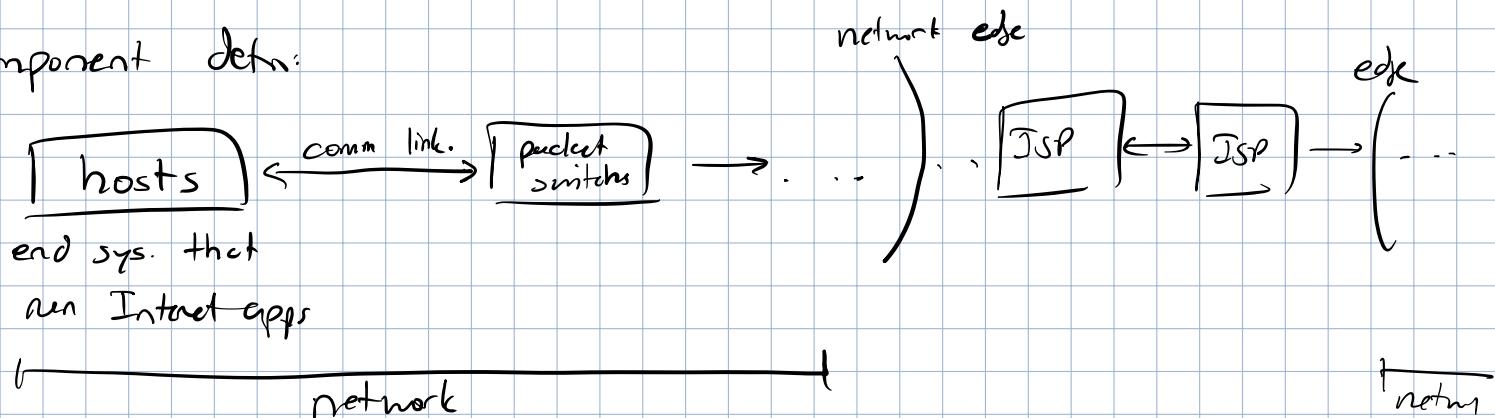


INTRODUCTION

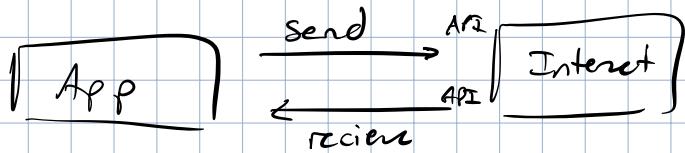
Definition of Internet

Component defn:



Internet: network of networks.

Service-based defn:



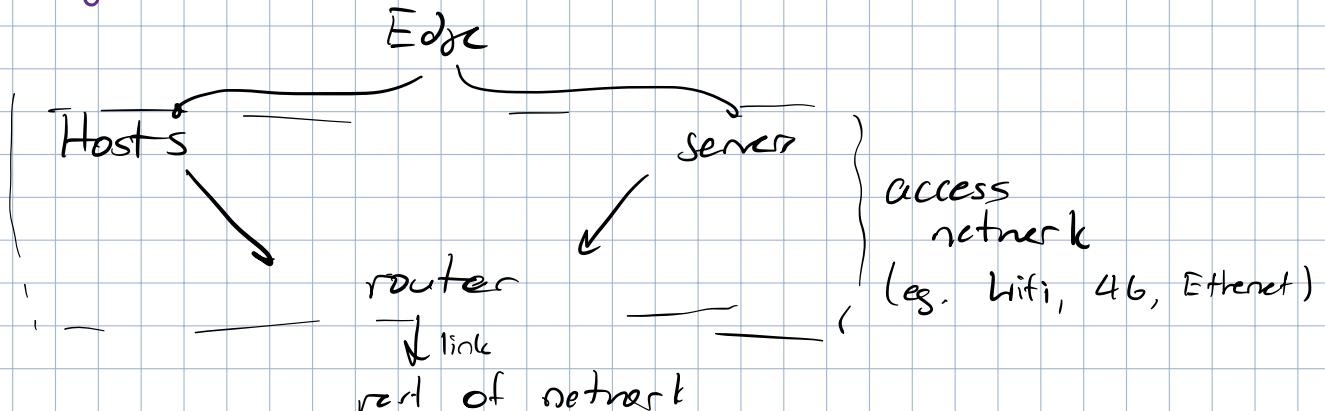
Protocols

Comm. standards that controls sending & receiving messages.

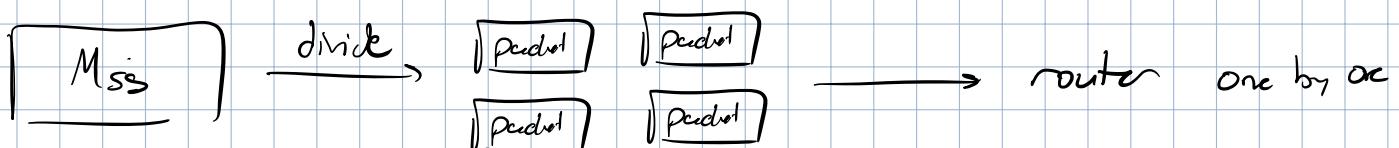
Defined by software

Reason: compatibility

Network Edge



Q: how do we send messages to router



Assume each packet is L bits, bandwidth of router is R b/sec.
↳ Time to send packet: L/R sec.

Links:

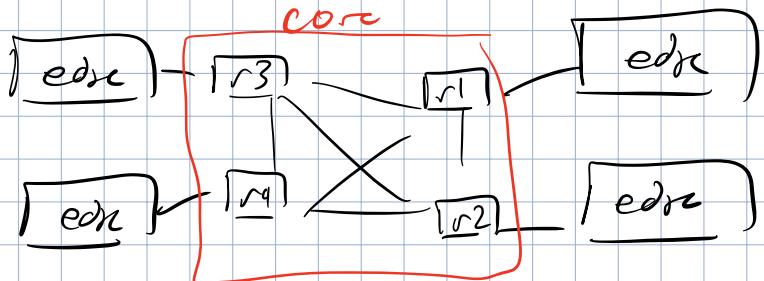
- ① Guided: transmission in solid media
 - Ethernet
 - Coaxial : 2 Cu wires \rightarrow bidirectional measg.
 - Fibre optic : 1 pulse \rightarrow 1 bit. Low error, high bandwidth

- ② Unguided: propagate freely

- Radio
 - Microwave
- } Cons:
1. Environment can have effect
 2. Half-duplex: can only send & receive for particular spectrum 1 at time.

Network Core

Dfn: interconnected routers:



Role of core:

1. Figure out best route from src to dest
2. Forward packets to dest.

2 methods of forwarding.

- ① Store-and-forward packet switching



entire packet must arrive at r1 b4t moving to r2

Q: How long to move 1 packet from src \rightarrow dest w/ n

packet switches w/ same transmission rate:



$$\# \text{ of links} = n + 1 \quad ? \quad \text{transmissn / link} = W/R \quad \text{total time} = (n+1) \cdot W/R$$

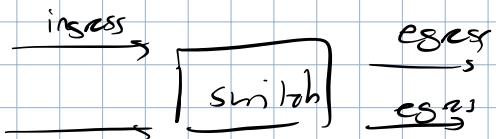
Q: time delta b/w packages?

$W/R \leftarrow$ wait for package to be uploaded

Q: Total time to send p packets?

$$\begin{aligned} \text{total time} &= \text{time to send 1 packet} + \text{time to send } (p-1) \text{ packets} \\ &= (n+1) \cdot W/R + \sum_{p-1} W/R \\ &= \boxed{W/R (p+n)} \end{aligned}$$

Con: queuing

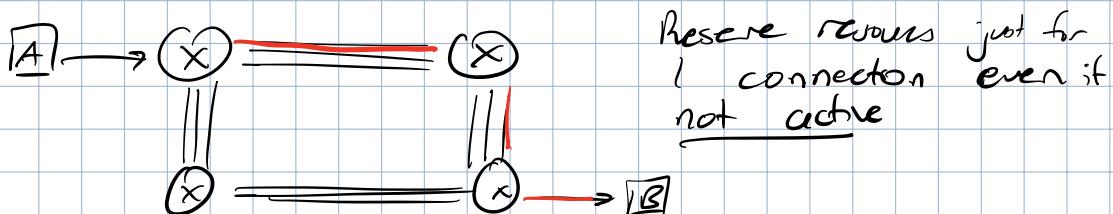


Ingress rate > egress rate for some time \rightarrow queuing

Implications:

1. Packet loss \leftarrow internal mem. can't store
2. Latency \leftarrow switch needs to prioritize packets

② Circuit switching:



How?

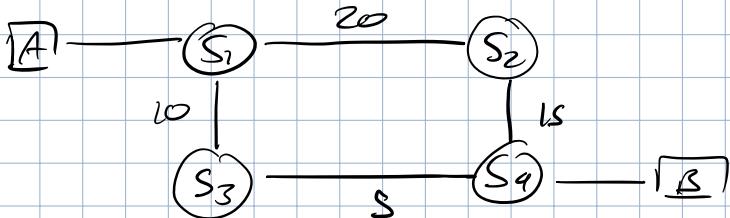
- ① Frequency div. multiplex. (FDM)
 - Reserve freq. / call
- ② Time div. mux. (TDM)
 - Each call has own tr.

Pros: interference w/ other comm. zone, reliable for comm.

Cons:

1. Suboptimal resource alloc.
 2. Takes time to alloc. resource
 3. Not resilient to network changes.
 4. Max # of conn.

Ex: //



(Q: Total # of calls from A → B?)

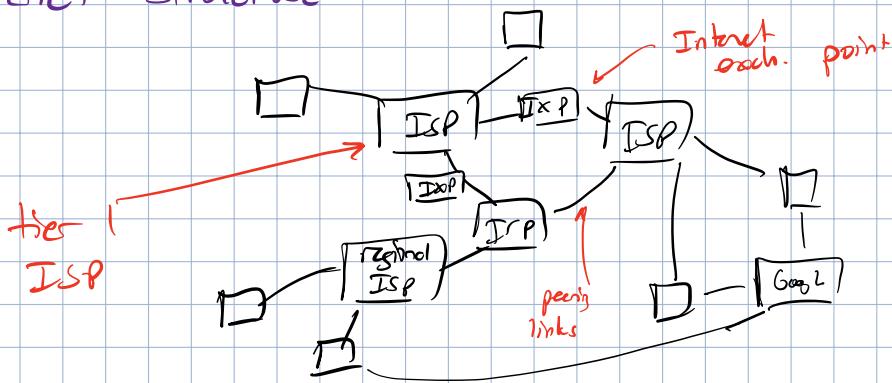
Procedure: find all paths & take min. # of conn on each path + sum.

2 factors?

- ⑦ $S_1 \rightarrow S_2 \rightarrow S_3$: 1S conn. } 20 conn. max.
 ⑧ $S_1 \rightarrow S_3 \rightarrow S_4$: S conn

Benefits of packet switching: far more efficient in bursty traffic.

Internet Structure



Alt:

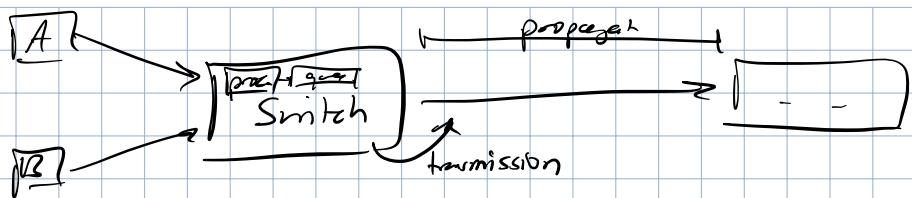
- ① Pairwise conn. b/w host
 - Not practical
 - ② Central system directs traffic
 - Too centralized
 - ③ Separate ISP.
 - Connect ISP together

Companies can run own network.

Performance

- ## ① Delay

Causes: queuing



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

delay / packet / node

d_{proc} : checking packet & output link
in usec.

d_{queue} : time at switch waiting for link

$$\text{Diagram: } \boxed{1} \xrightarrow{t_1} \boxed{2} \xrightarrow{t_2} \Rightarrow d_{\text{queue}} = t_2 - t_1$$

$$\text{Traffic intensity} = \frac{L \cdot \alpha}{R} \text{ vs. arrival rate}$$

Intensity $\left\{ \begin{array}{l} \sim 0 : \text{delay small} \\ \rightarrow 1 : \text{large delay} \\ > 1 : \infty \text{ delay (more work causing > scenario)} \end{array} \right.$

d_{trans} : bandwidth restricted

$$d_{\text{trans}} = L/R \quad \left\{ \begin{array}{l} \text{delay in upload} \\ \text{packet} \end{array} \right.$$

d_{prop} : round off length of link

$$d_{\text{prop}} = \frac{d}{s} \rightarrow \text{length of link} \quad \rightarrow \text{prop. spec. } (2 \times 10^8 \text{ m/sec})$$

Traceroute used to find real delay

- ① Send n packets
- ② On each switch, packets sent back
- ③ Repeat for all switches

⑦ Loss

Packet \rightarrow full queue \Rightarrow lose packet

Remediation:

- ① Retransmit from prev. node

- ② 11 send

- ③ Don't care } Protocol defines behavior

③ Throughput

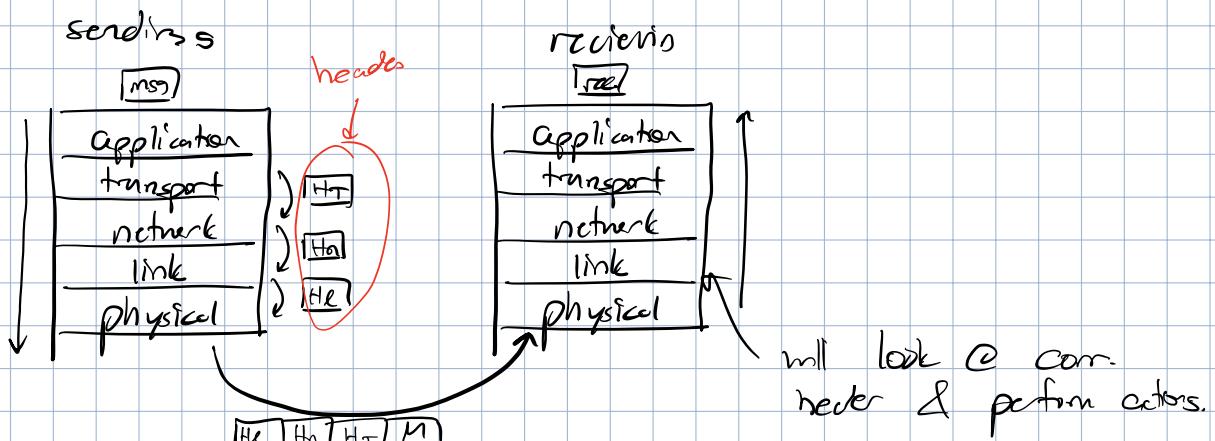
$$\frac{\text{bits sent}}{\text{time}}$$

instants: look @ specific t.

avg.: look @ time avg.

Throughput of path is limited by slowest link

layers

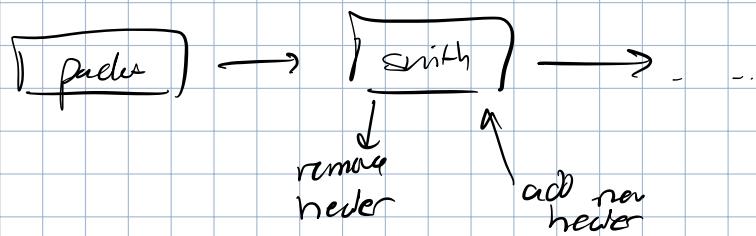


Layer i depends on Layer i + 1

Layer job:

- Application: network apps (HTTP, DNS, ...)
- Transport: how to transfer process → process (UDP, TCP)
- Network: routes datagrams from host → host
- Link: route datagram from neighbor
- Physical: physical movement of bits

Not all comp. of network implements all layers → only ones that matter



Dif. b/w switch & router: router implements network layer, switch doesn't

Pros: encapsulation

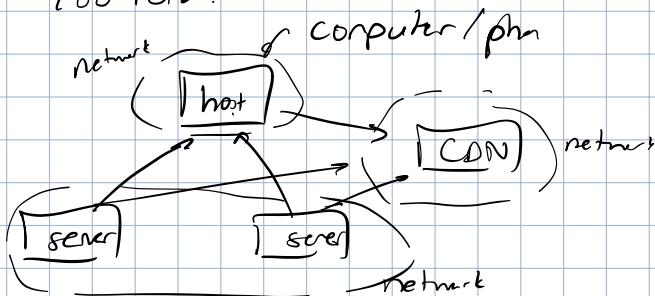
Cons: header overhead

APPLICATION LAYER

Principles

Network app: runs on diff. machines.

↳ Ex.: YouTube:



No software needed on apps for network core func. (routing, queuing)

Paradigms

Client-server

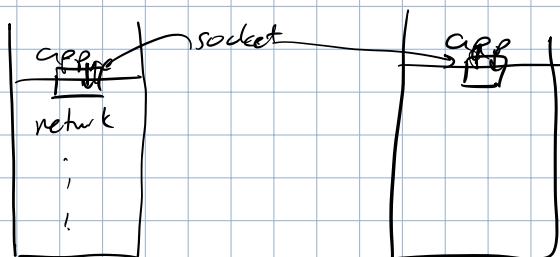
- Server is always-on
- Client contacts server for info, no comm. b/w each other
- Protocols: HTTP, IMAP, FTP

Peer-to-peer

-
- The diagram shows two 'host' nodes connected directly to each other, forming a peer-to-peer connection. A third 'host' node is shown below them, indicating they provide services to each other.
- Peers provide service to each other
 - Not always on / connected

Q: How to send messages from pro. in 1 machine → another machine?

A: sockets.



Q: How to identify proper process?

A: IP addr for machine + port no.

↳ corresponds to process (e.g. 80 for HTTP serv)

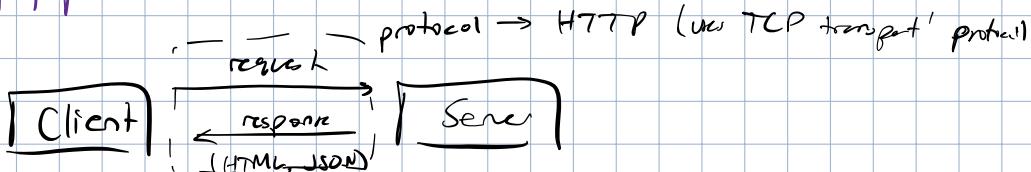
App layer protocol defns:

- ① Type of message: req./rep.?
- ② Syntax of message
- ③ Semantics of message

App needs following characteristics of transport service:

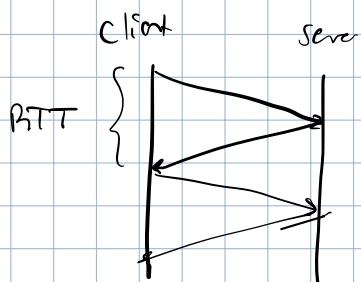
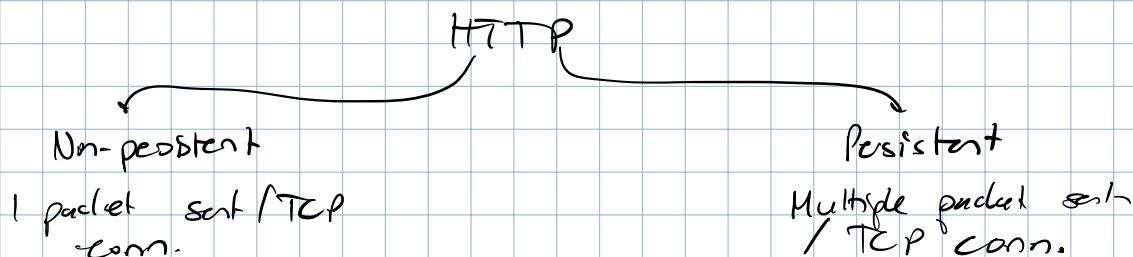
- (1) Data integrity
- (2) Timeliness
- (3) Throughput
- (4) Security

Web & HTTP



HTTP steps:

- (1) Client initiates TCP conn. to server (port 80)
 - (2) Server accepts conn.
 - (3) HTTP message is exchanged
 - (4) Conn. closed
- } HTTP is stateless, server shouldn't know about past requests.



B/c conn. opened, not as much RTT overhead (as little as 1)

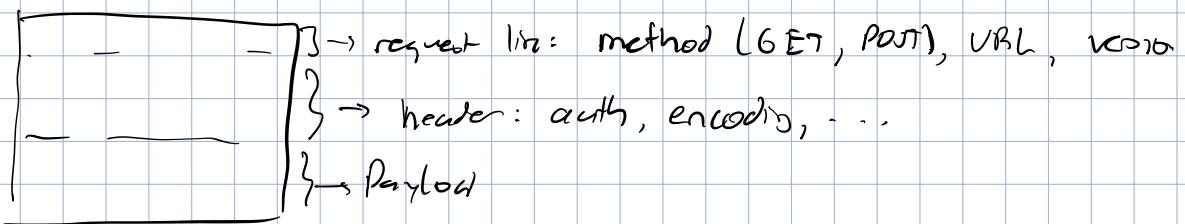
∴ Total for getting 1 object:

- + 1 RTT for conn.
- + 1 RTT for request resp.
- + File transmission

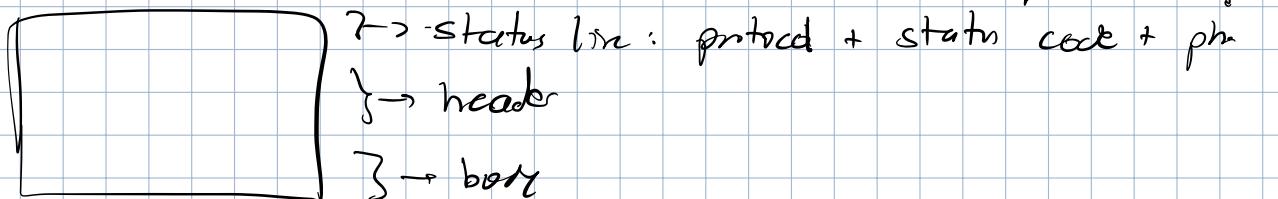
Conn.:

- (1) OS overhead
- (2) Not great for parallel conn.

HTTP request:



HTTP response:



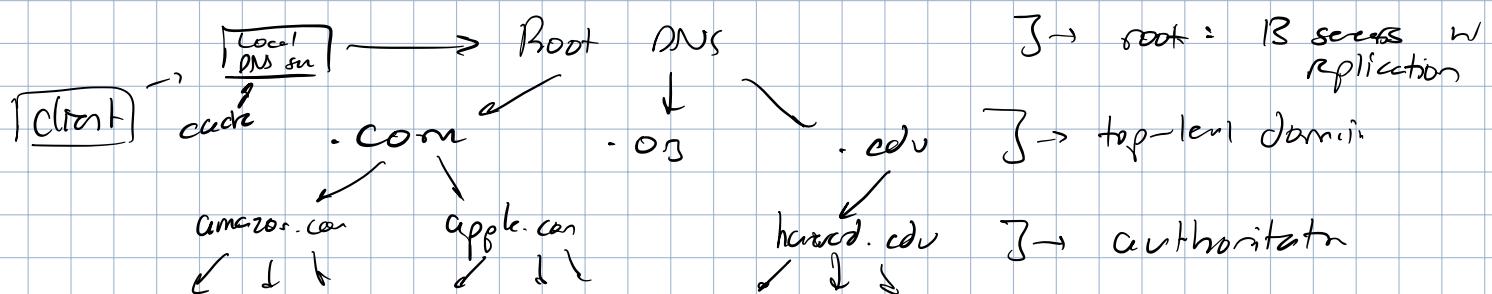
DNS



Other func.?

- ① Host aliasing (diff servers w/ same name)
- ② Mail aliasing (@univatloo.ca.)
- ③ Load balancing

Distributed & hierarchical \leftarrow resiliency & performance-sensitive



Traversal steps from local DNS server:

- ① Iterate: go through root \rightarrow top-level \rightarrow authoritative
- ② Precise: let root serve record & return.

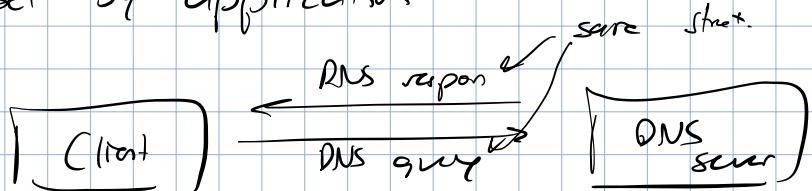
Caching is v. common

DNS records: (name, value, type, TTL)

A: (hostname, IP)
NS: (domain, authoritative server)
CNAME: (alias, actual, ...)
MX: (domain, envelope)

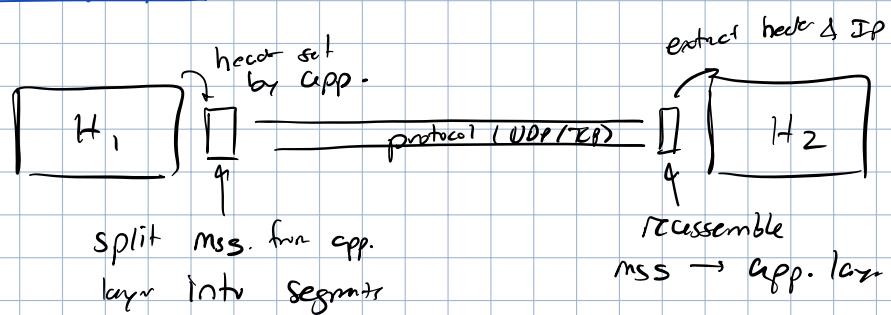
caching purpose

Set by application



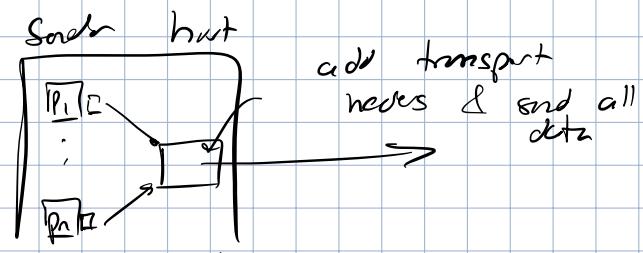
- ① ID for query
- ② Flags (eg. recursion)
- ③ Queries

TRANSPORT LAYER

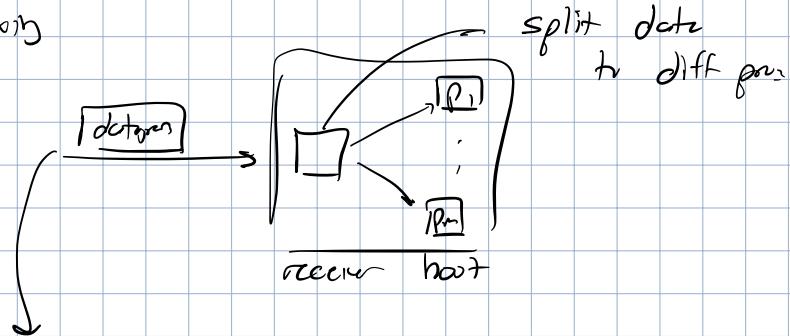


Multiplexing & Demultiplexing

Multiplexing:



Demultiplexing



use this to
find out who to send mess
to.

Port
Present by DS
Free

Types of demux:

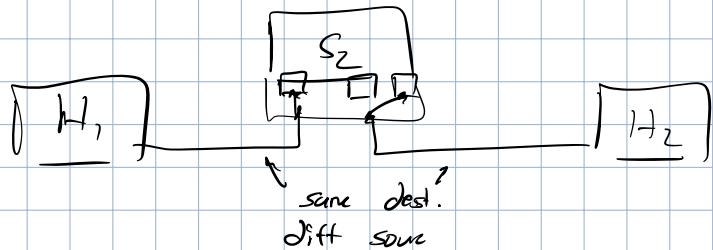
- ① Connectionless: only use dest IP & dest port to find correct socket
↳ Use in UDP

If multiple datagrams w/ same dest IP & port, diff. source IP & port
go to same socket

- ② Connection-oriented: use all 4 data-fields to demux

↳ Use in TCP

Allows multiple TCP sockets b/c not all conn. going to same socket



UDP: User Datagram Protocol

Bad bones

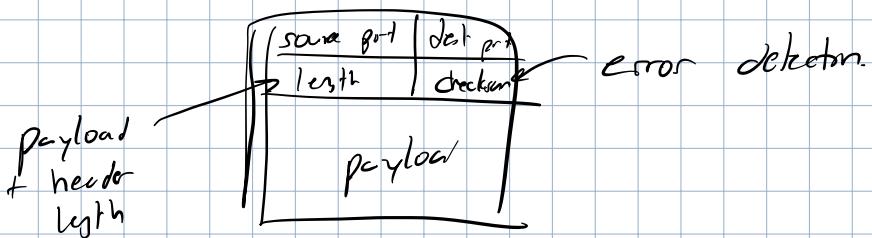
Best effort: segments can be lost / out-of-order

Connectionless: no handshaking but still 2-way comm.

Why?

- ① Little delay \rightarrow streaming, video games
- ② Simple
- ③ No congestion control

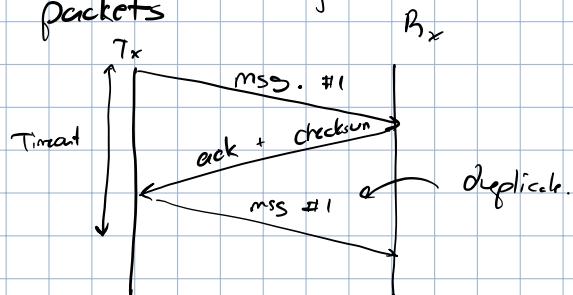
Add anything? Do it in app layer



Principles of reliable data transfer

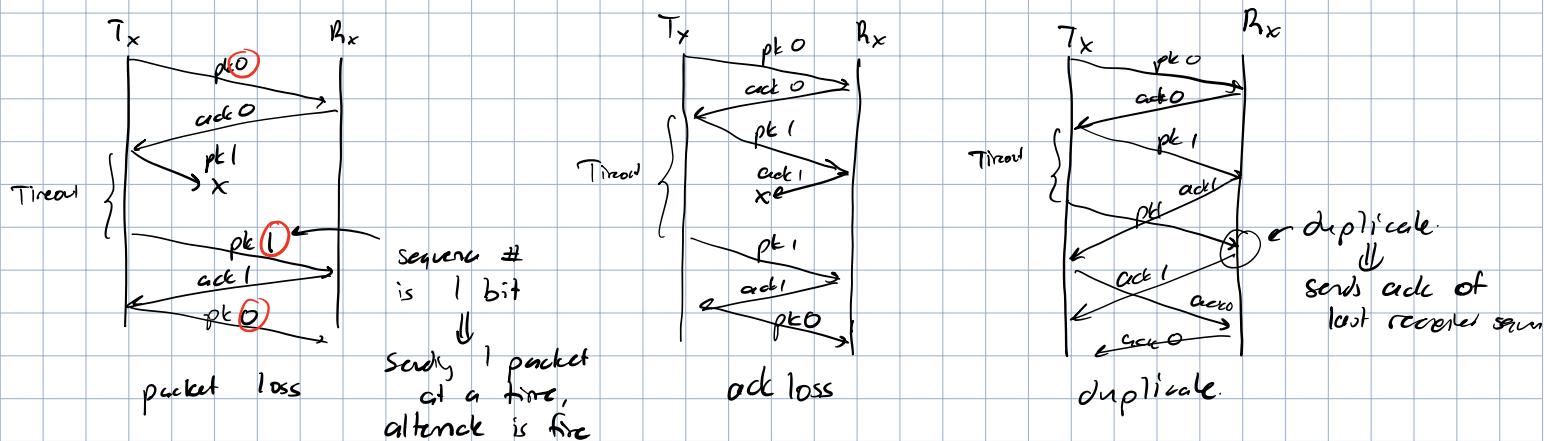
Unreliability of UDP:

- ① Packet errors \Rightarrow error det. mechanism needed
- ② Packet loss \Rightarrow feedback to sender & automatic repeat
- ③ Out-of-order packets } sequence #s for packet
- ④ Duplicate packets

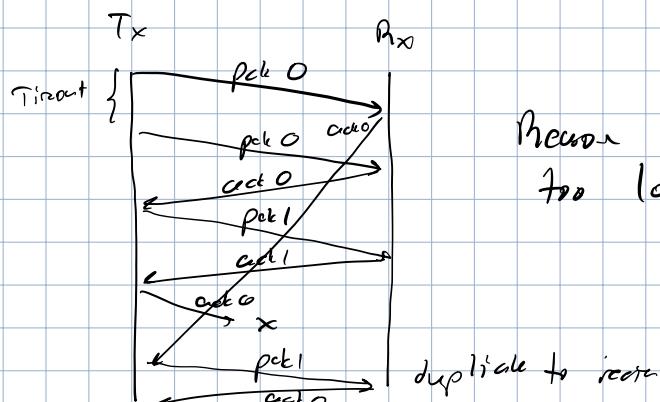


Stop-and-wait protocol

- ① Send packet w/ sequence #. If duplicate \Rightarrow discard
- ② Wait for ack / timer runs out
 - If ack \rightarrow next packet
 - Else \rightarrow retransmit



Stuck case:

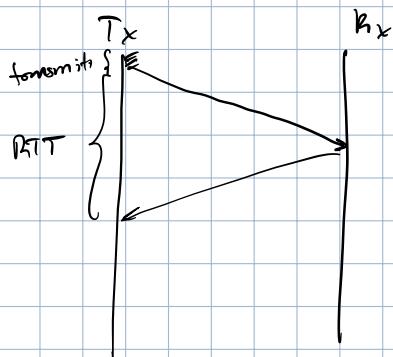


Reason for stuck: Timeout too long

duplicate to ack

Performance of stop-and-wait

$$U_{\text{sender}} = \frac{\text{time spent sending}}{\text{total time}}$$



$$\Rightarrow \text{Total time} = L/R + RTT$$

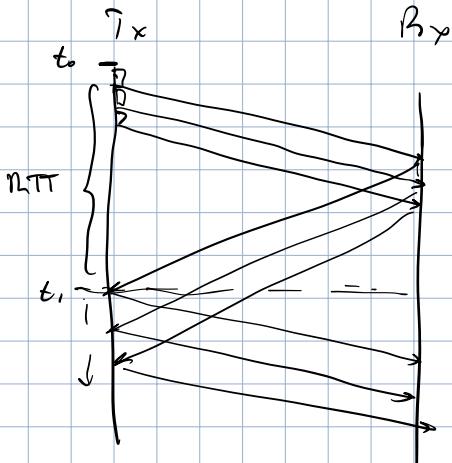
$$Tx \text{ transmit} = L/R$$

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

BAO!

Improve? Pipeline!

Defn: send multiple uncoded packets



$$\Rightarrow U_{\text{sender}} = \frac{3 \cdot L/R}{RTT + L/R} \Leftarrow t_1 - t_0 \quad (1^{\text{st}} \text{ packet})$$

↓

3x utilization!

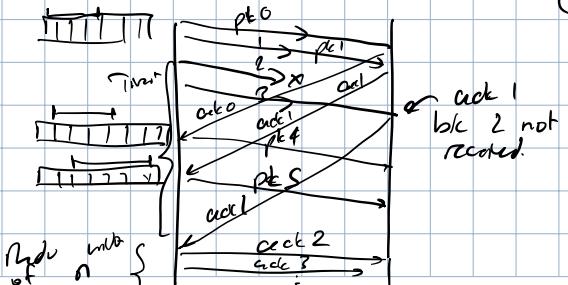
Pipelined protocols

Go-back n

- ① Window of n packets
 - ② Receiver sends cumulative ACK
 - ACK(p) → received all packets up to p (highest in-order seq #)
 - ③ If ACK(p) → update cont on Tx & transmit
 - ④ If timeout for oldest ^{uncode} packet expires ⇒ p+1
 - send next n packets starting from uncoded packet
- n=3
-
- The diagram shows a window of 3 packets. A timeout occurs for the oldest uncoded packet (blk 1), so the sender retransmits blk 1 and then blk 2. The receiver has received blk 1 and blk 2, but blk 2 was not received. The receiver sends an ACK for blk 1, but the ACK for blk 2 is lost.

Selective repeat

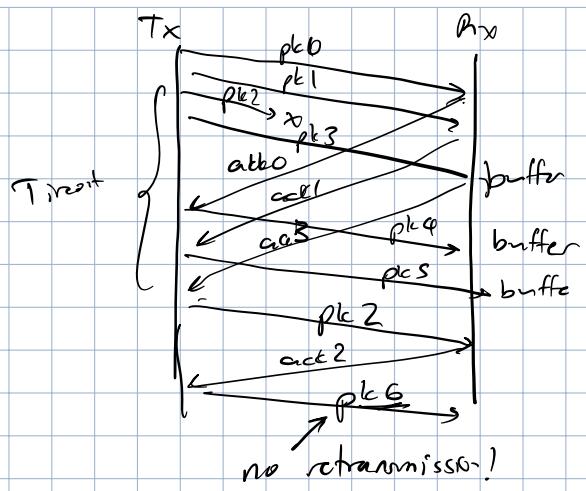
- ① Packet sent w/ own timer on window n
- ② Sender sends ack / packet
- ③ If packet / ack lost:
 - Sender repeats smallest uncoded packet
 - No duplicate packets
 - ↳ Out-of-order buffer by receiver



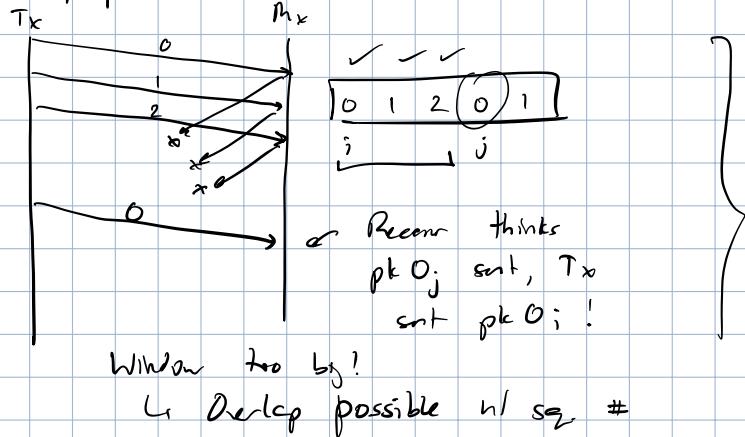
Optimize:

① Buffer packets sent in prev. windows

② Don't wait for the to retransmit



Selective repeat dilemma:



$$\text{Soln: window size} \leq \frac{\text{seq \# space}}{2}$$

Pipelined protocol perf.:

Window size W (decided on handshake)

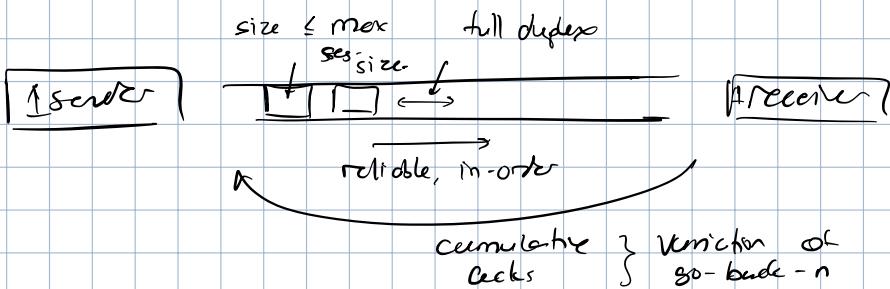
$$t_T = L/R + RTT \Rightarrow \text{time to send 1 packet \& Ack}$$

$$\underbrace{W \cdot L/R}_{\substack{\text{time to transmit} \\ W \text{ frames}}} \leq \underbrace{t_T}_{\substack{\text{1st frame} \\ \text{time}}} \Rightarrow \text{send } W \text{ packets but 1st arrives}$$

$$\text{Optimal window size} = \frac{\text{time for 1 packet round trip}}{\frac{\text{time to transmit 1 packet}}{t_i}} = \frac{t_T}{\frac{L/R}{t_i}}$$

$$\therefore \text{Utilization} = \min(1, \frac{W \cdot t_i}{t_T})$$

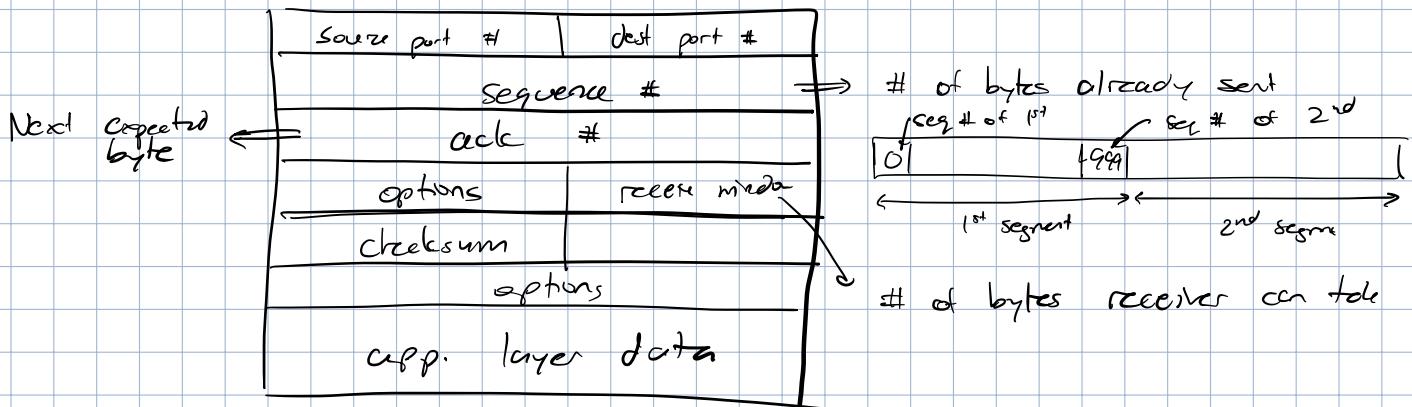
TCP



Other properties:

1. Pipelined
2. Handshaking to establish conn.
3. Flow control
4. Congestion control

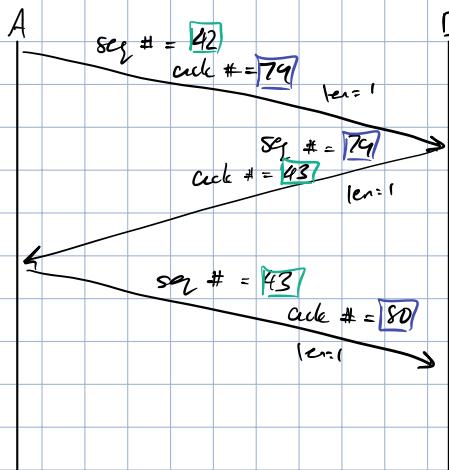
Segment structure:



If out-of-order segments: implementation can decide

- Can discard
- Can buffer (most common)

Full duplex:



TCP sender:

- Ⓐ Data from app. layer:

1. Create segment w/ seq # & receive window
2. Start timer if not running (oldest unACKed segmt)

- Ⓑ Timeout

1. Retransmit segment

2. Timer restart

(C) ACK receive

1. Update oldest unACKed

2. Update timer

TCP receiver:

(A) In-order segment arrives w/ expected seq. #. Everything up to seq # is ACKed
prev. receiver response ACK #

- Ideal

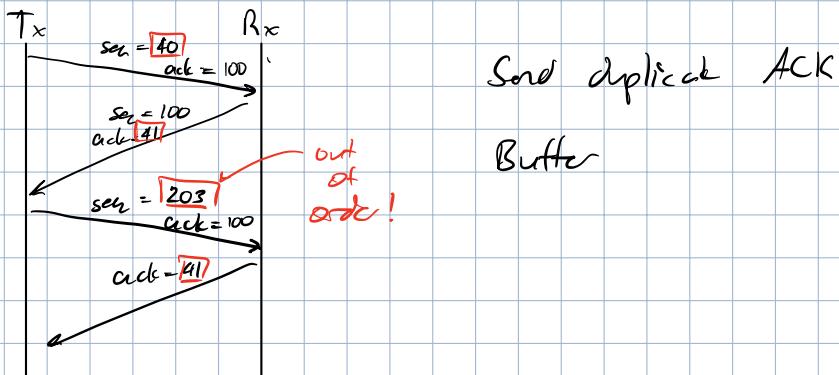
- Delay for 500ms → if no segt. send ACK

↳ Minimizes ACKs b/c cumulative

(B) In-order segment arrives w/ expected seq. #. Some other segment has not sent ACK (③)

Send cumulative ACK ⇒ ACK every 2nd packet

(C) Out-of-order segment w/ higher-than-expected seq #



(D) Segment completes gap

Send ACK if segment completely fills gap

Q: What is diff. b/w TCP & go-back-n?

TCP retransmits oldest un-ACKed packet rather than entire window

Even a scenario where it will not retransmit if cumulative ACK covers lost ACKs

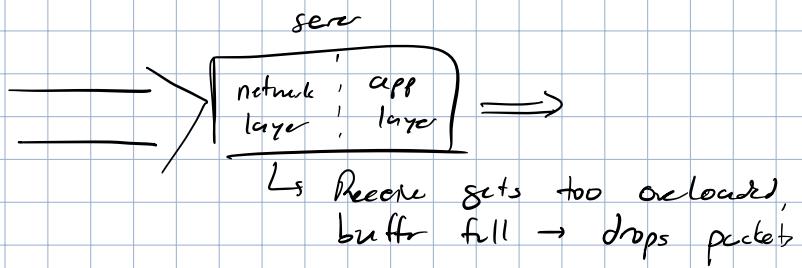
When to retransmit?

① Timeout

② Triple duplicate ACK (fast retransmit)

↳ Usually faster than timeout

Flow control for TCP

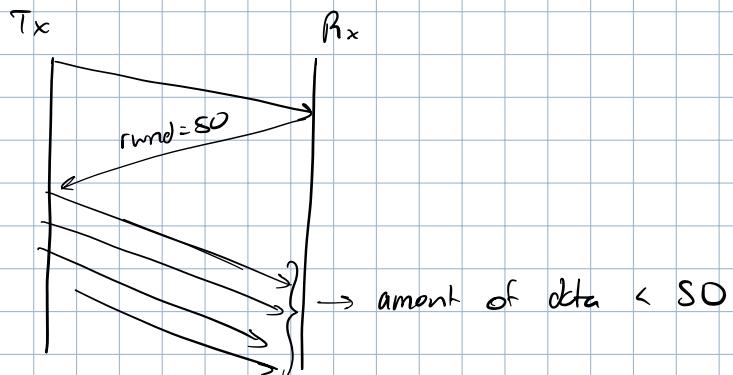


Soln:

① Receiver sets rwnd field in TCP header

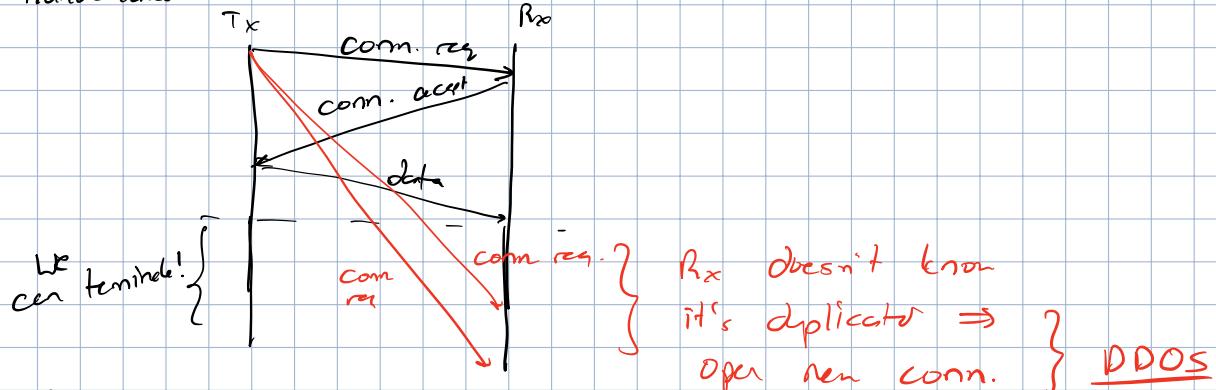
↓
amount of free space
in buffer

② Sender sets rwnd ⇒ limit amount of unACKed data to rwnd bytes.

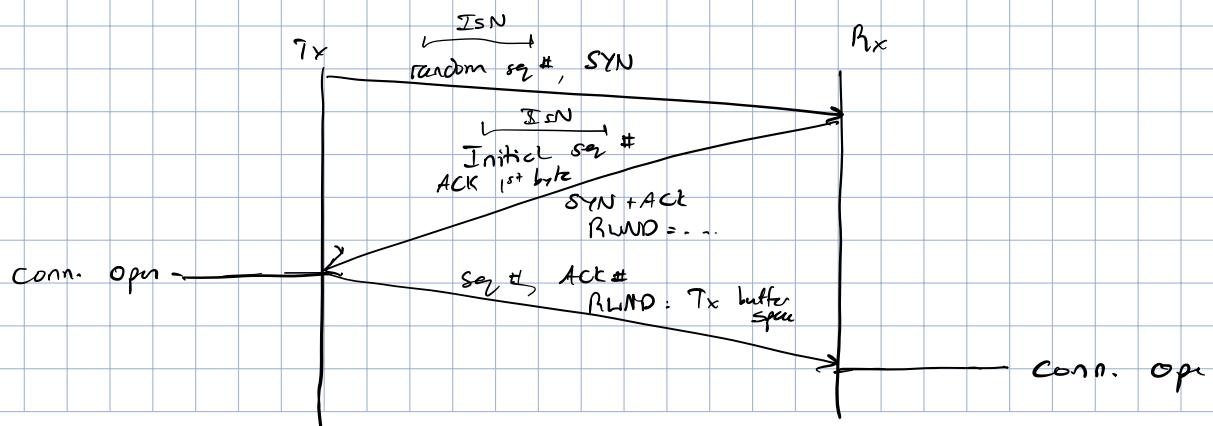


TCP Conn. Establishment

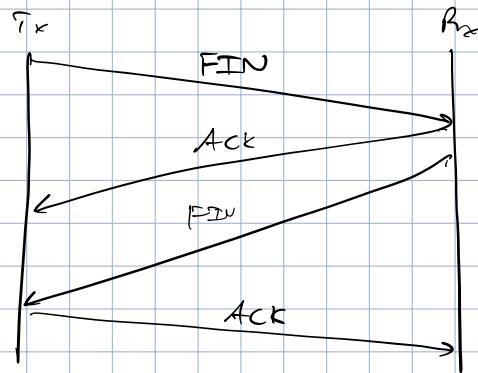
2-way handshake



3-way handshake:



TCP Closing Conn.



Congestion Control Principles

Congestion: too many senders sending traffic \rightarrow overwhelming
 ↳ Flow control: 1 sender D too fast

Approaches

E2E control

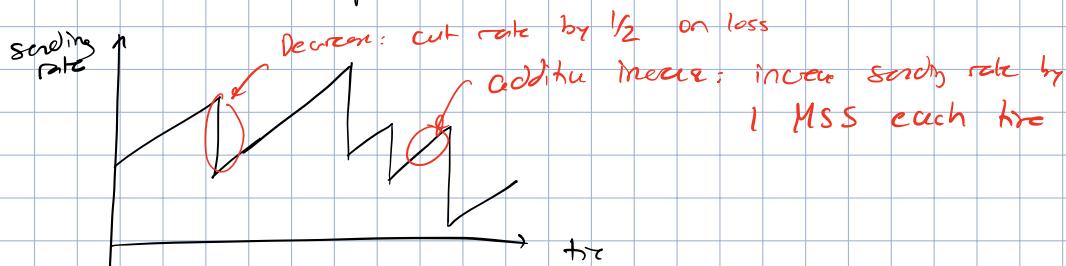
Info congestion from delays \rightarrow modify rates

Network-assisted control

Routers comm. flow rate & explicit congestion \rightarrow sources modify

TCP Congestion

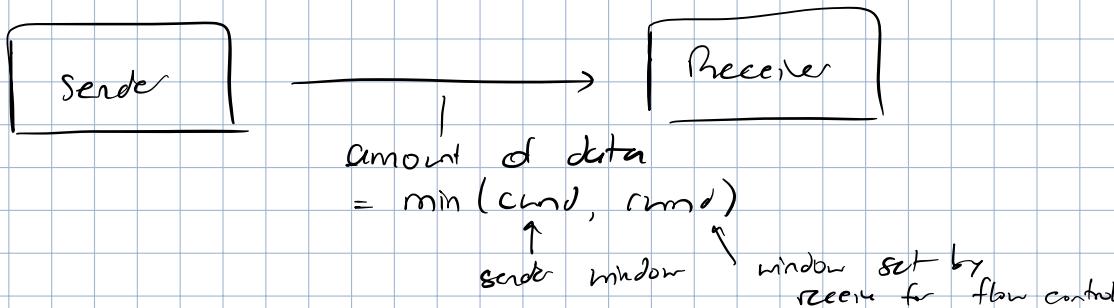
Additive increase, multiplicative decrease



} Stable & reduces congestion

Cases:

i) Ideal case



cwnd dynamically adjusted as messages received & acked

2) Timeout loss:

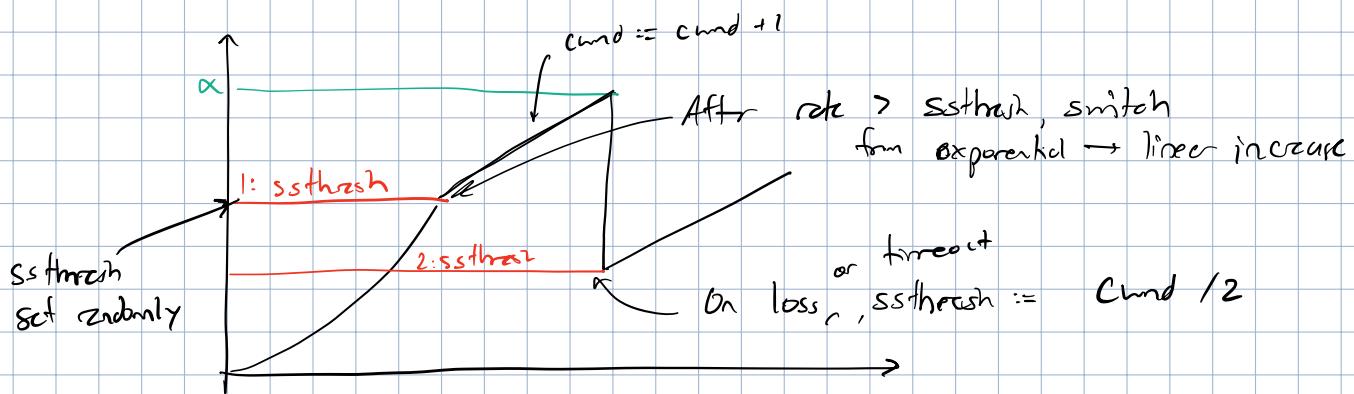
$$Cwnd = 1 \text{ MSS}$$

3) Lost packet (3 duplicate ACKs):

↪ TCP Reno: $Cwnd := Cwnd / 2$

↪ TCP Tahoe: $Cwnd := 1$

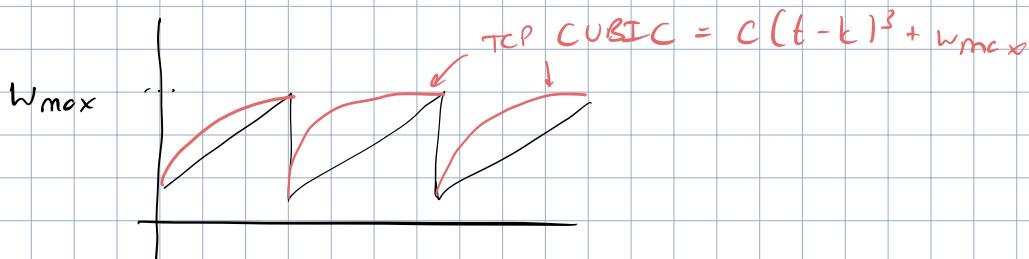
Additive increase details:



TCP throughput: $\frac{\text{avg. window size}}{\text{RTT}} = \frac{3}{4} \cdot \frac{\text{max window}}{\text{RTT}}$

Q: Rather than doing linear increase, better way?

↪ Ans: TCP CUBIC



V. popular TCP variation!

NETWORK LAYER

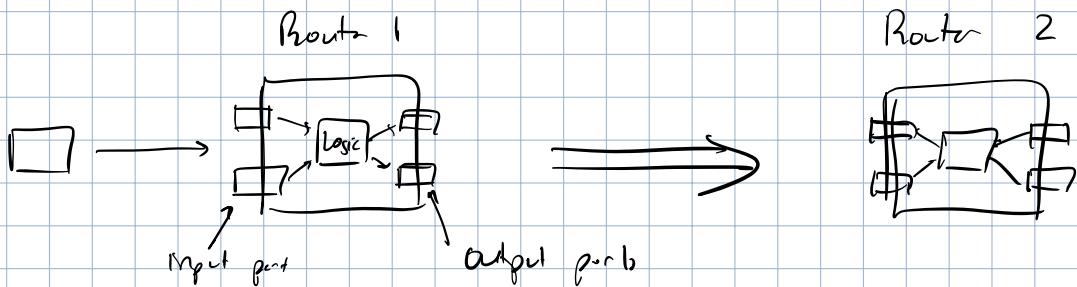
Overview

Objective:

- ① Route packets from host \rightarrow host (software algos)
 - ② Forward packets from input \rightarrow output ports (hardware)
- } Packet delivery across network

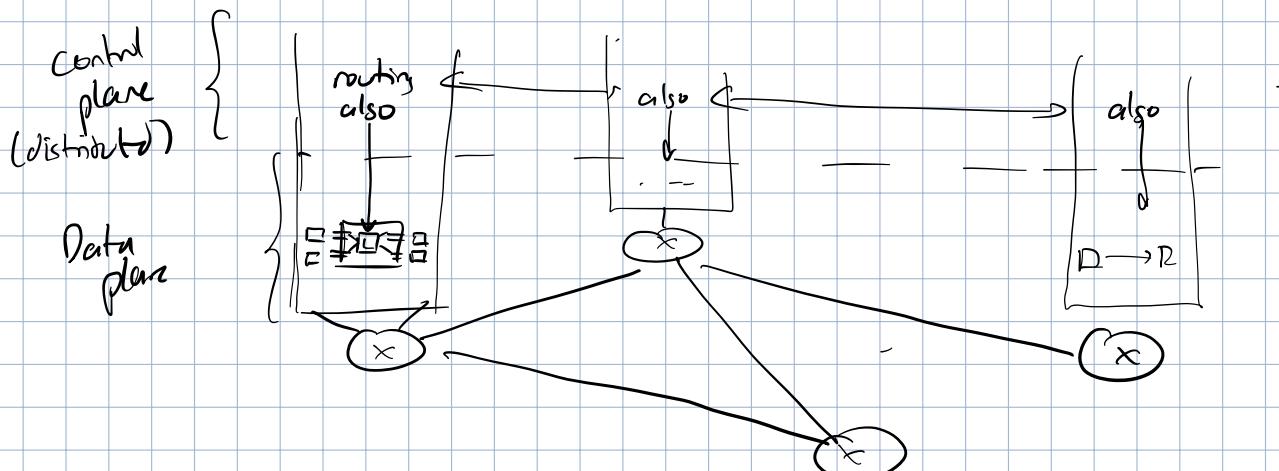
Every Internet-connected device has network layer

Core component: router

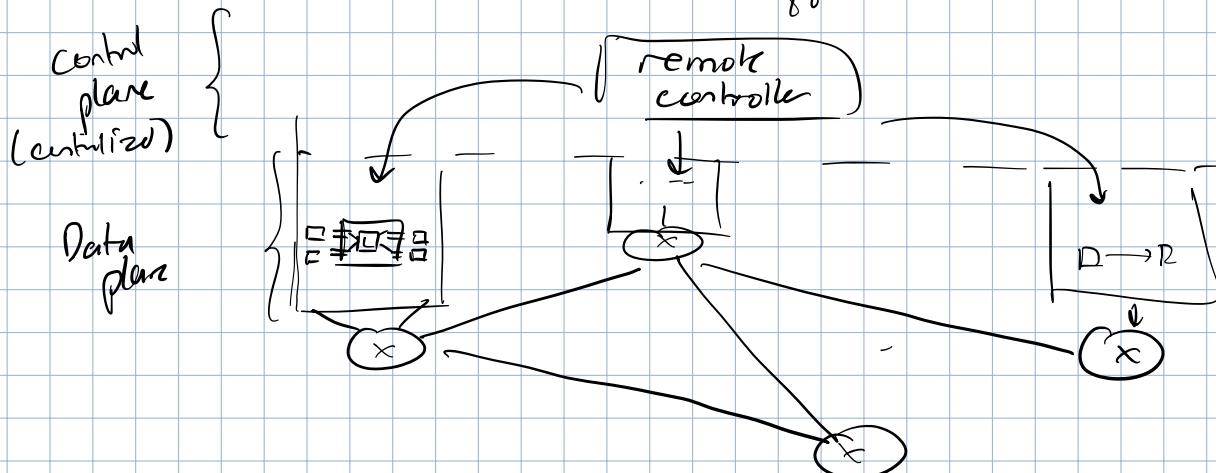


Control plane vs. data plane

① Traditional routing algo approach



② Software defined networking approach:

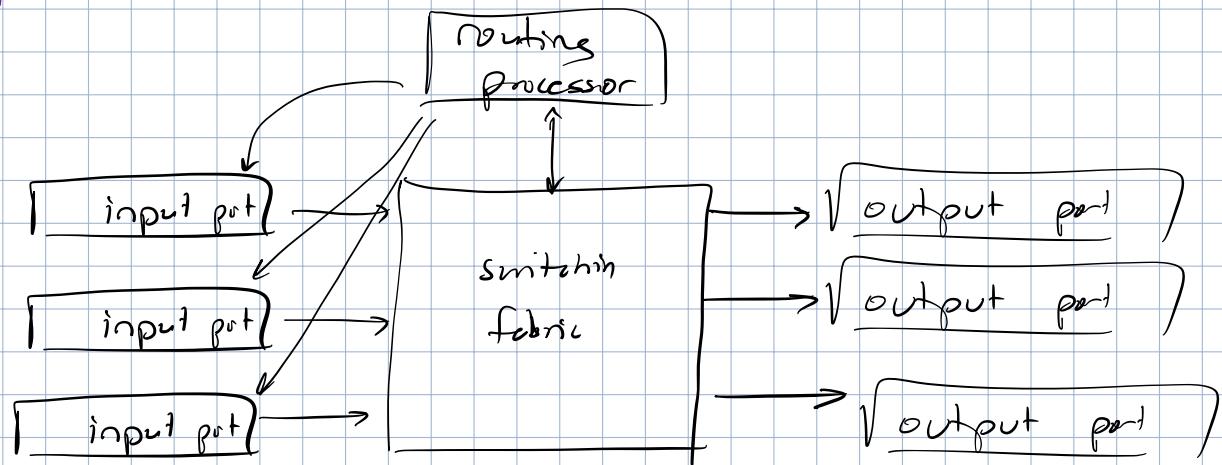


Data plane uses local forwarding table

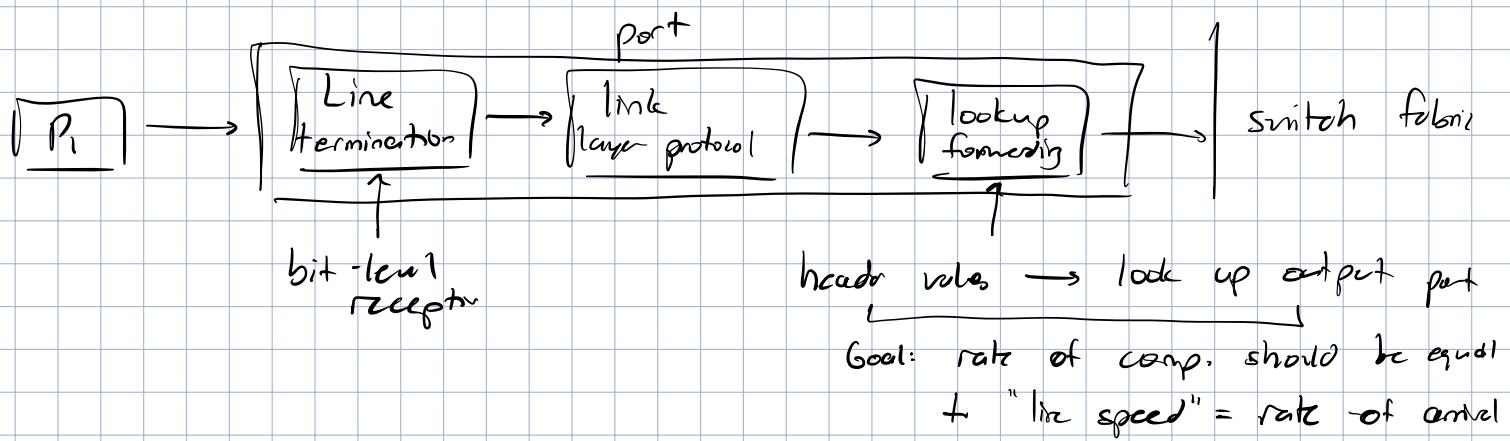
header bits	output
0111	4
1011	3

determined by control plane

Router



Input port functionality



If processing time \gg arrival rate \Rightarrow queuing

Dest. based forwarding: only use dest. IP addr. for forwarding

Switching fabric

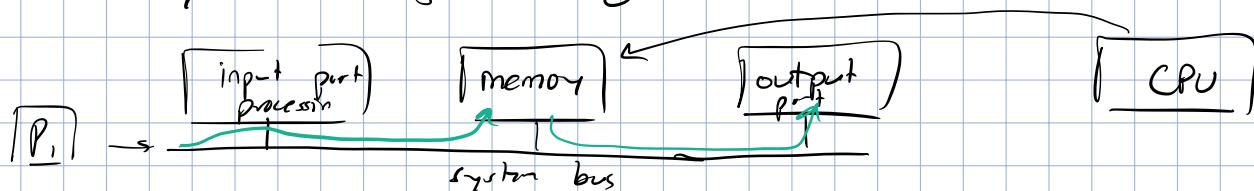
Obj: move packets from input port \rightarrow output port

Switching rate: rate of packet movement

\hookrightarrow Ideal case: NR (N : # of input ports, R : arrival rate)
otherwise queuing can occur

Types:

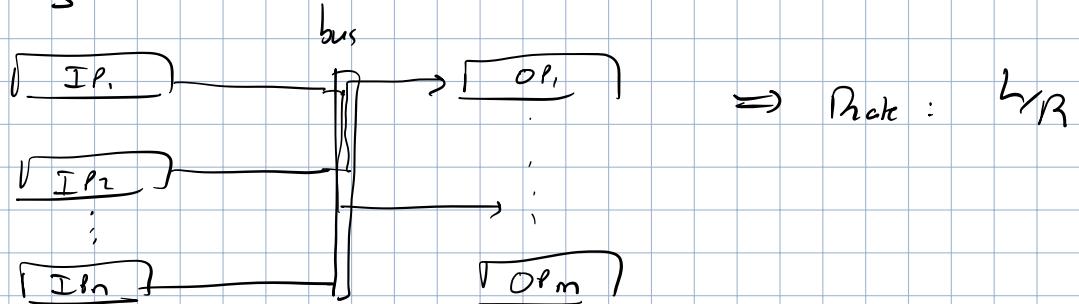
① Memory switching (1st gen routers)



Limitation:

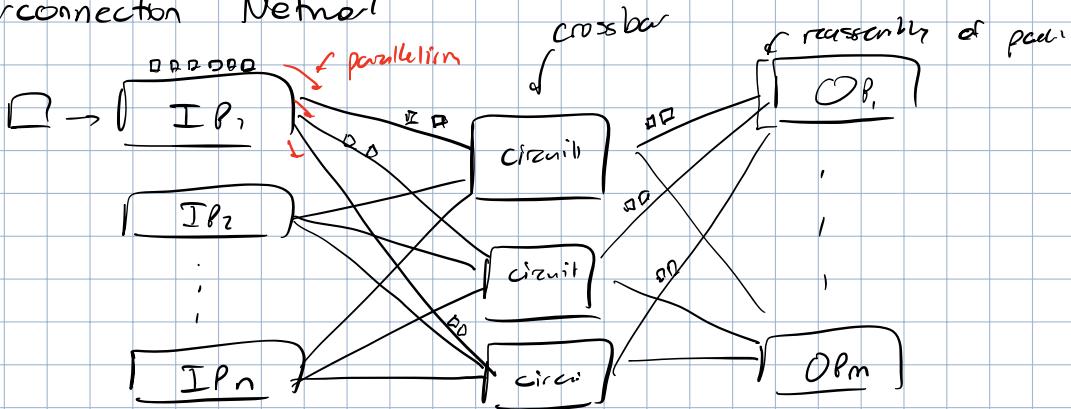
1. Packet moved 1 at the b/c packet blocks sys-bus
2. 2 bus crossings \Leftarrow limited by memory bandwidth
 \hookrightarrow Time to send: $L/R \times 2$

② Bus switching



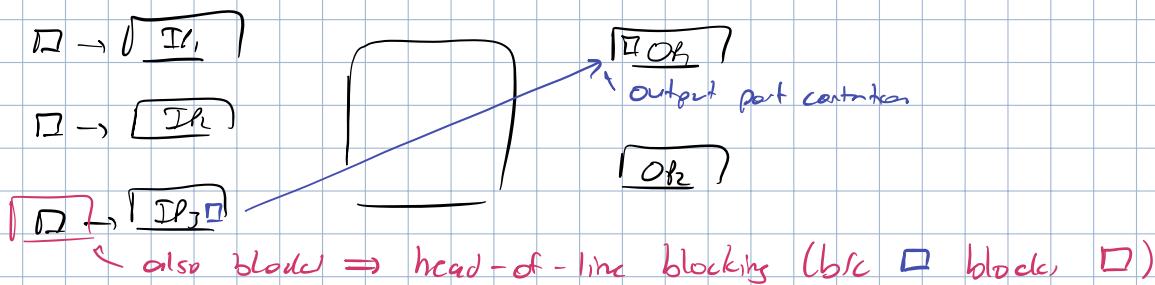
Limitation: bus contention (multiple packets go to same OP.)

③ Interconnection Network

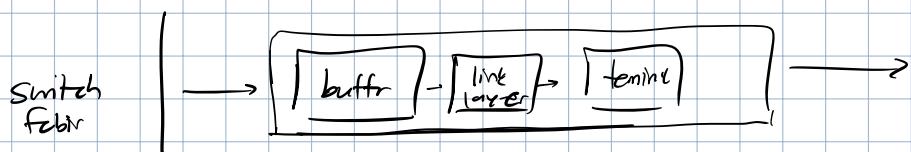


(Can add more crossbar circuits to form multistage switch)

If not fast enough \rightarrow queuing



Output port



{ Q: Why do we need buffer in OP?

} A: Link layer rate < switch fabric rate

{ Q: What if too many packets

} A: Drop policy

{ Q: FIFO service from buffer

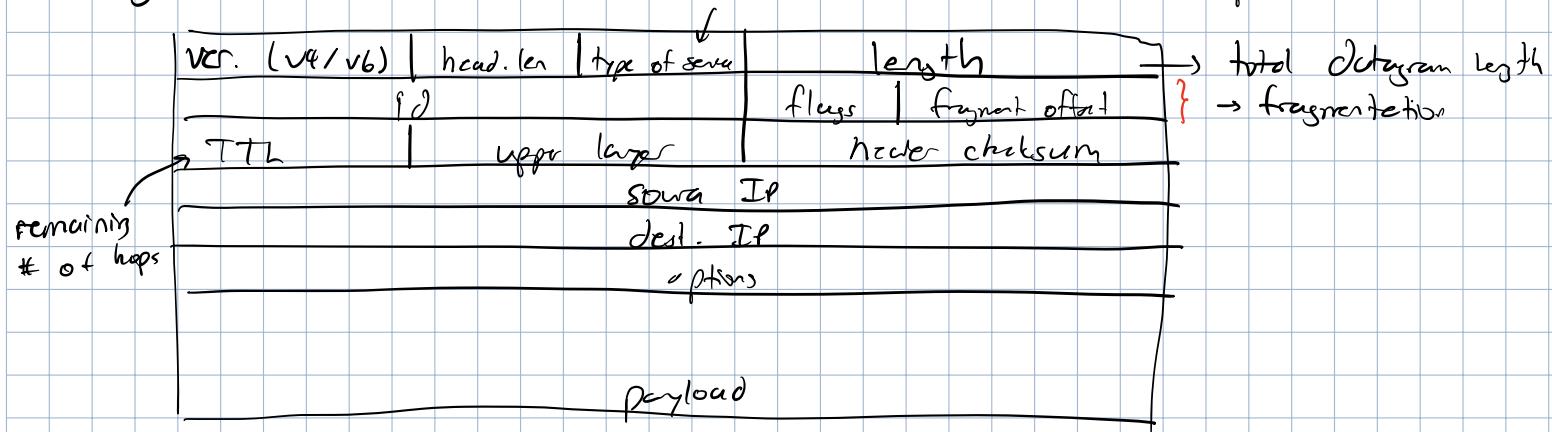
} A: No! Priority scheduling discipline

Internet Protocol (IP)

Resp.: Datagram format, addresses, packet handling

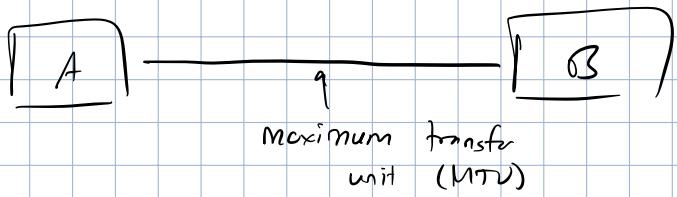
Datagram format:

congestion notifications / error type (flow control)



$$\text{Length} = 20 \text{ bytes of TCP header} + 20 \text{ bytes of IP header} + \text{payload}$$

Fragmentation:



If IP datagram $>$ MTU \Rightarrow split datagram into fragments

↳ Reassemble datagrams at destination, not at every node

If nofrag flag = 1 \Rightarrow no fragmentation, send error if $\| \text{datagram} \| > \text{MTU}$

How to reassemble fragments? Heads.

↳ Each fragment has length field, offset field (8 byte multiples)

Ex. 11

Datagram: $\boxed{\text{length} = 4000 \mid ID = x \mid \text{flag} = 0 \mid \text{offset} = 0}$, 20 byte header
 $\therefore \text{Payload} = 4000 - 20$
 $= 3980 \text{ bytes}$
 MTU: 1500 bytes

Fragments: \downarrow 1480 bytes of payload, 20 bytes of header

A: length = 1500 ID = x flag = 1 offset = 0

B: length = 1500 ID = x flag = 1 offset = 185 \downarrow 1480 bytes in 1st datagram / 8

C: length = 1040 ID = x flag = 0, offset = 870 \downarrow 1480 / 8

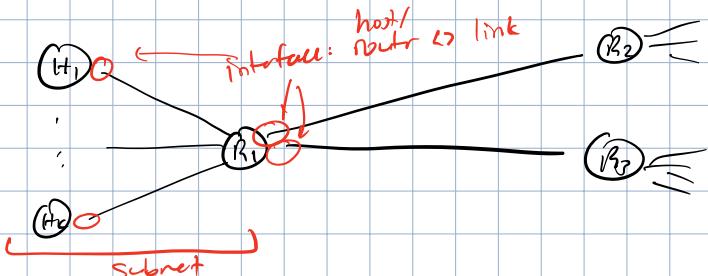
$$\# \text{ of datagrams} = \left\lceil \frac{\text{payload size}}{\text{MTU} - \text{header size}} \right\rceil$$

$$\begin{aligned} \text{len}_1 &= \text{remaining} + 20 \text{ bytes} \\ &= (3980 - 2 \times 1480) + 20 \\ &= 1040 \end{aligned}$$

\therefore MTU \geq header size \Rightarrow if not, non-header data cannot be sent

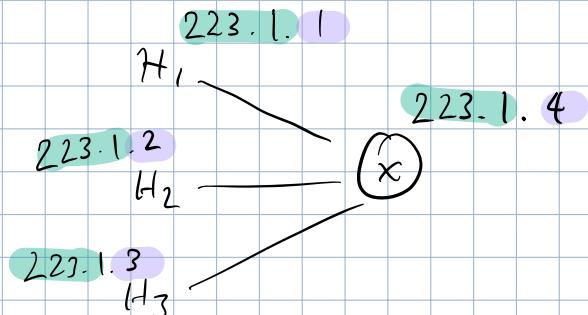
IP Addressing

IP address: 32-bit ID for a host or router interface



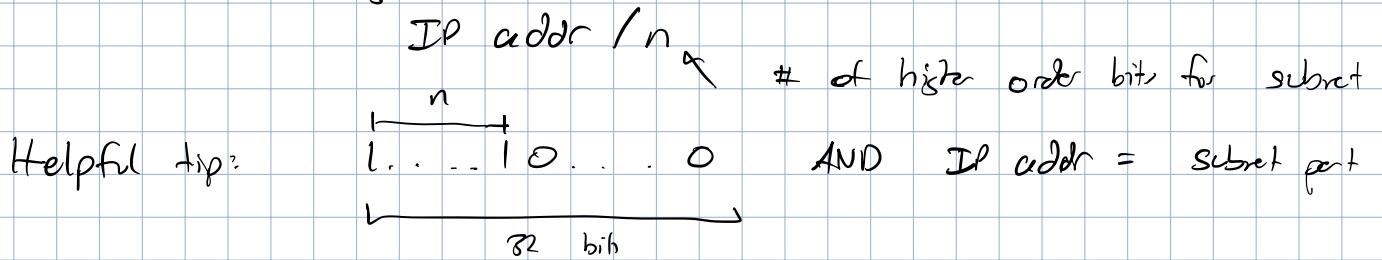
Subnet: Collection of devices that can communicate w/out route

- IP address
- Subnet part
 - Higher order bits shared by all devices in subnet
 - Host part
 - Unique IDs / host



Q: how to know what part is subnet & host?

A: Subnet masking



Q: How to assign IP addresses to subnet?

① Classful addressing

Fixed length prefixes

A: Prefix of 8 bits, starts w/ 0 $\Rightarrow 2^{24}$ hosts

B: Prefix of 16 bits, starts w/ 10 $\Rightarrow 2^{16}$ hosts

C: Prefix of 24 bits, starts w/ 110 $\Rightarrow 2^8$ hosts

Pros: simple

Con: fixed # of addresses

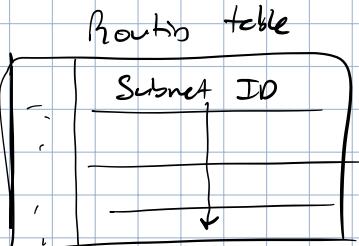
② Classless

CIDR: classless inter-domain routing

Subnet masking: a.b.c.d /n

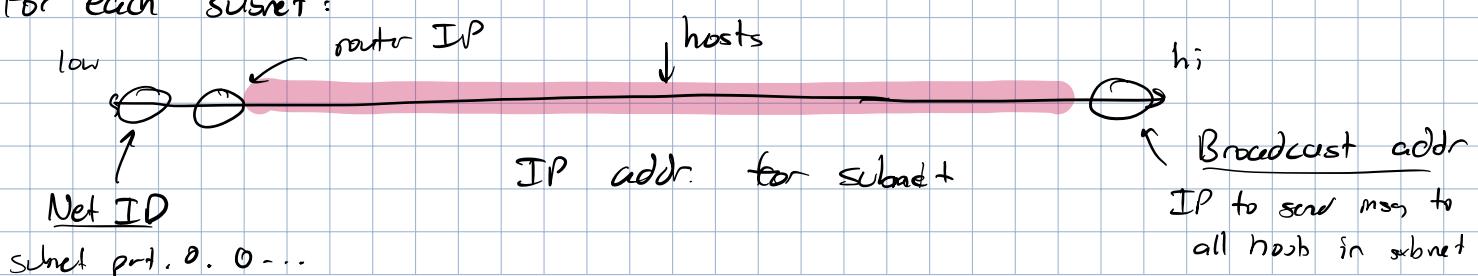
* Important: Subnet mask is consec.. Is followed by consec. 0s *

Q: How to send messages w/ IP addresses & subnet masks



\Rightarrow Does not have host part

For each subnet:



Ex:// Network has IP 10.2.5.16 /28. Router IP? Broadcast IP?

① Subnet mask:

11111111. 11111111. 11111111. 11110000
255. 255. 255. 240

128
64
32
16

② Network ID:

$\begin{array}{r} 10 \cdot 2 \cdot 5 \cdot \\ \hline 255 \cdot 255 \cdot 255 \cdot 240 \end{array} \rightarrow \begin{array}{l} 16 \rightarrow 00010000 \\ \hline 11110000 \\ 00010000 \end{array}$ ↗ No 1's
↓
10.2.5.16

③ Broadcast IP:

Make host part all 1's

10.2.5 | 00011111 ↗ 31 ⇒ 10.2.5.31

Ex:// Given IP addresses 192.168.16.0 /24. We need to create 4 subnets w/ following # of hosts.

- o LAN A: 120 hosts
- o LAN B: 50 hosts
- o LAN C: 12 hosts
- LAN D: 28 hosts

① Find # of addresses needed / LAN

$$A: 120 + 2 = 122$$

↖ network + broadcast

$$B: 52$$

$$C: 14$$

$$D: 30$$

② Find least # of bits necessary per network to distinguish host.

$$A: \lfloor \log_2 122 \rfloor + 1 = 7 \text{ bits}$$

$$B: \lfloor \log_2 52 \rfloor + 1 = 6 \text{ bits}$$

$$C: \lfloor \log_2 14 \rfloor + 1 = 4 \text{ bits}$$

$$D: \lfloor \log_2 20 \rfloor + 1 = 5 \text{ bits}$$

③ Subnet mask / network

Formula: $32 - \text{host bit } \#$

$$A: 32 - 7 = 25 \text{ bits}$$

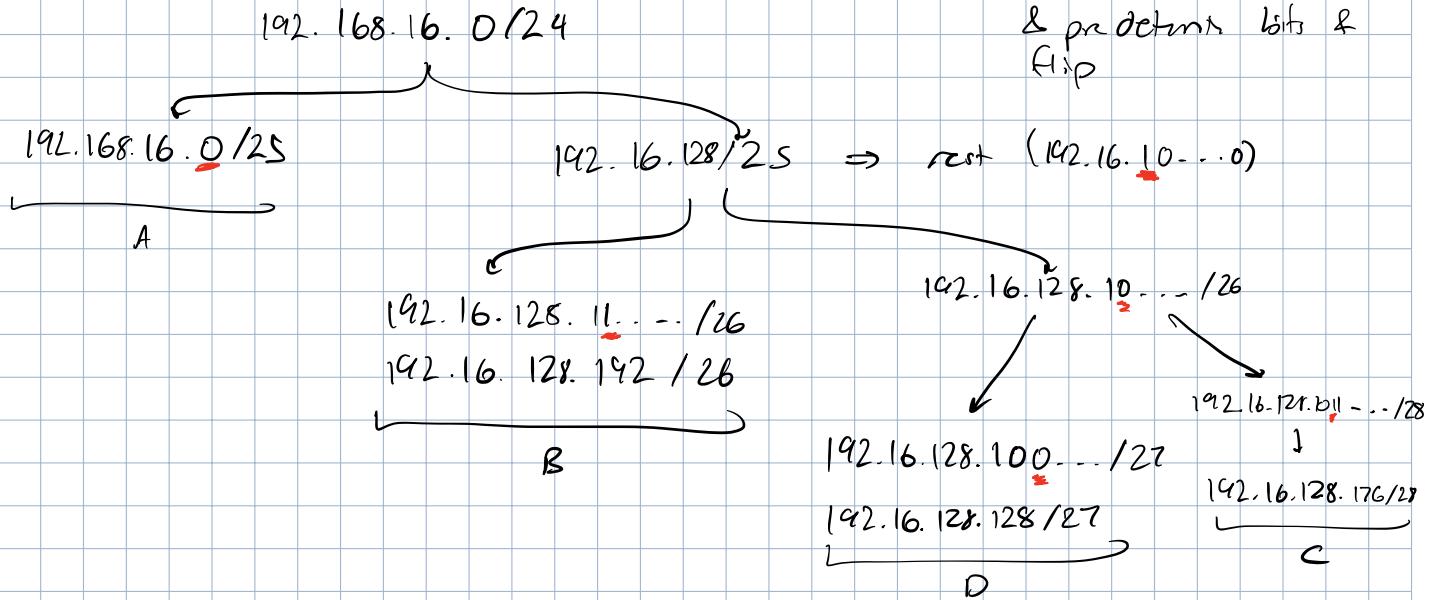
$$B: 32 - 6 = 26 \text{ bits}$$

$$C: 32 - 4 = 28 \text{ bits}$$

$$D: 32 - 5 = 27 \text{ bits}$$

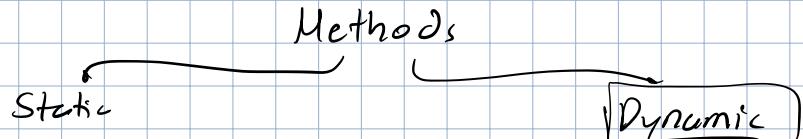
④ Split

Idea: work from large \rightarrow small
& pre-determine bits &
flip

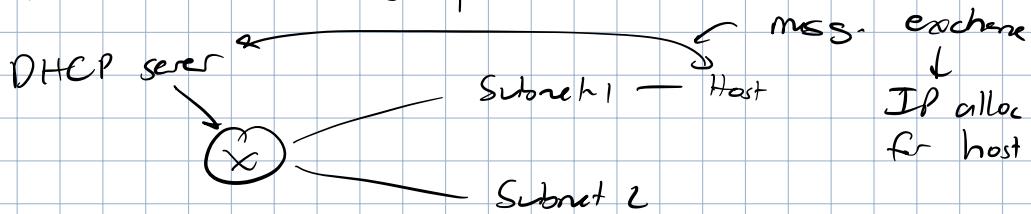


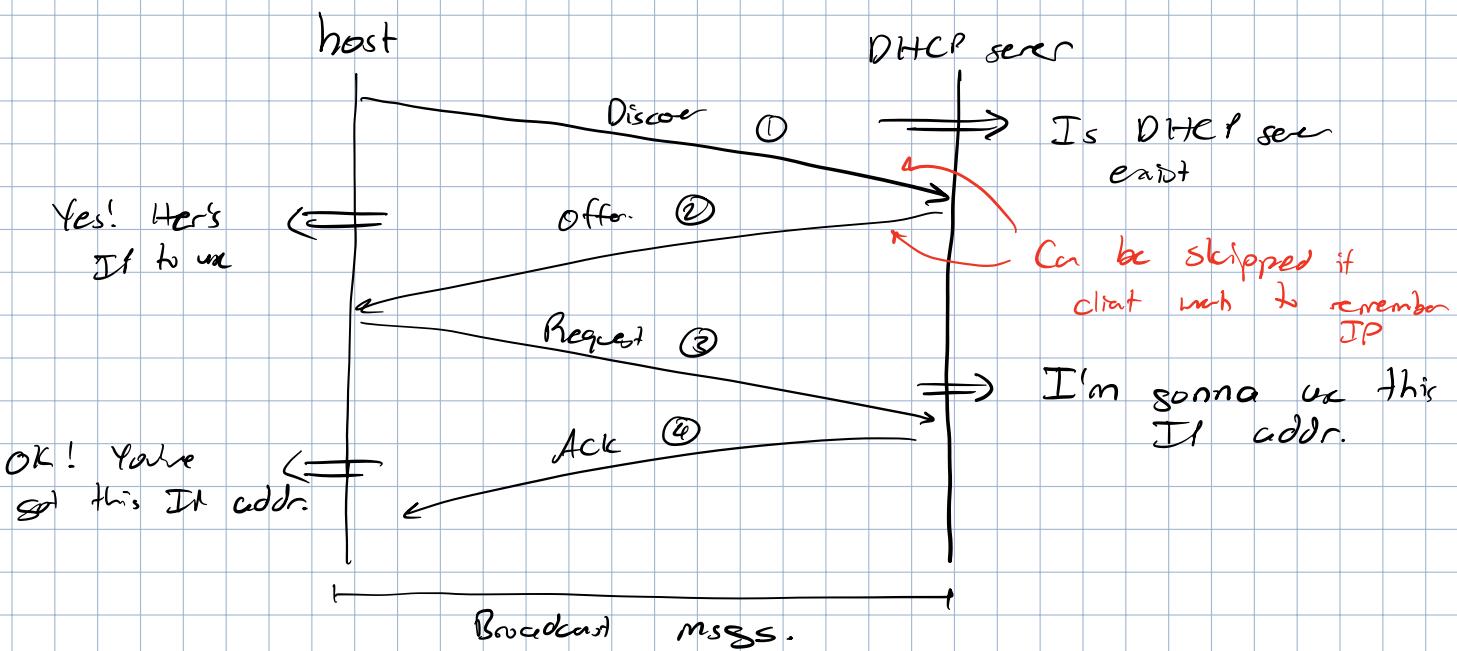
⑤ *Optional but useful* Calculate network & broadcast IPs to ensure no overlap!

Q: How can host get own IP address?



A: DHCP (dynamic host config. proto.)





① src : 0.0.0.0 , src port #
 dest : 255.255.255.255, dest port #
 yaddr : 0.0.0.0
 transaction ID: ID for entire process.

② src: DHCP IP addr., port #
 dest: dest of ①
 yaddr: host IP addr. it can use
 ID:
 lifetime: TTL for IP addr. \Rightarrow request for newer.

③ src: src of ①
 dest: broadcast IP
 yaddr: host IP addr. in ①
 ID:
 lifetime:

④ src: src of ②
 dest: dest of ③
 yaddr:
 ID:
 lifetime:

VDP

DHCP relay passes router, DNS IP + network mask.

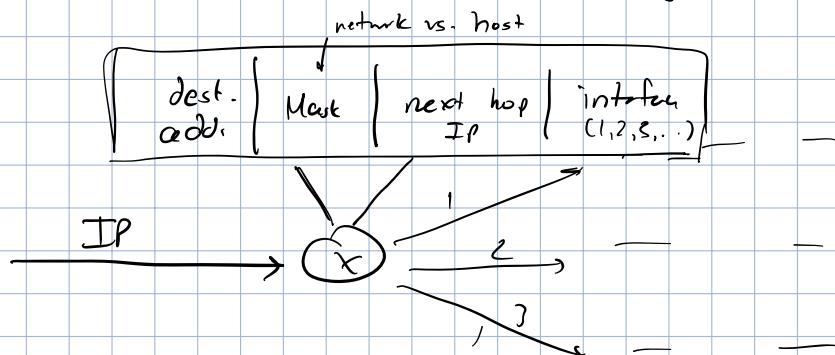
Q: How does network get IP

A: ICANN \rightarrow IP addresses w/ regional registries

Not enough IP addr in IPv4 (32 bit) \rightarrow IPv6 (128 bit addr.)

Packet forwarding

Most routers use table-driven routing:

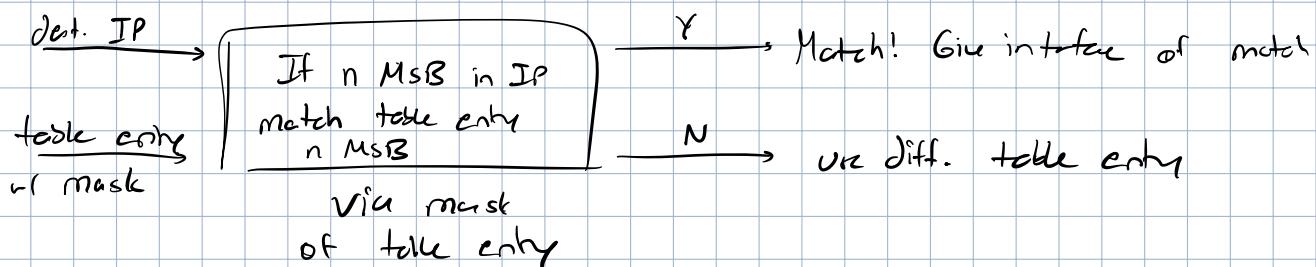


Each table has default entry if no matches

Routing protocols det. table content

Q: How to figure out right interface given dest. IP?

A:



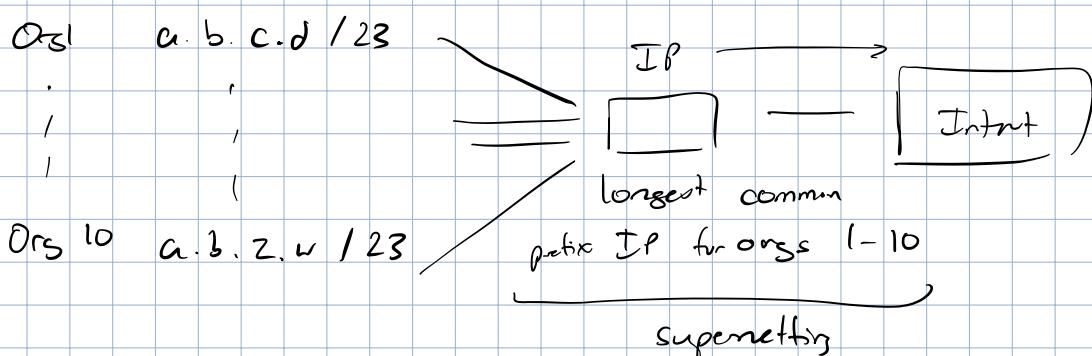
If multiple matches happen, choose entry w/ longest match bits

Ex:// If IP1 matches w/ entry 1 for 5 bits, matches w/ entry 2 for 10 bits, choose entry 2 interface.

Hierarchical addressing

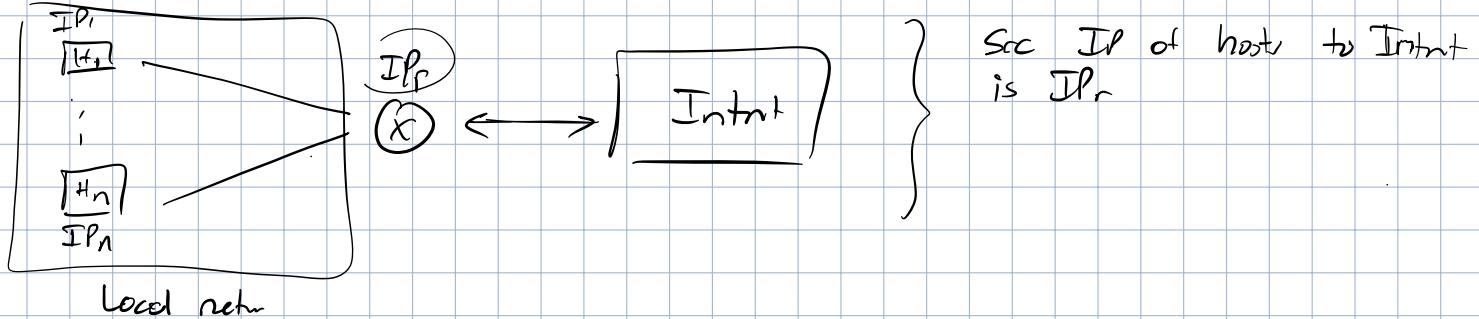
Problem: we don't want to publish all IPs of networks → tables get bloated.

Sol: hierarchical addressing



Network address translation (NAT)

Idea: do hierarchical addressing at smaller scale



IP_1, \dots, IP_n have same prefix: 10/8, 172.16/12, 192.168/16

↳ For any local network. Only IP_x has to be unique

Pros:

- ① Reduces # of unique IP addresses needed
- ② Can change IP of host w/out entire Internet knowing } decoupling
- ③ Can change ISP IP w/out changing host IP
- ④ Security b/c local host not directly addressable

Cons:

- ① Routers shouldn't process below layer 3
- ② Violates port-to-port agreement ← Router is changing port #
- ③ How to connect to server directly if behind NAT?

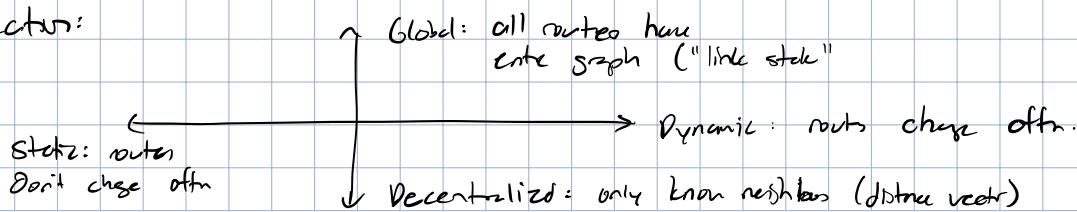
Responsibilities:

- ① Outgoing datagrams: replace src IP w NAT IP, port #
- ② Remember (NAT IP, port #) ↔ (src IP, port #)
- ③ Replace NAT IP w src IP for incoming datagrams

Routing protocols

Resp.: figure out best path from $H_1 \rightarrow H_2$

Classifications:



Link state

Centralizes \leftarrow entire network known to each host

Finds shortest path from node s to all nodes in network

↳ After k iterations of algo, we know least cost path up to k dest.

Algorithm (Dijkstra's algo):

Initialization:

$$N' = \{u\}$$

\forall nodes v :

if v is adj to u :

$$D(v) = c_{u,v}$$

else:

$$D(v) = \infty$$

Loop:

find node w not in N' s.t. $D(w)$ is minimum

$$N' := N' + \{w\}$$

for all nodes v adj. to w :

$$D(v) = \min(D(v), D(w) + c_{v,w})$$

until $|N'| = \#$ of total nodes.

Do algo for each node to figure out least cost path \Rightarrow routing table.

Distance vector

Based off Bellman-Ford DP algo.

Intuition:

① $D_x(y)$: least cost path from $x \rightarrow y$ \rightarrow distance vector

② $D_x(y) = \min_v \{c_{x,v} + D_v(y)\}$

$\nwarrow v$: all neighbors of x

Whichever node is min is next hop on way to dest.

In practice:

① Host send own distance vector to neighbors occasionally

② Neighbor nodes update own distance vector via DP algo

③ Should converge to actual least cost paths

④ Notify neighbors if DV changes

Properties: decentralized, async, self-stopping

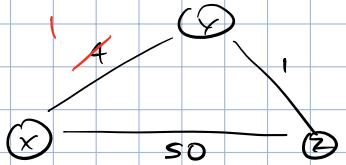
Initially: DV of node is only costs to neighbors.

Fast to link cost changes if good, slow if not good

Ex://

At t_0 : y detects change \rightarrow notify neighbor

$$D_y = \begin{bmatrix} D_x \\ D_z \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$



At t_1 : z gets update \rightarrow update

$$D_{z_bthr} = \begin{bmatrix} D_x \\ D_y \end{bmatrix} = \begin{bmatrix} 50 \\ 1 \end{bmatrix}$$

$$D_{z_aft} = \begin{bmatrix} D_x \\ D_y \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

At t_2 : y gets z 's update:

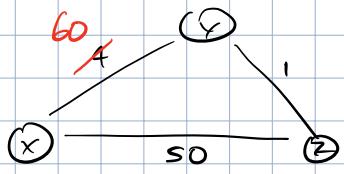
$$D_y \text{ before} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \xleftarrow{\text{Equal}} \quad \xrightarrow{\text{Equal}}$$

$$D_y \text{ aft} = \begin{bmatrix} D_x \\ D_z \end{bmatrix} = \begin{bmatrix} \min(1, 1 + D_{z,x} = 3) \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Stop notif.

Ex://

to: y detects and updates



$$D_y = \begin{bmatrix} D_x \\ D_z \end{bmatrix} = \begin{bmatrix} 4 \\ 1 \end{bmatrix}$$

$$D_{y,a} = \begin{bmatrix} D_x \\ D_z \end{bmatrix} = \begin{bmatrix} \min(60, D_{z,x} + 1 = 5 + 1 = 6) = 6 \\ 1 \end{bmatrix}$$

t_1 : z detects & update

$$D_{z,b} = \begin{bmatrix} D_x \\ D_y \end{bmatrix} = \begin{bmatrix} 5 \\ 1 \end{bmatrix}$$

$$D_{z,a} = \begin{bmatrix} \min(50, 1 + D_{y,x} = 1 + 6 = 7) = 7 \\ 1 \end{bmatrix}$$

t_2 : y detects & update

$$D_{y,a} = \begin{bmatrix} \min(60, 1 + D_{z,x} = 1 + 7 = 8) = 8 \\ 1 \end{bmatrix}$$

Count to infinity problem

2-node instability problem

(x) — (y) — (x) \rightarrow takes long time for both nodes to realize infinite cost to get to x

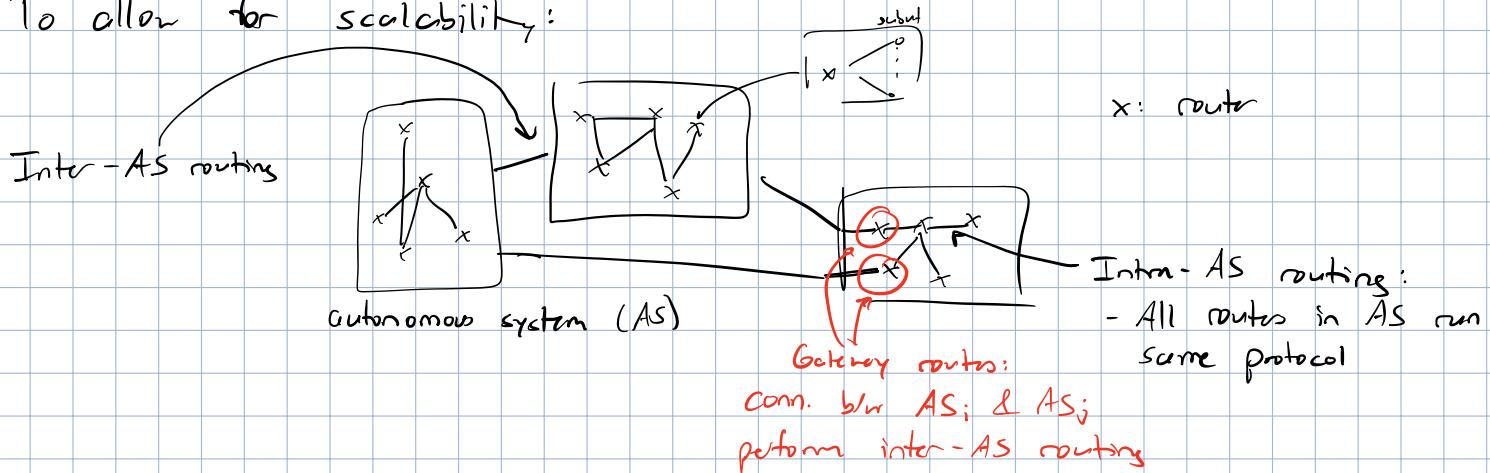
Soln:

- ① Split horizon: B doesn't tell A about path $X \leftarrow A$, assumes A knows.
- ② Poisoned reverse: A tells B about infinite cost when getting back to $X \rightarrow B$ steps.

OSPF

We want to scale up soln. to routing & allow for autonomy in networks

To allow for scalability:



These routing mechanisms work together to build forwarding tables

Responsibilities of inter-domain routing:

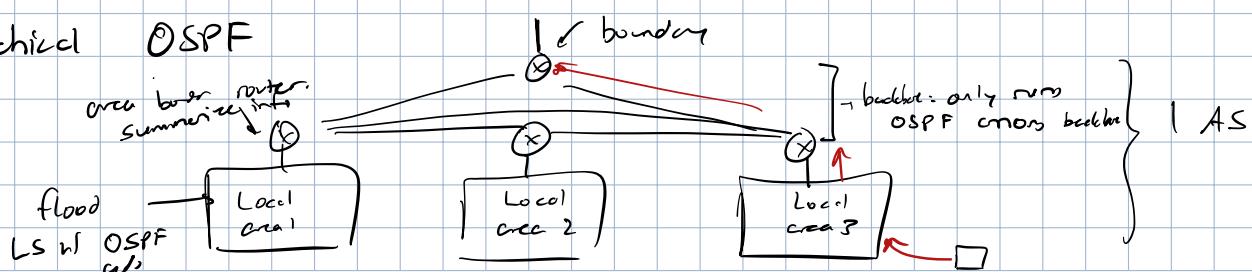
- ① Learn which dest. reachable through neighboring AS
- ② Propagate reachability info to all routers in AS

Intra-domain protocols to route within AS:

- ① RIP (routing info. proto.): exchange DVs within AS, not used
- ② EIGRP (enhanced interior gateway routing proto.): uses DV, used to be Cisco prop.
- ③ OSPF (open shortest path first): link-state

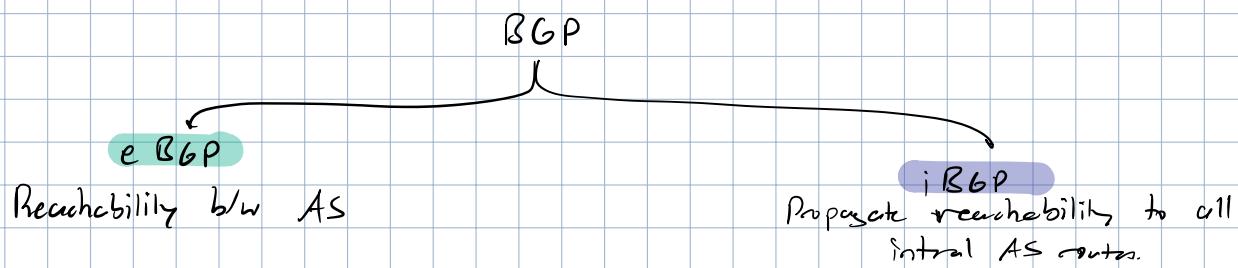
↳ Each router floods OSPF adverts to all routers in AS w/ costs,
↳ use Dijkstra's to figure out best path to each node. Uses IP, not TCP/UDP
↳ Authentication → secure.

Hierarchical OSPF

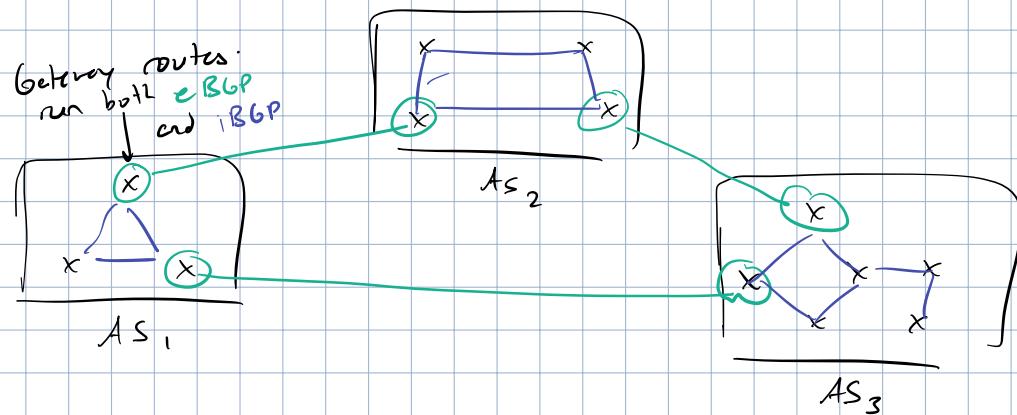


BGP

Border Gateway Protocol: standard inter-AS routing proto.

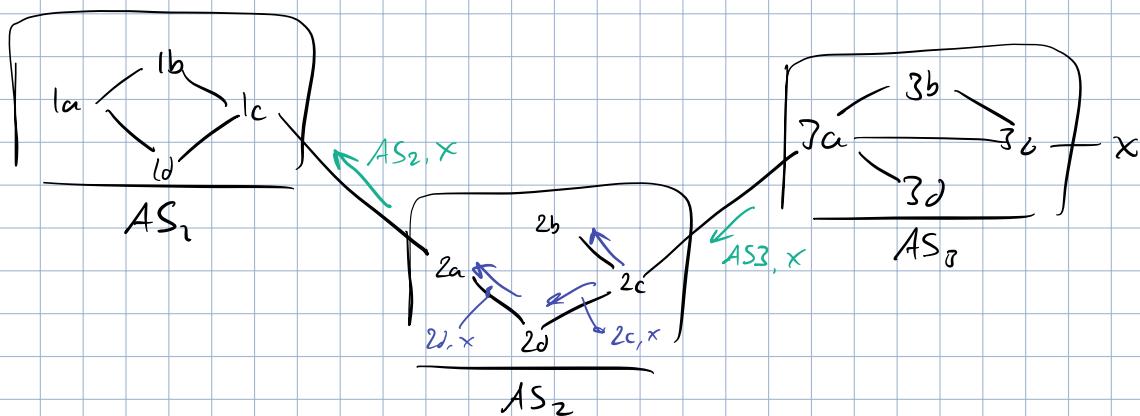


Each AS uses policies to determine best route based off reachability



When new router added to AS, BGP will announce it can reach router to neighboring AS

↳ Across semi-persistent TCP conn.



BGP msg components

Prefix
Destination address

Attributes

① AS-PATH: list of AS which ad has passed

② NEXT-HOP: most recent router

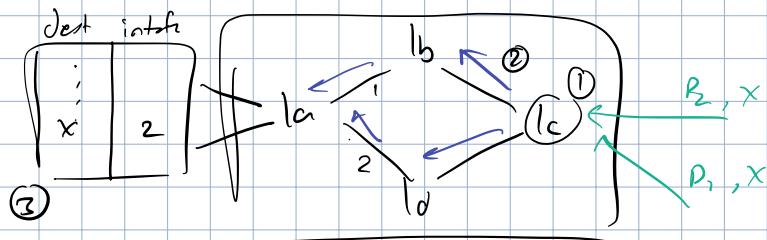
BGP msg types

<u>OPEN</u>	<u>UPDATE</u>	<u>KEEPALIVE</u>	<u>NOTIF.</u>
Open TCP conn. W/ BGP peers	Advertise new path	Keep conn. alive for fast updates. Also an ACK for OPEN	Report mistake / close conn.

When ad received by gateway, policy in action

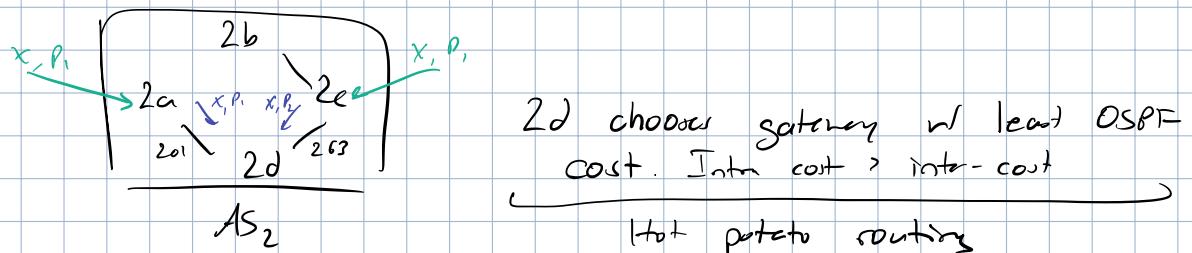
- ① Choose to reject ad (import policy)
- ② Choose to advertise new path to other AS (e.g. don't want other providers to know)
- ③ Choose to keep existing route (e.g. shortest AS-PATH, closest NEXT-HOP, other policy)

Connection w/ intra-AS routing:



- ① 1c makes decision to use P2 via eBGP
- ② iBGP propagates to all routers about X and to use 1c
- ③ OSPF/route: tells router which interface to use to get to 1c if X is dest.

What happens if internal router gets 2 iBGP msgs?

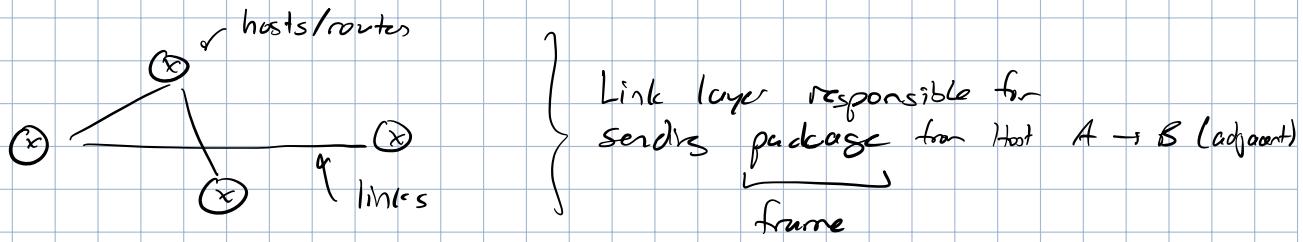


Why do we need distinction b/w intra vs. inter routing?

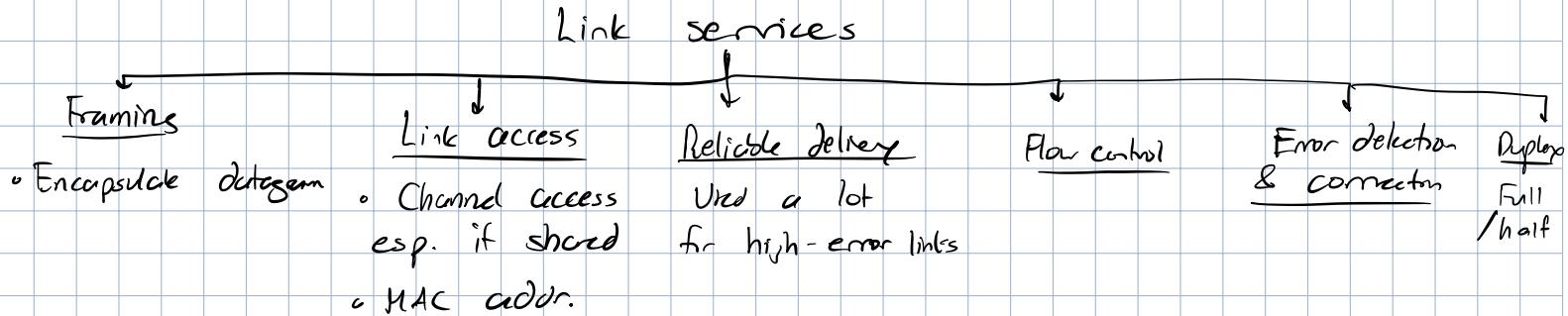
- ① Policy: inter-domain vs. more policy considerations
- ② Scale: hierarchical tables
- ③ Performance: intra focuses on perf., inter focuses on policy

LINK LAYER

Introduction



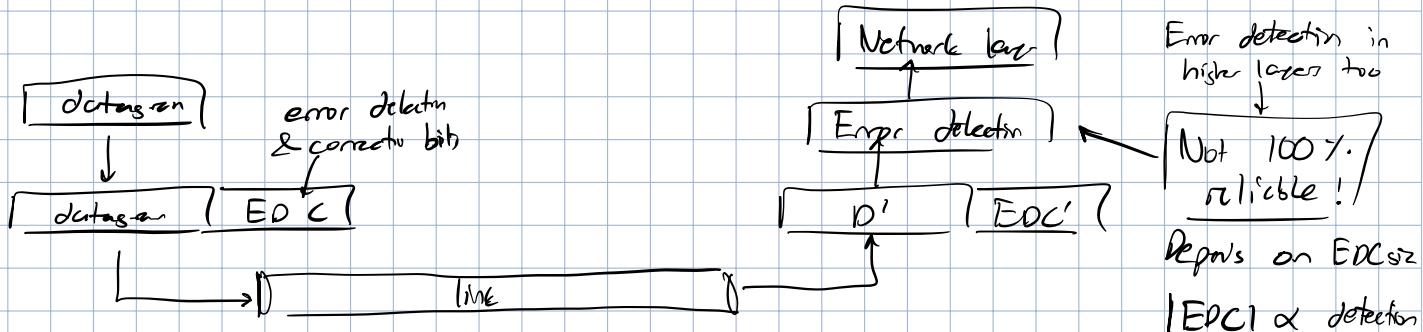
Link protocol depends on link type



Implementation: network interface card (NIC) on chip in host

↳ Combo of SW, HW & firmware

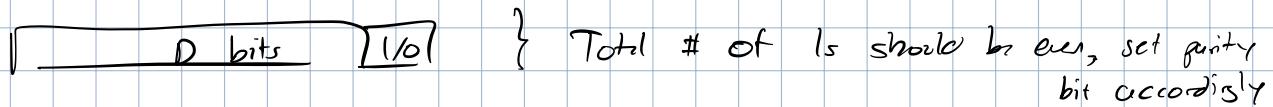
Error detection & correction



Methods:

① Parity checking

A: Single bit:



Check on receiving side if parity of 1 bit is same

Con: some parity can be achieved even if errors exist

↳ Eg: 0 → 1 & 1 → 0

Hard to correct! Which bit flipped?

B: 2D bit parity

Ex:// Msg: 10101 11110 01110

Sending side

① Table

b_1	[1 0 1 0 1]	
b_2	[1 1 1 1 0]	
b_3	[0 1 1 1 0]	

② Row & col parities

	1 0 1 0 1 1	row
	1 1 1 1 0 0	
	0 1 1 1 0 1	
col	[0 0 1 0 1 0]	

parity of parity,
choose either row or col

③ Add bits to msg

10101 111100 011101 001010

Add row parity bit to each group

Make one group w/ col parity.

Error detection:

Ans:

1 0 1 0 1 1	
1 1 1 1 0 0	
0 1 1 1 0 1	
0 0 1 0 1 0	

Col & row error detected

↳ Correct!

Can only detect up to all 3-bit errors.

Receiving side

① Reconstruct table from bits

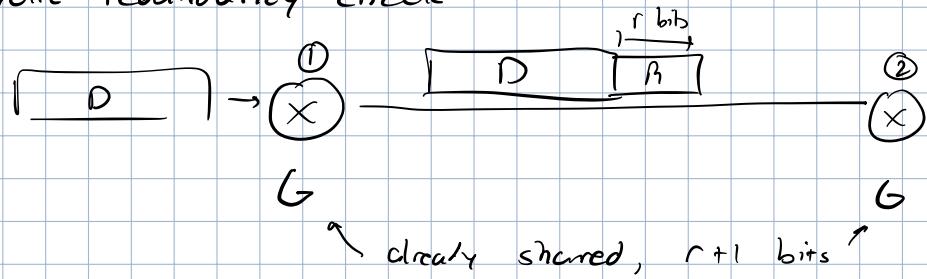
1 0 1 0 1 1	
1 1 1 1 0 0	
0 1 1 1 0 1	
0 0 1 0 1 0	

② Reconstruct table to verify
parities → error detection.

1 0 1 0 1 1	
1 1 1 1 0 0	
0 1 1 1 0 1	
0 0 1 0 1 0	

Error detected but not correctable

② Cyclic redundancy check



$$\textcircled{1} \quad \langle D, R \rangle = D \times 2^r \text{ XOR } R$$

Choose r CRC bits R s.t. $\langle D, R \rangle \bmod (G \setminus 2) = 0$

- $\textcircled{2}$ Create $\langle D, R \rangle$ & check if $\langle D, R \rangle$ divisible by $G \bmod 2$. If not \rightarrow error

} Detect errors
 $\langle r+1 \text{ bits} \rangle$

Ex: $D = 101110$. $G = 1001$. Find CRC bits

- ① Find r (# of bits of CRC)

G has 4 bits $\rightarrow r = 3$ bits

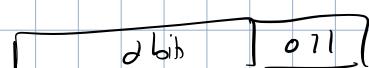
- ② Bit shift D by r bit

101110 \longrightarrow 101110 000

- ③ 'Divide' (XOR) $D \times 2^r$ by $G \rightarrow$ remainder of r bits

Handwritten long division diagram for CRC. The divisor G is 1001. The dividend is 101110 followed by three zeros. The quotient is 10101 and the remainder is 001. Annotations include: "Cannot match" with arrows pointing to the first two digits of the dividend; "3 bits w/ 4-bit G" written below the divisor; "down" indicating the direction of subtraction; and "r bits $\rightarrow R$ " at the bottom right.

- ④ Put this in msg



Widely used method! Receiver will do same division of $D+R$ by $G \Rightarrow$ should be 0 if no error

Multiple access protocols

Q. If we have shared link medium, how do we share link efficiently



A: use multiple access protocol. Det.

1. When should node transmit
 2. Coordination via channel/link only
- } distributor also

Ideally:

If channel has bandwidth R bps

1. If 1 node wants to transmit \rightarrow transmit at R
2. If M nodes $\rightarrow R/M$ rate per node
3. Fully decentralized (no special coord.)
4. Simple

MAC protocols

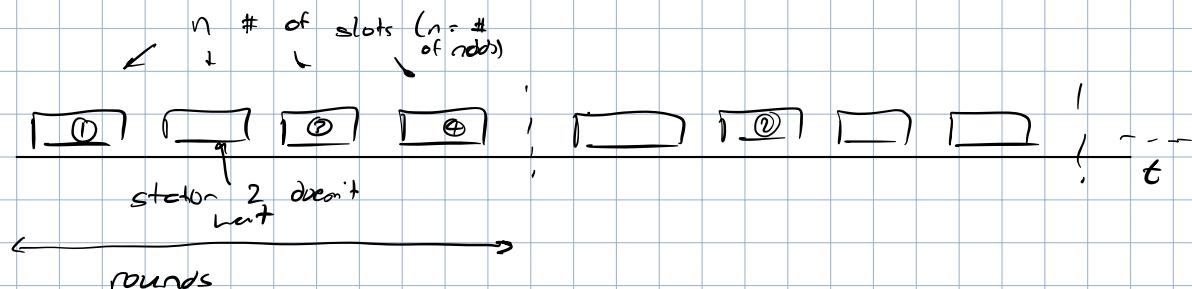
Channel partitioning
Divide channel &
alloc. each piece \rightarrow node

Random access
Allow msg. collisions
& retrans.

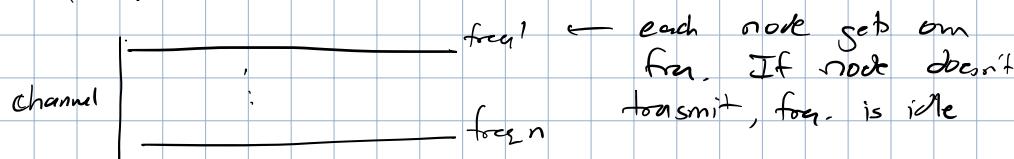
Taking turns
Nodes have full access
& switch depending on load
of msg.

① Channel partitioning

A: TDMA (time division multiple access)



B: FDMA (freq. DMA)



Con: not efficient & not dynamic, needs some coord.

Pro: simple

② Random access protocol

Each node sends @ R bits \rightarrow protocol defines error detection & recovery

A: Slotted ALOHA

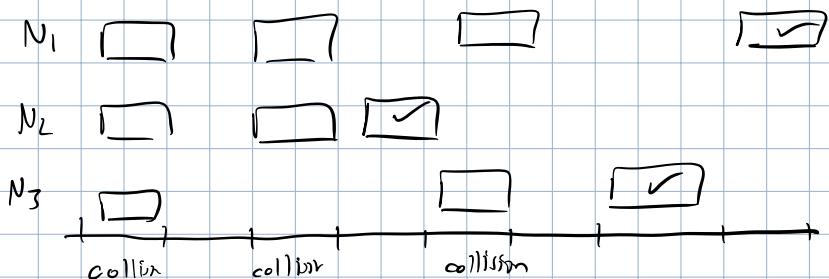
Assumptions:

- ① All frames equal size (L bits)
- ② Time divided into equal slots ($\frac{L}{R}$ sec.)
- ③ Nodes synchronized & all detect collisions.

Operation:

- ① Node sees slot \rightarrow transmits
- ② If no collision: send next frame @ next slot
- ③ If collision: node transmits frame @ next slot w/ $\text{prob } p$ prob all nodes from retransmits @ next frame \rightarrow collision

Ex://



Pros:

- ① Continuously transmit at R bps
- ② Decentralized (only slot syncs)
- ③ Simple

Cons:

- ① Collisions \rightarrow wasted data
- ② Accidentally cause more collision if cannot detect collisions fast enough
- ③ Synchronization required.

Efficiency: 37% (at best prob. $p^* = \frac{1}{N} (\# \text{ of nodes})$)

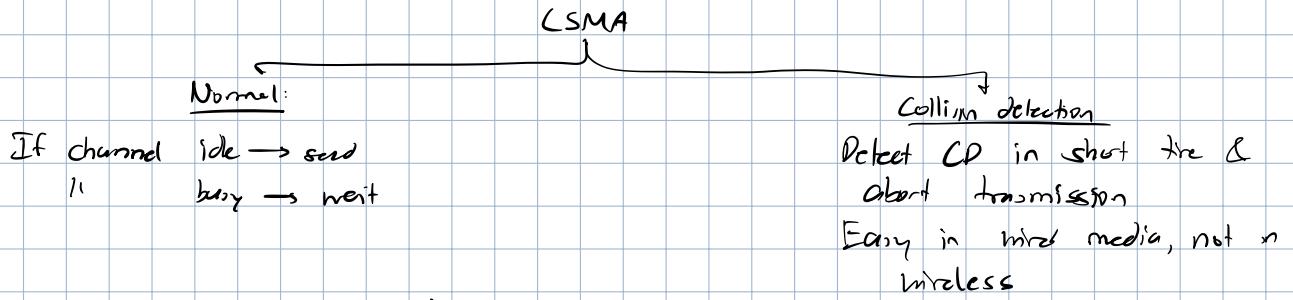
B: Pure ALOHA

No slotting \rightarrow more chaos & collisions

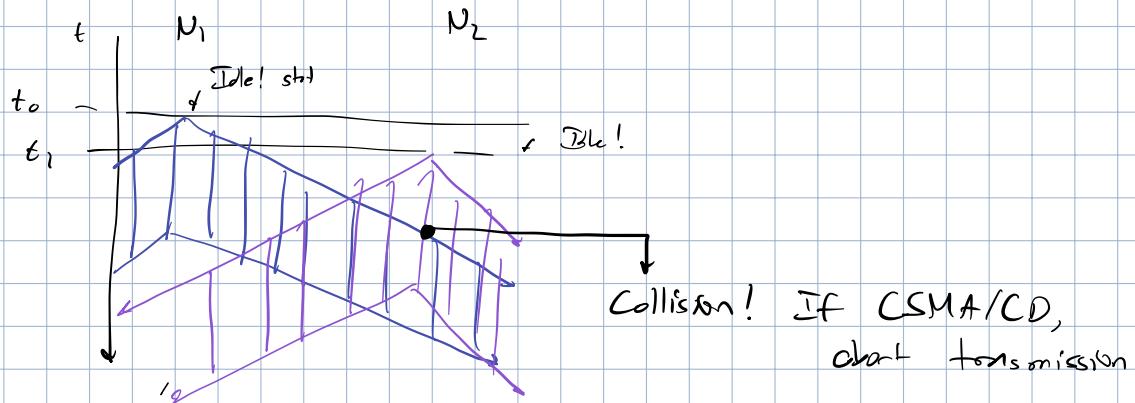
Efficiency: 15%. (at best prob. $p^* = \frac{1}{2N}$)

C: CSMA

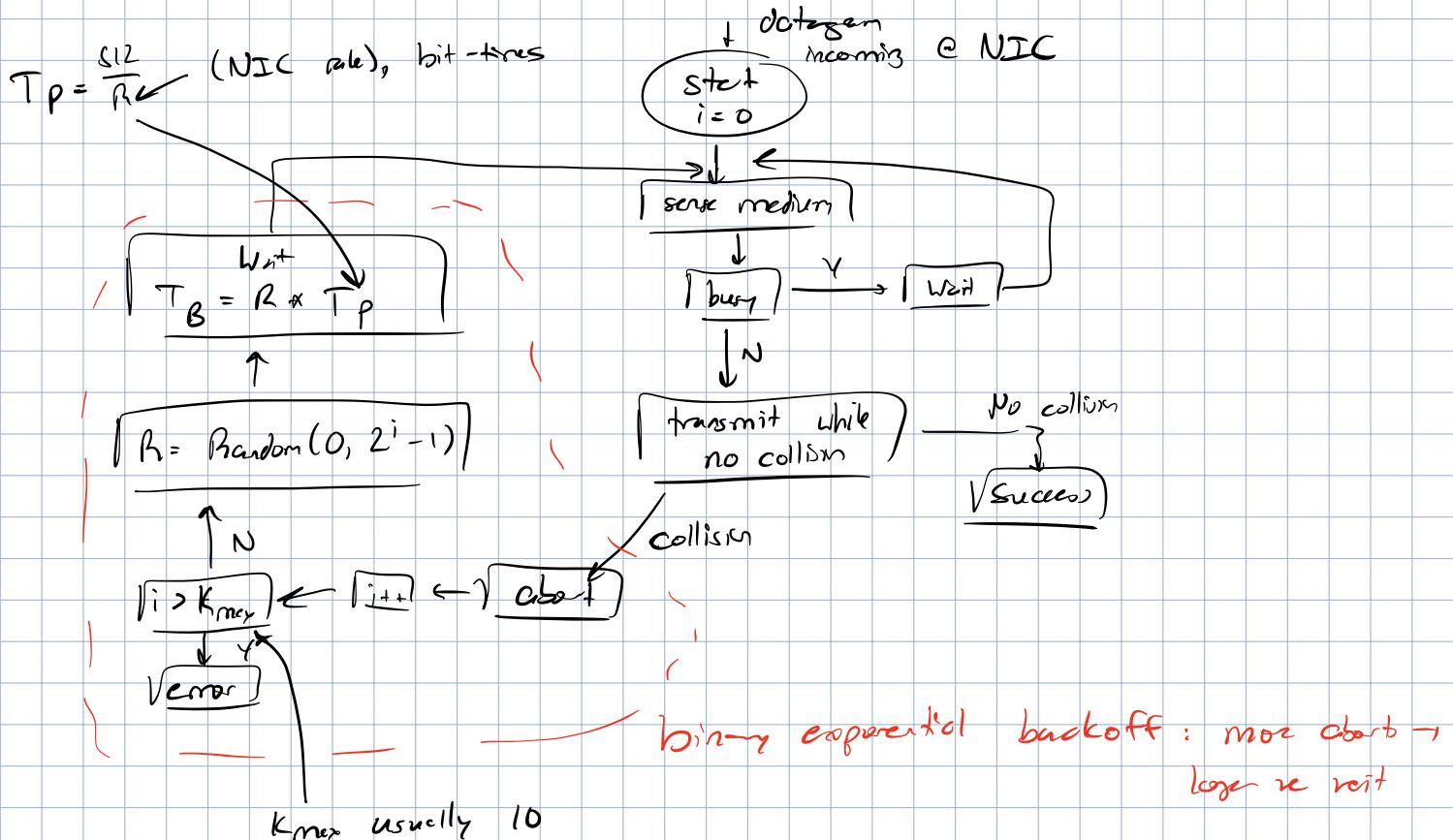
Carrier sense multiple access



You can still have collisions!



Ethernet CSMA/CD also:



$$\text{Efficiency} = \frac{1}{1 + \frac{S + t_{prop}}{t_{trans}}} \rightarrow$$

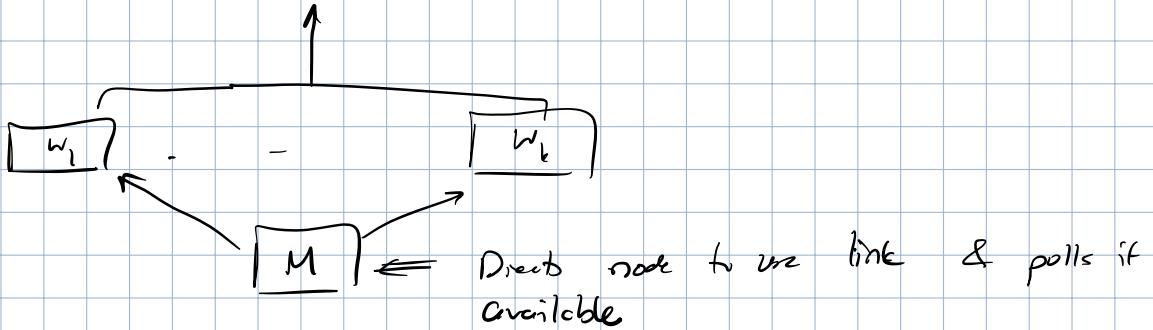
t_{prop} : max delay b/w 2 nodes in LAN
 t_{trans} : time to transmit max-size frame.

If $t_{prop} \rightarrow 0$ or $t_{trans} \rightarrow \infty \Rightarrow$ efficiency $\rightarrow 1$

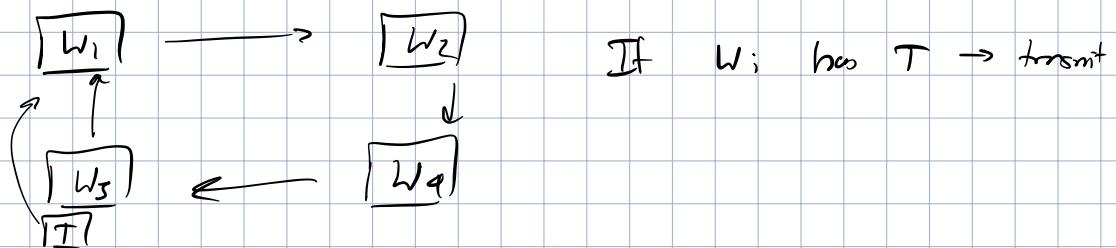
Not the fairest but good perf.

③ Taking turns

A: Polling



B: Token passing



Taking turns and mac for wireless (Bluetooth)

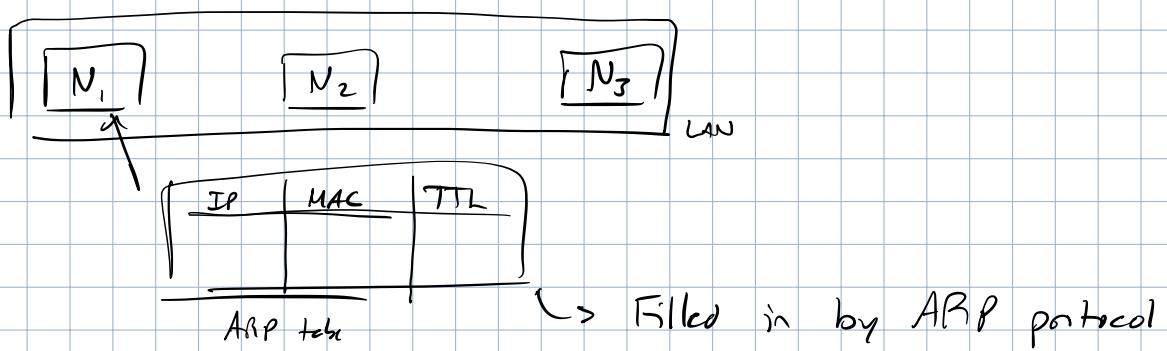
Addressing

Obj: physical address to communicate w/in adjacent nodes

Ans: MAC address (48 bits, in NIC, unique w/in subnet)

MAC addr. allocated by IEEE & portable (unlike IP addr.)

Q: How do we use MAC address in LAN?



ARP protocol: N_1 wants to transmit to N_2

- ① N_1 broadcasts ARP query w/ N_2 's IP

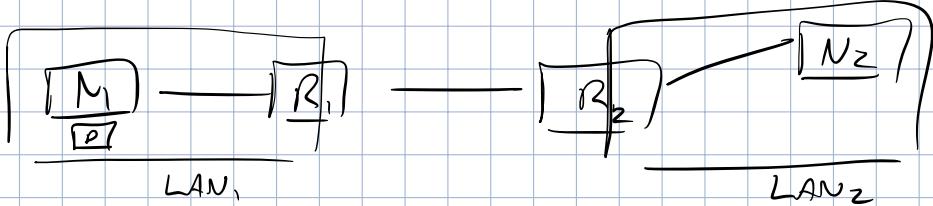
Src MAC: - -
Dest MAC: FF .. FF
Src IP: - -
Dest IP: IP_{N2}

- ② N_2 receives & sends ARP response

Src MAC: MAC_{N2}
Dest MAC: MAC_{N1}
Src IP: IP_{N2}
Dest IP: IP_{N1}

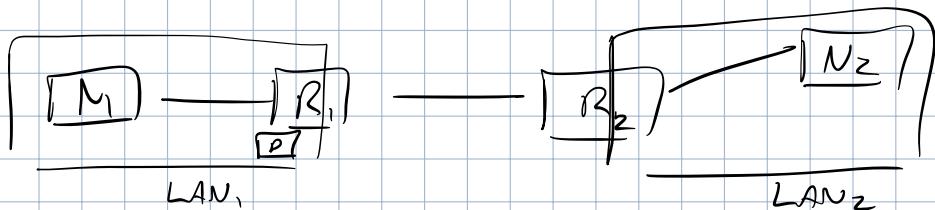
- ③ N_1 receives & updates ARP table

Q: What if we need to send to diff subnet?



- ① N_1 will set link frame w/ MAC addr. of R_1 , not N_2

- ② Broadcast frame to LAN₁, denies; if device MAC matches, remove datagram → sum + IP (only R_1 responds)



- ④ IP layer gets interface → link-layer gets MAC to next hop & sends

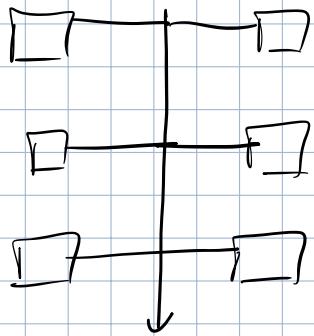
- ⑤ Repeats until router attached to N_2 & replaces MAC w/ MAC of N_2

Ethernet

Dominant wired LAN link type (fast, cheap)

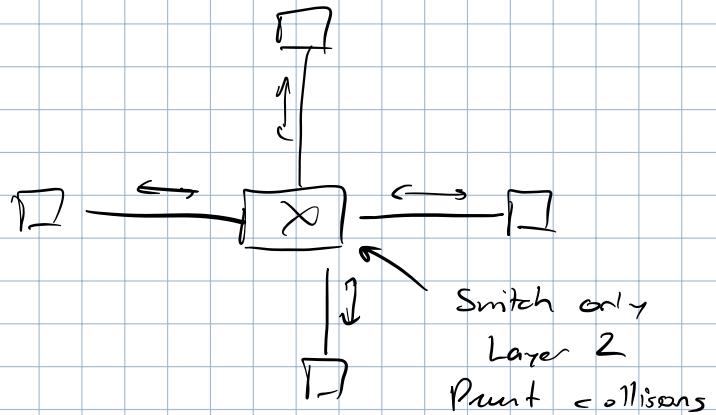
Types of topology:

① Bus:

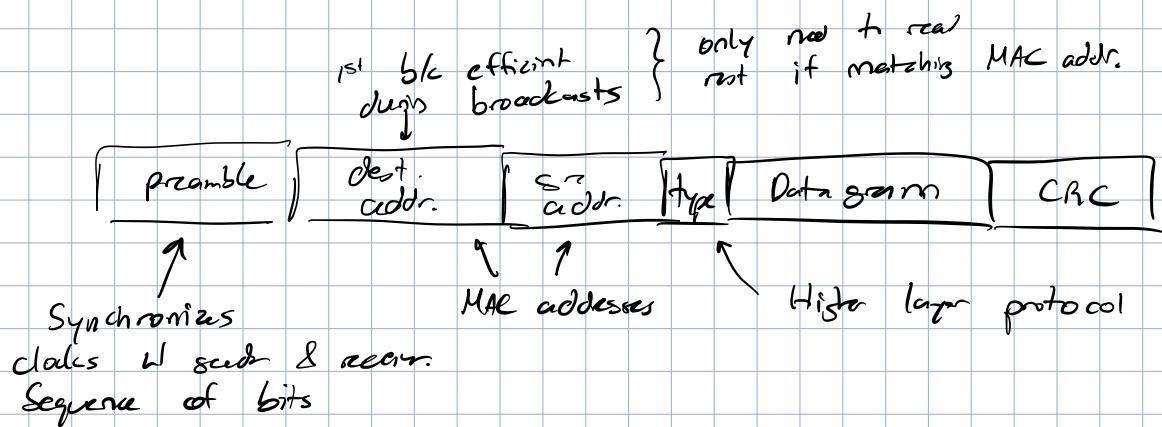


Collisions possible
Popular in 90s

② Switch



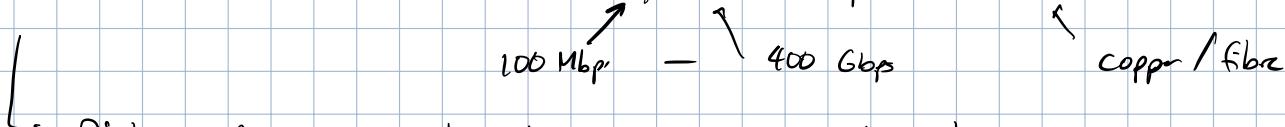
Frame structure for Ethernet



Properties:

- ① Connectionless: no handshaking
- ② Unreliable: no ACKs, no NACK (if packet lost, GL)
- ③ MAC protocol: CSMA/CD w/ binary exp. backoff.

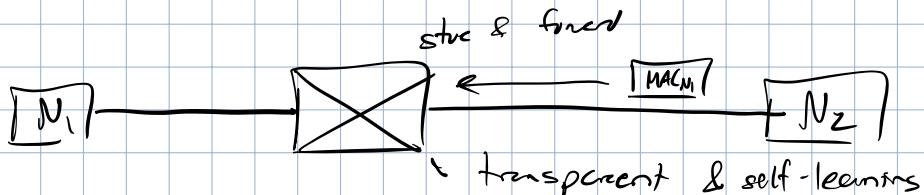
Ethernet standards: Differ in speed & physical layer



Distinguish: <rate in Mbps> BASE - <medium types>

100, 40, 600
if \geq Mbps \rightarrow appear \sim 6 (40G)
TX \Rightarrow copper, X = # of mtrs
else: fiber optic

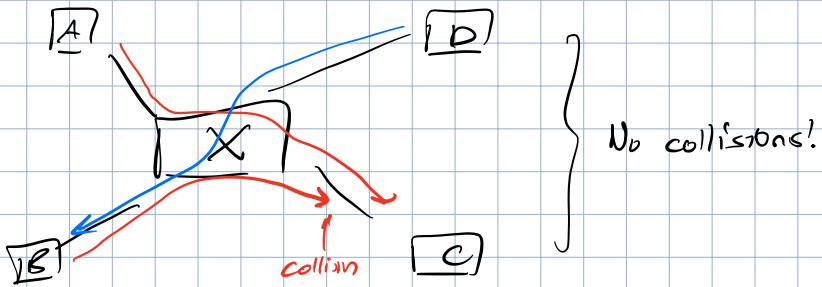
Switch



- ① Look @ packet MAC
- ② Forward to interface w/ corresponding MAC

Each host has conn. to switch, each run Ethernet protocol (CSMA/CD)

- ↳ Collisions can still happen if multiple nodes sending on same link
- ↳ Each link is full duplex



Q: How does switch figure out which interface to forward packets?

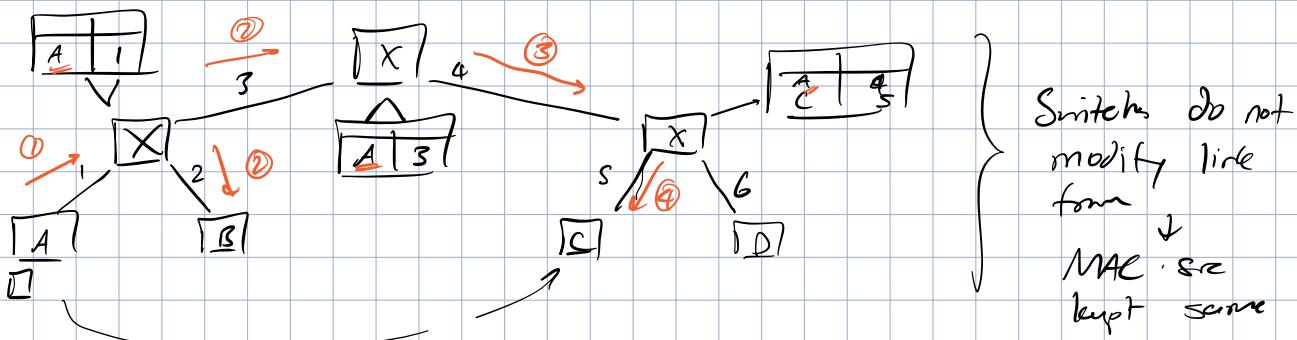
A: Switch table

MAC addr.	Interface

Learning process:

- ① Frame comes in → enter MAC addr & interface
- ② Check table if dest. MAC in switch table:
 - i) If match found:
 - A: Match is on same interface (going to take it here) \Rightarrow drop
 - B: Else: forward
 - ii) Else: broadcast frame to all hosts

Switches can be put in a hierarchy, works



Cycles are possible \rightarrow protocols can deal w/ it

Q. Diff. b/w hubs, switches & routers

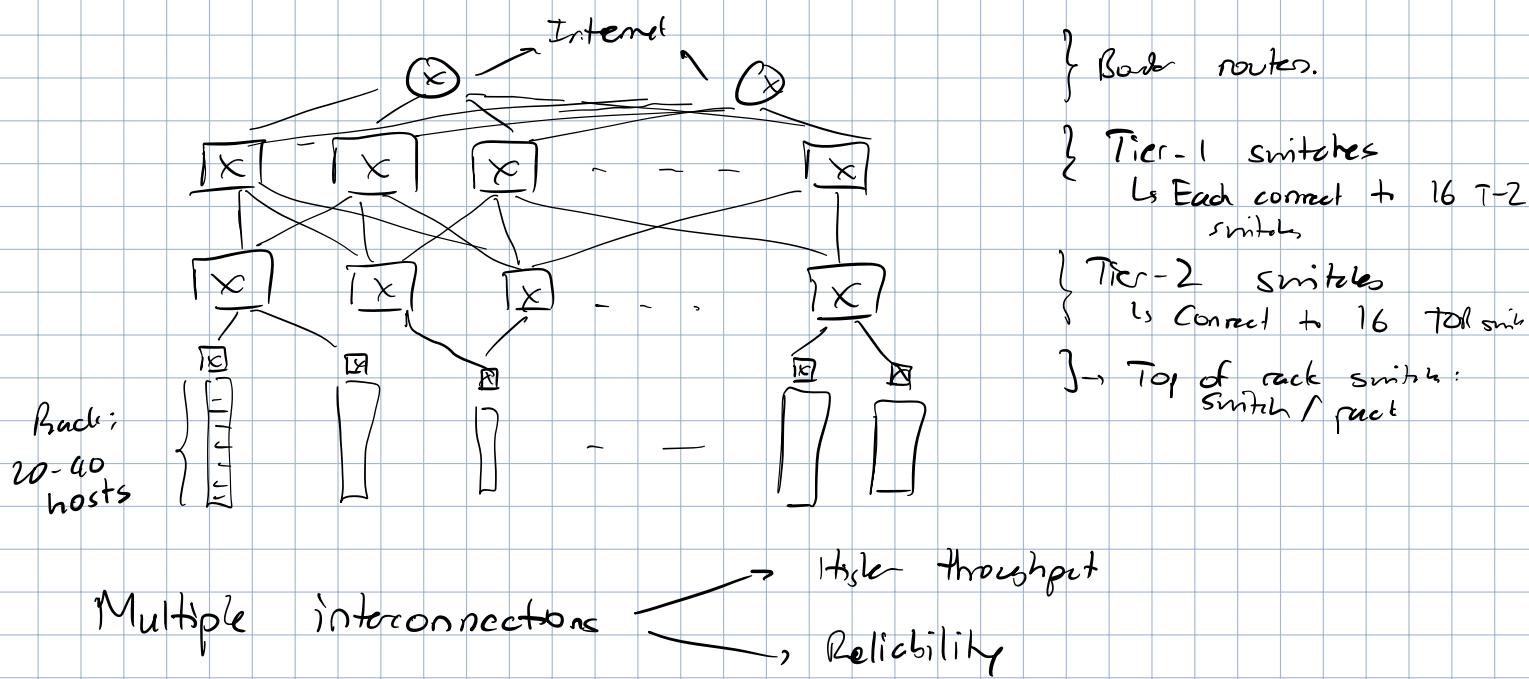
repeaters, flood network w/ messages

	Hubs	Switches	Routers
Layer	Physical layer	Link	Network
Plug-and-play	✓	✓	✗ (configs for routing also, OSPF weights, etc.)
Collisions	✓	✗	✗
Security (only participants hear packets)	✗	✓	✓
Optimal forwarding (always use shortest path)	✗ (flooding means we could forward packet in non-optimal path)	✗ (cycles)	✓

Datacenter networking

100 000+ hosts working together \Rightarrow multiple apps, reliable, balance load

General network struct.



Innovations:

① Load balancer routing (app. layer): directs workload to hosts \rightarrow local

② Data center TCP (DCTCP):

- DC OCN
↑
quantized
congestion
notification
- A: Explicit congestion notify: intermediate nodes send info to routers via packet headers about congestion
 \hookrightarrow no wait for dropped packets
- B: Priority flow control

③ Hop-to-hop backpressure detection & rate-limiting

④ RoCE: remote DMA over converged Ethernet

Interact w/ hosts as if 1 CPU \rightarrow direct memory access

⑤ Software-defined networking \rightarrow routing algo