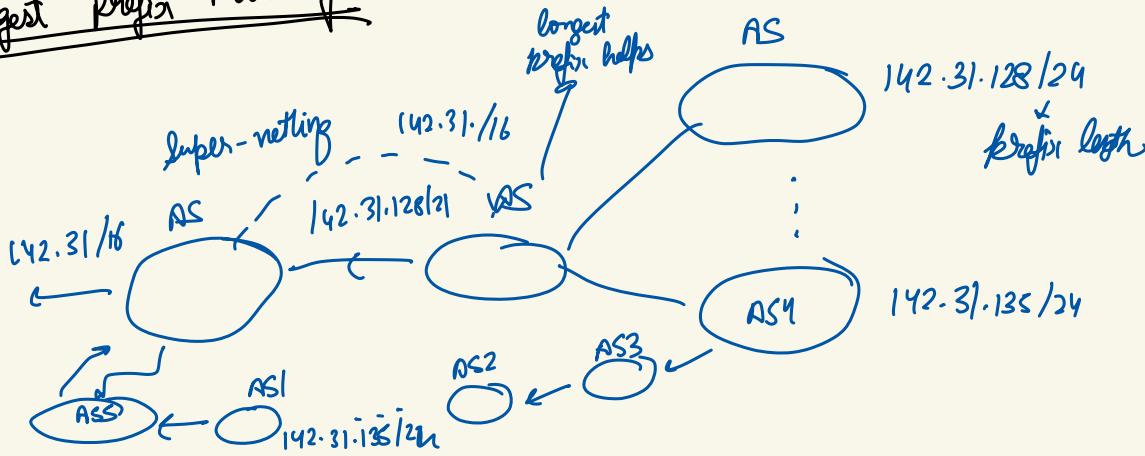
DST	Next	Tag	Cost
R <sub>1</sub>	R <sub>11</sub>		5
R <sub>2</sub>	R <sub>12</sub>		12
⋮			
142.32.6/24		R <sub>1</sub>	
142.32.6/24		R <sub>2</sub>	

Yours

$R_3 \rightarrow R_3$ : SRC ID

lookup, find matches and forwards to closest "tag"  
→ hot potato routing

### Longest Prefix Matching



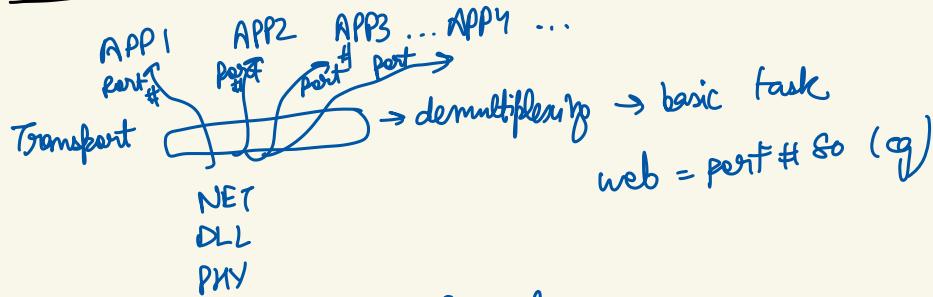
ASS gets a packet with dst IP = 142.31.135.20

Prefix	Next hop
142.31.1 /16	-----
142.31.135 /24	....

Rule = choose the longest prefix to route

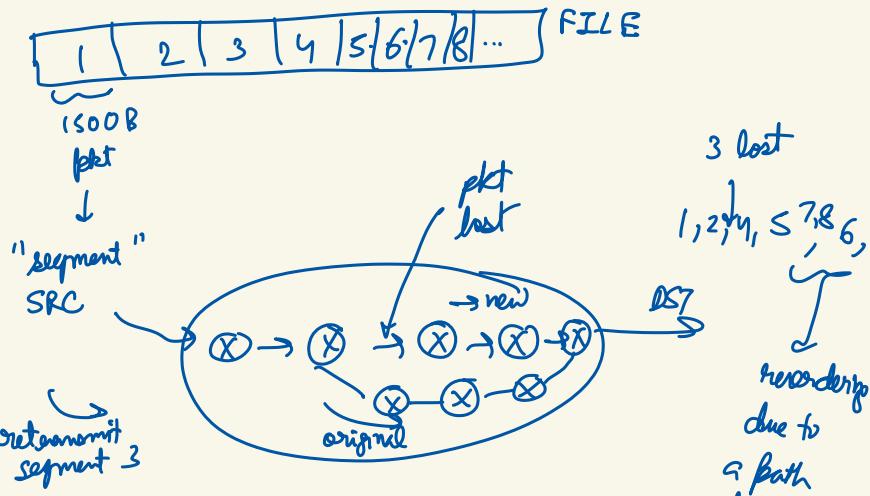
CAM: Content addressable memory  
for easy lookup.

## Layer-4 Transport layers



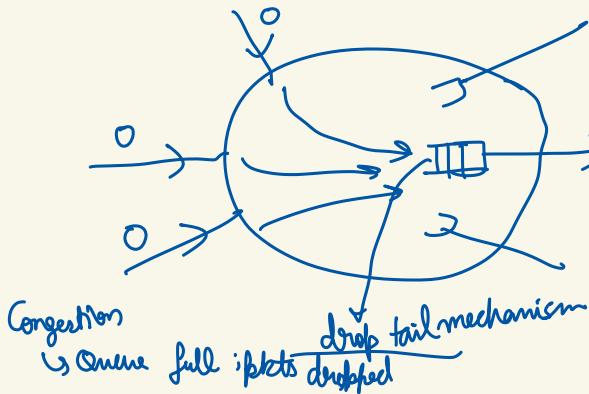
UDP: User Datagram Protocol

File Transfer:



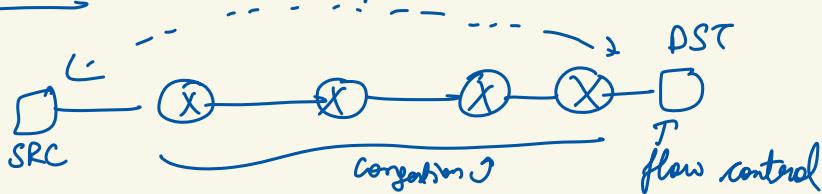
Transport (TCP):

reorders segments  
Reliable data transfer: retransmitting lost segments  
→ Congestion and flow control  
↳ after packet loss



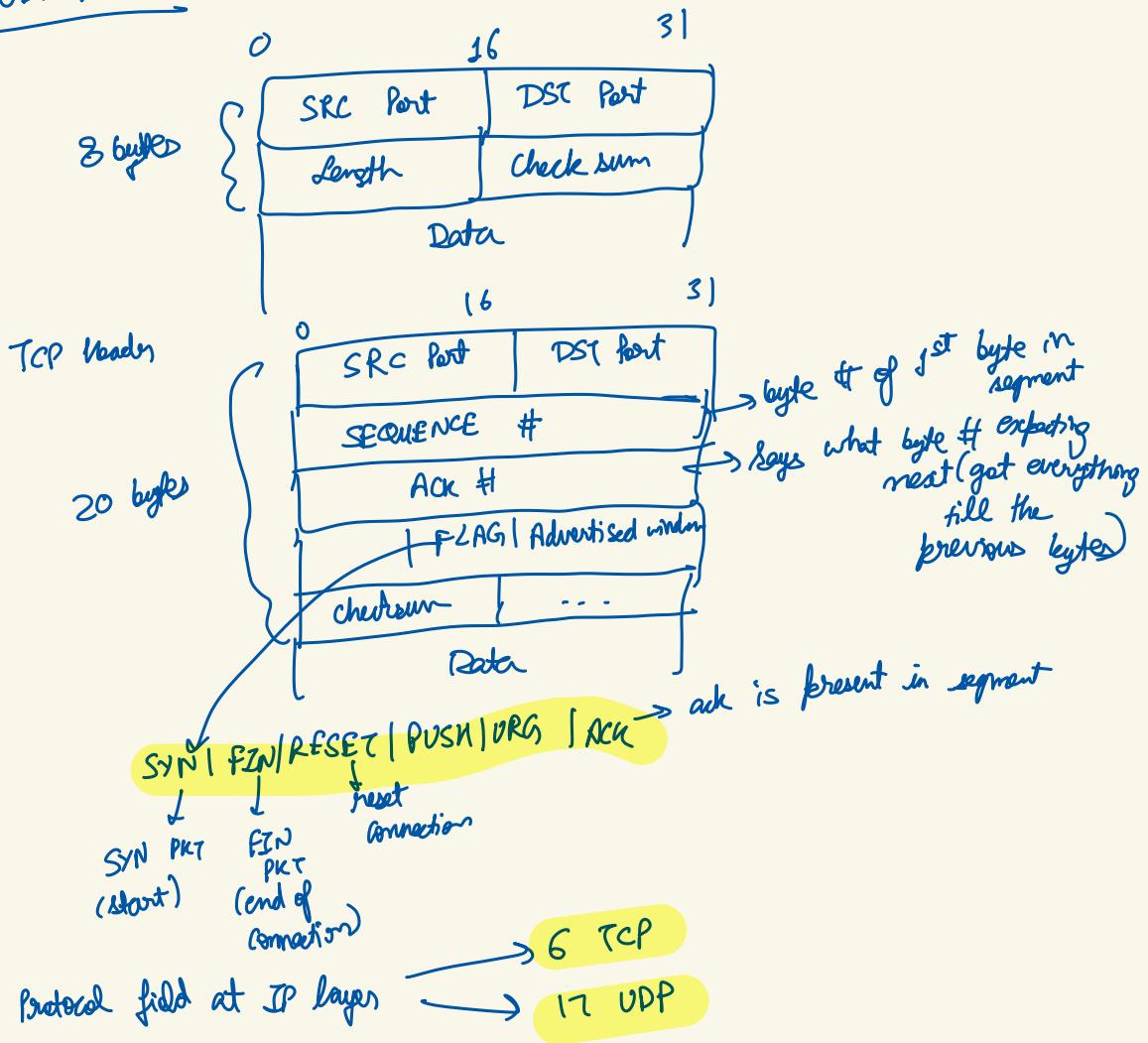
TCP inflates  
congestion  
does congestion  
control  
reduce segment  
sending rate

## flow control



Van Jacobson - 1980s

## UDP Header

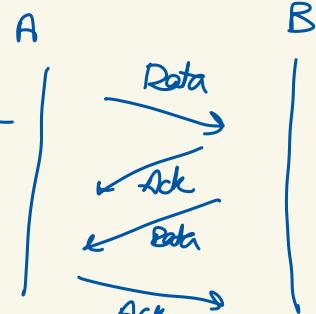


# TCP

SEQ #	ACK #
-------	-------

Connection	Establishment
------------	---------------

we can have  
data and ack  
in the same  
packet too



Passive Participant  
(server)

Active open (client)

Listening

SYN packet (SEQ=x)  
SYN flag set in  
header

(ACK=y)  
SYN+ACK  
both flags set

starting seq #

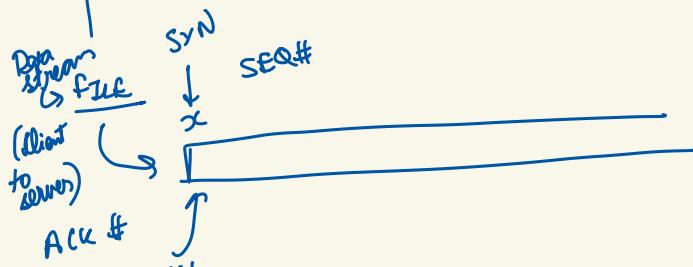
SYN = synchronization

3-way handshake

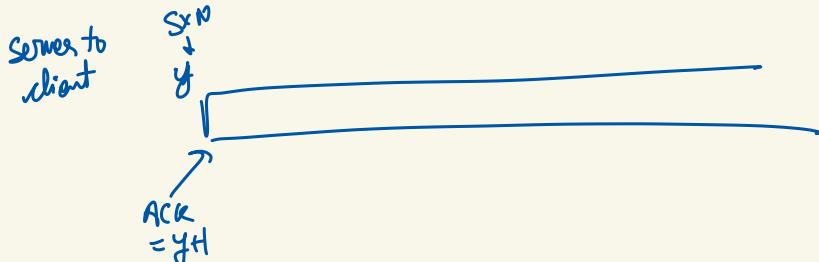
SYN, FIN segments  
don't have data

But considered to communicate  
1 byte.

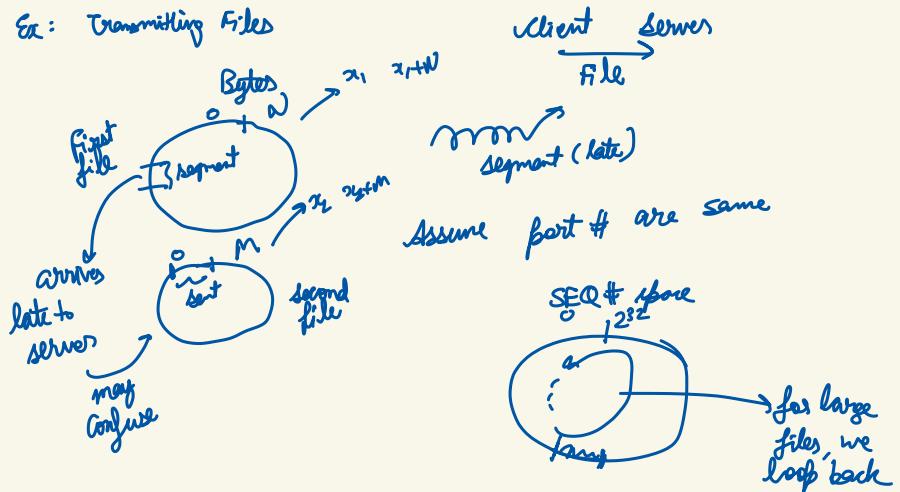
x → randomly selected



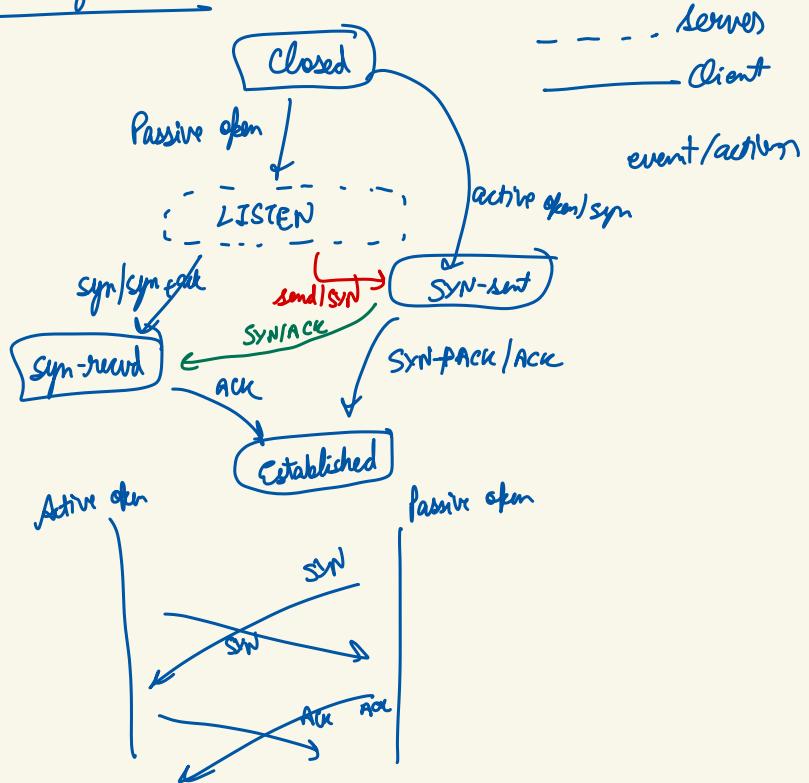
I got everything from start till one less than the ACK#



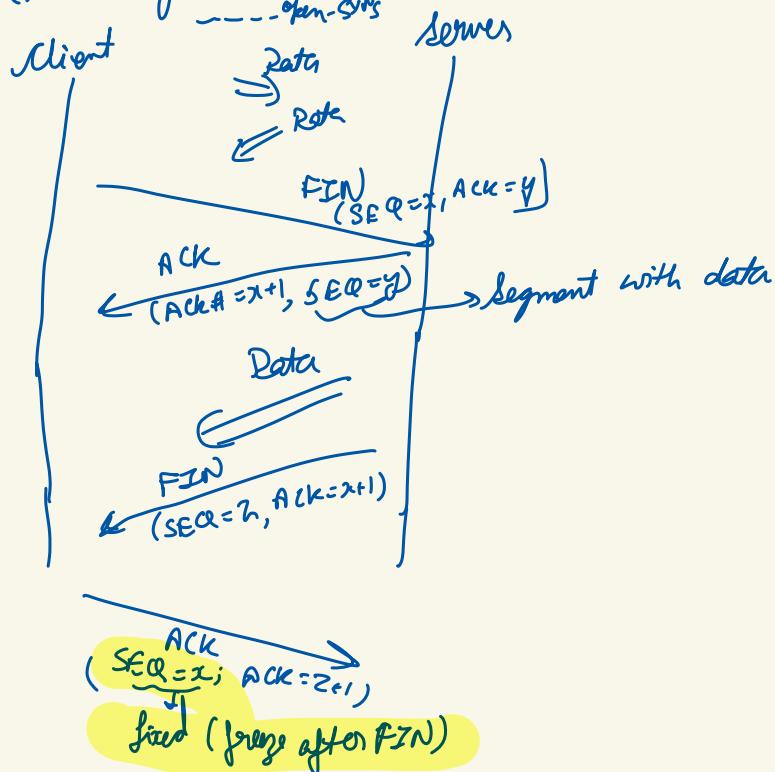
Ex: Transmitting Files



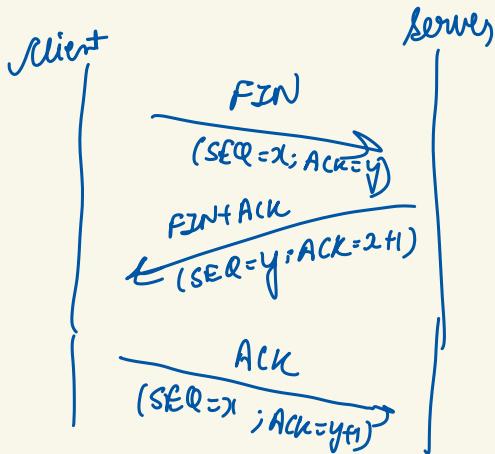
### State Diagram



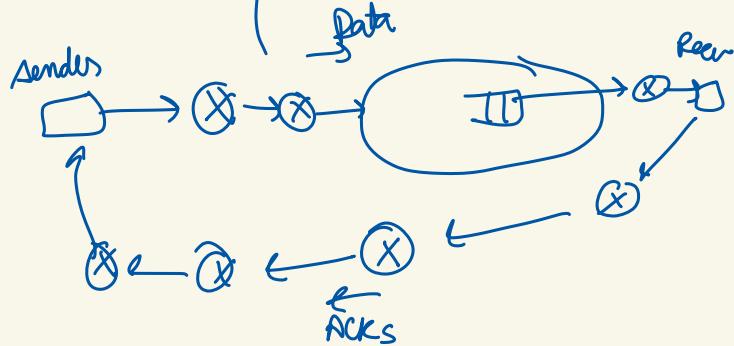
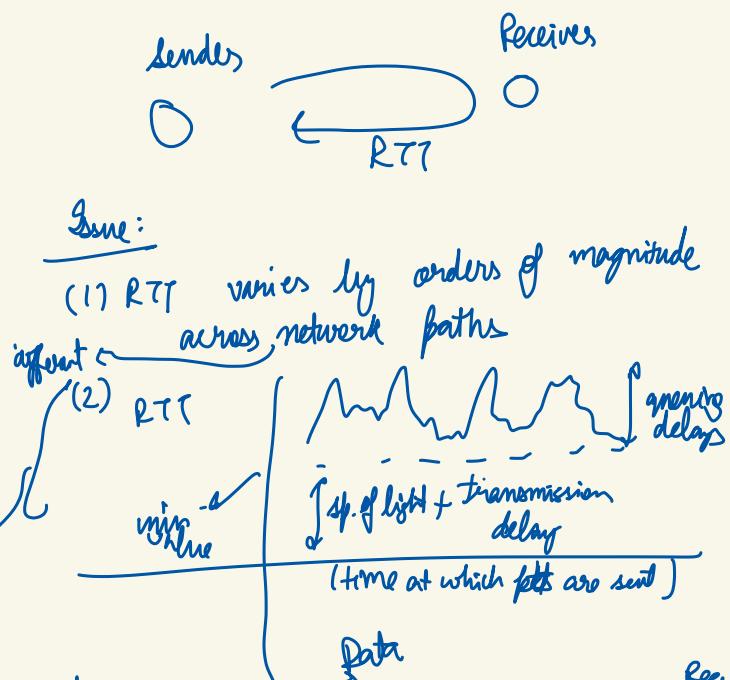
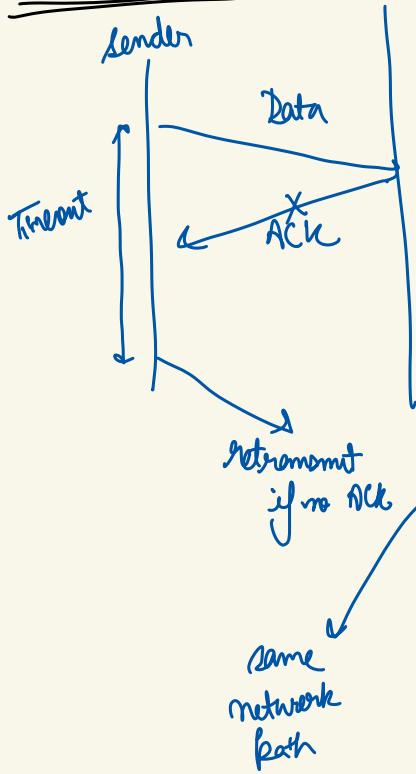
Connection Termination : (Use FIN)  
 (1) Half-close (closed by one end at a time)



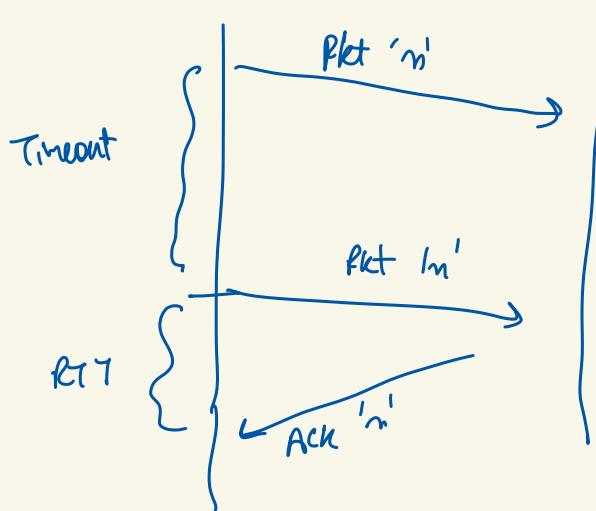
(2) 3-way handshake termination (both closed at the same time)



## TCP Timeout



## 3. Packet loss



## Old Algorithm

Sample RTT = most recent RTT measurement

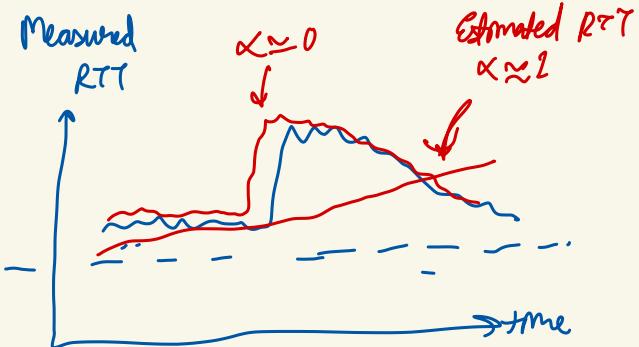
Estimated RTT

$$= \alpha \text{ Estimated RTT} + (1-\alpha) \text{ Sample RTT}$$

like an average

moving average

$$\text{Timeout} = 2 \times \text{Estim. RTT}$$



$$\alpha \in (0, 1)$$

$\alpha \approx 1 \Rightarrow$  less importance to the current measurement

$\alpha \approx 0 \Rightarrow$  lot of importance to recent measurement  
not good for spikes

## New Algorithm

Gaussian Distribution:

Idea: Set Timeout = mean + m \* std

measure: Estim RTT as before

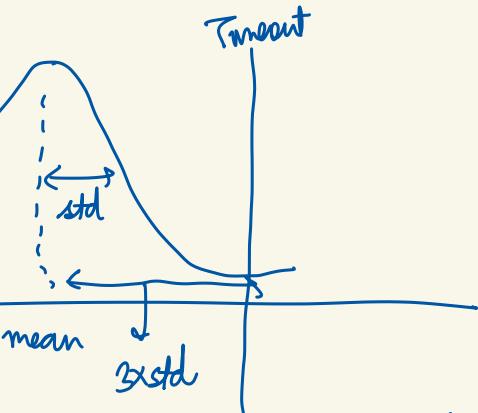
$$\text{Diff} = \text{sample RTT} - \text{estim. RTT}$$

assume mean

Suppose  $x_1, x_2, \dots, x_n$  are some values from random variable  $X$ .  
(i.i.d.)

$$\text{Estim mean } \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\text{Mng. std.} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2}$$



$$\text{Mean deviation} = \frac{1}{N} \sum_{i=1}^N |x_i - \bar{x}|$$

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

$$\text{Dev} = \text{Dev} (1-\beta) + \beta |\text{Diff}|$$

Deviation

most recent deviation

$$\text{Timeout} = \mu \times \text{estimated RTT} + \phi \times \text{Dev}$$

$$\mu = 1$$

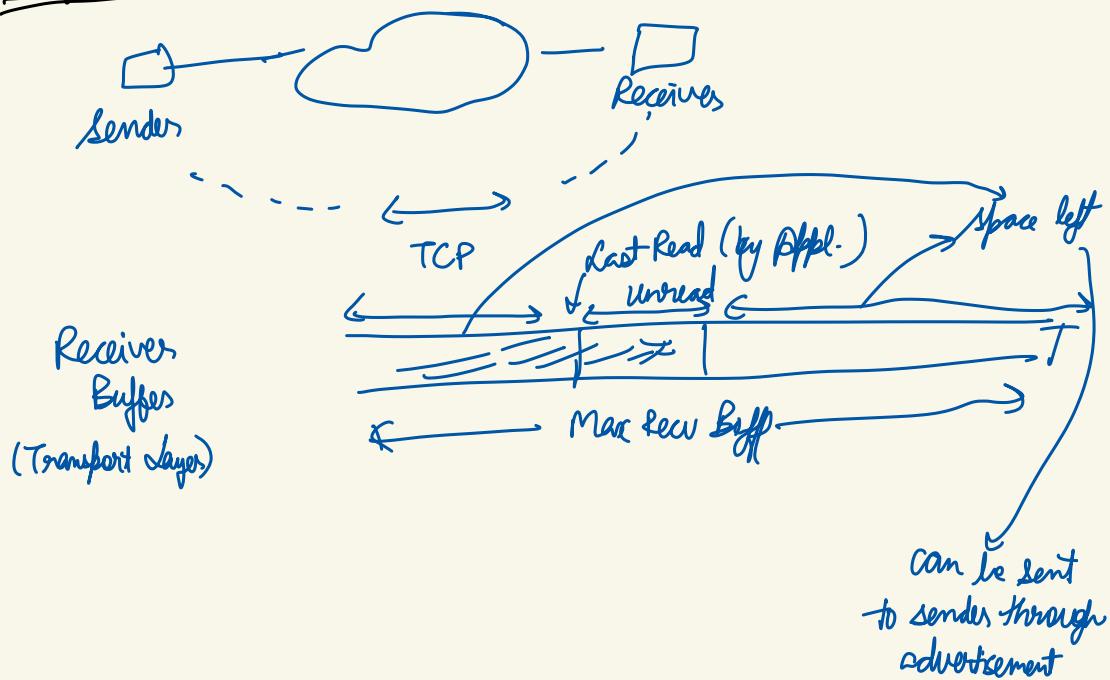
$$\phi = 4$$

$$\text{Recall } \text{estimated RTT} = \alpha \text{ estimated} + (1-\alpha) \text{ Sample RTT}$$

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

$\text{dev}^m$  = Do not use RTT measurement for retransmitted pkt.

## Congestion Control



Flow-control  $\Rightarrow$  Recls with Congestion at receivers

Adv. window (field in TCP header)

$\Rightarrow$  space left in the Recv Buffer

Sender has a "window"

window = max amount of data (in bytes) which is outstanding  
(sent into network but not ACKed)

Suppose RTT is const.

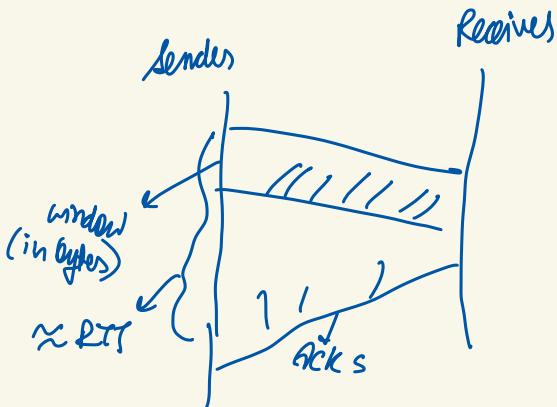
$$\text{Data rate} = \frac{\text{window}}{\text{RTT}}$$

| .

want

$$\text{window} < \text{Adv. window}$$

Ensured to take care of Flow control

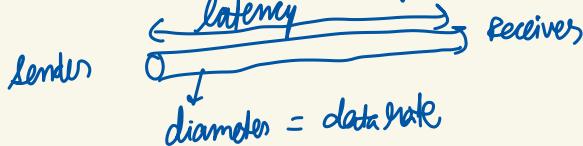


How to deal with Congestion at receivers?

Calculate Cong. window

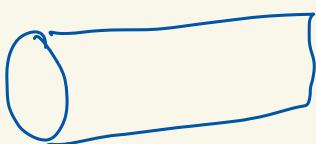
$$\text{window} = \min(\text{Adv. window}, \text{Cong. window})$$

$$\text{window} = \frac{\text{Data rate} \times \text{RTT}}{\text{Delay - Bandwidth product}}$$

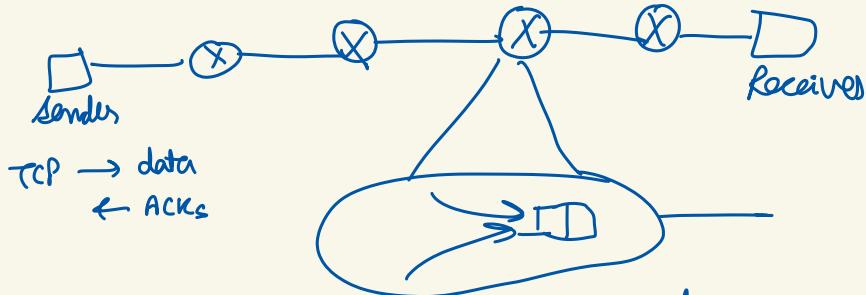


pipes

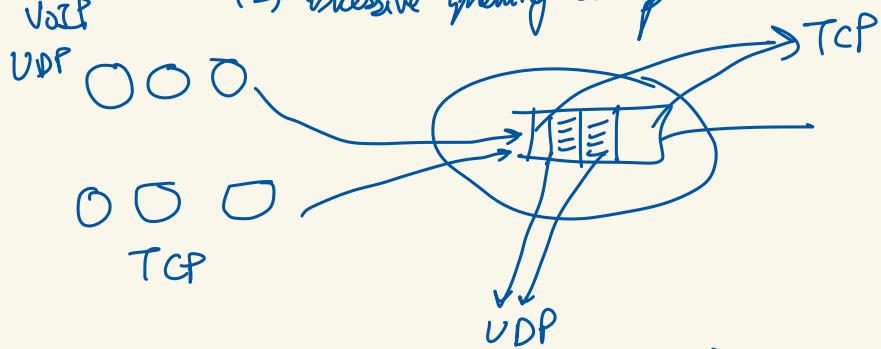
volume  $\rightarrow$  analogous to window size.



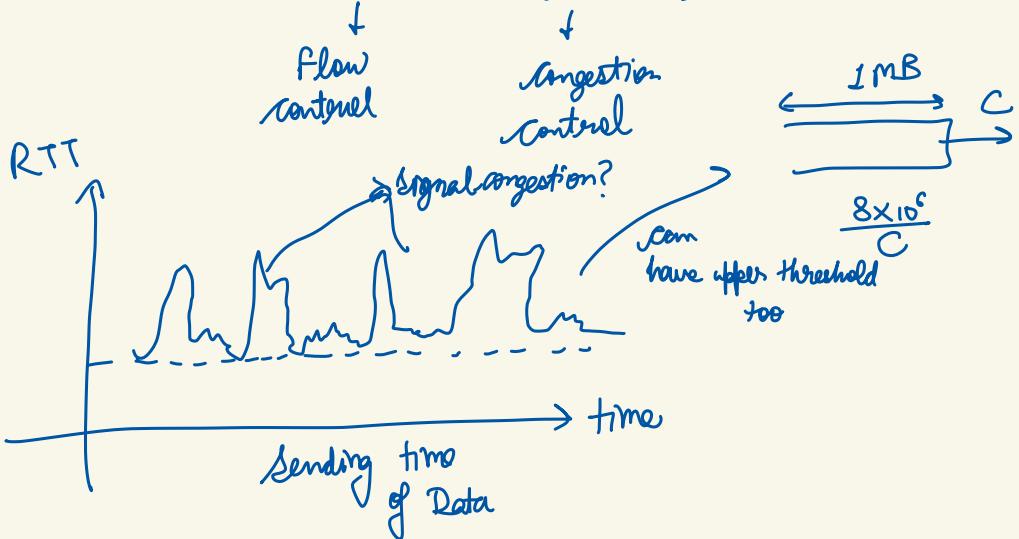
## TCP Congestion Control



Congestion Problem : (1) packets get dropped  
 (2) Excessive queuing delay



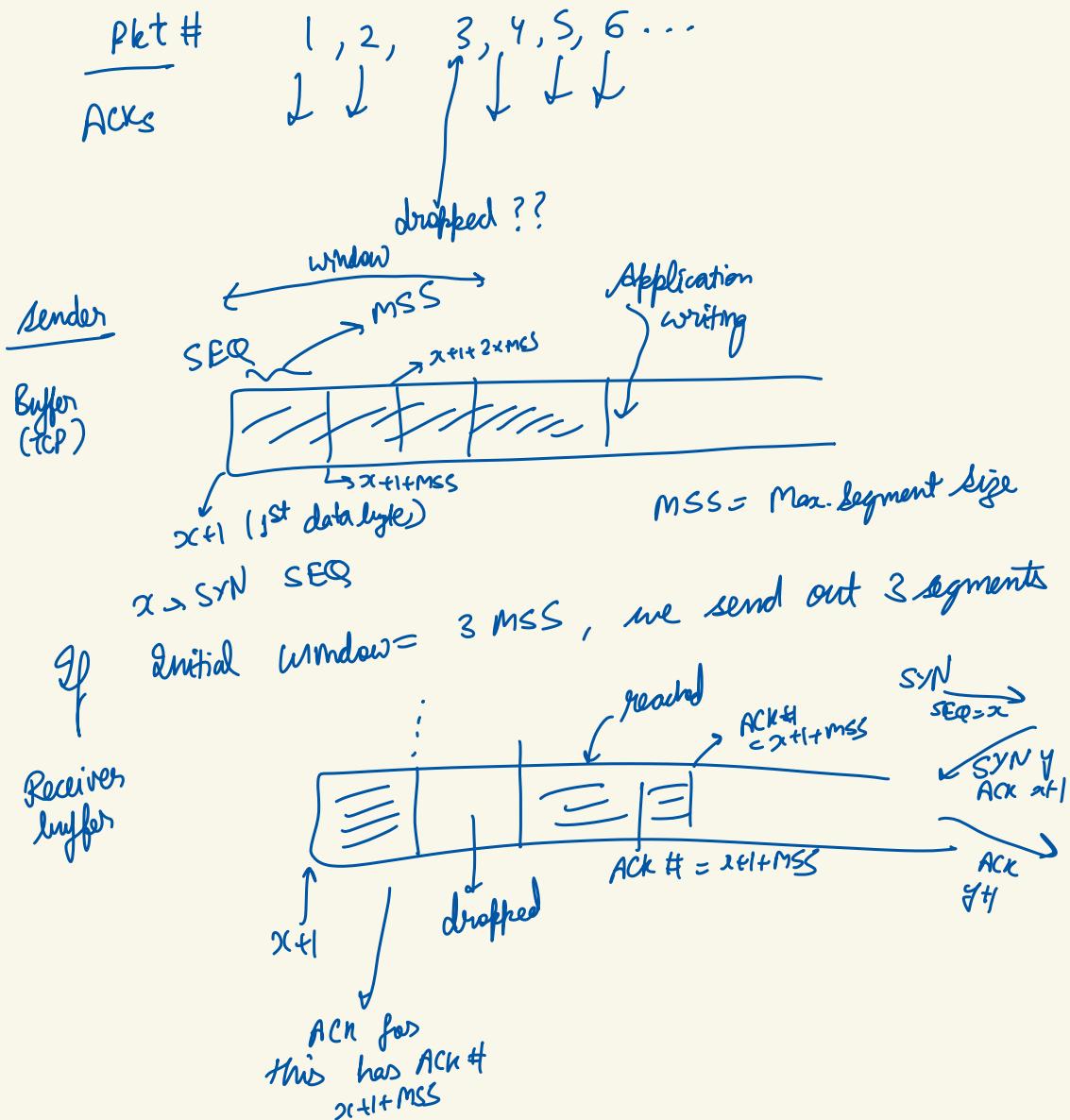
$$\text{Window} = \min(\text{Adv. window}, \text{Cong. window})$$



## Detect packet loss

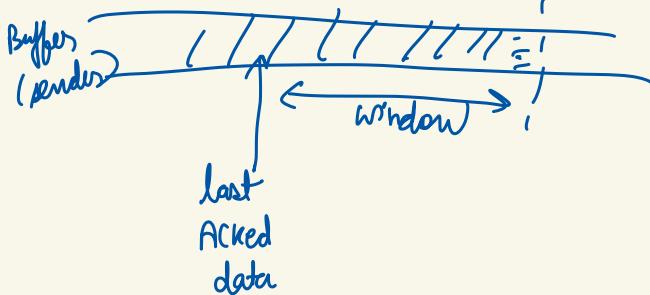
One idea: Timeout

Second idea: Use ACK feedback to infer losses



Cumulative Ack  $ACK\#_z(\text{sent}) \rightarrow$  Everything from start till byte  $(z-1)$  has been received

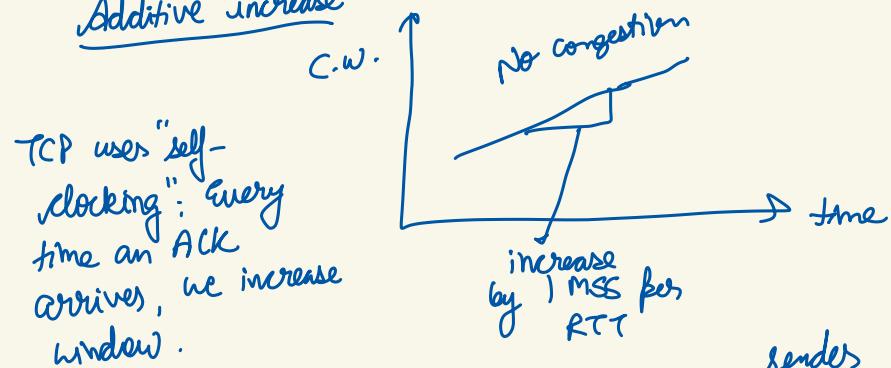
Duplicate ACKs  $\rightarrow$  ACK # of some previous ACK sent out  
 ↳ Each dup ACK says one more segment received.



### Principles of Congestion Control

1. If no congestion, increase CW (congestion window) conservatively.  
 Rotate  $\approx \frac{\text{Window}}{\text{RTT}}$   $\left[ \text{Window} = \min(\text{Adv. window}, \text{cw}) \right]$   
 For now assume window = CW  
Congestion window

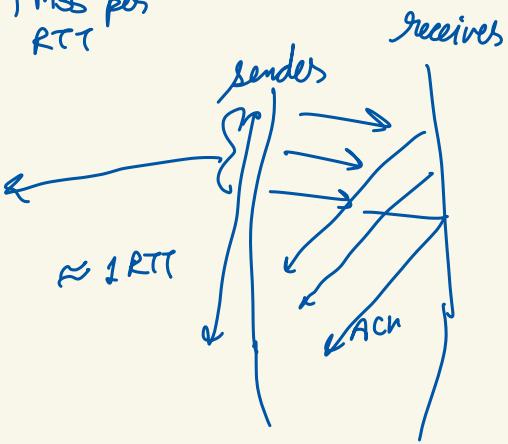
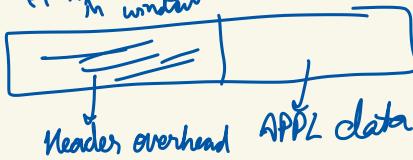
#### Additive increase



CW, MSS

$$\# \text{ segments in window} = \frac{\text{CW}}{\text{MSS}}$$

$$\# \text{ ACKs in window} \rightarrow \text{CW/MSS}$$



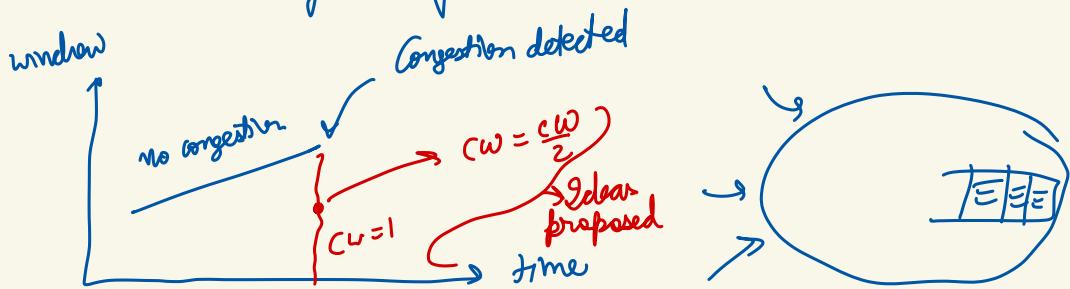
$$\text{Per ACKs window increase} = \frac{\text{MSS}}{\# \text{ACKs}} = \frac{\text{MSS}}{cW/\text{MSS}} = \frac{(\text{MSS})^2}{cW}$$

Additive increase

On receiving an ACK (Assuming no congestion)

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

(2) If congestion is detected, decrease window size aggressively.



TCP Tahoe  $\rightarrow cW=1$  on detecting congestion

TCP Reno  $\rightarrow cW = cW/2$  Multiplicative decrease

TCP Vegas  $\rightarrow$  study later

TCP Africa  $\rightarrow$  not study

AIMD: Add Incr Multi. Decrease  
Increase

(3) Initially

set  $cW$  small size  
we do not know the appropriate datarate

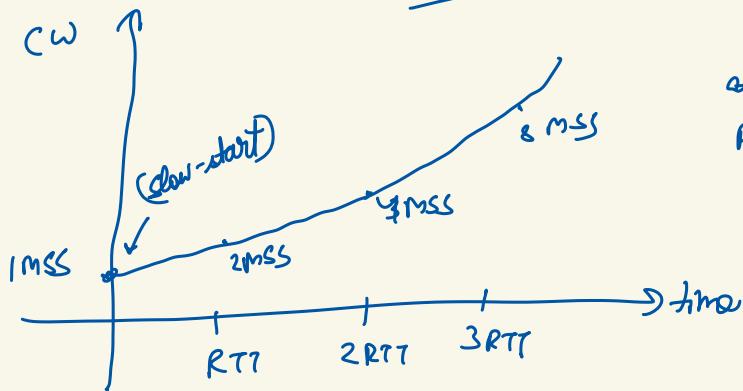
AI may be too slow

Initial window  $\approx 10,000$  bits  
 $\text{data rate} \approx \frac{10^4}{1} = 10^4 \text{ bits/sec}$

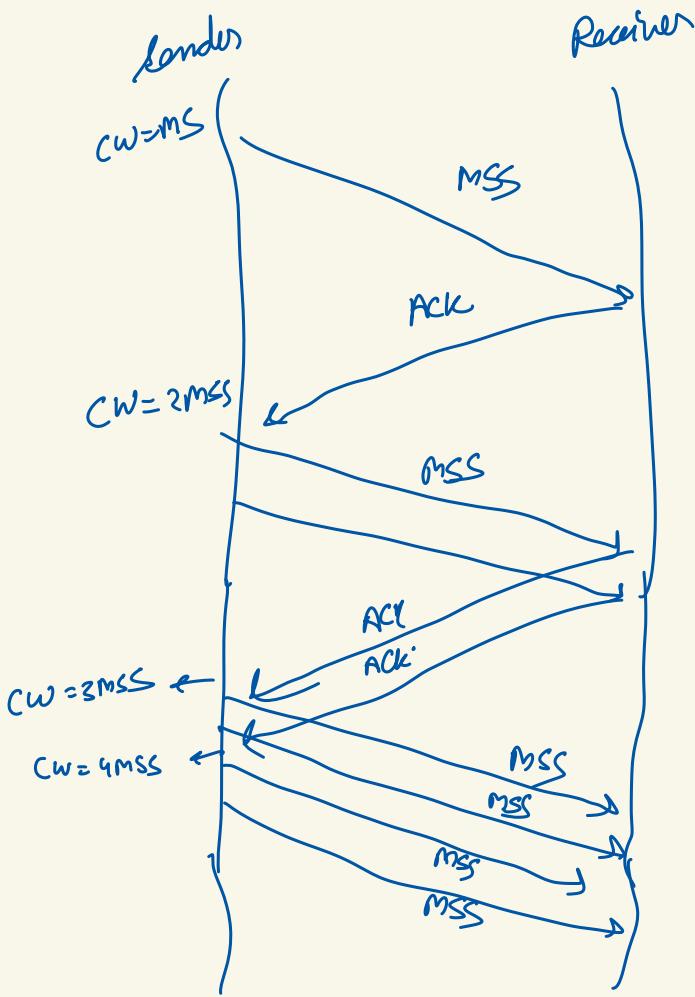
Be aggressive to increase window (initially)

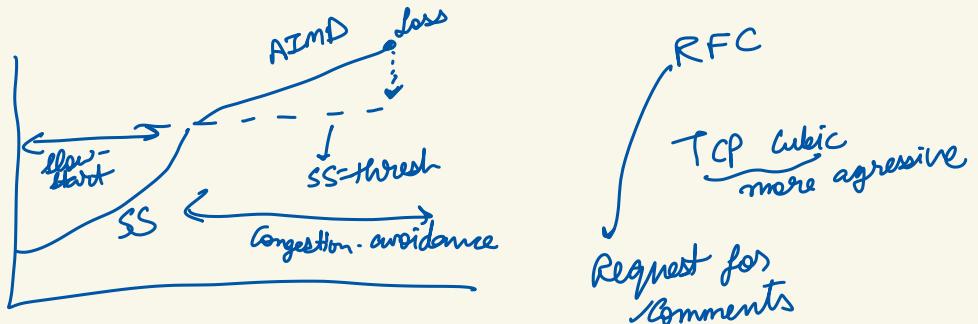
## Exponential Increase

Implementation: # ACKs (in window)  
 $= \frac{Cw}{MSS}$



# senders  
PER ACK  
 $Cw += 1MSS$   
(total increase  
for  $\frac{Cw}{MSS}$  ACKs)  
 $\approx \frac{Cw \times MSS}{MSS} \approx Cw$

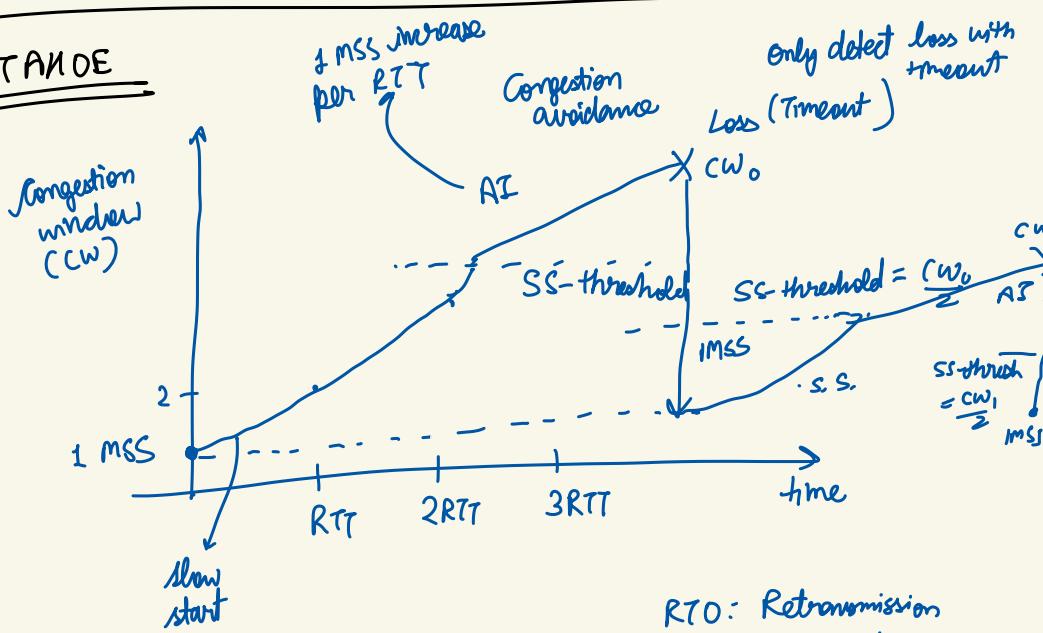




TCP TAHOE  
 ↓  
 RENO  
 +  
 VEGAS

CUBIC → latest

TCP TAHOE

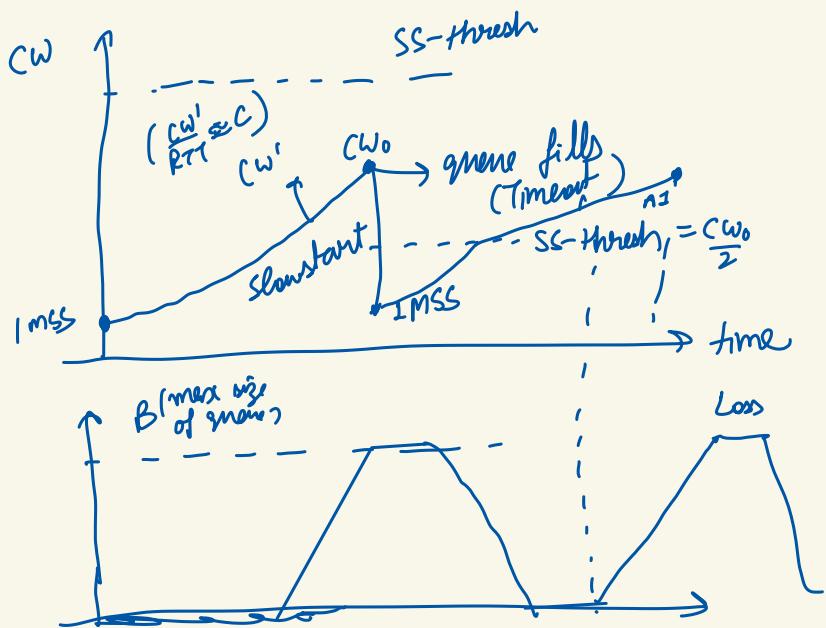
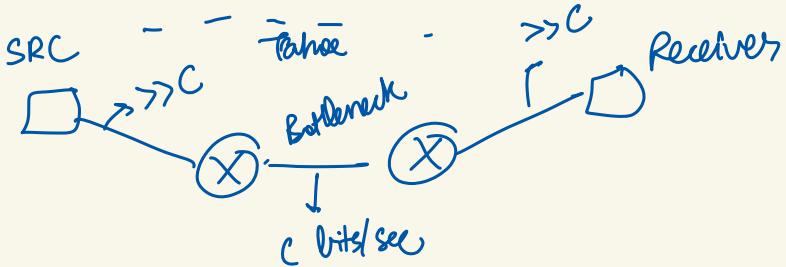


RTO: Retransmission Timeout

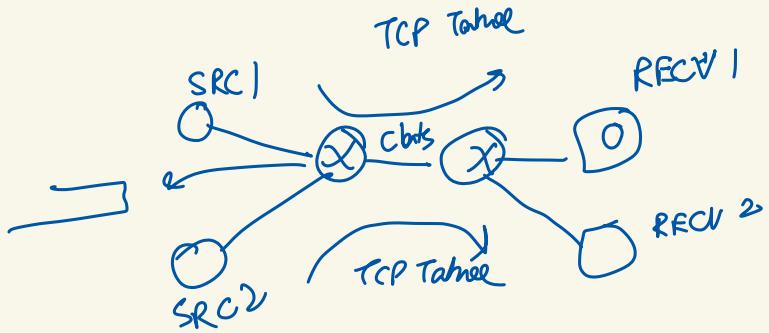
Tahoe Rules

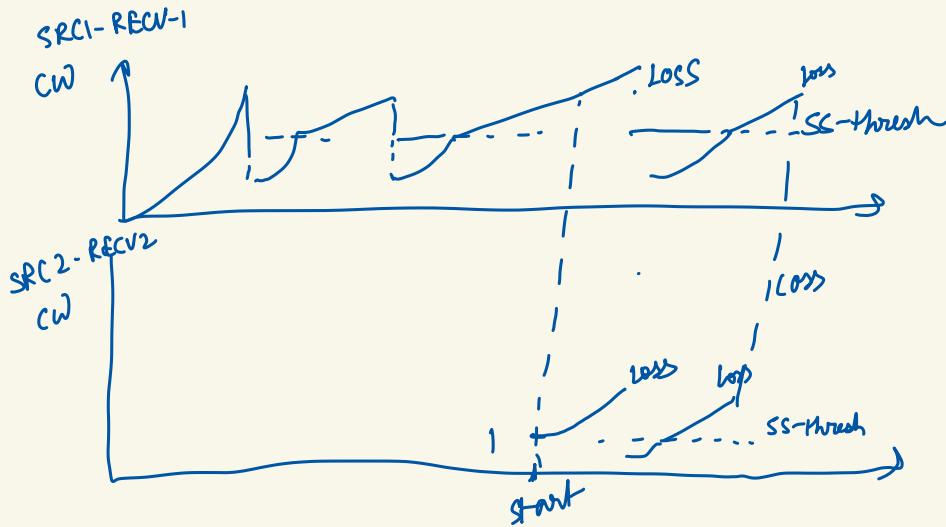
- $CW = 1 \text{ MSS}$  (Initialization)
- Loss  $\begin{cases} CW = \frac{SS-thresh}{2} & \\ CW + = 1 \text{ MSS} & \text{for each received ACK (Slow start)} \\ \downarrow \text{if } CW > SS-thresh & \end{cases}$

$CW += \frac{(MSS)^2}{CW}$  for ACK received (congestion avoidance)



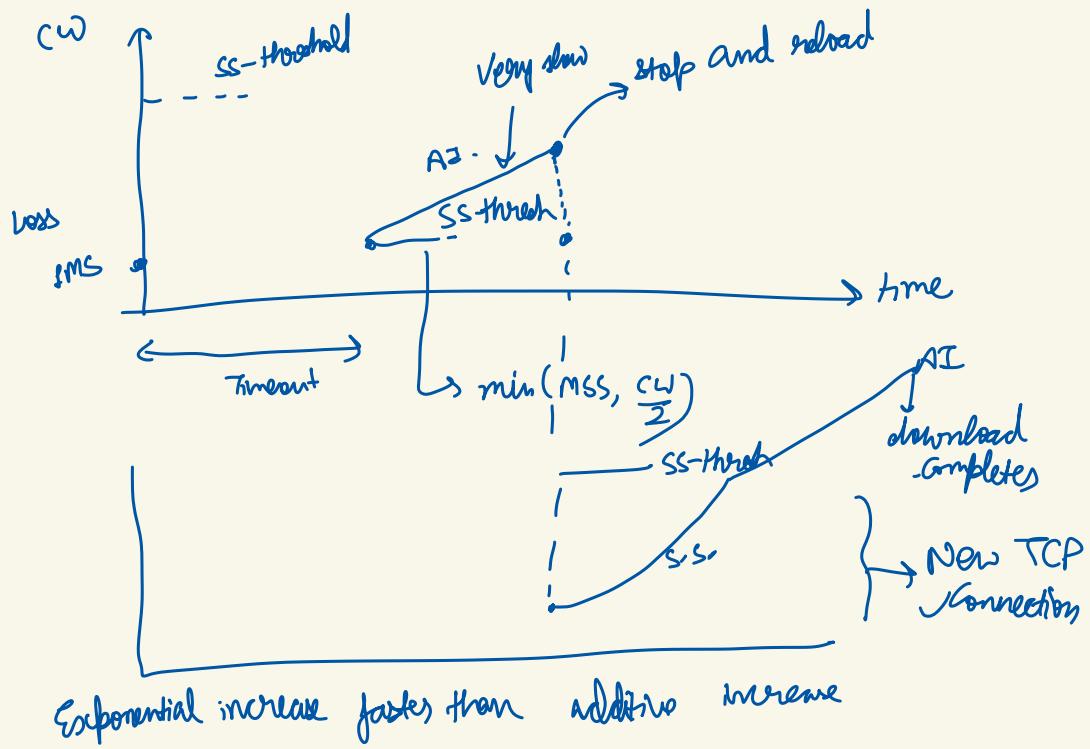
### Two TCP Tahoe Flows





Stop and Reload:

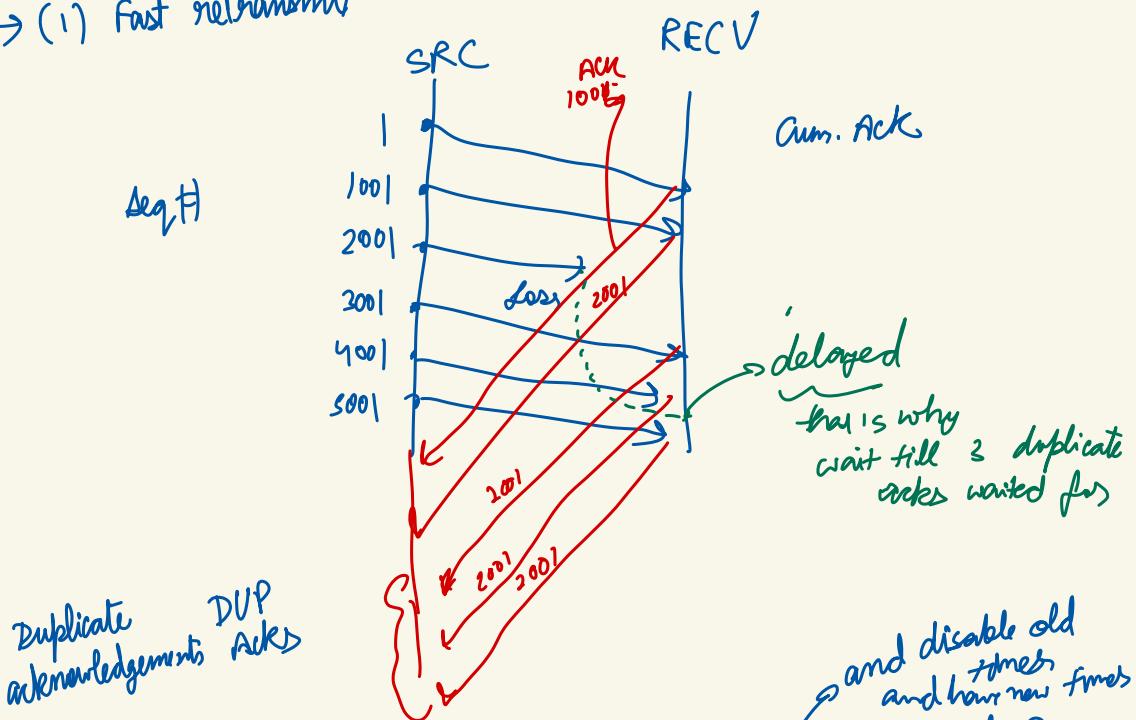
webpage download slow, hit stop and reload, downloads fast



- Question : (1) Should  $CW = 1MSS$  after loss  
 (2) Can we detect loss in other ways?

### Enhancements

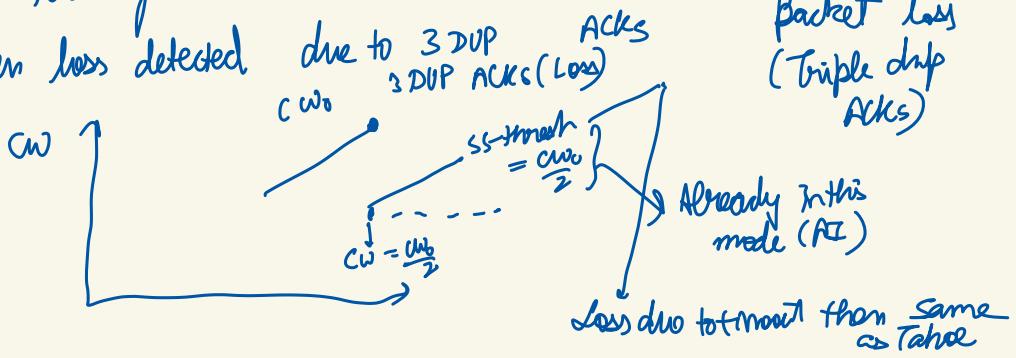
- (1) Fast retransmit



Fast Retransmit : Retransmit segments with SEQ # 2001 if 3 DUP ACKs received with SEQ # 2001

- (2) fast recovery

when loss detected



Loss due to timeout  $\rightarrow$  assumed worse than 3 DUP ACKs

TCP Reno: Includes Fast Retransmit, Fast Recovery.

### Miscellaneous

RFC 5681  $\rightarrow$  Congestion Control

RFC 6298  $\rightarrow$  timer

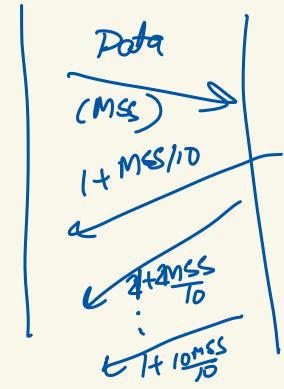
In slow start, suppose we receive an ACK which acknowledges  $N$  bytes of new data  
then  $CW + = \min(N, MSS)$

In congestion-avoidance  
( $CW \geq SS-thresh$ )

for every ACK which  
acknowledges New data

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

So this is done to  
prevent attack  
If we  
increase  
by  $MSS$   
each time, then  
we can Explode CW



RFC 6298: Initial RTO = 1 sec or higher

Loss detected by timeout  $\Rightarrow$  SS-thresh = max  $\left( \frac{\text{window}}{2}, 2MSS \right)$

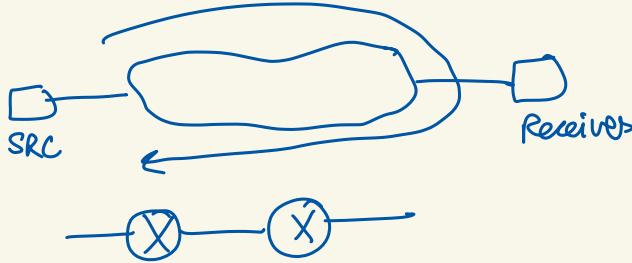
$$\Rightarrow RTO < \min(2 \times RTO, \text{max-RTO})$$

$\downarrow$   
max (cw, Advertised window)

(To avoid high data rate in  
case of congestion)

TCP Vegas → 1994-95

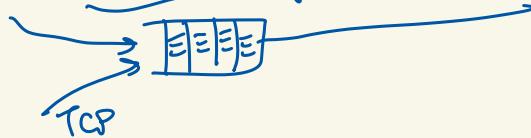
RTT detect Congestion



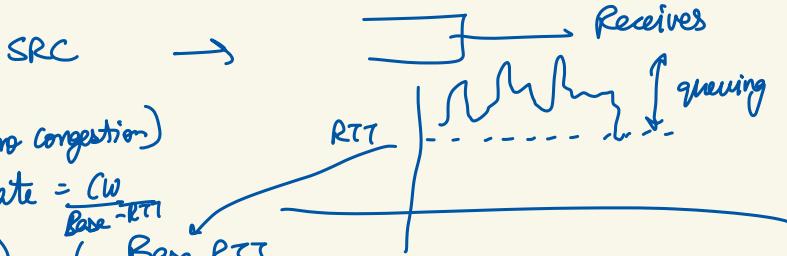
Detect loss by T.O.  
and 3 DUP ACKs  
and modify CW  
and SS-thresh as  
in TCP-Reno.

Does RTT for congestion help?

VoIP (over UDP) latency sensitive



Model path as a single queue



If queue empty (no congestion)

then expt rate =  $\frac{Cw}{Base-RTT}$

( $RTT = Base RTT$ )

$$\text{Actual rate} = \frac{Cw}{RTT} \leq \text{Expt rate} \quad (\text{since } RTT \geq Base-RTT)$$

↳ smooth out the measured RTTs in the recent past

$$\text{Diff} = \text{Expt. rate} - \text{Actual rate}$$

In congestion avoidance ( $CW > SS\text{-thresh}$ )

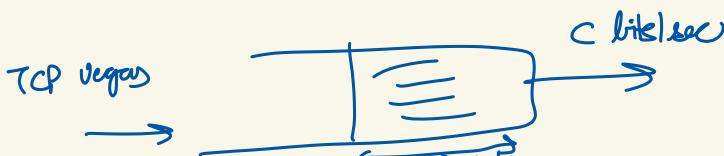
$$\text{Diff} < \alpha \Rightarrow \underset{\downarrow \text{parameters}}{CW} \rightarrow$$

is increased by 1MSS per RTT  
(eg:  $CW += \frac{(mss)^2}{CW}$ ) at every ACK for new data

$\beta < \text{Diff} \Rightarrow \text{CW decreased by 1 MSS per RTT}$

$\alpha < \text{Diff} < \beta \Rightarrow \text{CW is not modified.}$

$$\text{Diff} = \frac{\text{CW}}{\text{Base-RTT}} - \frac{\text{CW}}{\text{RTT}} = \text{CW} \left( \frac{\frac{\text{RTT} - \text{Base RTT}}{\text{Base RTT} \times \text{RTT}}}{\text{Base RTT} \times \text{RTT}} \right)$$



$$\text{Diff} = \text{CW} \left( \frac{B \text{ bits}}{\text{RTT} \times \text{Base-RTT}} \right) \xrightarrow{\text{queuing delay}}$$

$$\underbrace{\frac{\text{CW}}{\text{RTT}}}_{\text{Actual rate}} \approx C$$

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

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

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

$$\text{Base-RTT} = 100 \text{ ms}$$

$$\underbrace{\alpha < \text{Diff} < \beta}_{\text{targeting}}$$

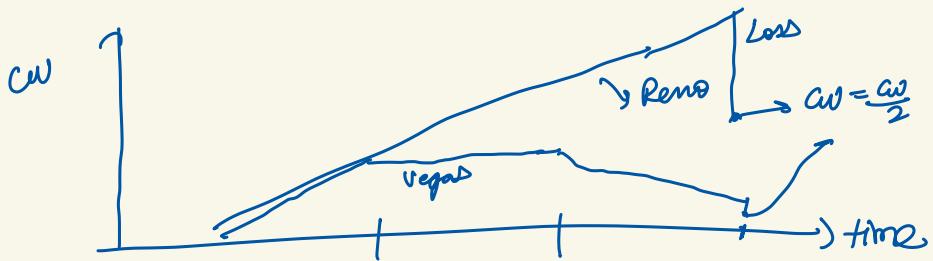
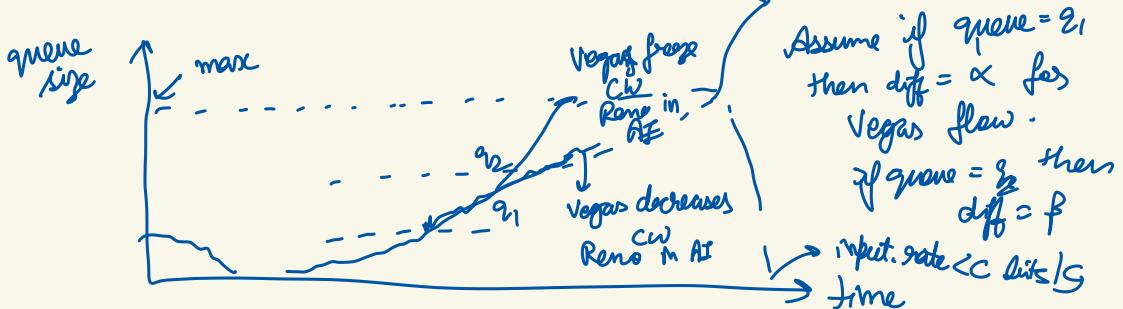
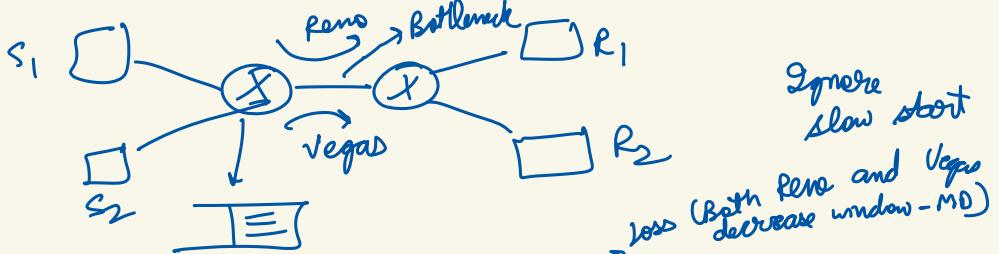
$$\alpha < \frac{B}{\text{Base-RTT}} < \beta$$

$$\alpha \cdot \text{Base-RTT} < \beta < \frac{\beta \cdot \text{Base-RTT}}{= 6 kB}$$

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

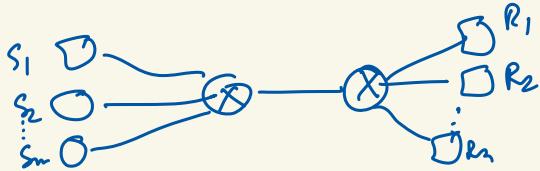
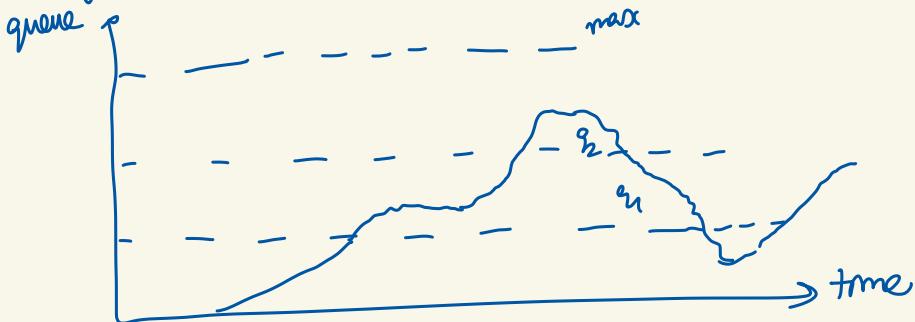
2) What if we replace all Reno flows in Internet with Vegas?

1) Reno vs Vegas



Reno is aggressive compared to Vegas

2) Only Vegas flows



$S_i \rightarrow R_i$  is a Vegas flow

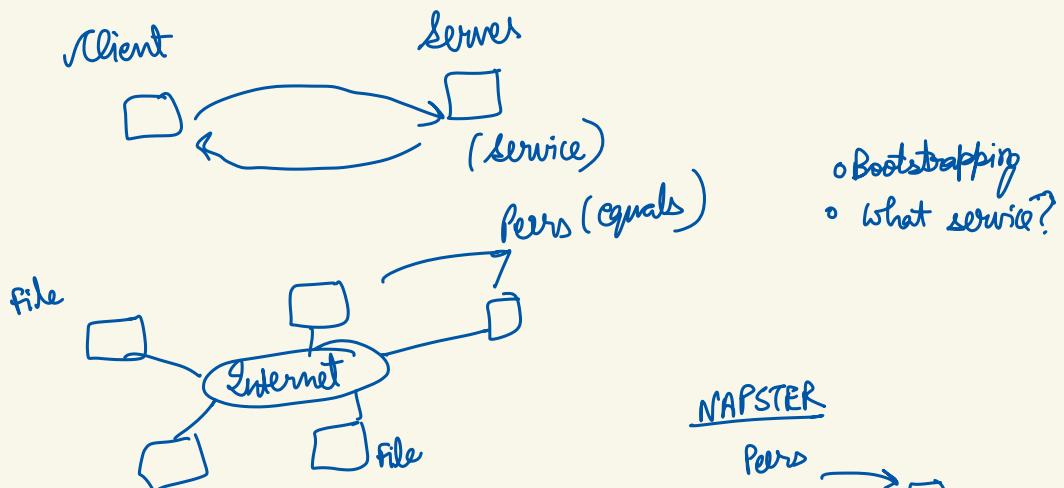
TCP - S tuple  
(SRC & DST IP, Port, Protocol)

Assume all have  
some RTT

↓  
IP header  
(says TCP is banjaxed)

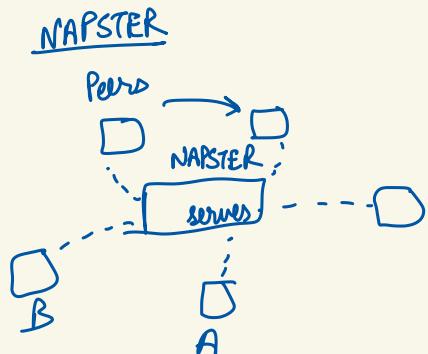
- (a) No packet-loss
  - (b) queuing delays between  $\frac{q_1}{C}$  and  $\frac{q_2}{C}$   
where  $C = \text{output rate of the queue}$
  - (c) Throughput can be higher than for the all-Reno case, because here the queues never empty unlike the all-Reno case
- Claim : all Vegas Case give 50% higher throughput than all Reno Case

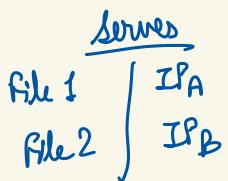
## P2P Peer-to-Peer Networks



1990's  
share digital content

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

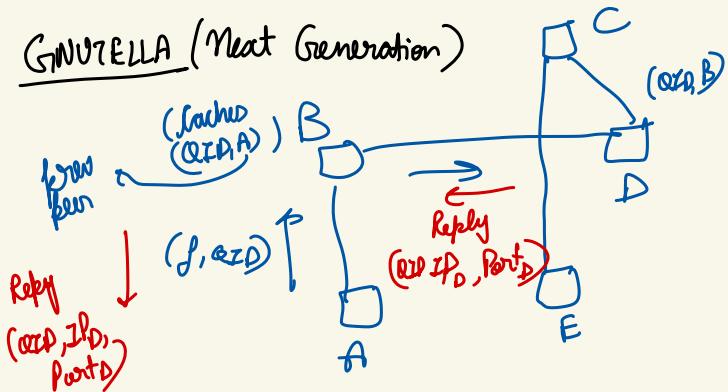




- C logs in
- search for song
- Match with file 1, ( $\text{file 1} \leftrightarrow \text{IP}_A$ ) given to C
- C can connect to A and get file

Drawbacks : Centralised  $\rightarrow$  single point of failure (not fault tolerant)  
 $\rightarrow$  legally easy to take down

### Gnutella (Peer Generation)



### Bootstrapping

- (1) Application may have IPs of some other peers
- (2) Lookups IPs from one or more websites

Suppose A wants to search for "f"

- Limited Broadcast ( $\text{TTL} = n$ ,  $\text{TTL} --$  at every hop  
 $\text{TTL} = 0 \Rightarrow$  don't forward query)

Suppose  $n=2$

First method: Query had "f", QID (Unique ID for query)  
 Reply over the path which the query traveled.

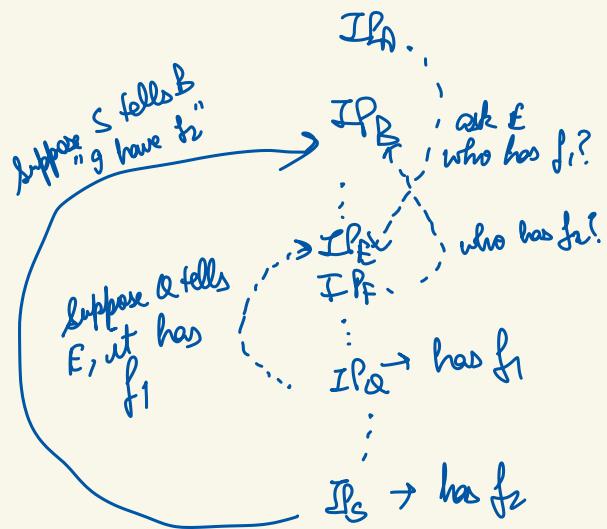
Second Method: Query has IP<sub>A</sub> as well  
 Reply directly to A

$$V0.4 \rightarrow \text{max TTL} = 7$$

$$V0.6 \rightarrow \text{max TTL} = 4$$

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

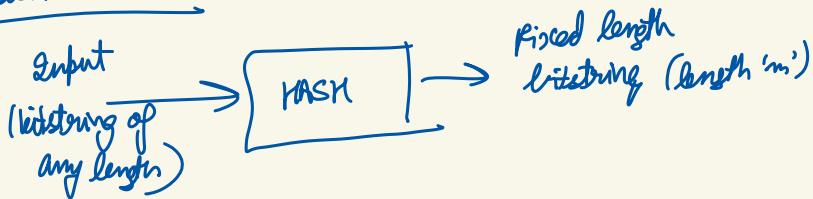
### IP Address



### Files

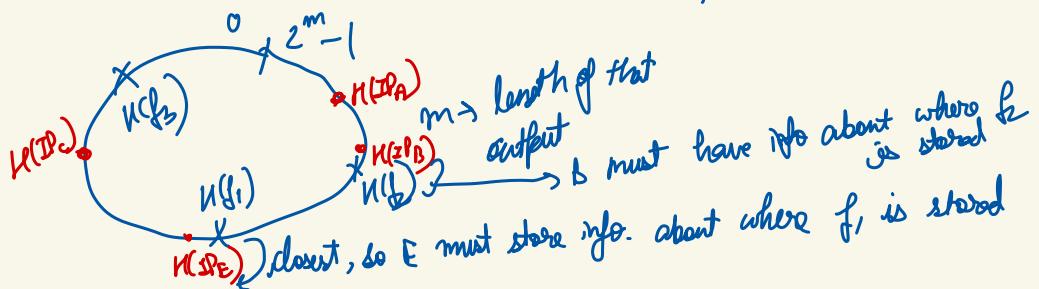
$f_1 \rightarrow$  mapped to  $IP_E$   
 $f_2 \rightarrow$  mapped to  $IP_B$   
 $\vdots$

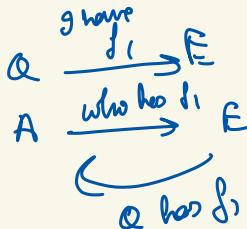
### Use Hash function



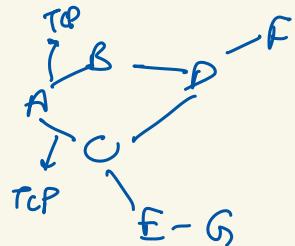
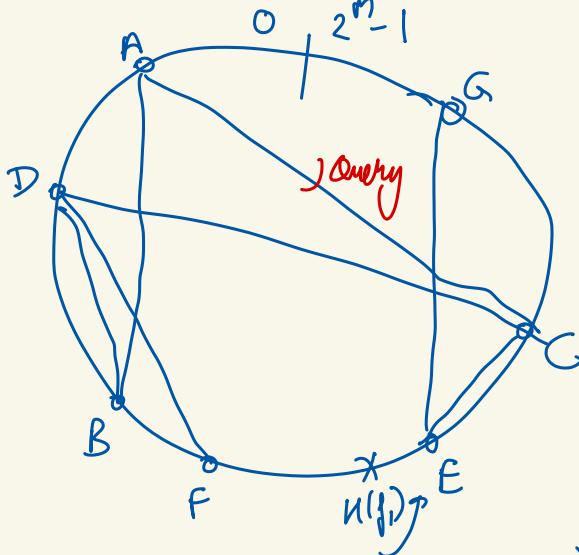
$$\begin{aligned} f_1 &\rightarrow [H] \rightarrow H(f_1) \\ f_2 &\rightarrow [H] \rightarrow H(f_2) \end{aligned}$$

Some IPs might get loaded because of this





Want to store info about where an object is stored at node whose node ID is closest to objID of the object.



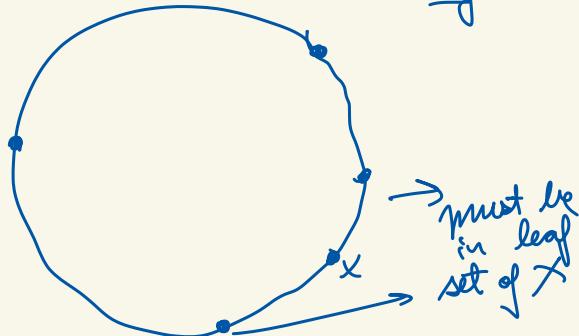
Idea: Route query messages close to  $H(f_i)$  in this virtual space till we reach 'E'.

- Suppose A wants to find out where  $f_i$  is stored
- A finds which of its  $\underbrace{\text{Connected peers}}_{\text{leaf set of } A}$  is closer to  $E$
- This repeats till the query reaches a node who has no connected peer closer to  $H(f_i)$

## PASTRY

In leaf node of  $X$ , include  $\frac{L}{2}$  closest nodes on either side of  $X$  in the virtual space.

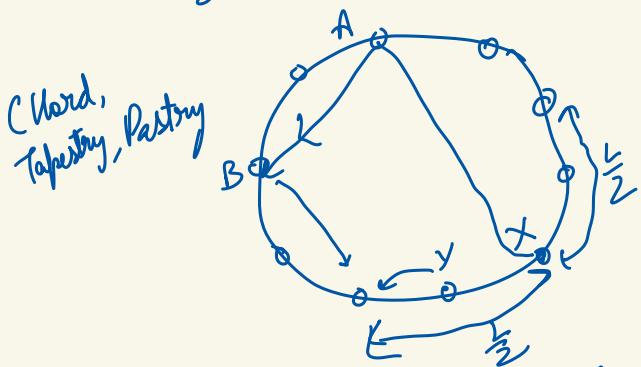
Eg  $L=2$



This ensures that the routing protocol gets to node closest to  $H(f_i)$ .

PASTRY ensures that leaf set of any node is of size  $L + O(\log N)$   
 $\# \text{peers in network}$

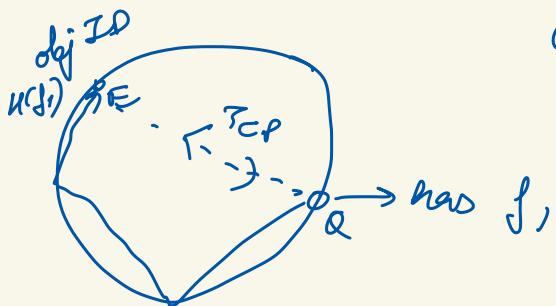
Suppose P2P network already has the property that leaf set of all contains  $\frac{L}{2}$  closest nodes on either side of the node.



( $L=4$ )

New node  $X$  routes a message to find mode whose node IP is closer to own IP

From  $Y$ ,  $X$  determines  $\frac{L}{2}$  nodes who are closer to it on left and right.

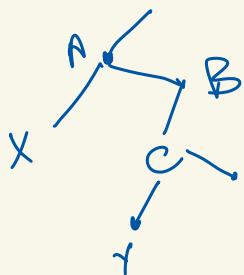


Query has Q's IP address

- E contacts Q
- Q tells E that it has  $f_i$

similar procedure to find where  $f_i$  is stored.

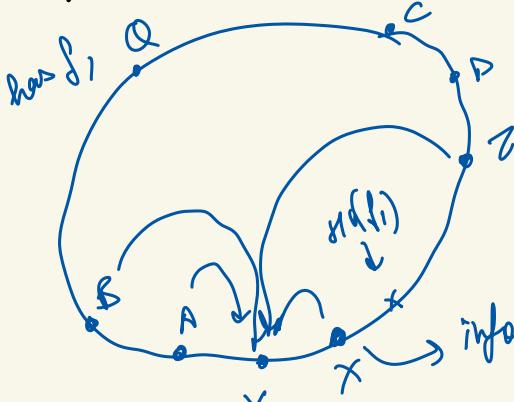
DHT  $\rightarrow$  Distributed Hash Table



X can add A, B, C, Y to its leaf set

$O(\log n)$   
this ensures that  
leaf set is of  
size  
 $O(\log n) + L$

Q: what if a node fails?



info about where  $f_i$  is stored

Y stores information which is present in 2 nodes on either side.  
All peers send keep alive messages to peers they are connected to (leaf set).

(- if any node 2 nodes to right or left fails then it finds out about nodes further away and ensures

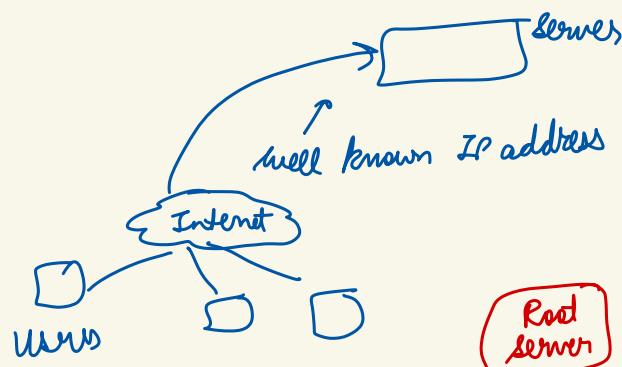
## Domain Name System

Originally → IP Addresses

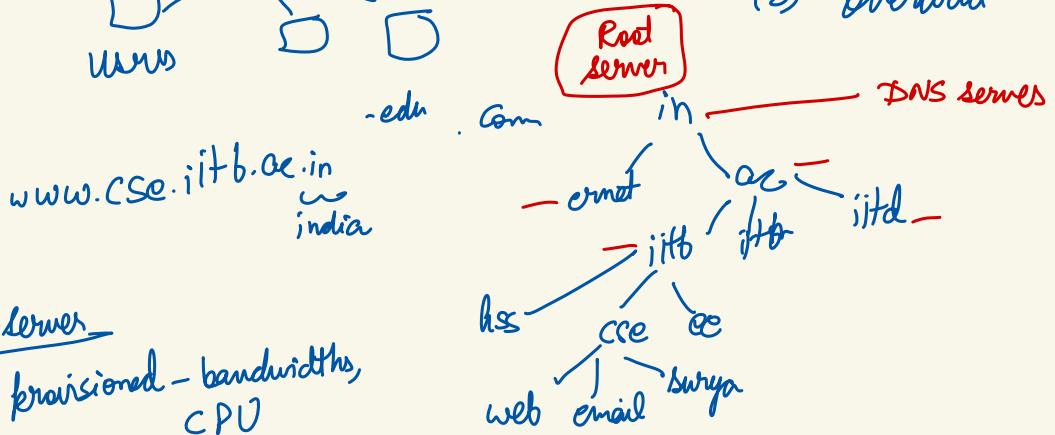
121.5.7.32

www.google.com  
www.cse.iitb.ac.in } User friendly → DNS → IP addresses  
has to be robust

Idea :-



- (1) single point of failure
- (2) overload



## Root server

well provisioned - bandwidths, CPU

## DOS - Denial of service

- Redundancy

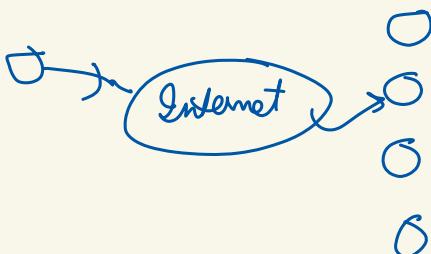
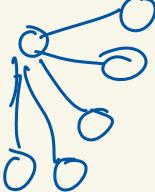
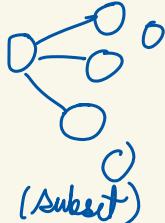
13 Root servers

A. root-servers.net  
B. root-servers.net  
C. root-servers.net

DNS packets size of 512 bytes  
has space for 13 addresses at most

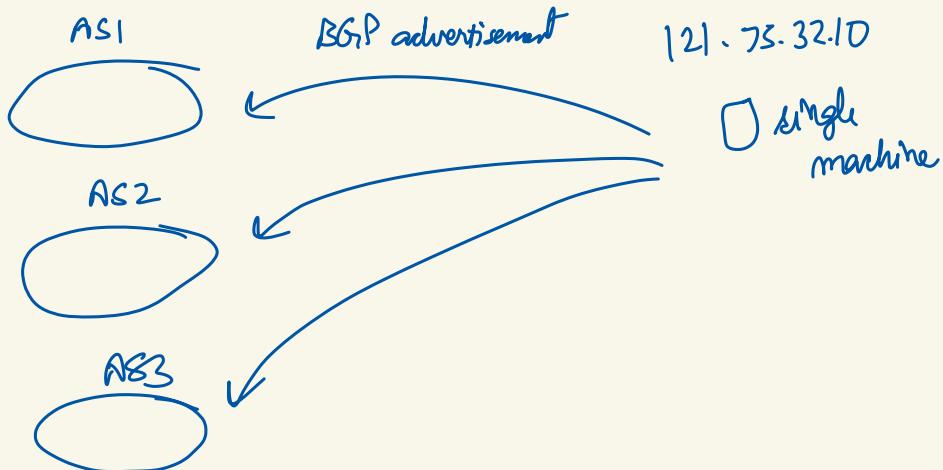
- Uses anycast for more redundancy

unicast, multicast, broadcast



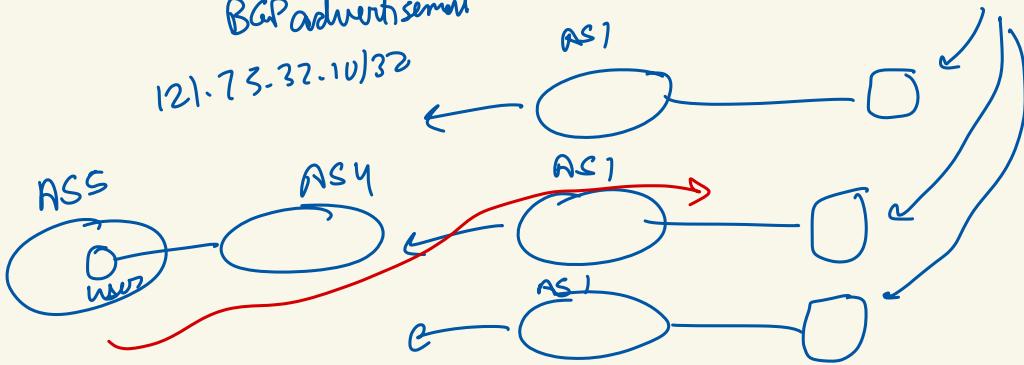
pkt reaches any one server

A. root-servers.net → mapped to many physical machines geographically distributed



121.75.32.10

BGP advertisement  
121.73.32.10/32



## Resource Record

$\langle \text{Name}, \text{Value}, \text{Type}, \text{Class}, \text{TTL} \rangle$   
says how  
value should  
be interpreted

Type

A

NS

CNAME

MX

Value

IP address

Name server

(host running the  
DNS service  
corresponding to "Name")

(alias of "Name"; canonical  
name of the host specified in  
"Name")

(name of host running  
mail server  
specified in "Name")

always for the  
internet

Linux : dig www-google.com  
↳ number decreasing

### Examples

Root server has  
(edu, a3.nstld.com, NS)  
↓  
Name of DNS  
servers of "edu" domain  
(a3.nstld.com, 192.5.6.32, A)  
↓  
IP of "Name"

Server a3.nstld.com has  
(princeton.edu, dns.princeton.edu, NS)  
(dns.princeton.edu, 128.112.185.5, A)

Server dns.princeton.edu  
(dns.cs.princeton.edu ...)

Server dns.cs.princeton.edu  
(penguin.cs.princeton.edu, 128.112.15.6, A)

(www.cs.princeton.edu, catwods.cs.princeton.edu, CNAME)

(cs.princeton.edu, mail.cs.princeton.edu, MX)

(mail.cs... , 128.12.1.5, A)

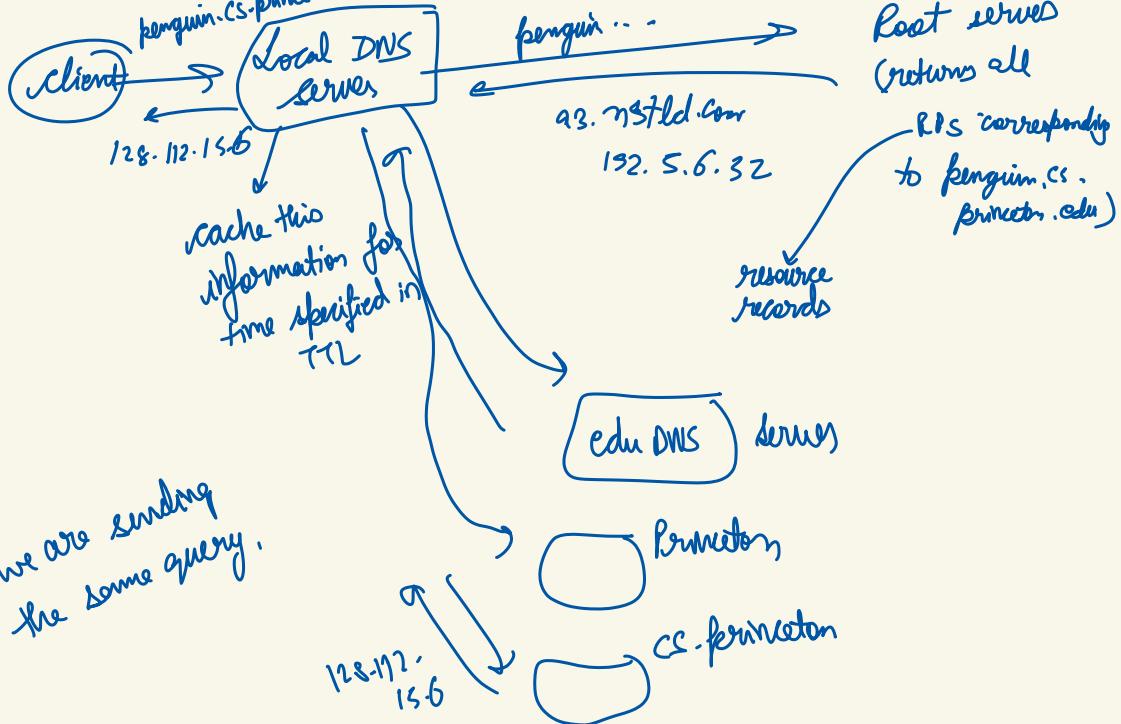
Look up:

DNS runs on UDP (port 53 for servers)

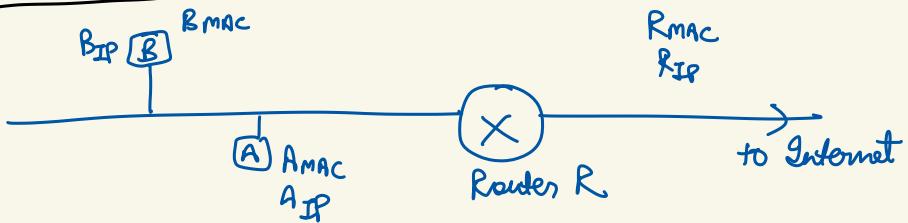
Manually config

or DHCP get local DNS servers

DNS query

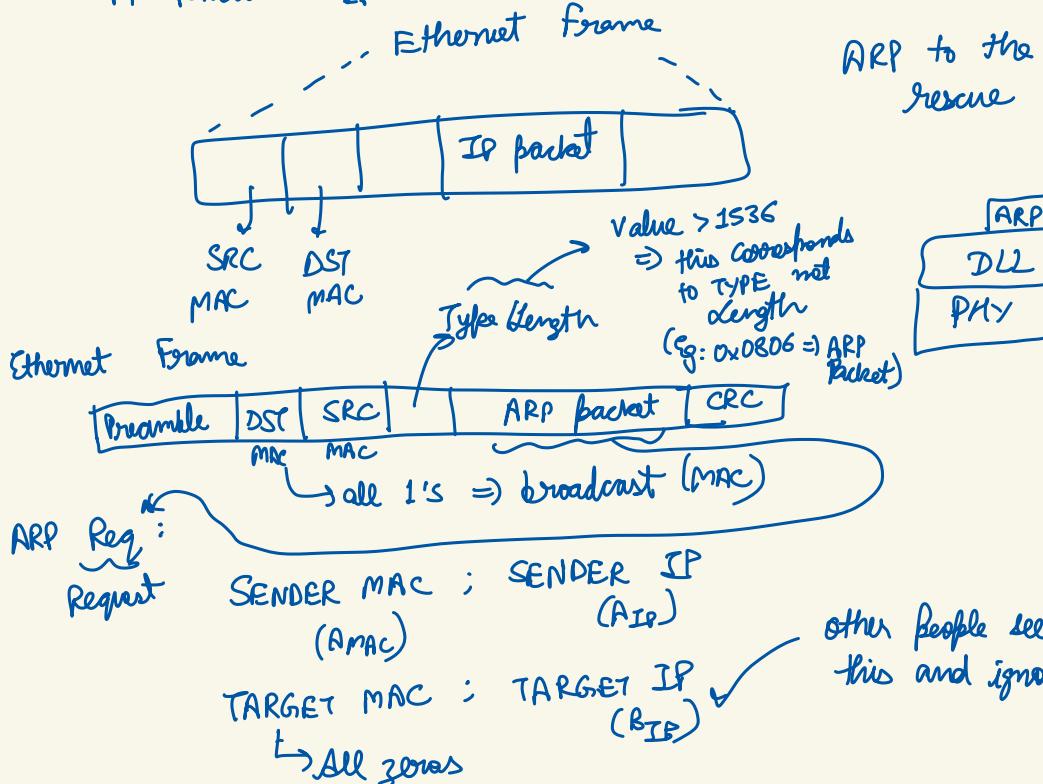


## ARP: Address Resolution Protocol



A wants to send an IP packet to B

A knows B IP; A does not know B MAC



ARP Reply:  
(from B to A)

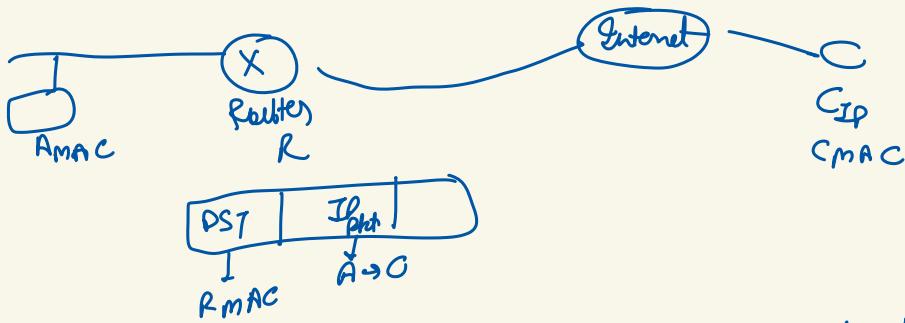
↓  
Unicast  
Ethernet frame

SENDER MAC : SENDER IP  
(BMAC)  
(BIP)

TARGET MAC ; TARGET IP  
(AMAC)  
(AIP)

A → Preamble DST SRC (B) ARP Reply [CRC]

Information about RMAC stored in ARP cache at A with a timeout.



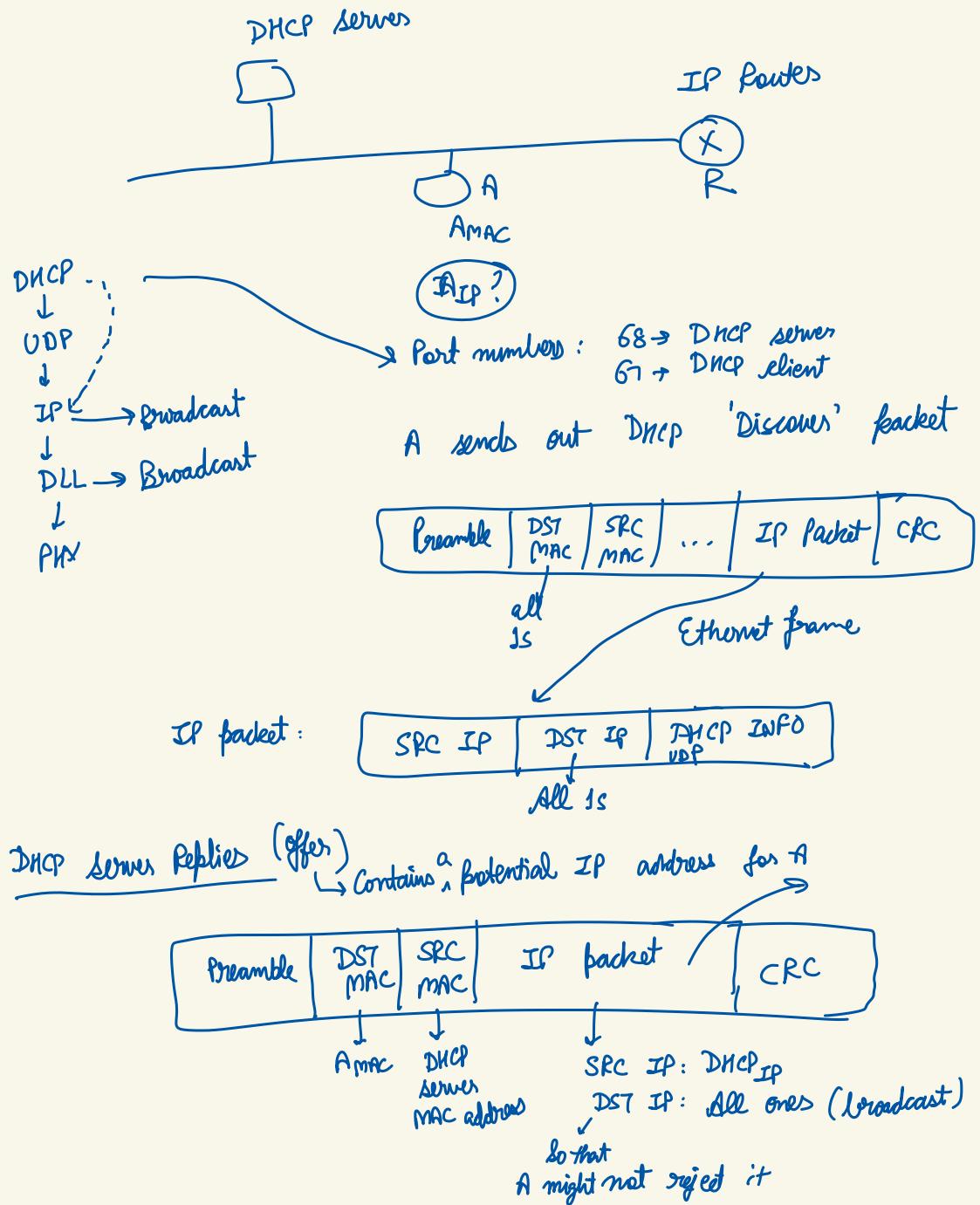
- (1) How does A know if DST-IP belongs to own network or not?  
(2) If DST-IP is not in own network, how to know RIP, RMAC

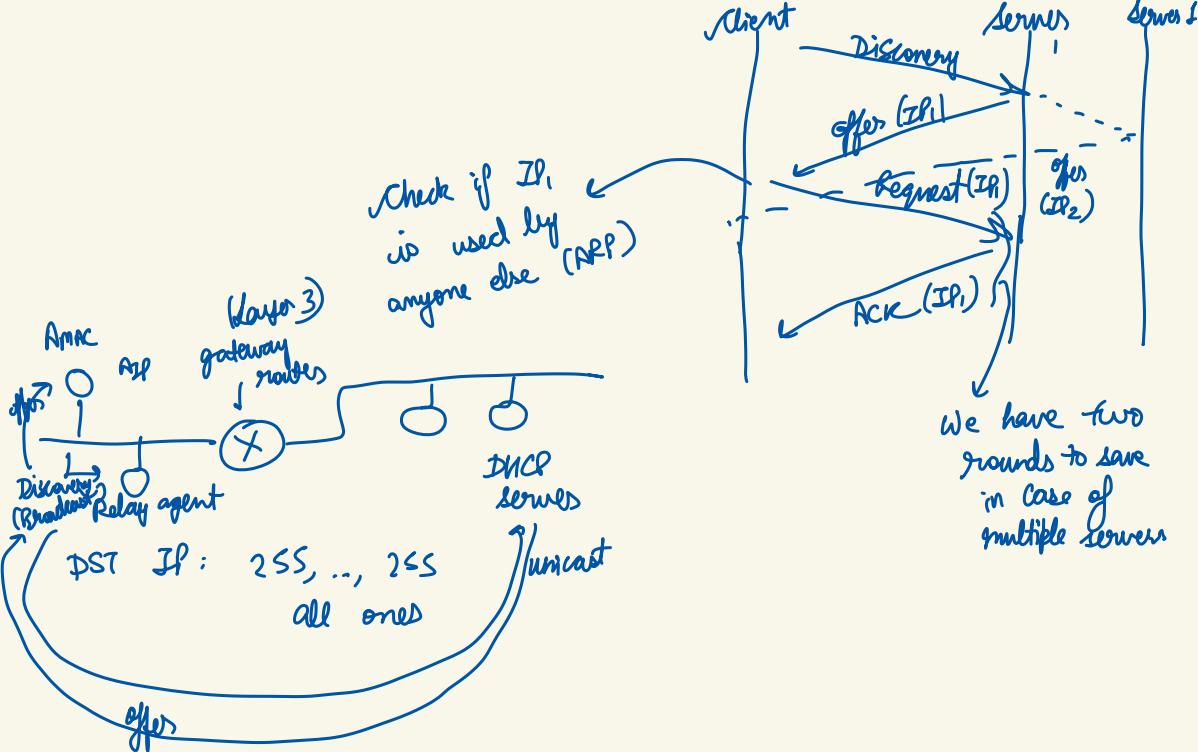
(1) AIP :  $a_1 \cdot a_2 \cdot a_3 \cdot a_n$   
subnet mask : Enc:  $\underbrace{255.255.255.0}_{\text{all 1s}}$

go (DST IP) AND (MASK) = (AIP) AND (MASK)  
then DST IP is in my network

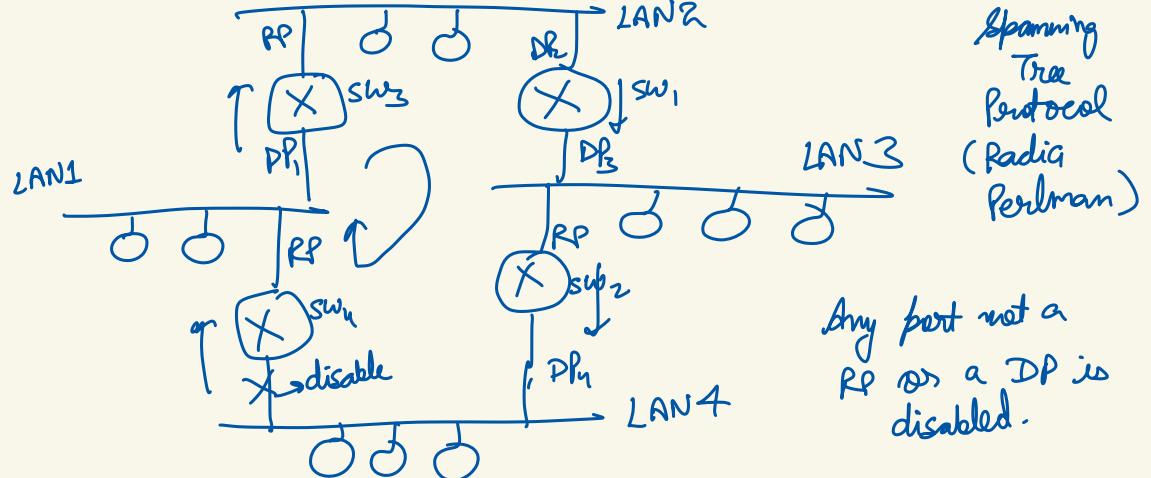
(2) Suppose know RIP, how to find RMAC  
Use ARP

# DHCP : Dynamic Host Configuration Protocol





\* Ethernet switches are sometimes called bridges.



### SPT Protocol

1. Elect root bridge
2. Each bridge finds which port is closest to root and assigns this port as root port.  
(Tie breaking rule)
3. All bridges connected to a LAN, elect one among them to forward frames on that LAN (Designated port)

### Details:

1) Elect root

Bridge ID :

↓  
lowest becomes root.

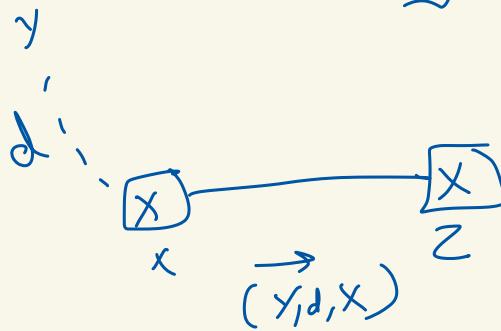
def : 32768  
 $0 - 61440$  (multiples of 4096)  
Configurable port (2 bytes)  
MAC Address  
6 bytes  
(smallest MAC of all ports)

Each bridge tells its neighbours  
(y, d, x) → distance to my ID  
the smallest ID heard till now

SW; → ID;  
SW1 → (1, 0, 1)  
SW4 → (4, 0, 4)

$SW_2 : (1, 1, 2)$   
 $\vdots$   
 $SW_4 : (2, 1, 4)$

after hearing  
 $SW_2, SW_3$

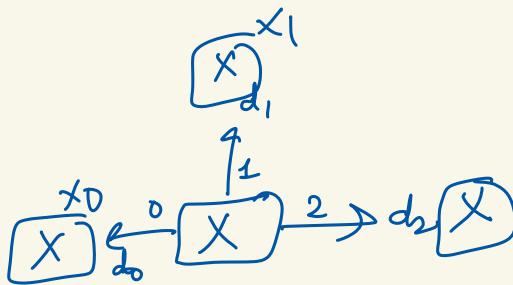


$(y_2, d_2, z)$

If  $y < y_2$  then  
 $y_2 = y$  and  $d_2 = d$   
 $+ \text{dist}(x, z)$

If  $y = y_2$  but  $d + \text{dist}(x, z) < d_2$   
then  $d_2 = d + \text{dist}(x, z)$

2) Root Port

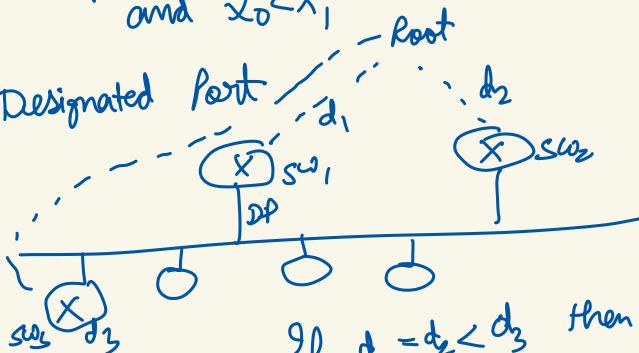


If  $d_0 = d_1 < d_2$

and  $x_0 < x_1$

if more than one port has the smallest distance to root, then tie-break based on IP (smallest) of neighbours on ports

3) Designated Port



$d_1 < d_2$   
 $d_1 < d_3$

If  $d_1 = d_2 < d_3$  then tie break  
based on IP of  $SW_1, SW_2$

Ports not RP or DP are disabled