

---

---

---

---

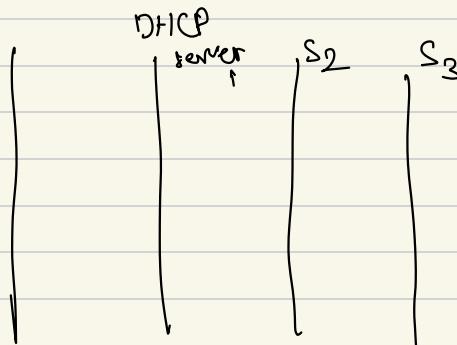
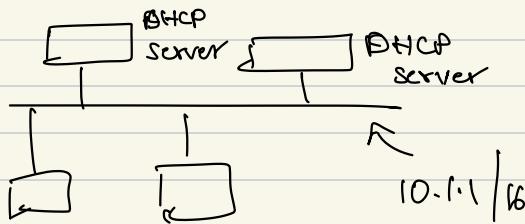
---



# 8/10/2024 CS348 - Computer Networks :-

ARP :-

DHCP :-



UDP      TCP      → If you care about latency      → If you care guarantee than if it stays for 10sec/a minute.

APPL  
layer

DHCP (Port = 68 for client A) (Port = 67 for DHCP Server)

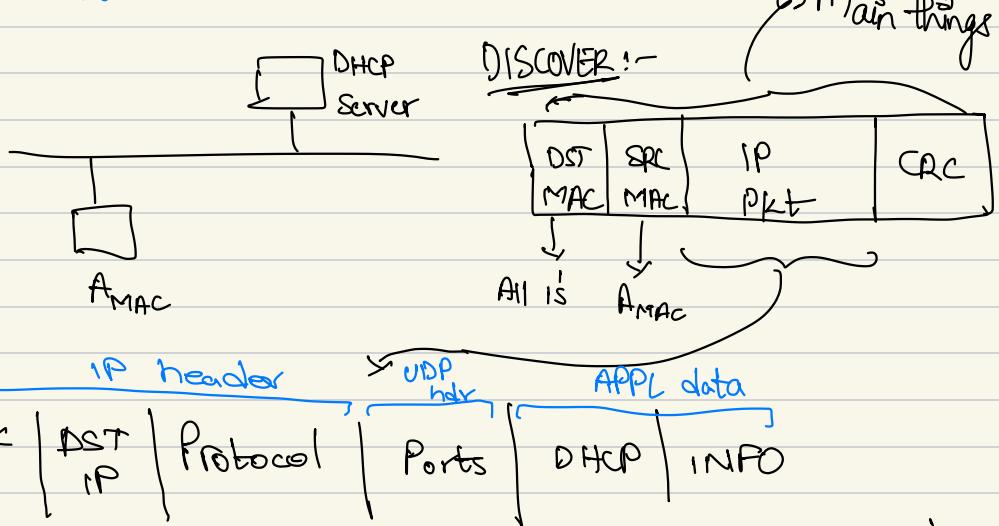
UDP      TCP

IP  
DLL  
PHY

when a pkt received, go up from DLL, IP. In IP field will have UDP / TCP



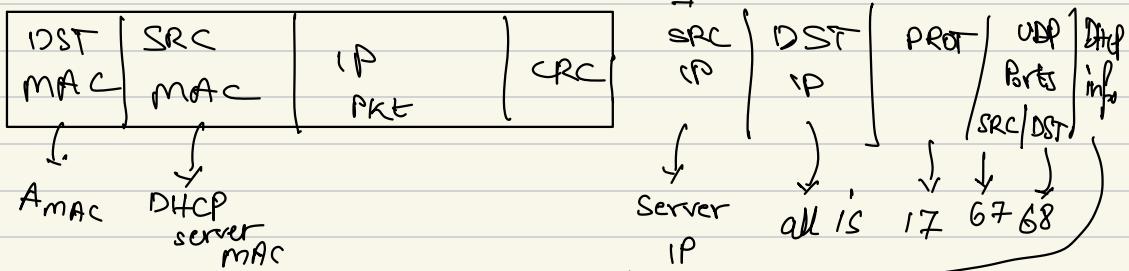
\*



(didn't listen explanation/default values of fields)

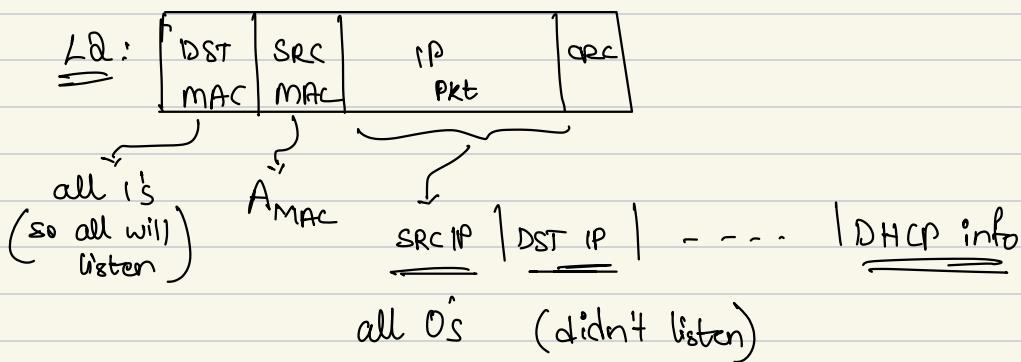
\*

OFFER!-

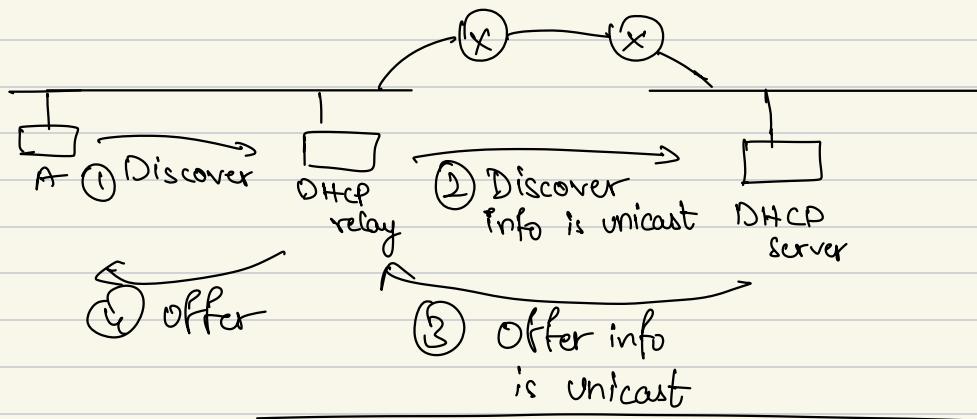


offered IP, server IP (which server is offering)

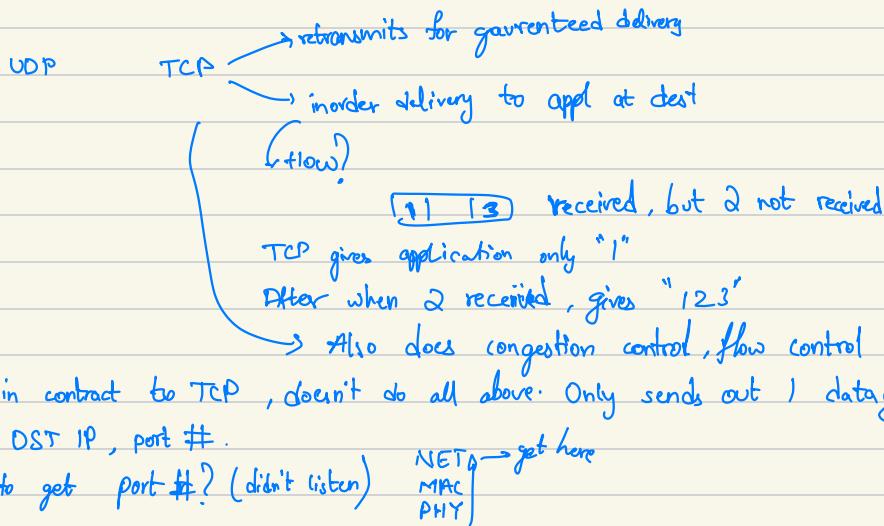
\* Request ( $A \rightarrow$  sent to all)



\* Relay: sort of a proxy



(1) relay CS 268 :- (didn't listen 30 min)  
\* APP1 APP2 APP3 ... -



### When to use UDP? :-

(1) If msg fits in 1 pkt



\*

Like TCP in your App written above UDP  
(But TCP is a tricky & large code. You'll not want to write it in your app! -)

Like TCP in your App written above UDP  
(But TCP is a tricky & large code. You'll not want to write it in your app! -)

\* UDP header :- (didn't listen)

\* TCP header :-

0	16	82
SRC port	DST port	
sequence numbers.		
Acknowledgment number		
FLAG	Advertisement window	
Check sum	Urgent PTR	
DATA		



How to decide 1,3 come  
but 2 didn't

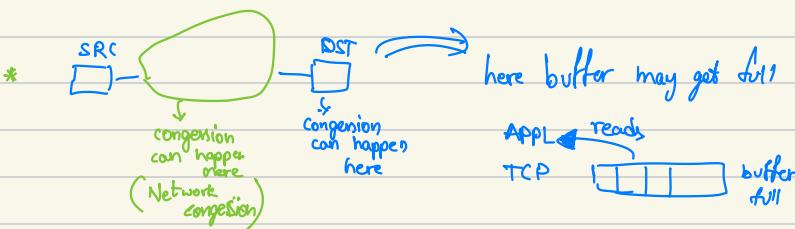
\* TCP is 2 way communication



\* FLAGS bits:

SYN	FIN	RESET	PUSH	URG	ACK
BIT	BIT				BIT

( { } )  
Sent to setup connection      Used to close connection

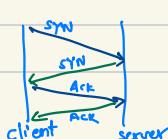


Then DST can send explicit message that I'm running out of buffer  
↳ Advertisement window used

14/10/2024

clam :- (didn't listen 30min)

\* 4 way handshake :- (little bit diff from 3-way handshake)



How to identify a particular TCP connection?

Ans:- IP : SRC IP, DST IP  
 Ports : SRC Port, DST ports

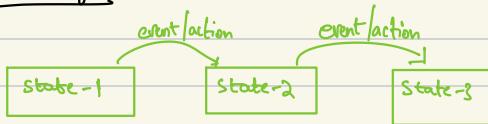
Though TCP is a 2-way connection, for now let us call 1 SRC & other DST

There are 16 UDP also.

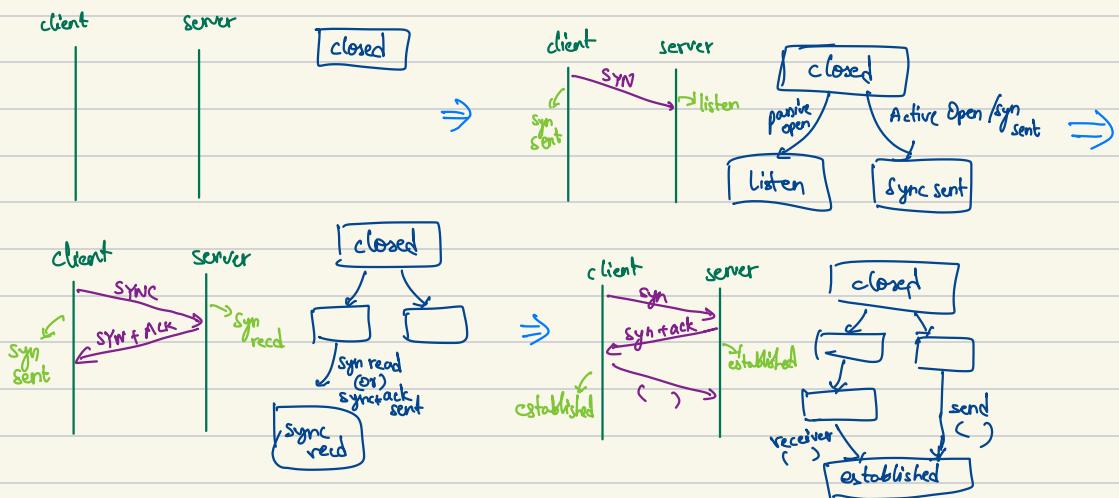
These aren't  
 didn't listen.

\* (said why we need state diagram)

State diagram :-



→ Now consider state diagram for establishing a 3-way handshake

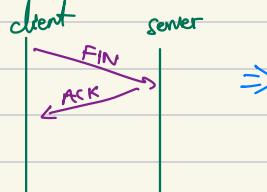


\* Connection termination :-

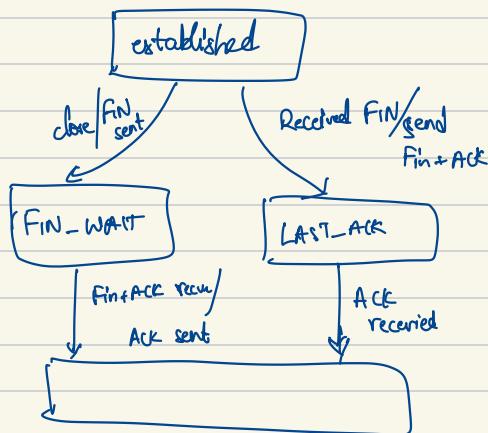
1<sup>st</sup> way :- 3-way handshake



2<sup>nd</sup> way :- only 1 terminate at a time



State diagram :- (1<sup>st</sup> way)



Problems with this:-

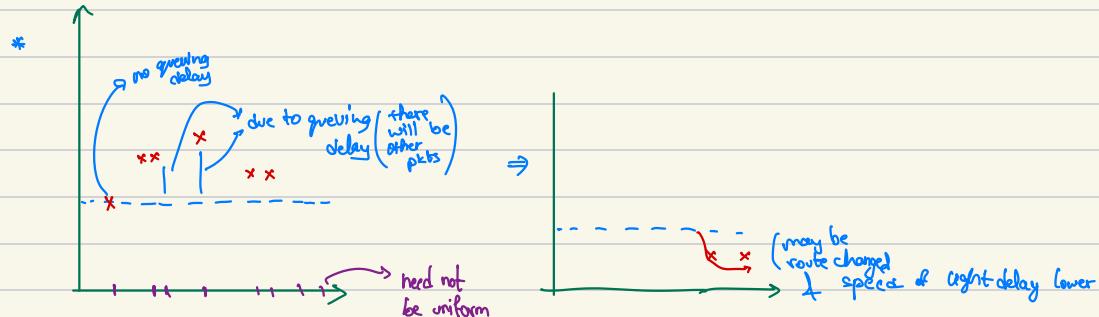
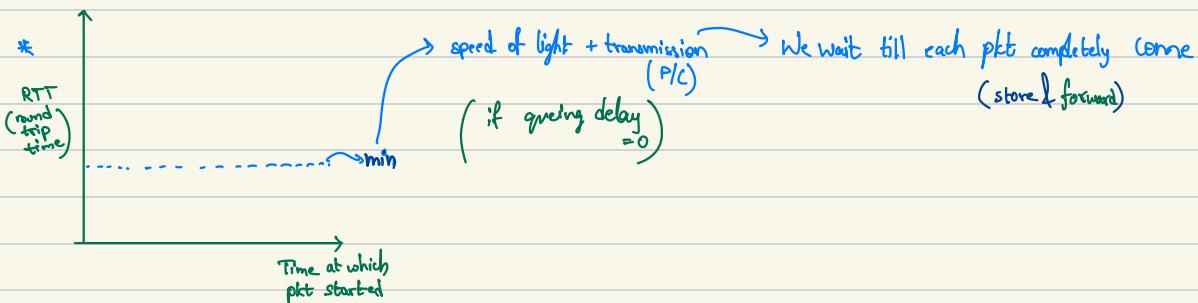
- \* What if (Fin + ACK) not received by client?  
(i.e., lost)  
client stays in FIN-LAST
- \* What if ACK not received by server?  
server stays in LAST-ACK

Soln:-

\* Time wait (didn't listen)

Nxt class congestion control.

15/10/04 clam- :- (missed 20 min)



\* Plot histogram of . . . (didn't listen)

\* Random numbers:  $x_1, x_2, \dots, x_N$ ,  $M = \frac{1}{N} \sum_{i=1}^N x_i$ , Estimate of std. dev. =  $\sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - M)^2}$

But there are 1000's of pkts every sec & many TCP ports data going. Everytime you receive a pkt, if you want to square, it is a overhead.

\* Used Mean deviation instead of standard deviation

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

\* Algo used by TCP :- (for RTO)

- latest RTT estimate : Sample RTT
- Current estimate of mean estimate : Estimate RTT
- Update step:-

$$\text{Estimate RTT} = (1-\alpha) \text{Estimate RTT} + \alpha \text{Sample RTT} \quad (\text{Exponentially moving average})$$

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

- difference : Sample RTT - Estimate RTT  
Like  $x_i - M$
- deviation =  $(1-f)$  deviation +  $f$  [difference]

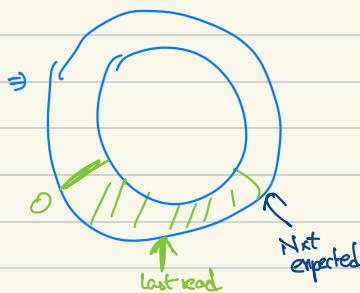
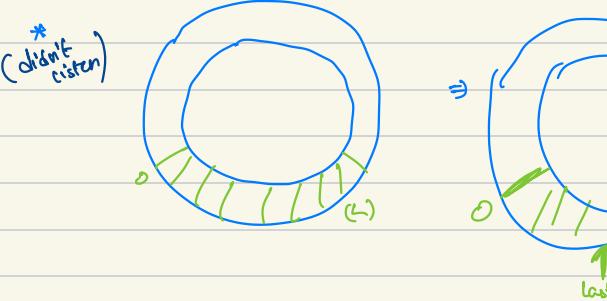
$$\bullet \quad \text{Timeout} = M + \phi$$

$$\mu = f \quad \phi = u \quad \alpha = \frac{1}{f} \quad \beta =$$

(TCP vega,  
didn't give values.  
Hence wasn't accepted)

Congestion & flow control :-

• TCP have to make a guess whether there is a congestion in network or not.



\* Not clear  
Windows = min  
of max bytes  
in added src  
can send out

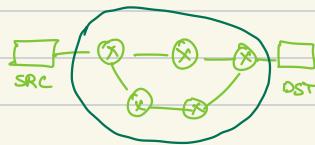
(Congestion window, Adv window)  
depends on router  
Depends on network congestion

17/10/2024

class - 32 :-

\* We've learnt RTT estimation

TCP congestion control :



(didn't listen 5 min)

The middle nodes do not have power. They just forward pkt (so no one wanted that) (More efficient possible if this is not case)

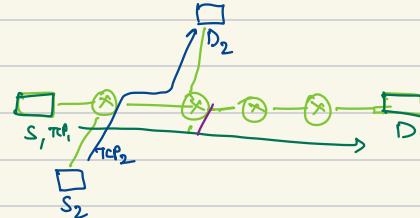
Power at Ends → TCP

End hosts running TCP :-

- They don't know link speeds
- They don't know link utilizations
- They don't --- (didn't listen)

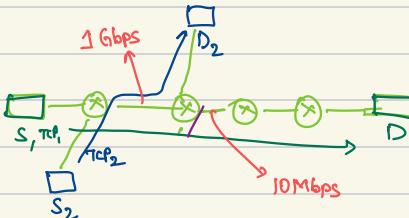
Here it is not rule that both should get same

$$\left\{ \begin{array}{l} \text{TCP}_1 \rightarrow 10 \text{ Mbps} \\ \text{TCP}_2 \rightarrow 990 \text{ Mbps} \end{array} \right.$$



Here to be fair, we divide equally

$$\left\{ \begin{array}{l} \text{TCP}_1 \rightarrow 5 \text{ Mbps} \\ \text{TCP}_3 \rightarrow 5 \text{ Mbps} \end{array} \right.$$



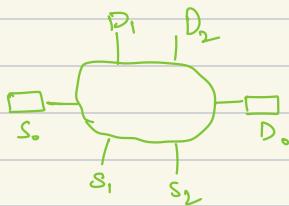
Congestion control issues :-

TCP connection

- wants to use bandwidth resources efficiently
- do not want to cause congestion (pkt loss, Queues filling up)
- fairness

→ one TCP connection should not have most of bandwidth at expense of other connection (A very high level idea. Actual fairness is much complicated to define)

These are our requirements. How do we design such protocol?  
It's non-trivial, one sees feels .....



a) How to set data rate of TCP?

Idea 1:-

1 sec      1 sec

• Each sec, send  $10^3$  pkts,  
each of size  $10^4$  bits

10 Mbps

|||||  
1 sec      1 sec      1 sec

what if in 1 sec, we send everything in first  
0.01 sec  $\Rightarrow$  we get 1 Gbps & be calm for  
0.99 sec. Congestion still takes place bcz of  
sudden burst into network.

\* so maybe 10ms is better window than 1sec

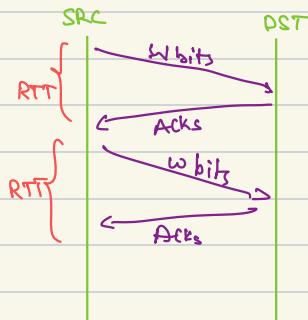
Idea 2:-

window based       $W = \text{max amount of un-acked data in flight of } \dots$   
data rate control :



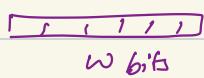
how does this do "data rate control"? There is no kind of timer as  
there before.

$\Rightarrow$



$$\text{Data rate} = \frac{w}{RTT} \rightarrow \text{can be } \uparrow \text{ or } \downarrow \text{ to vary data rate}$$

This is what TCP does.



$W$  can also affect congestion at queues also. If  $W \downarrow$ , it is expected to have less traffic at each queue

\* Algo to decide  $W$ :

$$W = \min(\text{Congestion Window}, \text{Adv window})$$

(Cw)

left out ----, can be sent to SRC via ACK

↓

This is the one which require algo

\*

Initial value of  $w = ?$

SRC

DST

Let  $10^3 \text{ pkts} = 10^7 \text{ bits}$  ( $1 \text{ pkt} = 10^4 \text{ bits}$ )

& RTT be 1ms

$$\text{data rate} = \frac{10^7 \text{ bits}}{10^{-3} \text{ sec}} = 10 \text{ Gbps}$$

Will definitely congest the network

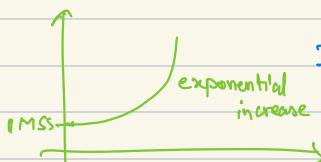
\* We think conservatively & initially send

$$W = 1 \text{ MSS} \text{ (maximum segment size)}$$

i.e., 1 pkt  $\sim 10^4 \text{ bits}$

} sending 1 pkt  
will definitely not congest network

Slow start: Start with  $W = 1 \text{ MSS}$  (since initially not utilising resource)



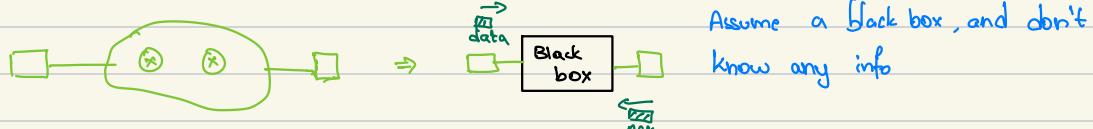
If cause congestion then slowdown

- $Q_1$ : How to know congestion happened
- $Q_2$ : How to slow down

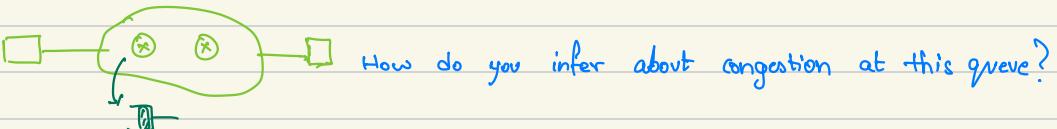
21/10/2021

class - 33 :-

### TCP Congestion Control :-

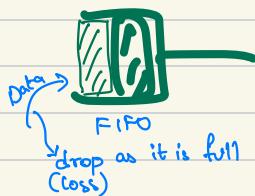


\*



Idea 1 :- See ↑ or ↓ in RTT

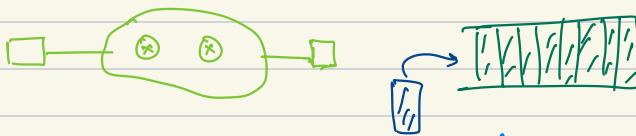
Idea 2 :- Packet loss



Q) How do we infer that packet loss occurred?  
(said one way. didn't listen)

Idea 3 :- ECN: Explicit Congestion Notification

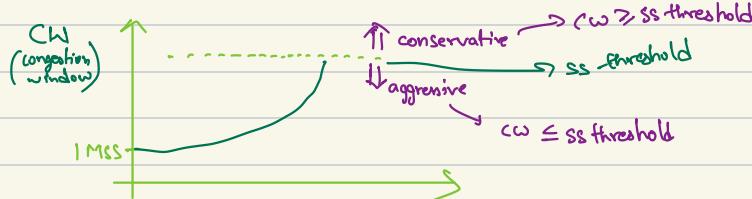
(was research interest in 80's)



(Sort of a hack)  
In b/w router take part in TCP

Before congestion happens itself,  
we inform DST/RCV by setting a bit in  
TCP hdr of packet

\* Another issue discussed: Range of available bandwidth is very high (10 kbps - 10 Gbps)



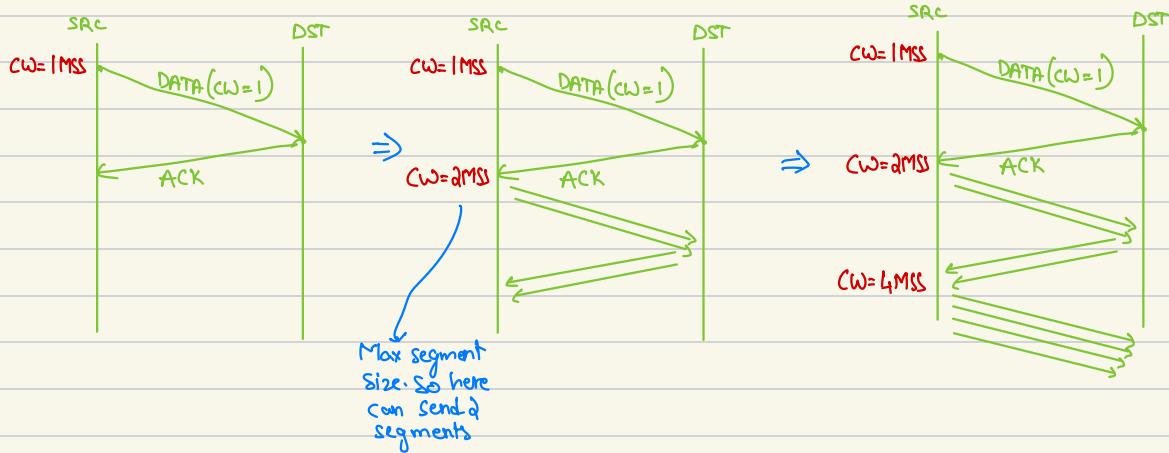


both are needed,

if we only do linear, it might take forever to find true bandwidth.

### \* How to practically increase CW:-

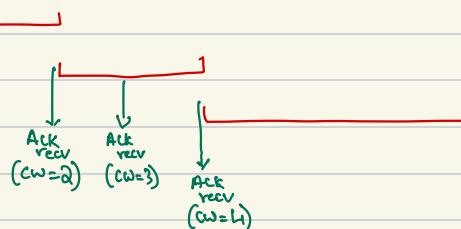
#### ① Slow start: Double CW for every ACK



\* On receiving an ACK,

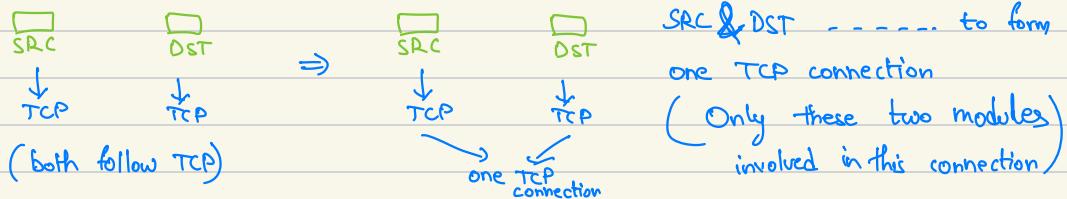
$$CW += 1 \text{ MSS} \quad \left( 1 + 1 + 2 + 4 + 8 + \dots \right)$$

exponential ↑

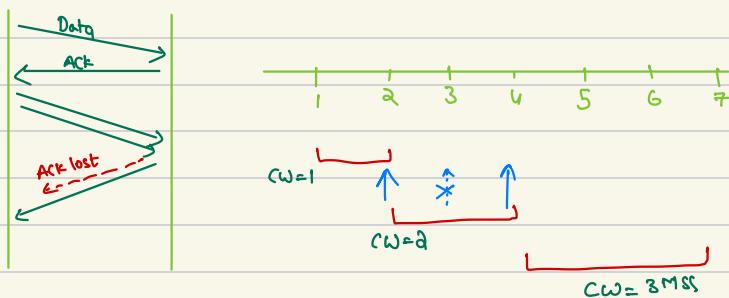


- Moving the window to right
  - Increment CW by 1
- Two things happening here

Note:- TCP isn't a Network level thing all have same parameters



RFC 5681: (A TCP protocol)



② Additive Increase:-

\*  $CW += 1$  per RTT (in additive inverse)  
↓

$CW += x$  per ACK

How much should  $x$  be?

In CW, how many segments are in each RTT?  $n = \frac{CW}{MSS}$

$$n \cdot x = 1 \cdot MSS$$

$$n = \frac{MSS}{\left(\frac{CW}{MSS}\right)} \Rightarrow x = \frac{(MSS)^2}{CW}$$

Note:-

Whenever  $1 \cdot MSS$  QMSS is written, it is like  $1 \cdot MSS$ ,  $2 \cdot MSS$

RFC 5681 :-

22/10/24 Class - 34 :-

TCP congestion control :-

- \* (Recap) slow start :  $CW += MSS$  per Ack , Additive increase :  $CW += \frac{(MSS)^2}{CW}$  per Ack
- \* What should we do when pkts dropped?



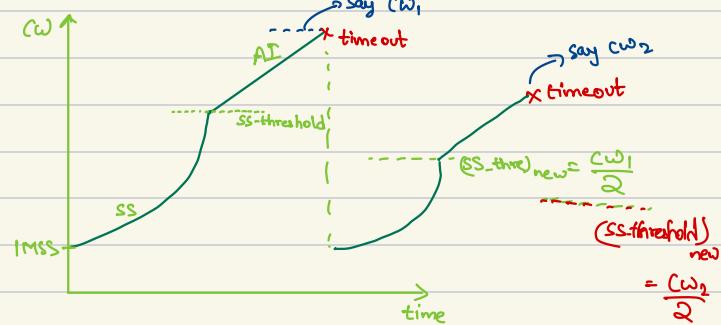
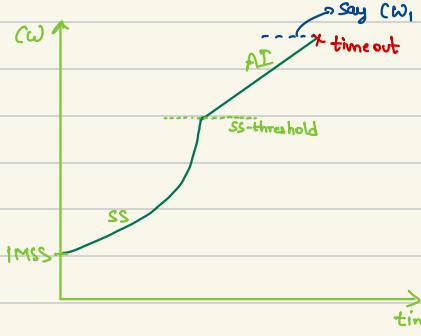
\* Congestion avoidance

AI-M.. :- Additive  $\uparrow$  & when a pkt drops, exponential  $\downarrow$

(drastic change when even little load occurs)

\* TCP TAHOE :- (one early version)

(very aggressive in  $\downarrow$  window size) (congestion avoidance)

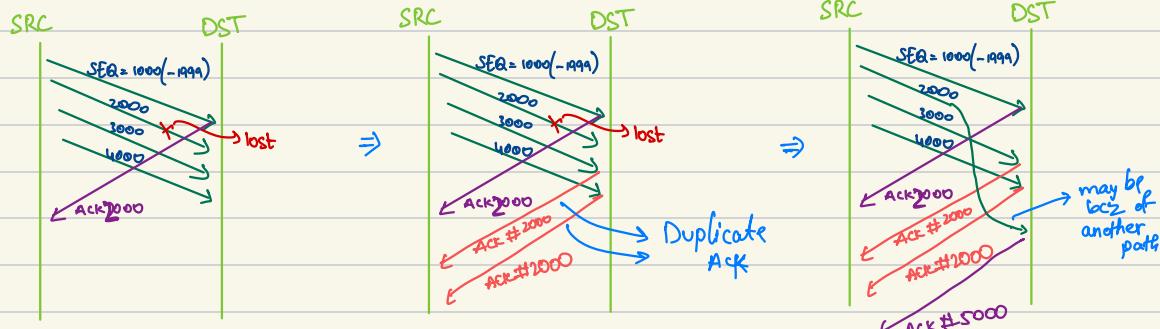


timeout : (if Ack not received in  $RTT$  time after sending a pkt)

- \* This was not default . It is drastic to wait until timeout  
Default is TCP Reno

\* Before TCP reno, we'll see what & why & how are ideas diff from TCP Tahoe

Can we assume pkt lost when ACK # 2000 received multiple times?



\* Triple duplicate loss :-

If 3 duplicates received  $\Rightarrow$  assume pkt loss

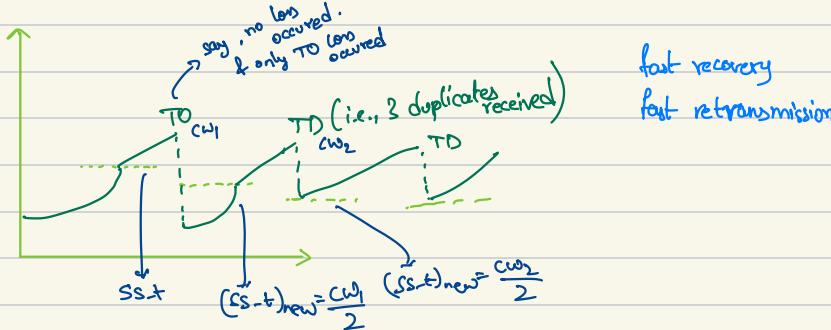
\* Should we bring CW to 1 MSS after loss?

No

$\rightarrow$  TO (time out) loss is on drastic congestion

$\rightarrow$  TD (triple dup) loss is not that drastic

TCP Reno :-



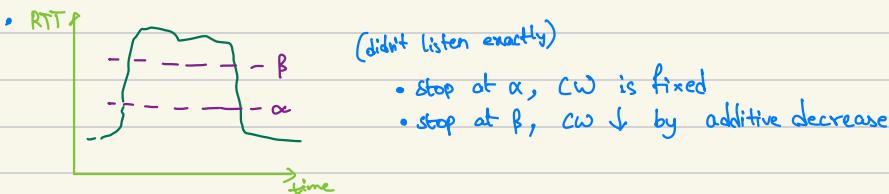
TCP Vegas :- Parameters, not known how to set so researchers didn't feel good

why class-35 :- (missed 10 min)

TCP Vegas :-

- If queue getting full & if congestion is increasing, Reno waits until pkt loss occurs.  
But Vegas doesn't.

- Slow start, TO loss, TD loss  $\Rightarrow$  same as Reno  
But when  $cw > ss\text{-threshold}$ , it does "Congestion Avoidance"



Vegas Rules :- (For congestion avoidance)

- \* Base RTT - min observed RTT in some recent time window (say  $W$ )  
RTT - current (smoothed) estimate of RTT

- \* Suppose no congestion, then  $RTT = \text{BaseRTT}$

- \* But if they are different,

$$\text{Expected rate} = \frac{W}{\text{BaseRTT}}, \text{ Actual rate} = \frac{W}{RTT} \quad \left. \right\} \text{diff} = W \left[ \frac{1}{\text{BaseRTT}} - \frac{1}{RTT} \right]$$

- \* If  $\text{Diff} < \alpha \Rightarrow$  A.I (as in Reno)

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

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

- \* Suggested  $\alpha = 30 \text{ kbps}$ ,  $\beta = 60 \text{ kbps}$

- \* Other issue than  $\alpha, \beta$  parameters, Question was can Vegas be used when some other in Network uses other TCP's?

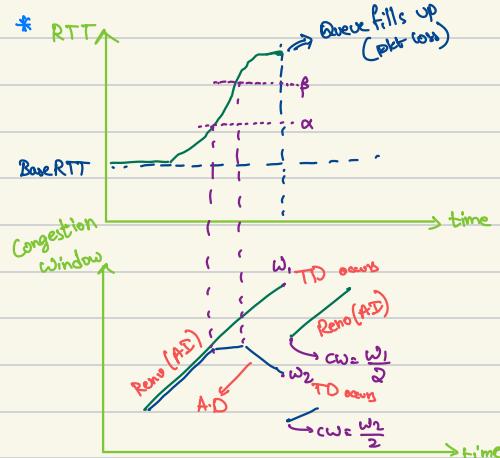
plan to change window as



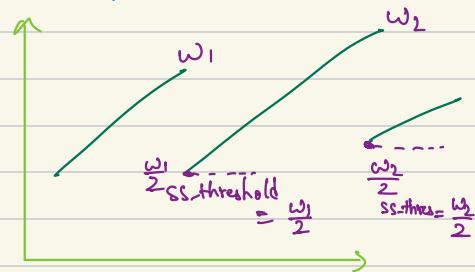
Vegas  $\Rightarrow$



Even decrease it if  $\text{diff} \geq \beta$



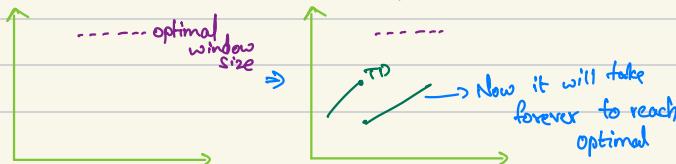
It is not that ss-threshold always decreases



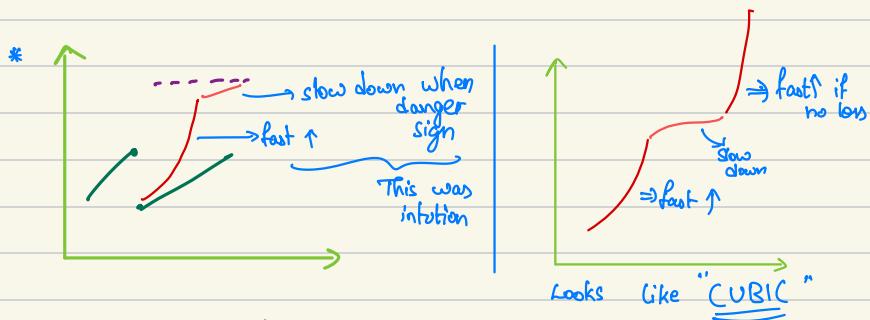
- So overall we see, when both are running  
VEGAS get unfair

- TCP Reno was default for long time (1988-2012) by many OS
- TCP Cubic  $\Rightarrow$  now default (Linux, MAC os)

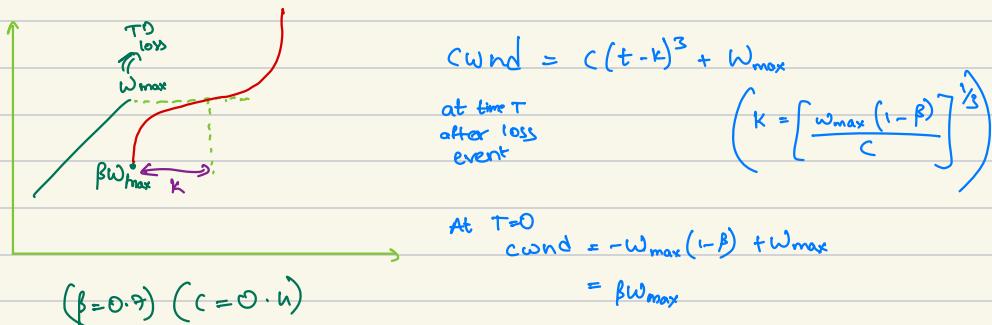
- Q) Why TCP Reno was good enough? We have much speed now. Why?



What if +5 MSS instead of 1 MSS? (didn't listen)



\* So researchers actually tried CUBIC



\* Next 2 weeks Application Layer

application

class-36 :-

- \* APPL
  - DLL, PHY many companies. But some survived. LiCh < WiMAX
  - WiMAX was ahead of LiCh. But due to economics (companies said upgrading to LiCh easy)
  - NET is stuck with IP
  - Transport has UDP, TCP → does nothing → does more than required

If you want diff from transport,

can build in App layer. Ex:- Qui-, google's. use TCP cubic

Others no big innovation. But App layer have huge innovation

Ex:- (Have to be in right place in right time)

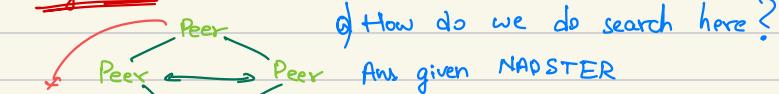
### I) Client-Server:-

- Smartpix (10 yrs ago) (IITD sir student made)
- OLA (IITB alum)
- HTTP (Web) → by high energy physicists (CERN)
- Facebook, . . . .

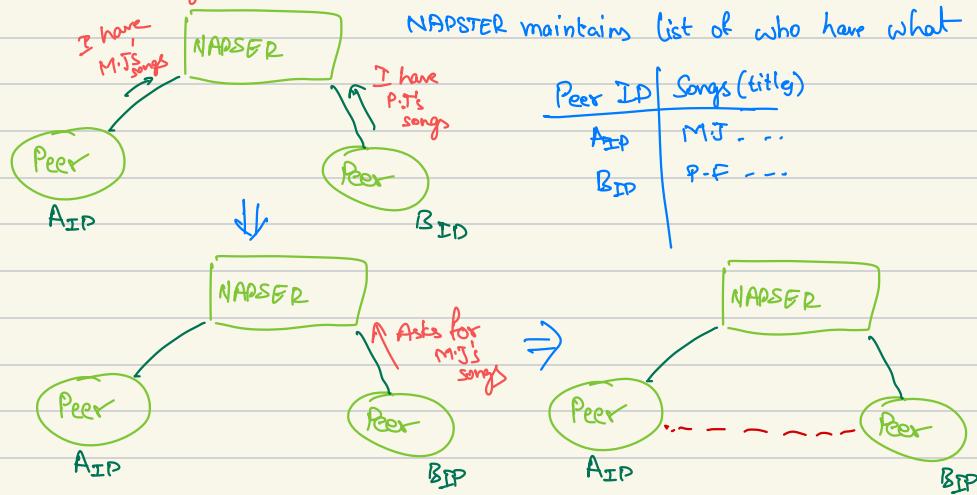
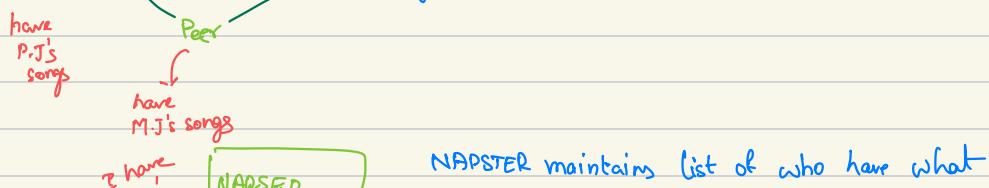
Another technology,

### II) P2P :- (peer to peer)

#### 1<sup>st</sup> generation:- NAPSTER



(NAPSTER was before google)



\* Bcoz of some legal issues, NAPSTER closed

\* It has some issues as centralized server

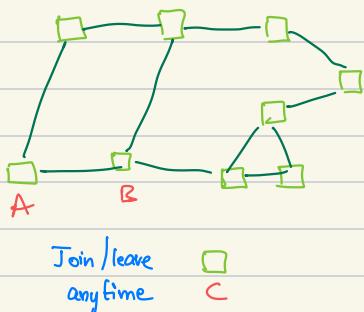
#### Problems with Centralised server:-

##### 1. Single Point of failure:-

Two aspects • Technical • Legal

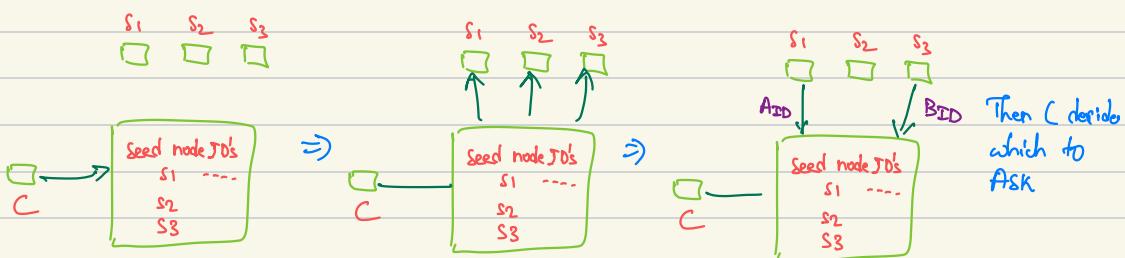
- Technical: i) Power cut    ii) DDos attack (Just sending many requests to single server to get down)
- Legal: If it had to shut down, just go to 1 place

## 2<sup>nd</sup> generation:- GNUTELLA



- There are some seed nodes  $S_1, S_2, S_3$   
They'll keep track of recent transfers  
Also keeps contacting periodically, to know if they're still in network (hello messages)
- In the application software of every node,  
Now C requests  $S_1, S_2, S_3$

Seed node J's
$S_1$ ---
$S_2$
$S_3$

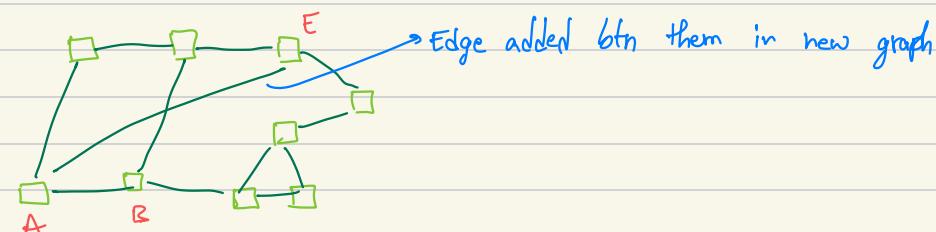
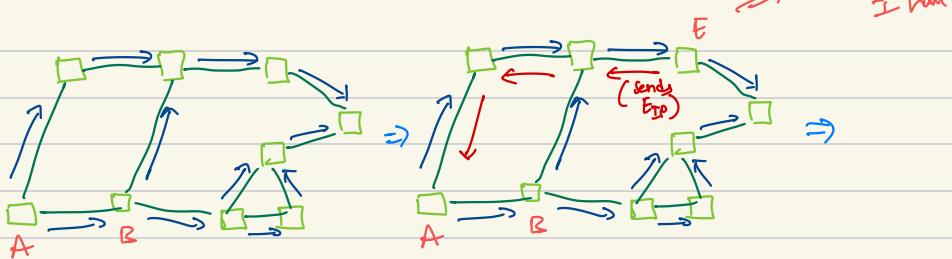


\* Seed doesn't have much info, it doesn't want to get more involvement.  
This is where innovation came in  
BitTorrent.

Idea 1 :- when a file comes, I'll broadcast that I have. But then each node become knapster server.  
Good for small networks. Not scalable

Idea 2 :- (Flip Idea 1) Broadcast the request

Note:- Assume everyone honest



\* Brunella did more advanced.

Expanded ring broadcast :-

29/10/14 Class-37:-

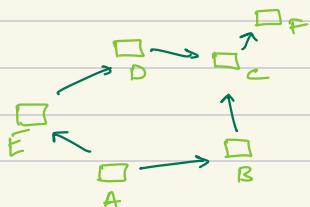
\* 15 min Recap (didn't listen)

=7 initially

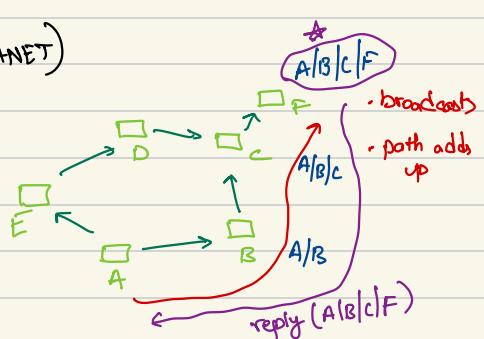
\* For file request, A can local broadcast (i.e., limited depth) (global can make bandwidth full at all instants)

→ Digression

Similar to P2P, we have "Mobile ad hoc Network" (MANET)



- Tanks - Moving
- If want to search for F



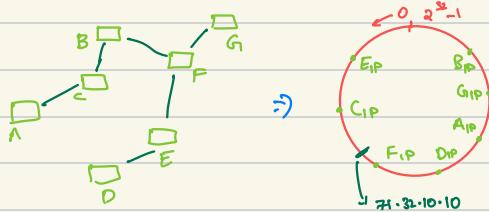
\* This is called "Source routing".

- broadcasts for path
- DST replies with 1 or 2 possible path
- Source decides which path data goes

Note:- Why not use D.V? All moving so weights change. Q) What risk is there normally also. After A knowing that path A/B/c/F, if C moves out? Ans:- Have to do some hacks ex:- B does source routing

\* Coming back, P2P wanted a ~~better~~<sup>efficient</sup> way of data ...

### 3<sup>rd</sup> Generation :- Distributed hash table (DHT)



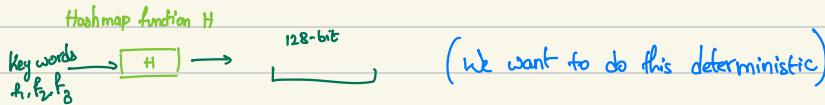
Since diff ports world connected  
via P2P, can assume uniformly  
distributed IP .

If A have to sent to,  
A want to know which  
of nodes in network has  
(closest IP to 71.32.10.10)

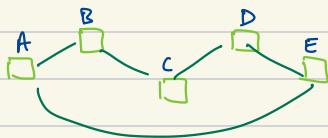


- First A sends "search" to C i.e., send to closest IP neighbour
- Then C sends "search" to F

\* Having above algo as alternate for broadcast, we want to bring a mapping from file name/key words  $\rightarrow$  IP



\* Consider a network as below



(didn't listen)

\* High Level :-

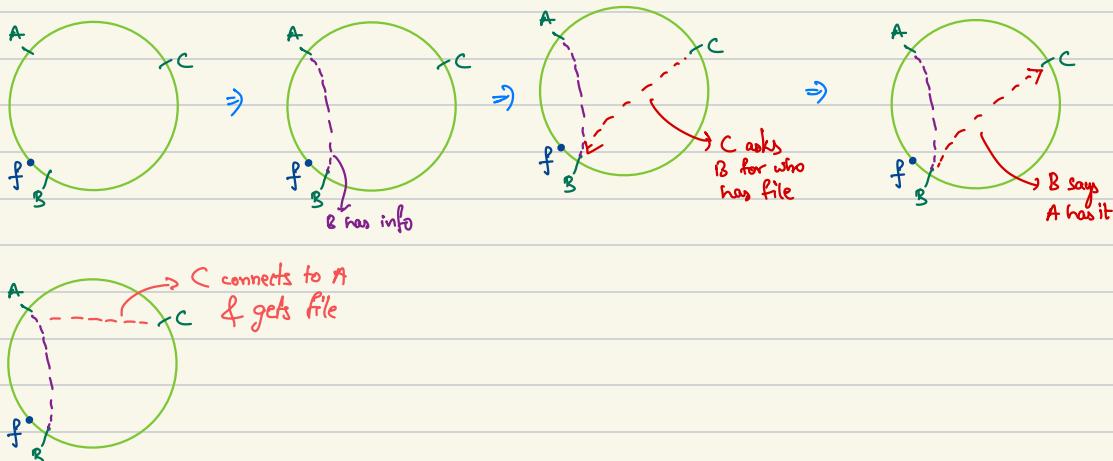
Suppose node A has a file 'f'. A finds node (say B) where IP is closest to f in common space i.e.,

$$\text{dist}(H(f), H(B_{ip})) \leq \text{dist}(H(f), H(X_{ip})) \quad \forall \text{ nodes } X \text{ in network}$$

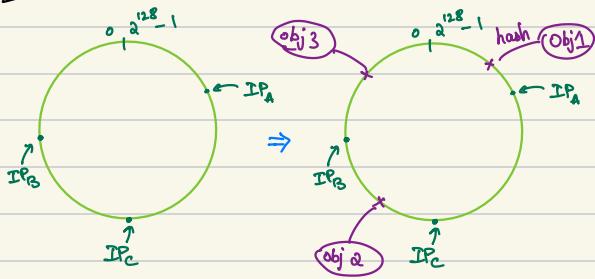
Then B maintains in a table that

file	node
f	A_ip

- C finds node closest to 'f' in common space
- C asks B: "Who is f?"
- B gives A\_ip to C
- C connects to A and downloads file

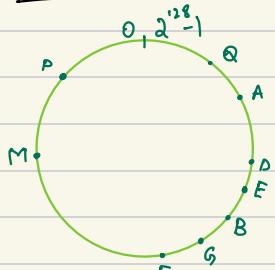


### Wifay Class-38:-



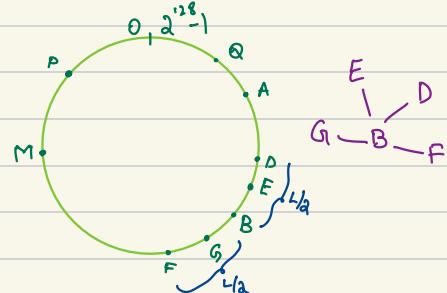
- \* Suppose Node C has Obj1
- It doesn't want to say everyone that it has Obj1
- \* Finds that A has IP closest to Obj1
- \* Tells A that, I have Obj1
- \* When B wants Obj1, asks A (closest to Obj1)
- Then A says C has it.
- \* B connects to C & gets Obj1

### Poorni :-

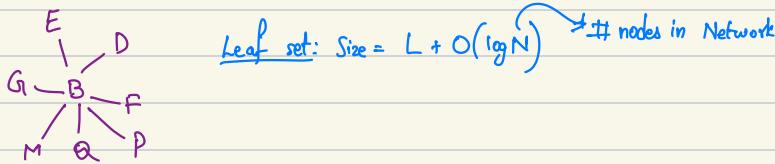


\* Leaf Set: Neighbors in P2P network of any node

\* if L=4, d on right, d on left (according to IP)  
L ⇒ set by administrator of network



\* Apart from  $\gamma_2, \gamma_2$  B could be connected to other nodes in network



\* If M wants an Obj1 which is in  $(IP_{B-\gamma_2}, IP_{B+\gamma_2})$

Not sure

\* For M, closest neighbour to Obj1 is B,  
M asks B for Obj1

\* Then B knows D is closest to Obj1, so asks D to see about it  $\Rightarrow$  Not sure

\* D knows  $-\gamma_2, +\gamma_2$  neighbours to itself &  
so knows that it is closest to Obj1. Hence replies to M about who has obj1.

\* Node X wants to join

\* clear on board

\* bcz of <sup>longest</sup> prefix match, in  $O(\log N)$  steps we reach closest

\* Node A gets fail

\* Say A,B on right of X, C,D on left of X \* Maintain some heart beat messages so X knows  
A failed

\* Now X knows A failed. But want to connect to another on right

\* It knows B is rightmost & asks for info of its immediate right

\* What about info in A? Is it lost?

The copy of info at A is always sent to its  $\frac{L}{2} \frac{L}{2}$  neighbours.

This is how node failure handled

Now-a-days BITTorent : simultaneous download. (May be "K-" has has bandwidth)

more complex. Seed & now to know which all nodes have file

## Domain Name System :- (DNS)

\* We don't prefer 32.75.5.9 IP . We want google.com , iitb.ac.in

URL  
Transport layer

But these can't be understood by lower L layers. They understand IP, MAC (or) port number for x  
So Appl layer should take care

\* Another adv:

Today google may have 8.8.7.5 , if it want to change 8.8.7.3 . It has to be changed by DNS . & one can still use "google.com"

\* Flexibility similar to IP ( machine's IP can change anytime)

\* Has to be robust.

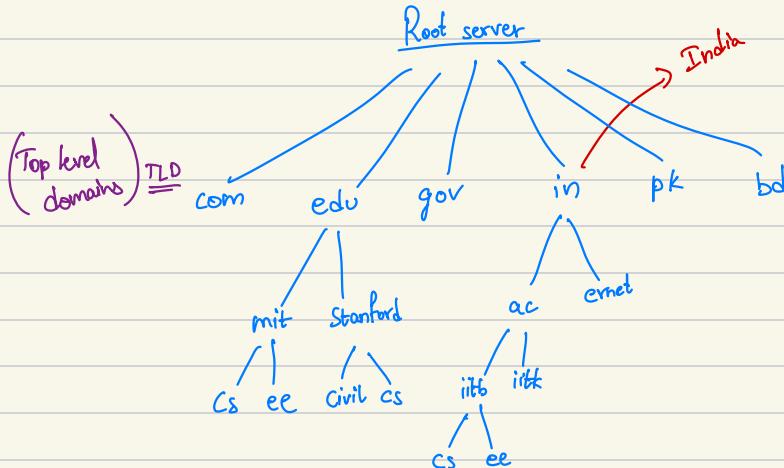
Naive :-

Maintain single server , which has URL $\leftrightarrow$ IP mapping . One queries , What is IP of a URL.

Problem: Single point failure , Not scalable

Ideal :-

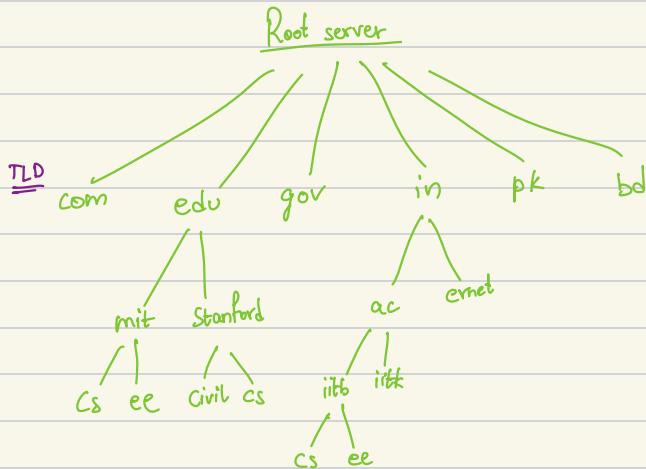
We have hierarchy



Root server, in , ac , iitb , cs  
all have a "DNS server"  
↳ Has mapping

5/11/24

class-39 :- (didn't listen 10min)



\* DNS server of parent domain should know the DNS servers of immediate children

- i) Name
- ii) IP address } necessary info

(But actually, in can remember cs)  
But not necessary

\* Why this is scalable?

Say we want to add `surya.cse.iitm.ac.in`. We only need to change at cse DNS server.

If centralized server  $\Rightarrow$  changing difficult

\* (didn't listen 5min) 13 servers there. What if all made down?

\* A. root-server.net = 20.32.5.23

We keep multiple physical machines with same name. We don't care where exactly it goes

(advertisements --- BGP routes ---)

↳ called ANYCAST

UNICAST = single destination

BROADCAST = All nodes are destination

MULTICAST = subset of nodes are destination

ANYCAST = Any node of a subset can be destination

↳ Had economic issues. Why not broadcast have that?

broadcast done locally. They planned to do multicast over all network.

(like cricket live, send some plot to many at a time)

Resource records:

<NAME, VALUE, TYPE, CLASS, TTL>

in  
(internet) ↳ how long ..

## TYPE:

.

- CNAME → Canonical name (alias)

(MX → Name of email server of domain specified in NAME field      cs.mit.edu      ee.mit.edu)

July class-#0: (missed 20 min)

## \* HTTP:-

(missed format)

### Requests:-

For example, GET <url>

### → Request operations:-

GET → Retrieve document

HEAD → Retrieve Meta Info

(e.g.: last-modified, length..)

OPTIONS → Available Options (HTTP version at server)

POST → Give new information to server

PUT → Store/modify info at specified URL on server

DELETE → Delete specified URL

giving info to server

### Responses:-

For ex, (startline) HTTP /1.1 200 ACCEPTED <CRLF>

(HDR) Content-length:

Expiry time:

<CRLF>

(Body)

### Request codes:

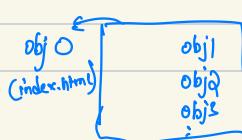
1xx = informational

2xx = success

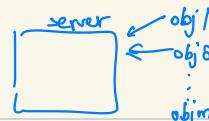
3xx =

## HTTP & TCP:-

Say a webpage (this itself is an object) contain some objects



Browser → HTTP → TCP



objects come after another  
is not a good idea

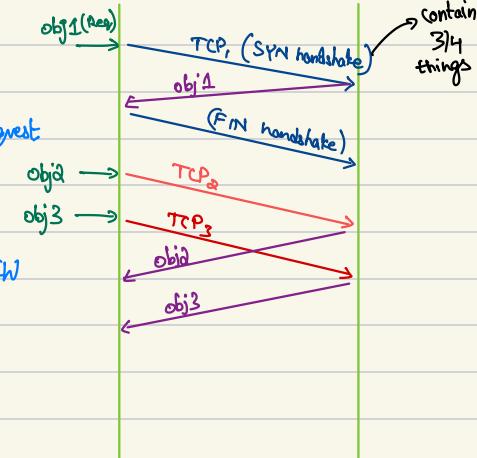
Q) Send all requests at once / sequentially?  
How many TCP connections?

HTTP does a naive approach

\* New TCP connection for every GET Request

Browser - client (HTTP)

Server

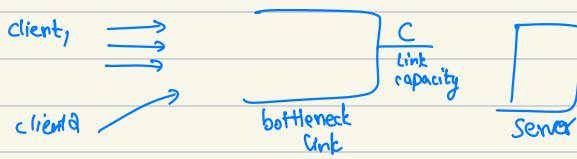


Issues with this :-

- 1) Each TCP connection takes time to learn optimal CW
- 2) Overheads (hand shake)
- 3) Server state is large due to multiple TCP connections being open
- 4) Fairness

(Say client<sub>1</sub> has m TCP connections

client<sub>2</sub> has 1 TCP connection



$(m+1)$  TCP connections

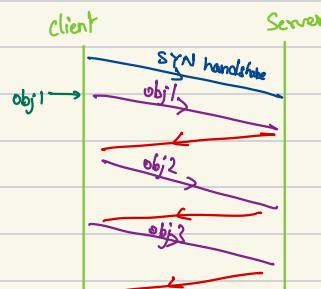
Each gets  $\frac{C}{m+1}$  (assuming same TCP version, RTT etc..)

$$\frac{C}{m+1} \rightarrow \text{client}_1$$

$$\frac{C}{m+1} \rightarrow \text{client}_2$$

\* So they invented HTTP/1.1 - (around year 2000)

All requests in single TCP connection,  
called "Persistent TCP connection"



We have a queue  
& all are done sequentially

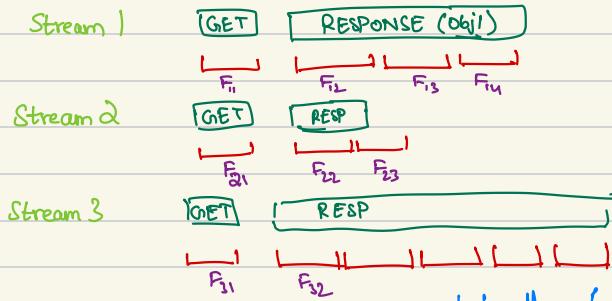
\* HTTP 2 :- (around 2015) (best of both 1.0, 1.1) issues in 1.1 :- What if obj1 takes more time but 2 takes less time

→ Head of line blocking :-

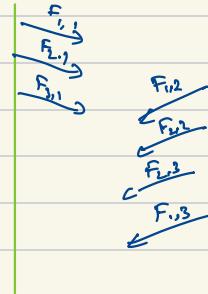
obj 'i' gets delayed ...

- Single TCP
- divide requests into streams

Say we have 3 objects with 3 streams



Same TCP



but allows frames  
from diff streams in the TCP

\* HTTP 3 more optimization