

**Internet:** infrastructure that ties together computing devices | no notificator,

↳ Decentralized (↳ Best effort service model): no guarantee on if data will be delivered

↳ Dumb infrast. w/ smart end hosts, e2e, layering, "human waist"

**Interoperability:** entities can connect & comm w/ other entities easily

**Protocol:** specification of msgs that comm. entities exchange (syntax & semantics)

**Bandwidth:** num bits sent/received per time (bps)

↳ transmission delay: size of packet / bandwidth

↳ Propagation delay: len of link  $\frac{\text{in } m}{\text{speed of light } (\text{cm/s})}$ ; non-zero (use smallest bandwidth given)

**Bandwidth-delay product (BDP):** bandwidth  $\times$  prop delay (bits)

Latency = packet delay = trans delay + prop delay + queuing delay

Utilization: usage / maximum  $\downarrow 10^3$  prop delay  $\rightarrow$  1st bit:  $(1/10^6) + 1/10^3$

(ex) 100B over 1 Mbps  $\rightarrow 10^6 \text{ bps} \rightarrow 8 \times 100 = 800/10^6 \rightarrow$  last bit:  $(800/10^6) + 1/10^3$

**Forwarding table:** next hop; network addr (where) vs. port name (which host)

Control plane	Data plane	Management plane	Packet switching	Circuit switching
Compute func tables (routing algo/protocols), global, per network event	actually find packets (forward), local (arriving packets & rt table) per packet arrival	interact w/ systems/humans to config & monitor device	Best effort model packet by packet basis Better efficiency, fast start, easy recovery, simple implementation, efficient, best for bursty	Reservation model Predictable/understandable, good for smooth apps, bad for small data Cons: tear-down failure, setup time, bad handling, complex, endhosts detect failure

Autonomous System (AS): set of nets all managed/supr by single entity/func

Internet service provider (ISP): set of packet switches & comm links, it ASes

Statistically multiplexed: do not provision for worst case; hope all peaks don't happen at the same time; peak of avg demand << avg of peak demands

Smooth: low ratio btw app peak to avg transmission rate (voice: bursty: high)

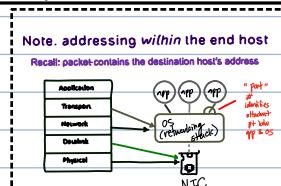
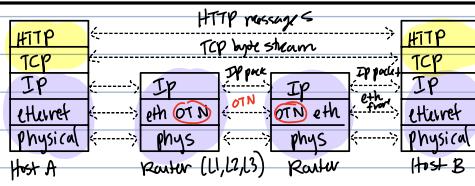
Transient overload: queue absorbs bursts, drain queue

Persistent overload: queue force to drop packets

**Layering:** well-def interfaces; layer only interacts w/ one above & one below

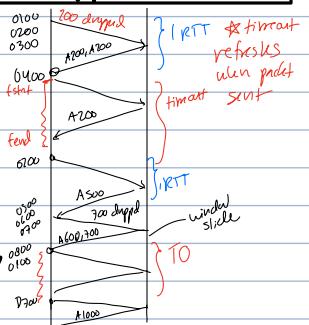
\* depends on layer below \* supports above \* independent  $\rightarrow$  parallel innovation

Layer	Desc	Protocols
L7 application	Browser, mail/client, only at host, exactly-once delivery	SMTP HTTP DNS NTP
L4 transport	Part of OS (only at host), at-least-once, reliable data delivery, Linux	TCP UDP
L3 Network	Best eff global packet delivery, best effort delivery, usually part of OS	IP
L2 Data link	Best eff local packet deliv, hardware/drivers ex. ethernet faeries	Ethernet FDDI PPP
L1 Physical	Physical transfer of bits; hardware/drivers/wire	optical copper radio PSTN



Routers: L1, L2, L3

Switches: L1, L2



**Network stack:** app(data)  $\leftarrow$  stack  $\leftarrow$  NIC (packets) func as a feat.

**Physical part:** link  $\leftarrow$  switches (logical); app  $\leftarrow$  OS stack

**Socket:** app  $\rightarrow$  gen socket  $\leftarrow$  logical part (OS uses port #  $\leftrightarrow$  socket)

**Routers:** look up where to forward/send packet; topology management, not host/host

↳ Spanning tree: undirected graph, symmetric, connect all nodes, least. based routing (User traffic: packets find according to shortest routes)

\* forwarding path == fast path, Control == slow (passing)

Control plane traffic: packets to this router

(@ BGP/IGP packets), send to controller and send to controller card

Router capacity =  $N$  (# of external ports)  $\times$  R (specif "line rate" of port)

project 1

- don't recharge poison'd routes if not in

- don't add entry for poison route

- periodically resend poison'd routes

- poison spread doesn't depend on resending

$$2^3 = 8 \quad 2^7 = 128$$

$$2^4 = 16 \quad 2^8 = 256$$

$$2^5 = 32 \quad 2^9 = 512$$

$$2^6 = 64 \quad 2^{10} = 1024$$

**Distance-vector:** tell neighbors my least cost entries to a dest. **①** If ad is next hop, replace. **②** If cur-route > ad route + dist, replace entry. **is my info**

**Interval:** resend every  $\times$  sec  $\rightarrow$  reliable vs. triggered: send on change, optimization

**Split horizon:** don't tell ur next hop if it's the dest (loops can still happen)

**Position reverse:** tell ur next hop infinity (loop state exist until next ad vs. explicit route poisoning)  $\rightarrow$  if down (faster propagation than timer) **SRP faster**

**Count to inf:** bouncing & max  $\rightarrow$  solve w/ a max cap  $\star$  SRP or PR doesn't solve.

**Routing state validity:** (1) no loops  
(2) no deadends  $\rightarrow$  state is valid

**Local cannot be evaluated, only global**

**Corrected / directed routes:** manual exp. duplication static: want them, no routing proto to discover if 2+ loops  $\star$  router expire

**Link-state:** global flooding (reliable: periodic resend); No center to  $\infty$ , can have loops; delay: time to detect failure, time to flood info, time to recompute

**CIDR:** classless inter domain routing (IP space) vs. Classful; A: 8, B: 16, C:  $2^{11}$  bits

**Aggr!** **(1x)**  $12.1.0.0/17 = 12.1.0.0$  to  $12.10.127.255$   $\star$  0 out if held for prefix **not aggr.**

**LPM:** longest prefix match; take longest route, overlapping possible

**Inter-domain:** between ASes **Transit:** carry on behalf of other AS **Stub**  $\star$  on behalf of dir. hosts

**Scalability via aggregation vs. Multihoming:** more than 1 provider, must announce

**BGP:** extend DV (autonomous, policy, privacy); export (which can enter) import (leaves)

**①** Many aggr dest **②** prioritize policy **③** DV  $\rightarrow$  PV: send entire route (avoids loops + flexible policy)

**④** Selective route ads  $\rightarrow$  reachability not guaranteed even if connected graph

**G-R** dest prefix ad by export to cust. | everyone | peer | cust | provider | cust | IP AS learns a route from a customer, AS will export route to anyone

**Assumptions for steady state guarantees** **①** cust-provider acyclic (peer is ok) **②** chain leads to TI

**Then guarantee:** **①** reachability (any 2 ASes can talk) **②** convergence (routers agree on paths)

**eBGP:** learn routes to different ASes **iBGP:** distr. knowledge internally **IGP:** internal routing **OSPF / RIP**

#	Rule	Next Hop	Description	May Conflict
1	LOCAL PREF	Pick highest #		
2	ASPATH	Pick shortest len		
3	IGP path	Pick lowest cost		
4	MED	Pick lowest		
5	Router ID	Smallest next-hop addr in local		

Version	4	Src IP	TTL
Harlen	size	Dst IP	options
	Type of service	Priority/delay	ECCN + DSCP
	Total len	size	Checksum
	ID	uniquely ID group of fragments	
	Flags	No frag (1) None frag (2)	1..2
	Frag offset		
	Protocol	which L4 demux	= TCP = UDP

**ctrl plane** **data plane** **Checksum:** updated at every router; entire header + TTL  $\rightarrow$  new calc

**Frag-off:** 8 byte boundary: add bytes before frag, div by 8  $\star$   $1380 + 1220 + 1200 - 2(10) = 0$  G size  $\star$  don't forget hdr!

**Hop limit == TTL** **VIS:** IPv6 (128): no checksum, frag, hdr len (const), + flow label

**Reliable transport:** W based: allow W packets in flight, allow sender to transmit for entire KTT size

$W \times \text{Packet size} = \text{RTT} \times B$  (bottleneck)  $\hookrightarrow$  As long as src retransmits unack'd packets  $\rightarrow$  reliable  $\neq$  in order

**Indiv ACK:**  $I \rightarrow \text{ack}(1)$  **(1x)** 1,2,3,4 sent, acks sent, ack 2 lost  $\rightarrow$  ack 2 retransmit (loss = retransmission, waste)

**Full info ACK:**  $\text{cum ack}$  highest + additional **(ex)** 1,2,3,4  $\rightarrow$  acks ( $\leq 1$ )... ( $\leq 4$ ) **(ex)** 1,2 lost, 3,4 lost  $\rightarrow$  ack( $\leq 1$ ), ack( $\leq 1+3$ ) (resilient + sizable overhead)

**Cum. ACK:** highest seq. for all received **(ex)** 1,2 lost, 3  $\rightarrow$  ack( $\leq 1$ ), ack( $\leq 1$ ).. (compact, more resilient, but incomplete info)  $\star$   $W=2BL$

**TCP:** cum ACKs + byte streams + seq num (byte offsets) + timers  $\star$  kth byte  $\star$  ISN + k + ACK = len(data)

**CWND**  $\leq RTT \times B$  **MSS (seq size)** = MTU - IP hdr - TCP hdr  $\star$  seq # == ISN + k  $\star$  ACK = len(data)

src/dst: demux			
seq #: start byte offset			
ACK: rest expected byte			
Checksum: hdr + payload			
hdr len: 4 bytes, max 52 bytes, min 20 bytes			
0 (reserved)			
Flags: SYN ACK FIN RST PSH URG			
Vigint. pfr: max window options (ignore)			

**(1x)** X, len B  $\rightarrow$  ACK = X + B back  $\rightarrow$  X + B len B scnt  $\rightarrow$  ACK = X + 2B  $\star$  no longer send but can receive

**Establish:** SYN SYN ACK ACK **(1x)** ISN = 19  $\rightarrow$  W = 10  $\rightarrow$  L send 19+1  $\rightarrow$  ACK = ISN + MSS  $\times$  (2+1)  $\rightarrow$  W. **VIS:** FIN ACK FIN vs. RST

**AIMD (most fair)** **AIAD (unfair)** **MIMD (unfair)** **MIAD (max. unfair)**  $\frac{x_1 + x_2}{2} > C$  then congested

**Flow ctrl: sender RWND (advertised)** **Congestion: avoid** **RWND (how much to send to links)**  $\star$  RWND  $>>$  (WND  $\star$  CWND max + CWND min)  $\star$  Avg num pack in flight =  $\frac{\text{CWND max} + \text{CWND min}}{2}$

**Avg throughput =**  $\frac{3 \times \text{MSS}}{2 \times \text{RTT}} = \frac{\text{avg # bytes inflight}}{\text{RTT}}$

**Slow start**  $\star$  MSS = packet size

**CWND init** = small const  $\times$  MSS

**Each ACK**  $\rightarrow$  (WND) = (WND + 1)  $\star$  same as  $\frac{2 \times \text{CWND}}{\text{RTT}}$

**End:** timeout: back to slow start

**3 dup ACKs:** fast recovery

**Fast recovery:** retrans single packet

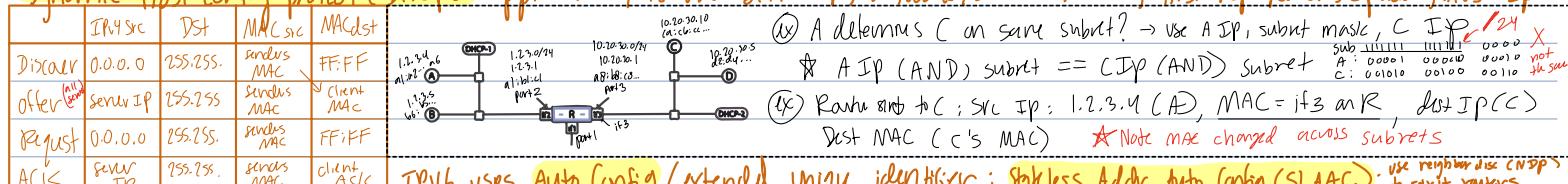
**End:** timeout  $\rightarrow$  back to slow start

**New ACK**  $\rightarrow$  congest avoidance



Ethernet L2: local area network (LAN) needs Media Access control (MAC); 48 bits, burned in, globally unique back off  
 Original: polling (Coordinator), token-passing; carrier sense. Mult. access w/ collision detection (CSMA/CD): listen while talk, stop if hear  
 Unicast: 1 recipient, use frame (Ethernet packet); checksum, preamble for seq. Dest/seq MAC; Broadcast: FF! FF Multicast: 01: 00:00:00 set bit to all in a group  
 ARP (for same subnet): Broadcast request, unicast reply → ARP table IPv4, MAC, interface, TTL (IP addr might change over time)  
 \* L2 may change in hops, but not L3 (IP) ★ ARP across LANs; check if addr in network (subnet mask) → first hop router ARP & send

**Dynamic Host config protocol (DHCP):** application; learn self IP, network/subnet mask, first hop router's IPaddr, DNS servers IP



E2E Network Addr translation (NAT): part of network translation; L4 (TCP/UDP) flow awareness, many hosts share single public address  
 Retranslate: Src IP → public IP, modify TCP src port to distinguish 2 flows ★ IP src addr (private) TCP src port (2 diff hosts contact) same remote addr

DHCP	ARP for MAC (to use our DNS)	DNS for request B IP	ARP req for B's MAC	TCP connect	HTTP req (TCP)	Tcp transient	A wants to C	DNS req 1st-hop	ARP for handshake	TCP handshake	TCP (HTTP) data trans	TCP teardown
Discover (send)	ARP req (send)	DNS req (send)	ARP req	TCP SYN	HTTP req (send)	Fin (send)	initially have ARP req (recv)	DNS req (send)	ARP req	ICP SYN	HTTP req (send)	Fin (send)
Offer (recv)	ARP reply (recv)	DNS reply (recv)	ARP reply	TCP SYN+ACK	ACK to req (recv)	ACK (recv)	ACK to reply (recv)	DNS reply (recv)	ARP reply	TCP SYN+ACK	ACK to req (recv)	ACK (recv)
Request (send)				TCP ACK (no reply req)	HTTP reply (recv)	Fin (recv)	Fin (recv)	DNS reply (recv)	ARP reply	TCP ACK	HTTP reply (recv)	Fin (recv)
Ack (recv)				ACK to reply (recv)	ACK to reply (recv)	ACK (recv)	ACK (recv)			ACK to reply (recv)	ACK (recv)	

**Host Networking** ★ implement functions @ host to support abstraction of network as fast, reliable, secure, ordered byte stream  
 ★ Traditional OS handle networking = hard (copies, development, CPU resources) Kernel / OS bypass: data copy App → OS NIC vs. store data into shared memory; offload copy work; packet processing & network protocols in user space ★ still need to reduce CPU cycles, offload processing, move to NIC, free CPU → ① Simple stateless: checksum, segment / tx/rx queue select

② simple stateful / accelerating data plane: BW mgmt, forwarding ③ Protocol offloads: TCP / RDMA, programmable data plane

**Remote Direct Mem Access (RDMA):** minimize CPU transfers; mem-to-mem bw servers; CPU just init, while NIC & network responsible for pair; send + receive; ptr to mem to transfer/receive data; NJC notifies CPU via completion queue elements (CQEs)

Pros: high performance (low latency by moving SW var), CPU efficiency | Cons: mem complex, limited protocol comp.

CC: delay based CC; Signals ↓ Swift: delay based, measure RTT as time, AIMD on packet add. high packet delay in loss CC

**Loss** BW ECN Delay Load balancing: ECMP to balance; dynamic bursts still imbalance → protection "mitigation"  
 Pros: high performance (low latency by moving SW var), CPU efficiency | Cons: mem complex, limited protocol comp.

**Traffic shaping:** spreading traffic over time to constrain to constraint BW use by users; based on policies/business priorities (via cc, no policy)

**Traffic policing:** spreading traffic over time; reduce bursts, queue/delay, packet loss; insert delays + window clocking/sliding to prevent bursts | Time when / Single time indexed queue: timestamp biffs; & determines rate biffs empty into Q on next; spin some | Timestamp: Earliest departure time via policy

**Quality of Service (QoS):** priority of this traffic for allocating buffers & BW at each hops  $\leq$  RTT time scales?

Issue: Cong at link  $\xrightarrow{\text{weighted fair q}} \text{how alloc vsrc?}$  Rule: Sender app  $\geq$  priority class (latency, throughput)  $\Rightarrow$  QoS (enclosed in hole)  $\geq$  1: Buff z; BW z

\* QoS only matters when switch port is 100% utilized when packets arrive & QoS priority policy under cong. i MUXing priorities

④ Timing wheel: class A rate 1 Gbps, B: 0.4 Gbps. Last packet A  $\Delta t = 100\ \mu s$  = B. Packet size 1500B; traffic class alloted packet in 8 ms

Horizon (time to cycle thru Qs) = 4, "nav" = 3rd Q, time count = 100 μs, "nav" last mark at  $\frac{100}{100} \mu s$   $\xrightarrow{\text{horizon}} 8\ \mu s$  bytes

① Max rate 1.5 Gbps for any traffic Q → time granularity (amount of time caused by each q)?  $\frac{\text{packet}}{\text{horizon}} = m = 1.5 \times 10^9 \times (8 \times 1500) \text{ min. bw} = 1.5 \text{ Gbps}$

② Min rate supported?  $Q 1 \text{ packet per horizon, min rate} = (1500 \times 8) / 4 = 3 \text{ kbps}$  ★ slowest we can send out, num Qs unknown

③ How many queues are there? # queues = horizon / granularity =  $4 \times 10^6 / 8\ \mu s = 0.5 \times 10^6$  107 + time granularity = time when next q selected

④ Packet ready in A; timestamp will it get? which queue?  $\text{Timestamp} = 100\ \mu s + \frac{\text{pkt\_size}}{\text{rate}} = 100\ \mu s + \frac{1500 \times 8}{1.5 \text{ Gbps}} = 112\ \mu s$  vs. 115ms, 112 falls in 3rd

**Wireless** Shared, attenuating signals w/ distance, env. change rapidly, packet collision hard to detect; fading, shadowing, multipath, path loss

Medium access controls: multi. devices sharing transmission CSMA (listen for others, don't transmit if they are) retransmission, obstacles

Hidden terminal problem: A → B ← C, A & C out of range to each other; A can sense C collison → CSMA fails  $\xrightarrow{\text{RTS/CTS}}$

RTS/CTS (req/clear to send): broadcast medium for all A  $\xrightarrow{\text{RTS}}$  B  $\xrightarrow{\text{CTS}}$  C (carrier sense wait for transmission to finish)

Exposed terminal problem: B → A, C → D, but B & C cannot transmit at same time  $\xrightarrow{\text{A (B C) D}}$  MACA (back off unfair) → MACA

**Cellular** Introduce ① Discovery: which tower to connect to ② Authentication: provide service? ③ seamless conn: no disrupt ④ Accountability: roaming plan

**Radio Access Network (RAN):** collection of towers (given set of ph) ; cellular core: RAN  $\leftrightarrow$  Radio gateway  $\leftrightarrow$  mesh  $\leftrightarrow$  Mobility mgmt (control functions)  $\leftrightarrow$  DB (subscriber info)  $\leftrightarrow$  packet gateway (cellular network & internet boundary) ★ IMSI: unique ID for subscriber

① Registration (DB stores IMSI, security) ② Discovery (faster conn. broadcast, pick tower via bmb scan + best) ★ IMEI: physical device ★ MSISDN: phone number

③ Attachment (IMSI + attach req  $\xrightarrow{\text{radio GW}}$  mobility manager → EMM data plane → record in DB) ④ Roaming (home routing vs. local breakout & surr)

⑤ Handover: closer to new tower → disconnect; co-exp bw tower, + serving device mobility manager, radio gateway ★ want highest RSSI (closest distance)

TCP throughput =  $\sqrt{2} \times (MSS / RTT \times \sqrt{J})$