

CS2105 Notes.

Date

No.

I Overview

Hosts

End systems running applications.

Packet

A package of messages sent between hosts

Store-And-Forward

A packet can only be sent to the next packet switch after it is fully received. Each packet is transmitted at full link capacity (link transmission rate or bandwidth)

Routers VS. Switches

Routers are often at the network core (with high traffic) while switches are often at the network access (boundaries).

Data Centers

A concentration of physical servers that aggregate as a (logical) server (primarily to handle high traffic).

Connecting End Systems with Edge Routers

- ① Residential Networks.
- ② Institutional Networks.
- ③ Mobile Networks

Delays

① Transmission delay

How fast the link can eat/handle.

$$d_{trans} = \frac{1}{R} \frac{\text{bits}}{\text{bps}}$$

$$\text{End-to-end delay} = \frac{d_{trans}}{R}$$

② Propagation delay

how fast the packet can travel (capped by the light speed).

$$d_{\text{prop}} = \frac{d}{s} \quad \begin{matrix} \text{distance} \\ \text{speed} \end{matrix}$$

③ (Nodal) Processing

Overheads incurred upon receipt, including checksum, destination (forwarding table) lookup.

$$d_{\text{proc}} < ms$$

④ Queuing Delay

Time spent waiting for turn to be sent out.

Depends on the traffic.

The total time needed for Nodal Processing.

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

Routing

Determines source-destination route to be taken by packets.

Circuit-Switching

Each link has N circuits exclusively for a pair of end systems (src, dest). Performance is guaranteed since there is no sharing.

Internet Service Provider (ISP)

Tier I: National/International coverage Tier II: Regional
Internet Structure: Network of Networks

ISP of both tiers.

IXP (Internet Exchange Point) between ISPs and content providers' private networks. Sometimes, a peering link communicates two ISPs.

Throughput

Usually measured in bps. The rate of bits actually transferred between 2 end systems.

Network Protocols

Define format, order of messages sent and received among network entities, and actions to be taken on message transmission, receipt.

Layers:

① Application: Supporting network applications
FTP, SMTP, HTTP.

② Transport: Process - Process data transfer
TCP, UDP.

③ Network: Routing of datagrams from source to destination

IP, routing protocols.

④ Link: Data transfer between neighbouring network elements.

Ethernet, 802.11 (WiFi), PPP.

⑤ Physical.

'traceroute'

- 3 probes each element along the path from src to dest to know how delays are distributed along that path.

- The 3 packets may be sent to different routers.

- * x n indicates no responses or refusal to respond within TTL. for n packets.

II Application Layer.

Server - Client Architecture

Server: Always-on host. Well-known, permanent IP.

Client: Might be intermittently connected and may have dynamic IP. Doesn't communicate with each other.

P2P Architecture

Self-salability: New peers bring in new capacity as well as new demands.

May intermittently connect and change IP.

Server Process & Client Process.

Who waits for connection invitation & who initiate a connection. (Can exist in P2P as well).

Addressing Process / Thread.

Identifier = IP + Port number (a 16-bit int).

Typically, one thread for a port number.

TCP Properties.

- Reliable Data Transfer (RDT).
- Flow control
- Congestion control
- Has connection

No guarantee of timing, min throughput and security

UDP Properties:

- Unreliable Data Transfer (UDT).

Url = Host Name + Path Name

www.sample.edu / someDept / pic.gif

HTTP: Hypertext Transfer Protocol

An App-layer protocol for client-server architecture.
 Uses TCP (has connection) (at port 80 for webs).

Stateless: servers know nothing about past requests.

Non-Persistent HTTP: (v.1.0 only has this).

Both use
TCP.

/ At most one object sent over a TCP connection
(At least 2x RTT for an object: including the RTT for establishing connection).

Persistent HTTP: (Default in v.1).

Multiple objects over a single connection.

Don't forget the first RTT to establish the connection!.

- When loading a website, one ^{request}RTT for getting the base html file, then one ^{request}RTT for each referenced objects.
- For HTTP, non-parallel and parallel flavours are offered.
- Aside from RTTs, also consider file transmission time.

? HTTP Methods.

- ① GET: Retrieve an object specified by a URI
- ② HEAD: Same as GET, but the response will contain the status line and the header section only.
- ③ POST: Submit data for processing
- ④ PUT: Update data. Need to include the entire body.
- ⑤ DELETE
- ⑥ CONNECT: Establishes a tunnel
- ⑦ OPTIONS: Describe communication options available.
- ⑧ TRACE: For debugging. Loop-back test.

? Methods for Form Uploading .

① POST: The submission is included in the message body.

② URL: Use GET to upload the object specified in the URL field of the request line. E.g.: <https://www.example.com/res?monkey?banana>

HTTP Response Status code.

200 OK

301 Moved permanently

400 Bad Request.

404 NOT FOUND

505 HTTP Version Not Supported

304 Not Modified.

Cookies

① Cookie header line of the response message (initial cookie identifier created by server).

② Cookie header line of subsequent request messages (so that the server can identify).

③ Cookie file managed by browser (so it knows what to send next time).

④ Cookie and past behaviors in database at the Web.

Web Caches (Proxy Server).

A server that acts as a cache.

Object in cache → return.

Else → get from the original server.

Conditional GET.

The request message includes the last modification date of an object. If up-to-date, the server returns 304 with an empty message body.

Domain-Name System (DNS). — An App-Layer Protocol
 IP address \leftrightarrow Domain name. (many-to-many)
 Over UDP. (but also has TCP options)

Local DNS Name Server

"Default name server" that each ISP has one.
 Acts as a proxy and may send outdated mapping.

Perform recursive query.

Root Name Servers (13 logical servers worldwide
 $\uparrow\downarrow$ with many minors)

.com, .edu, .org Servers.

Top-Level Domain/Tier I

.yahoo.com, .amazon.com, poly.edu, umass.edu Servers Tier II.
 $\uparrow\downarrow$
 (Could be authoritative),
 DNS servers

DNS Query: (Port No. 53 at server) (nslookup) (dig -t)

Iterative: The local DNS server asks every server along a path from the root down.

Recursive: Each lower server asks its parent server. The queries can be cached in servers for some time-to-live (TTL).

* Change ~~IP~~ may not known until all TTL expire

UDP Sockets

- No handshaking.
- Data can be lost or received out of order
- Sender explicitly attaches dest IP and port number to every packet.
- Uses `SOCK_DGRAM`
- Uses `clientSocket.recvfrom(bufsize)`
- Uses `bind()` for server to listen to a tuple (IP, port).

TCP Sockets.

- Each server has a 'welcoming' socket that listens to any request for connection.
- Exclusive socket is created for a single client upon connected. (error-checked, continuous)
- Reliable, in-order byte-stream transfer, There are no app message boundaries, and the app layer needs to figure them out.
- Different sockets do not imply different port numbers.
- Uses **SOCK_STREAM**
- Uses `bind()` and `listen()`.
- Uses `accept()` to create exclusive sockets
- Uses `recv(bufsize)`.

Demultiplexing:

TCP identifies the socket by (src IP, src port, dest IP, dest port).

UDP identifies the socket by only (dest IP, dest port).

TCP & UDP port number: 16-bit unsigned

III. Transport Layer

Transport VS. Network.

Transport: Logical communication between processes.

Network: Logical communication between hosts.

TCP UDP checksum

- including header = \sum checksum ≥ 0 while computing it.
- ① Add all 16-bit integers. Wrap around when there is a carry at the MSD.
 - ② Flip 0 and 1 (1's complement)

Complications of RDT

Date

No.

Problem	Fix
① Bit flips in body	Retransmission and acknowledgement
② Bit flips in ACK Resent duplicates are not recognised	sequence number
③ NAK is unnecessary	Repeat the last Ack as an NAK.
④ Packets might be lost	Timeout and resend
a) Duplicate ACK (means the one just sent is corrupted).	Ignore duplicate ACK. Resend upon timeout anyway. It is still a wait-and-stop protocol. The sender will ensure the current pkt is well received before sending another
⑤ Wait-and-stop is very inefficient.	($V_{\text{sender}} = \frac{D_{\text{trans}}}{RTT + D_{\text{trans}}}$ neglecting time taken by receiver) Introduce pipelined protocols.
⑥ In pipelining, packets may be lost or received out-of-order	Go-back-N and selective repeat.

Go-Back-N.

- Sender maintains a single clock for the earliest unacknowledged packet. P. (the send base).
- Receiver sends the highest continuous sequence number every time it receives a packet.
- Receiver discards not-in-order packets.
- Sender resends a whole sliding window starting from P upon timeout.

- Receiver acknowledges duplicate packets.
- Sender ignores duplicate acknowledgements. No.
- The sliding window moves as soon as another consecutive ACK comes. Moves = send new packets in the window.
-

Selective Repeat

- Receiver acknowledges individual correctly received pkt. and buffers out-of-order pkt for delivering in-order messages to the upper layer.
- Sender maintains a timer for each unacknowledged pkt. Upon timeout, resend.
- On receiver's side, duplicate (already buffered) pkt also need to be acknowledged. (The sender may not think the pkt is ^{not} acknowledged by the time it resends or the ACK might be lost as well.)

TCP

- Point-to-point.
- Has handshaking.
- Full duplex data (bi-directional data flow)
- In-order byte-stream (no app-level message boundaries).
- Pipelined (dynamic window size set by congestion/flow control).
- Has 2 buffers, one at sender and the other at receiver
- Sequence # = byte number of the first byte in the body.
- ACK # = the next expected byte # = the last seq # + the length of the last message.
- Cumulative ACK
- Initial seq # is randomly chosen.

Sender and receiver both need to give an ack # as well as a seq # in a message. The seq # A sends = ack # B sends to A last time.

Flow control:

- Receiver 'advertises' an acceptable rate of transmission.
- Sender controls the number of packets on the fly.
- Such that the receiver's buffer won't overflow

Connection Establishment: 3-way Handshake

① C → S SYN-bit = 1, seq# = x.

② S → C SYN-bit = 1, seq# = y, ack# = x+1

even if the segment is empty, ACK-bit = 1.

③ C → S ACK-bit = 1, ack# = y+1. This is an ack of ②

Connection Teardown: 4-way Handshake

① C → S FIN-bit = 1, seq# = x.

② S → C ACK-bit = 1, ack# = x+1.

③ S → C FIN-bit = 1, seq# = y

④ C → S ACK-bit = 1, ack# = y+1.

Client cannot send data, but can receive data.

Client waits for 2 max segment lifetime just in case server did

not receive ④ and keeps sending things.

Receiver Actions

Event	Action	Purpose
Ack'd: 3 4 5 "; "	Wait for 500ms If no second segment, send ack# = 4 (ACK'ing 3)	Not to send too many ack to overwhelm the net
Ack'd: 3 4 5 "; " pending segment arrived	Send ack# = 5 immediately	To avoid time out

Event	Action	Purpose
acked: 3 4 5 ... []	Send duplicate seg# = 3 immediately	For fast transmit, telling the sender there is probably a loss.
acked: 3 4 5 [] [] just arrived pending	Send seg# = 4 immediately	To avoid unnecessary timeout and retransmission
TCP doesn't specify how to handle out-of order segments.		

Sender Actions.

Pretty much the same as Go-Back-N, only that the sender only resends the segment with the lowest seq# upon timeout.

RTT measurement & setting

$$\text{SampleRTT} = \text{Ack-receipt} - \text{Tseg-trans} \quad \text{ignoring re-trans}$$

$$\text{EstRTT}_{k+1} = (1-\alpha) \text{EstRTT}_k + \alpha \text{SampleRTT}_k$$

$$\text{DevRTT}_{k+1} = (1-\beta) \text{DevRTT}_k + \beta |\text{SampleRTT}_k - \text{EstRTT}_k|$$

Typically, $\alpha = \frac{1}{8}$ and $\beta = \frac{1}{4}$. Timeout = EstRTT + 4 DevRTT

TCP Fast Transmit.

Don't wait for timeout.

Resend as soon as 4 ack's (3 duplicates) are received.

APP. Protocol	Tpt. Protocol	Port
HTTP	TCP	80 (default)
HTTPS	TCP	443 (default).
DNS	UDP	53
SMTP	TCP	25
DHCP	UDP	67 (svr) 68 (client).
RIP	UDP	520

CS2105 Notes. II.

IV Network Layer

- Network layer handles host-to-host communication

Functions

- Forwarding: determines output port on router based on the input port.
- Routing: determines the overall route.

Data Plane.

- Local, per-router
- Forwarding function

Control Plane

- Network-wide logic
- Two approaches:
 - Traditional Routing Algorithms: Embedded in ^{Router}
 - Software-defined networking (SDN): Implemented in remote servers.

IP addressing

- 32 bits for IPv4, 128 bits for IPv6.
 - Interface: connection between host/router.
 - Routers ^{typically} can have multiple.
 - Hosts can have multiple. (different media)
- Different interfaces have different IP addresses
 \Rightarrow A router can have different IP addr.

Subnets.

- Hosts in the same subnet can, logically, "directly" communicate, without routers.
- Have the same IP prefix.

CIDR Classless InterDomain Routing

- a.b.c.d/x.

Length of IP subnet Prefix.

+ [ID: 1
offset: 0
flag: 0
...]

[ID: 2
offset: 0
flag: 0
len: 1200]

[ID: 3
offset: 0
flag: 0
...]

[IP: 2 flag: 1 offset: 0 len: 500] [IP: 2 flag: 1 offset: 60 len: 500] [IP: 2 flag: 0 offset: 120 len: 240]

• Calculation Trick:

2^{32-x} public IP addr can be allocated to hosts in this subnet.

Assume $0 \leq x \leq 8$, then given a.b.c.d/x, IPs that are in the same subnet are from $a - (a \bmod 2^{8-x}).0.0.0$ to $(a - (a \bmod 2^{8-x}) + 2^{8-x} - 1).255.255.255$.

Remember to subtract ^{the} number of reserved addr if encountered in the range.

- Subnet mask is x' 1's followed by 0's.
- Special IP addresses.

0.0.0.0/8 Meta-addr.

127.0.0.0/8 Loopback

10.0.0.0/8]

172.16.0.0/12 } Private.

192.168.0.0/16

255.255.255.255 Broadcast.

- ICANN: Internet Corporation for Assigned Names and Numbers.

Allocates IPs to ISPs.

Largest Prefix Matching.

- Forward the packet to the router whose subnet IP has the longest shared prefix with the destination IP. \hookrightarrow guaranteed to be unique.

- Note: Two IP addr can have the longest common prefix even if their decimal representations look very different.

Getting an IP.

- Hardcoded
- DHCP.

▪ Plug-and-play; Dynamically gets an IP.

DHCP } server
| client } port number { 67

A yieldaddr has a lifetime -

Date

No.

- Can renew lease on addr in use.
 - Allows reuse of addr (only when connected).
 - Support mobile users.
- Execution of DHCP.
① Host broadcasts "DHCP discover" [optional].
② Server "DHCP offer" [opt].
③ Host "DHCP request".
④ Server "DHCP ack".

IP Datagram Fragmentation.

An arriving IP datagram can be decomposed into several smaller IP datagram at a router. They will only be reassembled at dest. +

On the length field at header.

UDP len. = header (8 bytes) + body, in bytes

TCP len = header (20 - 60 bytes) only.

4 bits $\Rightarrow 2^4$ values for 20 ~ 60 ?

Express $\frac{\text{len in bytes}}{4} \Rightarrow 5 \sim 15$.

i.e. in 4-byte words.

IP len = header (20 bytes) + data, in bytes.

IP offset = displ first byte in the same fragment).

in 8-bytes

IP Format.

• Frag: 1: not the last. 0: the last.

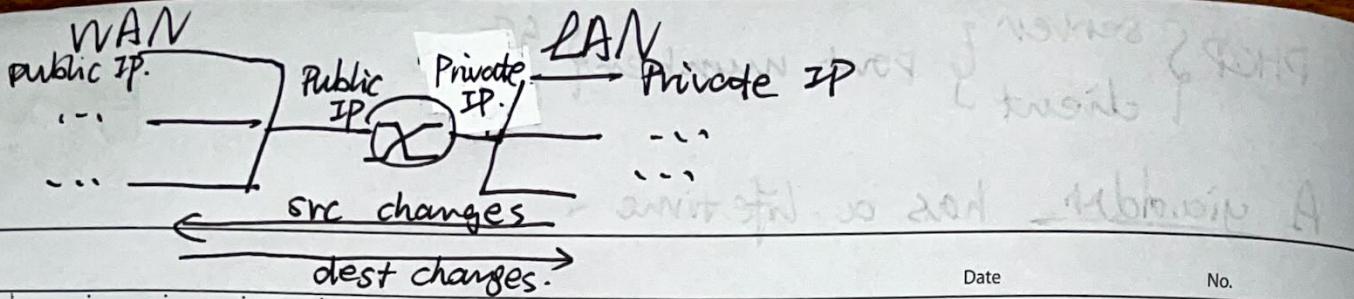
• Offset: relative position to the original beginning

$$\text{Offset}_{i+1} = \text{Offset}_i + (\text{length}_{i+1} - 20)/8$$

i.e. in 8-byte words.

NAT:

• Intuition: exploit the port number space to make up for shortage of IP addr space



- A "façade" router that sits between WAN and LAN has a public IP and a private IP.
- Implementation: Maintain a 1-to-1 mapping.
 - Replace (private IP, port) with (router IP, dedicated port) when sending out
 - Remember the mapping.
 - Replace (router IP, dedicated port) with (private IP, port) when sending in.
- Also a security plus in the sense that the outside world cannot directly address a private IP.
- Challenge: "NAT maze" Peer-to-peer apps don't work directly.

Routing.

- Intra-As: No policy decisions. Focus on performance.
- Inter-AS: Policy may dominate.

Cost:

- Could be constant.
- Or inversely related to bandwidth.
- Or related to congestion.
- Or distance
- Or financial.

$$c(x,y) = \text{cost}(x,y)$$

= ∞ if x and y don't neighbour.

$$D_x(y) = \text{cost of the least-cost path from } x \rightarrow y$$

Bellman - Ford Equation.

$$D_u(z) = \min_{a \in N} \{ c(u,a) + D_a(z) \}. \quad N: 1\text{-hop neighbours}$$

distance to every node.

Date

No.

Routing Algorithm.

① Receive distance vectors from neighbours

② Compute its own distance vector.

$$\begin{bmatrix} a_1 \\ \dots \\ a_m \end{bmatrix} = \begin{bmatrix} \min(C_{a,b} + b_1, \dots, C_{a,z} + z_1) \\ \dots \\ \min(C_{a,b} + b_m, \dots, C_{a,z} + z_m) \end{bmatrix}$$

③ Notify neighbours only when local distance vector changes, possibly due to:

- a) Link cost changes.
- b) Topology changes.
- c) DV update from neighbours.

Note: The above algorithm terminates/stabilizes in finite steps from $t=t_0$ if there are no further a) or b) changes from t_0 onwards.

However, the number of steps may be greater than the diameter of the graph.

RIP.

- Uses hop count as the cost \rightarrow not sensitive to traffic.
- Exchange routing table every 30 sec.
- "Self-repair": assume a neighbour has failed if no update for 3 min.
- Distributed, iterative, and asynchronous.

ICMP.

- Error reporting: unreachable host/net/port/protocol.
- Echo request/reply: ping. (uses a connection).
- Exchange info about net.

For example, a packet can be dropped due to TTL expired \Rightarrow send an ICMP to notify the sender.

V Link Layer.

Date

No.

Handles neighbour-to-neighbour communication over a single link.
Deals with addressing, protocol, and errors.

Concerns and Services.

- Framing.
- Link Access Control.
- Error Detection.
- Error Correction.
- Reliable Delivery: typically for error-prone links such as wireless links.

Link layer is implemented in "adapter" or a chip.
Adapters are semi-autonomous, implementing both link & physical layers.

(Even) Parity Checking (Single Bit)

- Number of 1's in the data bits is even.
- Can detect odd number of single bit errors but not even number of single bit errors.
- In practice, $\text{Pr}[\text{undetected error}] \sim 50\%$.
- Works well if errors are independent.

2.1 Parity checking.

- For an $m \times n$ matrix, there are $m+n+1$ parity bits.
- The $(m+1, n+1)$ bit should be the same regardless of whether it is evaluated by row or col.
- Can detect and correct single-bit errors.
- Can only detect a two-bit error.
- Can detect odd number of single-bit errors.

Cyclic Redundancy Check

r: length of CRC

R: CRC number

D: data, viewed as a binary

G: generator of $n+r$ bits, agreed by sender and receiver beforehand.

① Append r 0's to D, resulting D' .

② $R := D' \bmod G$ under Modulo 2 Arithmetic.

$$x+y = x-y = x \oplus y$$

③ Send and receive (D', R) , meaning D concat. R .

④ Receiver checks if $(D', R) \bmod G = 0$. If not, some error exists

Intuition: divisor | (dividend - remainder)

$$\begin{aligned} \text{dividend} - \text{remainder} &= \text{dividend} \oplus \text{remainder} \\ &= \text{dividend} + \text{remainder} = D' + R = (D, R). \\ \therefore G &\mid (D, R). \end{aligned}$$

- Can detect all odd number of single-bit errors.

- all burst errors of $\leq r$ bits

- all burst errors of $> r$ bits with prob $1 - (\frac{1}{2})^r$

- A.K.A. polynomial code. k-bit frame = k coefficients for a polynomial of degree $k-1$.

Types of Network Links

1. Point-to-Point Link: PPP, SLIP.

2. Broadcast Link (shared medium): Every node sharing the link receives a copy.

Multiple Access Protocols

- Random Access : Possible collisions \Rightarrow recovery is concerned
- Taking Turns :
- Channel Partitioning : Time/Frequency.

Ideal Multiple Access Protocol (given capacity R)

- 1. Collision Free (Mutual Exclusion).
- 2. Efficient: If only 1 node is transmitting, it can achieve R . (Progress).
- 3. Fairness. Each node has an average of R/N . (Starvation Free).
- 4. Fully Decentralized.
 - No out-of-band communication.

TDMA/FDMA.

- channel
partitioning
- Fixed timeslots/frequency bands that are pre-agreed / hardcoded.
 - For FDMA, the average is still $\frac{R}{N}$ because the difference in frequency on different bands are negligible.
 - Satisfies Property 1, 3, 4.

Polling.

- Taking Turns
- The master node polls every slave node in round-robin. Admitted slaves take turn to send at most m packets/frames.
 - Need to consider polling overhead.
 - Satisfies 1, 2, 3.
 \rightarrow not completely.

Token Passing.

- A special frame, called token, is passed along a (logical) ring topology (the underlying physical topology does not have to be a ring).
- When holding the token, can send something forward if finished.
- Satisfies 1, 2, 3, 4.
 \rightarrow not completely.
- Need to consider token loss and node/link failure.

data loss
system bug.

broken topology.

On Random Access: transmit at full rate R.
no a prior coordination.

On ALOHA collision detection: Listen while sending.

For some propagation delay too long for a collision to be detected while transmitting \rightarrow it relies on ACK.

Slotted ALOHA.

- All frames are L bits long.
 - (Time Quantum) = time for sending a frame.
= Timeslot.
- ① If there is a frame to send, wait for the beginning of the next slot.
 - ② No collision \Rightarrow success.
Collision \Rightarrow failure \Rightarrow retransmit in each following slot with prob p . until success.
- Satisfies 3, 4.
 - About 2:
 - a) Full rate once a node starts.
 - b) Only 37% ($\frac{1}{e}$) of R if there are many active nodes.
 - c) Slots are wasted due to collisions or being empty (the prob leads to this).

Pure/Unslotted ALOHA.

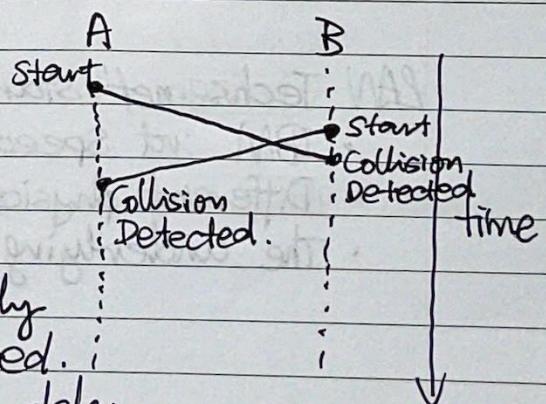
- No timeslots + No synchronisation.
- ① Transmit immediately.
 - ② For failure, wait for 1 unit time, retransmit the frame with prob p .
- Chance of collision increases, resulting in halved overall capacity (18%).
 - Satisfies 3, 4.

CSMA.

- Listen before speaking.
- Collisions are still possible.
- Unslotted.

CSMA/CD.

- Abort transmission immediately if there is a collision detected.
- Retransmit after a random delay.
- Satisfies 2, 3, 4.



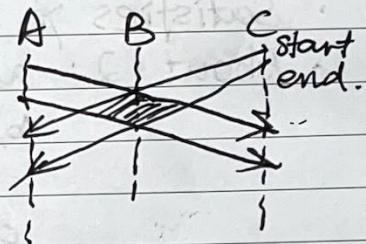
* Constant delay prob p only burdens the future. It can even cause more collisions.

Back-off Algorithms. (Binary Exponential Backoff).

- After the m-th collision:
 - choose $k \in \{0, \dots, 2^m - 1\}$ at random.
 - such that the prob of retransmitting is $p = \frac{1}{2^m}$.
- Intuition: More collisions imply heavier loads.
- Possible Improvement: Can start over again from $m=0$ after a period of time, otherwise old past history influences now too much.

Minimum Frame Size.

- Ethernet: 64 bytes.



Only B detects the collision.

MAC Address.

- A.k.a., physical/LAN address.
- Upon receipt of a frame, a device checks if the destination is its MAC addr
 - yes \rightarrow extract datagram and hand over it.
 - no \rightarrow discard.
- Typically 48 bits, burnt in ROM (sometimes software settable).
- Broadcast Address: FF-FF-FF-FF-FF-FF.

802.3 Ethernet Standards.

- Different speeds: (2M, 10M, 100M, 1G, 10G, 100G) bps.
- Different physical layer media: cable, fiber optics.
- The underlying protocol is CSMA/CD with backoff.

- Frame Structure:

8 bytes						6	6	2	46~1500	4	MTU
Preamble	Dest. Addr	Src Addr	Type	Data	CRC						
10101010 × 7 (AA) ₁₆ , 10101011 × 1 (AB) ₁₆ . used to synchronise clock rates. (the width of a bit), important when there are 19 or 20 0's.	6	6	2	46~1500	4						

T
 }
 10101010 × 7 (AA)₁₆,
 10101011 × 1 (AB)₁₆.
 used to synchronise
clock rates. (the
width of a bit),
important when
there are 19 or 20 0's.
 }
 The higher layer
protocol, e.g., Novell IPX,
AppleTalk, ARP.
 }
 No
 retransmission
 if corrupted.
 ; unreliable
 dropped
 frames will
 only recovered
 when the
 higher layer
 offers
 retransmission

Ethernet Topology.

• Bus

- All nodes can collide with each other.
- If the backbone cable fails, the entire network will fail.
- Very slow due to collisions → not ideal for large networks.

• Hub (Star)

- A physical layer device that acts on individual bits rather than frames.
- Re-create and boost the bit signal, send it to all other interfaces → logically still a bus.
- Cheap, easy maintenance, very slow, not ideal for large networks.

• Switch (Star)

- A link layer device.
- Acts on individual frames in a store-and-forward fashion.
- No collisions.

Q: How can A know if B is in the same subnet?

A: Via DHCP and subnet name and mask.

Date

No.

Ethernet Switch.

- Selectively forwards frame to one or more links
- Transparent; hosts are unaware of switches.
- Plug-and-play (self-learning).
- Multiple simultaneous transmissions with buffer but without collisions.
"Simultaneous" in the sense that $A \rightarrow A'$, $B \rightarrow B'$ can be done in parallel.

That's why nodes in LAN are logically one hop apart.

Frame Filtering/Forwarding.

- Forwarding table:
 $\langle \text{MAC addr, Interface, TTL} \rangle$ → renewed upon every receipt.
- Each tuple means: "I have seen MAC x from Interface y." This tuple is valid for TTL.
- Suppose the table has only one record $\langle x, y, 60 \rangle$.
 - On receipt of x on y: discard.
 - On receipt of x on z \neq y: send to y.
 - On receipt of w on z: send to every interface except z. "flood."

ARP

- A link layer protocol.
- Each IP node has an ARP table about hosts in the same subnet.
 $\langle \text{IP addr, MAC addr, TTL} \rangle$ → typically a few min.
- Plug-and-play.

Sending Frame in the Same Subnet. ($A \rightarrow B$)

- ① If A knows B's IP. in the ARP table → trivial.
- ② Otherwise, use ARP.
 - A broadcasts to $(FF \times 6)$ with B's IP.
 - B replies with B's MAC.
 - A caches B's MAC until TTL expires.

* LAN \neq subnet.

It may contain multiple subnets.

Sending Frame to Another Subnet.

- Cannot just set dest = B's MAC, otherwise it will be discarded by all nodes in the same subnet including (the link layer of) the router, due to MAC mismatch.
- Correct approach: set IP dest = B's IP, MAC dest = R's MAC. Upon receipt at R, R changes MAC dest = B's MAC, MAC scr = R's MAC.
 $\therefore <\text{MAC src, MAC dest}>$ changes when subnet changes

* What if a node does not even know about its gateway router?

- ① There is a default gateway in the ARP table, usually.
- ② Via DHCP, get the subnet name and mask, and hence the gateway router
- ③ Use ARP proxy reply. The ARP broadcast query with B's IP will now ^{be} handled by the router as a proxy. The router will reply A that it can reach B by its MAC. \Rightarrow From now, every frame from A to B will have MAC dest = R's MAC, IP dest = B's IP.

VI Security.

Malicious behaviors include:

- Eavesdrop: intercept
- Insert
- Impersonation
- Hijack : Man in the middle. Take over the connection
- DoS.

Aspects of Security

- Symmetric Key → Confidentiality: only someone can understand
- Signature & Certificate → Authentication: confirmed identity.
- Hash → Message Integrity: not altered without detection.
- MAC • Access & Availability.

Notation

$$K_A K_B(m) = m \rightarrow \text{the plaintext.}$$

Caesar's Cipher

- An example of substitution cipher.
- Key: the shift number.
right shift by n letters: every char is moved to the right by n letters in the alphabet.

plaintext	ciphertext for n=3
a	d
b	e
...	...

- Easy to do brute-force.

Monoalphabetic Cipher.

- A one-to-one mapping of letters.
- $26!$ possible keys, $26! - 1$ of which are meaningful
- Easy to break with statistical analysis. (frequency analysis). Some two-and-three-letter occurrences are often together (is, are, of, to, -ion ...).

Polyalphabetic Encryption

- Addresses the problem, also the fundamental weakness, of having only one mapping.
- Uses multiple mappings. Select n substitution ciphers and a pattern $C_1 C_3 C_2 \dots$
- $K(m_1 m_2 m_3 m_4 m_5) = C_1(m_1) C_3(m_2) C_2(m_3) \dots$

Block Cipher

- Processes m in blocks of K bits.
- Uses a 1-to-1 mapping. e.g., $001 \rightarrow 111, 010 \rightarrow 100$.
- Number of keys $(2^K)^K$
- Hard to break with frequency analysis if K is a good choice.

DES:

- 56-bit symmetric key, 64-bit block.
- Broken in < a day.
- 3DES: encrypt 3 times with 3 different keys.

AES

- 128-bit blocks, 128-, 192-, or 256-bit keys.
- If a computer can break DES in 1 sec, it will break AES in $14^9 \times 10^{12}$ yrs.

Setting of Attacks.

- Ciphertext Only.
- Known-plaintext. (may not be totally known)
- Chosen-plaintext (use the encryption algo as a black box arbitrarily).

Public Key (Asymmetric) Encryption.

- $K^+(\cdot) \neq K^-(\cdot)$, but $K^-(K^+(m)) = m$

Given K^+ , it is merely impossible to find K^-

- RSA is an example.

Modulo Arithmetic.

$$(a \stackrel{+}{\times} b) \bmod n = [(a \bmod n) \stackrel{+}{\times} (b \bmod n)] \bmod n.$$

$$\therefore (a \bmod n)^d \bmod n = a^d \bmod n.$$

RSA

① Choose 2 large prime numbers p, q .

$$\textcircled{2} \quad n = pq, \quad z = (p-1)(q-1)$$

③ choose $e < n$ that is co-prime with z

④ choose d such that $ed \bmod z = 1$.

$$\textcircled{5} \quad K^+ = (n, e), \quad K^- = (n, d).$$

Encryption.

$$c = m^e \bmod n.$$

Decryption.

$$m = c^d \bmod n = (m^e \bmod n)^d \bmod n$$

Note that K^+ and K^- are commutative, meaning
 $K^+(K^-(m)) = K^-(K^+(m))$

The underlying challenge: factorising a big number is hard.

Session Key.

- DES is $\times 100$ faster than RSA, but need K_s in advance.
- Use RSA to deliver K_s .
- Use K_s to en/decrypt messages for the session.

Motivation for Alteration Detection.

1. checksum only detect accidental errors not attacks
2. CRC clashes are common and easy to find.

Minor change in input \rightarrow minor change in output

Cryptographic Hash Function.

• Message of arbitrary length \rightarrow fixed-length digest.

• Computational infeasible to find $x \neq y$ s.t.

$$H(x) = H(y)$$

• MD5 $\rightarrow \mathcal{S}^{128} = \{0, 1\}^{128}$ Cryptographically broken.

• SHA-1 $\rightarrow \mathcal{S}^{160} = \{0, 1\}^{160}$ deprecate.

-fingerprint

Does not have confidentiality

Message Integrity.

• Use a shared "Authentication Key" s ✓

• Send $(m, H(m+s))$; Receiver checks if $H(m'+s) = H'$

Message Authentication Code.

Password.

- Saved in hash \rightarrow can only be reset

Digital Signature

- Desired properties

↳ Verifiable \rightarrow both parties can check

↳ Unforgeable \rightarrow only authorized parties can generate.

- Send $(m, k^-(m))$ or $(m, k^-(H(m)))$

signature = message encrypted with private key
fingerprint

- Receiver verifies $k^+(k^-(m)) = m'$ or

$$k^+(k^-(H(m)))' = H(m')$$

Certification Authorities.

- A list of CAs for receiver to query trusted public keys
- These CAs' messages are also signed. \Rightarrow their signatures are maintained as universal knowledge by Trusted Root CAs.
- Entity E provides CA with its proof of identity. \Rightarrow the certificate contains $(K_E^+, K_{CA}^-(K_E^+))$.

Firewalls

- Isolate an organisation's internal network from the outside world.
- Re-construct segments/datagrams/messages if needed.
- Three types.

1. Stateless packet filters *

2. Stateful packet filters

3. Application gateways

- Limitations: IP spoofing, performance bottleneck

Stateless Packet Filter.

- Decides to forward/drop based on addr, port #, ICMP message type, SYN/ACK bits.

Access Control List.

- Conditions are evaluated top-down.

VII Multimedia Networking.

At APP-layer and Transport layer.

3 Types of Applications.

1. Streaming stored audio, video.

Streaming : playout begins before downloading the entire file.

Stored (at servers/CDNs) : can transmit faster than it is played. \Rightarrow buffering.

2. Conversational

Interactive, delay-sensitive.

Intolerable with >400 ms

3. Streaming live

Typically done with CDNs.

Video.

- Frames played in sequence at constant rate.
 \downarrow
 Arrays of pixels \hookrightarrow bits.

- Exploits Redundancy
- Spatial Coding : repeated colors \rightarrow count the occurrences
 - Temporal Coding : similar frames \rightarrow send differences only
 - CBR
 - Not responsive to complexity change.
 - High rate to handle sudden increase in complexity.
 - Suitable for real-time encoding.
 - VBR
 - Changes according to how spatial/temporal coding changes
 - Suitable for on-demand video.

Audio.

- Sampling (quantized to bits)

- At constant rate:

 - Telephone: 8000 /sec.

 - CD 44100 /sec

- Bitrate = sample size \times sampling rate.

 - CD: 1.411 Mbps.

 - MP3: 96, 128, 160 kbps.

 - Internet telephony: ≥ 5.3 kbps

Continuous Playout Constraint.

- Motivates
- Once client starts playing, the speed must match
 - but there are jitters.

Other Challenges

- Client Interactivity: pause, fast-forward, rewind, jump etc.
- Lost/retransmitted packets.

Buffering.

- Initial fill of buffer until playout at fp.
- Playout begins.
- Buffer level = $x(t) - r$ varies.
fill rate varies \leftarrow \rightarrow constant playout rate.

Consider the average, $\bar{x}(t)$

$\bar{x}(t) > r \rightarrow$ eventually full

$\bar{x}(t) < r \rightarrow$ eventually empty. (buffer starvation)
less likely with larger delay

∴ The initial fill should be large enough to absorb the variance of $x(t)$.

Calculation Tricks:

- Measure everything in 'chunks'

Push-Based Streaming

- Uses UDP. \rightarrow no congestion control \rightarrow no max rate restrictions.
- Short playout delay. (2-5 sec).
 - Error recovery at app-level if time permits.
- Video chunks encapsulated using RTP (seq #, timestamp, encoding ...).
- Control Connection maintained by RTSP (usually over TCP)
 - Handles client command (pause, play ...).
 - Establishes/controls media sessions.
- Drawbacks
 - More cost/complexity due to a separate RTSP.
 - May not go thru firewalls.

Pull-Based Streaming

- Over HTTP/TCP.
- Easier to get thru firewalls.
- Infrastructure (CDNs, Routers) fine-tuned for it.
- Fill rate fluctuates.
- Greater playout delay.
- Congestion Control exists.
 - Additive increase V.S. Multiplicative decrease.

VoIP

- End-end-delay requirement.
- Delay Loss \rightarrow
- < 150 ms : good. includes app-level (packetization)
 - > 400 ms : bad. playout delay.

Network loss \rightarrow

- Loss over 10% \rightarrow unintelligible.
- Packets only generated during talk spurts.
- 20 ms chunks at 8 Kbps = 160-byte chunks
- $(\text{chunk} + \text{header}) \in \text{TCP or UDP}$

\downarrow
 { seq # }
 { timestamp }

\hookrightarrow more common.

- Receiver has a constant q_r . chunks generated at t will be played at $t+p$.
 - Large $q_r \rightarrow$ less loss (fewer "too-late" chunks)
 - Small $q_r \rightarrow$ better exp.
 - But no q_r guarantees optimal performance

Adaptive Playout Delay.

- Estimates net delay, adjusts playout delay (Yes there are many playout delays for a conversation, unlike videos). of each talk spurt.
- Silent periods are compressed/elongated.
- EWMA:

$$\underline{d_i} = (1-\alpha) \underline{d_{i-1}} + \alpha \left(\frac{\text{time received}}{\text{time sent}} - t_i \right)$$

delay est.
 after i-th
 chunk.

small.
e.g (0.1).

time received
time sent

Apply after
receiving
the i-th
chunk.

$$\underline{v_i} = (1-\beta) \underline{v_{i-1}} + \beta \left| \underline{r_i} - t_i - \underline{d_i} \right|$$

est. avg
 deviation
 of delay

small.

last actual
delay
predicted
next
delay.

$$(\text{playout time})_i = t_i + d_i + 4v_i$$

Only
applicable
for the
first packet
in a spurt!!

Remaining packets in the same spurt are played out periodically.

Forward Error Correction (FEC) \rightarrow VoIP

1. Simple FEC.

- 1 extra chunk XOR-ed by n chunks
- lower performance by multiplying $\frac{n}{n+1}$.
- longer playout delay \rightarrow $n+1$ must arrive before playing.

2. "Piggyback" FEC.

- Piggybacks low-quality stream of the last (and second last ...) chunk. \rightarrow conceal loss

3. Interleaving.

- chunks are split into n small units, disordered, and reassembled to n -unit chunks
- If packet lost, still have most of every chunk.
- longer playout delay because of the need to wait for different chunks to be generated.

DASH.

- Deals with 3 things.

- When to request
- What encoding rate (quality)
- Where to request.

- Data encoded into different qualities, cut into short segments (streamlets, chunks).

- Manifest file.

- ABR.

- Advantages

- Simple server

- No firewall problems

- standard (image) web caching works (due to use of HTTP).

- Disadvantages

- 2-10-sec media segments \rightarrow not suitable for two-way / interactive apps.

CDN

- Option ①: large "mega-server"

- does not scale (with other shortcomings)

- Option ②: Store copies on geographically distributed sites.

- Enter deep; get deep into access networks, close to end users.
- Bring home; larger clusters at IXPs near access networks.

Execution:

- ① CDN stores copies at different nodes
- ② Client requests for content and gets manifest.
- ③ Client retrieves the highest quality affordable.
- ④ May dynamically choose different rate/copy in case of high-latency/congestion/loss/failure

Appendix

About DHCP

1. Implemented at internet edge.
- * Gives the following:
 - IP Addr of first-hop router.
 - Name & IP addr of DNS server.

3. Network Mask

Subnet Mask Numbers: 128, 192, 224, 240, 248, 252, 254, 255.

Autonomous Systems

- Can have multiple routers and subnets.
- Belong to an organization/institution/entity, which can have multiple ASes.