

Transport Layer Protocols

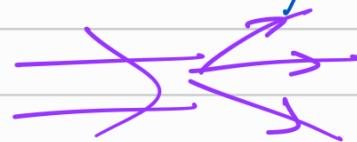
- inside the host.
- comm. b/w processes.

(Network Layer: Comm b/w hosts)

TCP: reliable, in order, congestion control, flow control, connection setup.



multiplexing (sender)



demultiplexing (recv)

- Conn-less Demux (UDP)

IP datagram with same dstn port no → same socket at recv.

- Conn"-Oriented Demux (TCP)

4 tuple: {src IP, src port, dest IP, dest port}

each socket of server associated with a diff connecting client

(same dstn port no then too into diff socket, based on tuple)

→ UDP:

best effort service

no RTT delay of handshake

Small header size.

No state maintd. at endsrecv (for +)

no state maintained at snes, recv. (hash)

(may add reliability, congestion control at application layer HTTP/3).

use: multimedia apps, DNS, SNMP, HTTP/3

Checksum:

treat contents of UDP segment (incl header fields & IP addr) as seq of 16 bit integers

add them, one's compl sum

$$\begin{array}{r} \text{1101} \\ \text{1001} \\ \hline \text{10110} \\ \text{0111} \\ \hline \text{1000} \end{array} \begin{array}{l} (\text{1 overflow}) \\ (\text{sum}) \\ (\text{checksum}) \end{array}$$

→ Reliable data transfer

RDT 1.0 : perfect channel, no errors / loss

RDT 2.0 : (with bit errors)
— checksum
— ACK / NAK

RDT 2.1 : seq. # (0, 1) for corrupted ACK / NAK

RDT 2.2 : NAK free

recv sends ACK of last pkt recvd OK.
(recv needs to incl seq no)

RDT 3.0 : (errors + loss) timeouts

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

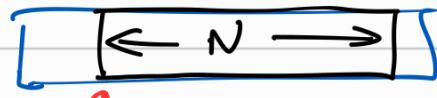
(do pipelining)

Go Back N

- window size = $N = (\text{transmitted unack. pkts})$
- $\text{ACK}(n)$: $\text{ack} \neq \text{pkt} \leq n$.
 \hookrightarrow send-base at $n+1$.
- timer for oldest in flight pkt.
- timeout (n): retransmit $n \rightarrow n + \text{nextseq.}$



Send-base nextseq no.



recv_base

recv: send ack for correctly recv highest seq no pkt.

Selective Repeat

recv buffers out of order pkts
 sender maintains timer for each unACK pkt
 retransmits only the pkt whose timeout

problem: window size and seq # magnitude
 should not be close
 $(\text{seg size} \gg \text{window size})$

TCP Overview

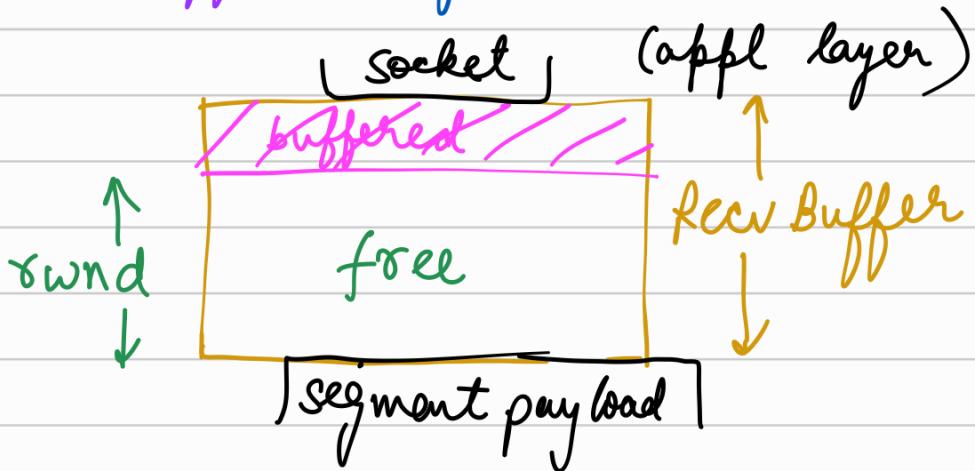
- * point to point (one sender one recv)
- * reliable in order byte stream.
- * bidirectional data flow in same dir

- * MSS : max seg size
- * cumul ack , window size fix .
- * TCP flow control + congestion control
- * handshaking initialises sndr, recv

→ Flow Control

Recv tells free buffer space in rwnd field in TCP header.

RecvBuffer : (def = 4096 B) - set using socket options



sndr : amt of unACK in flight data \leq rwnd

→ 2 way Handshake



→ End-End Cong Control : (TCP does this)
no explicit feedback from network
cong. inferred from observed loss / delay

→ Network Assisted Cong. Control
router sends direct feedback to sending / receiving host

recv host tells cong level / sets sending rate
TCP ECN, ATM, DEC

→ AIMD (probing for bandwidth)

inc send rate by 1 until pkt loss occurs (congestion), then dec send rate expoⁿ at each loss

TCP Reno - Cut in half on loss detected by triple duplicate ACK.

TCP Tahoe - Cut to 1 MSS when loss detected by timeout

- cwnd : sender window size.

Last Byte Send - Last Byte ACK \leq cwnd.

$$\text{TCP rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ B/sec}$$

- TCP slow start "ssthresh"

initially cwnd = 1 MSS.

double cwnd every RTT, until first loss event

new ssthresh = $\frac{1}{2}$ (cwnd before loss)

- TCP Cubic

w_{max} - sending rate at which congestion was detected

after cutting rate in half on loss, initially ramps to W_{max} faster, and approach W_{max} more slowly.

→ Delay based TCP Cong. Control

$$RTT_{min} = \text{min obs RTT (uncongested)}$$

$$\text{Best throughput} = \frac{cwnd}{RTT_{min}}$$

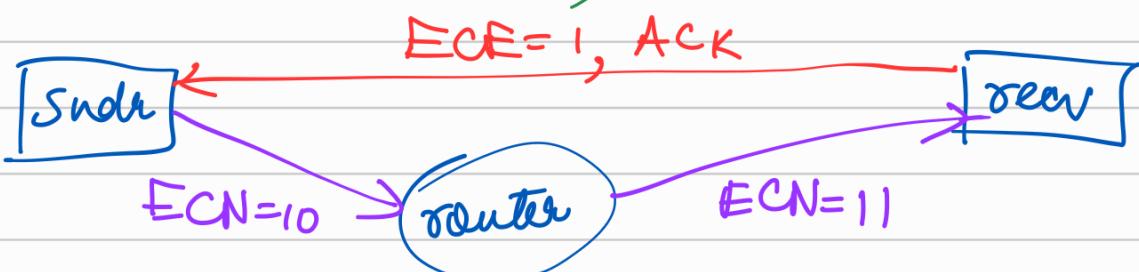
$$\text{measured } " = \frac{\# \text{Bytes sent in last RTT}}{RTT_{measured}}$$

if meas. close to Best \Rightarrow inc cwnd linearly
 if meas. far below Best \Rightarrow dec cwnd linearly

→ ECN (Explicit Congestion Notification) (Network Assisted Cong. Control)

- 2 bits in IP datagram header, marked by router to indicate congestion.
- destⁿ sets ECE bit on ACK segment to notify sndr of congestion.

both IP (ECN bit) & TCP (ECE bit)



→ TCP Fairness

K TCP connections share same bottleneck router of capacity R \Rightarrow each gets avg of R

Start at any points, hosts reach equal bw share in AIMD.

- Multimedia apps: use UDP, tolerate loss

long fat pipes: large data transfer.

Wireless: loss due to noise, mobility

Long delay link: Exts. long RTT

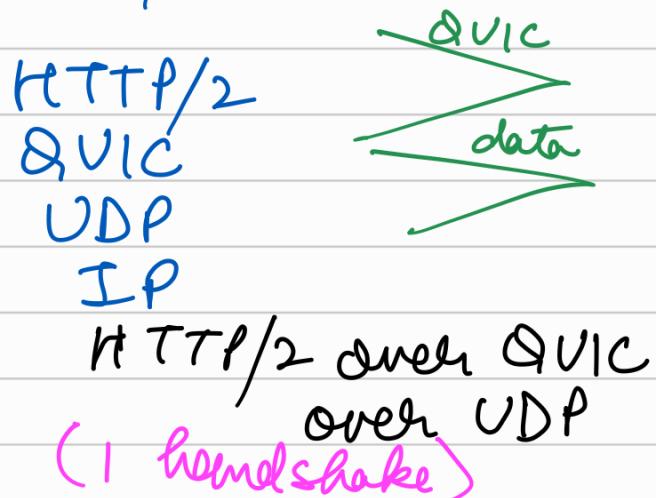
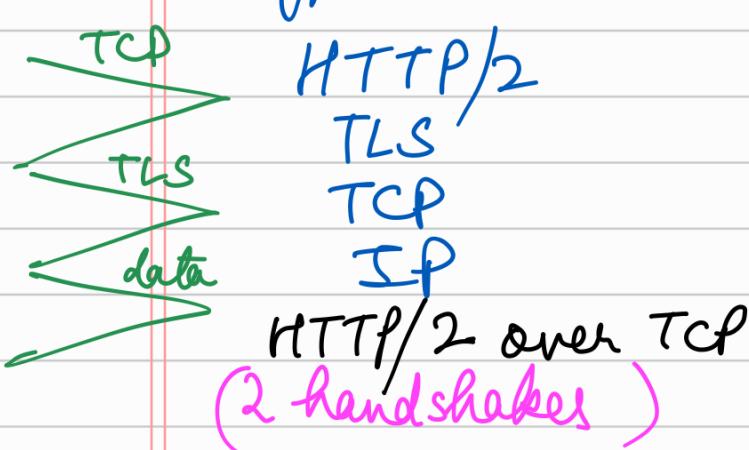
Data centre network: latency sensitive

Background traffic flows: low priority (in backlog)

HTTP/3: QUIC — moving transport layer func to appl layer on top of UDP
(Quick UDP Internet Connections)

QUIC gives:

reliability, cong. control, authentication, encryption, state estab, loss/error control.



QUIC - parallelism, no FOL blocking

→ TCP Throughput

W - window size (B) when loss occurs.
avg window size = $\frac{3}{2} W$

$$\text{avg throughput} = \frac{3 \cdot W}{4 \cdot RTT} \text{ bytes/sec}$$

$$\text{TCP throughput} = \frac{1.22 \text{ MSS}}{RTT \sqrt{L}} \quad (L = \text{segment loss prob.})$$

Ch-4

Data Plane :

local, per router function. $\text{inp port} \rightarrow \text{out port}$

Control Plane:

end-end path routing

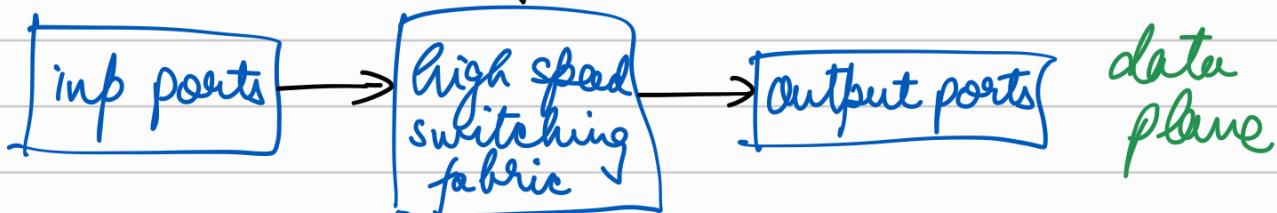
Fwdⁿ routing algos

software Defⁿ Network (SDN)

Routing algo \rightarrow local fwd table.

routing processor

control plane



Input Port func

lookup, forwarding, queuing (if datagrams arrive faster than fwd rate into switching fabric).

Match + Action:

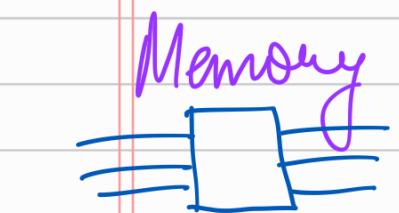
Dstn based fwd - based only on dest IP addrs.

Generalised fwd - based on any set of header field values

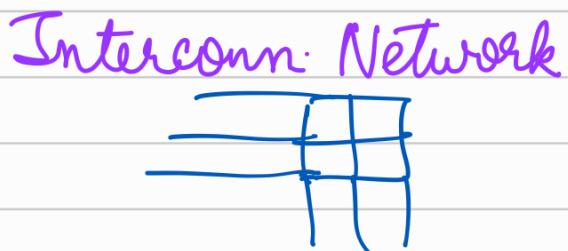
→ Longest Prefix Match for distn IP addrs.

- Switching Fabric

Inp port to output port switching rate \rightarrow NR ideally
(N ports, R pkts per port)



(2 bus crossings per datagram)



(multistage switch)

Inp port queuing - HOL blocking, if sw. fabric slower than NR.

Output port queuing -

Buffering: datagrams arrive from fabric faster than link transmission rate.

(Network neutrality: Drop policy, priority scheduling)

Buffer size:

① $RTT \times C$ ($C = \text{link capacity}$)
 $C = 10 \text{ Gbps}$, $RTT = 2 \text{ SD ms} \Rightarrow 2.5 \text{ Gbit buffer}$

② N flows $\Rightarrow \frac{RTT \times C}{\sqrt{N}}$

Packet Drop: tail drop
priority drop.
Marking (ECN bits, RED).

→ Pkt Scheduling

- ① FCFS / FIFO.
- ② Priority Sch: FCFS within priority class.
- ③ Round Robin (based on any header field)
- ④ Weighted Fair Queuing: $\frac{w_i}{\sum w_i}$
(min bw guarantee)
each class gets weighted amt of service in each cycle.

→ IP addressing

IP address: 32 bit identifier associated with each host or router interface

interface: conn b/w host/router and physical link

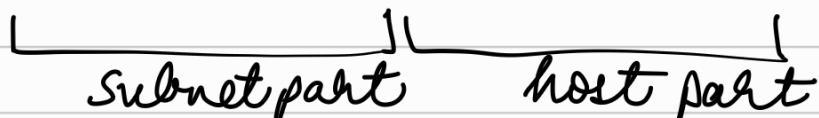
outers: multiple interfaces

hosts: generally one/two (wireless + ethernet)

→ Subnet

device interfaces that can physically reach each other, without passing thru a router. Have common high order bits.

32 bit IP addr



$223 \cdot 1 \cdot 3 \cdot 0 / 24 \Rightarrow 24 = \text{subnet mask}$.
high order 24 bits - subnet part of IP addr, 8 bits diff for host.

CIDR (Classless InterDomain Routing)

subnet portion of addr. of arbitrary len.

$a \cdot b \cdot c \cdot d / x \Rightarrow x$ bits in subnet

Host IP addr

- hard-coded by sysadmin in config file
- DHCP: dynamic host config "protocol (plug and play)

- i) host broadcasts DHCP discover msg (optional)
- ii) DHCP server gives DHCP offer msg (optional)
- iii) host req. IP addr : DHCP req msg
- iv) DHCP server sends addr : DHCP ack msg

DHCP colocated in router.

can give: addr. of first hop router

name + IP addr of DNS server

network mask (network v/s host portion of addr info)

Network IP addr

network gets allocated portion of its provider ISP's addr space. (ICANN gives)

ISP: 200.23.16.0/20

200.23.16.0/23

200.23.30.0/23

8 orgs.

→ NAT

all devices in local network share one IPv4 addr as far as outside world is concerned.

all devices have 32 bit private IP
addr: (10/8, 172.16/12, 192.168/16) that can only be used in local network

{ IPx local port no} → NAT IP x port no
diff

src IP, port # → NAT IP, port #
NAT IP, port # → src IP, port #

- only one IP addr needed from provider ISP for all devices
- can change addr of a host without telling outside world
- can change ISP without changing addr of devices in local network
- security

→ IPv6

"Classless"

fixed length header

flows → pKts of same flow: send directly, no need to intersect deeper.

Tunneling -

IPv6 datagram carried as payload in IPv4 datagram in IPv4 routers

[IPv4 headers, IPv4 src/dstn, IPv4 data]

[IPv6 headers, IPv6 src/dstn, data]

→ Generalised Fwd / Match + Action

flow: defined by header field values
(in link, network, transport layer)

- ① Match - pattern in table to header fields
- ② Action - drop / fwd / modify / send to ctrl.
- ③ Priority - disambiguate overlapping patterns
- ④ Counters - # Bytes , # packets

— Middle Boxes

NAT, Firewalls / IDS, Appl Specific,
Load Balancers, Caches

— Fragmentation

MTU: max transfer size (largest possible link level frame).

large datagram fragmented, reassembled only at dstn! IP header bits used to identify, order related fragments

offset in each datagram

TCP Reno

Triple Dup

$$ssthresh = \frac{cwnd}{2}$$

$$cwnd = \frac{cwnd}{2} + 3MSS$$

Timeout

$$ssthresh = \frac{cwnd}{2}$$

$$cwnd = 1$$

TCP Tahoe

TD

Time

$$ssthresh = \frac{cwnd}{2}$$

$$cwnd = 1$$