

# CS2105 Comp. Networks Notes

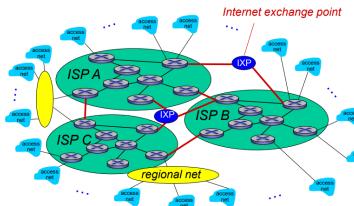
AY23/24 Sem 1. [github.com/gerteck](https://github.com/gerteck)

## 1. Computer Networks Introduction

- Fundamental concepts and principles behind computer networking, with internet as case study.
- Connected by communication links and packet switches.

### Internet

- The Internet is a network of connected computing devices (e.g. PC, server, laptop, smartphone).
- Such devices are known as hosts or end systems. Hosts run network applications (e.g. Tele, browser, Zoom)
- Packet switching network, users' packets share network resources that are used on demand. Excessive congestion is possible.
- **Network of networks** Hosts connect to Internet via access ISPs (Internet Service Providers), which themselves are interconnected.



### Network Edge (Access Network)

- The access network is the network that physically connects an end system to the first router on a path from that end system to any distant end system.
- E.g. Residential access networks, mobile access networks.

### Network Core

- A mesh of interconnected routers.
- Data is transmitted through network through:
- **Circuit switching:** dedicated circuit per call
- **Packet switching:** data sent through net in discrete "chunks"

### Circuit Switching

- End-end resources are allocated to and reserved for "call" between source & dest:

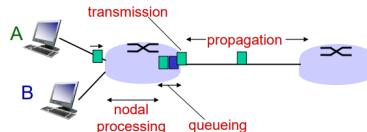
- call setup required, circuit-like (guaranteed) performance
- circuit segment idle **silent periods** if not used by call (no sharing)
- commonly used in traditional telephone network

### Packet Switching

- Host sending function: breaks application message into smaller chunks, known as packets, of length  $L$  bits. Then, transmits packets onto the link at transmission rate  $R$ .
- Packets are passed from one router to the next, across links on path from source to destination.
- link transmission rate is aka link capacity or **link bandwidth**
- **Packet transmission delay** (time taken to transmit  $L$ -bit packet into Link) =  $\frac{L}{R}$
- **Store-and-forward:** entire packet must arrive at a router before it can be transmitted on the next link.

### Delay, Loss, Throughput

Four sources of packet delay, suffered at each node, from source to destination. (**Total nodal delay**)



- **Nodal Processing Delay:** Time required to check for bit errors and determine output link. Typically  $<$  msec.
- **Queueing Delay:** Time waiting in queue for transmission, depends on earlier arrived packets. Typically  $<$  msec. Additionally, possibility of packet loss if router queue full (buffer has reached finite capacity), and drops packet.
- **Transmission Delay:** Time to push (last bit) of packet on wire.  $d_{trans} = \frac{L}{R}$ .
- **Propagation Delay:** Time to propagate to next router.  $d_{prop} = \frac{d}{s}$ , where  $d$  is length of physical link,  $s$  propagation speed.
- End to end packet delay is total time taken for packet to travel from source to destination, consisting of the 4 delays.
- **Throughput:** How many bits can be transmitted per unit time, usually measured for end-to-end communication (as opposed to bandwidth for specific link).

### Internet Protocol Stack

Protocols logically organised into 5 "layers" according to purpose. (Additionally presentation and session layers not included)

- **Application:** Where network applications and app-layer protocols reside. Packet here called message.  
Examples: HTTP, SMTP, FTP
- **Transport:** Transports app-layer messages between application endpoints. Packet here called segment.  
Examples: TCP, UDP
- **Network:** Moves packets (datagrams) from one host to another. Includes IP protocol and other routing protocols.
- **Link:** Moves packet from one node to another. Packet here called frame.  
Example: Ethernet, WiFi
- **Physical:** Moves individual bits within link-layer frame from one node to another. Link and transmission medium dependent.

## 2. Application Layer

### Principles of Network Applications

#### • Client-Server:

- Server waits for incoming requests and provides required services to client. Easy scalability.
- Client initiates contact with server and requests service.

#### • Peer-To-Peer (P2P):

- No dedicated server, instead it relies on direct communication between pairs of intermittently connected hosts called peers.
- Self-scalability, each peer generates workload but adds service capacity by distributing files.

### Process Communication

- **Socket:** A software interface that process uses to send messages and receive messages from network. Generally a combination of IP address and port number.
- **IP Address:** A 32-bit quantity, uniquely identifies host.
- **Port Number:** A 16-bit integer used to identify a receiving process running in a host.

### Requirements of Transport Service

- **Reliable Data Transfer:** Data to sent correctly and completely vs. loss-tolerant.
- **Throughput:** Bandwidth-sensitive apps may need guaranteed throughput of r bits/sec.
- **Timing/Delay:** Real-time applications generally require low delays to be effective.
- **Security:** Encryption, data integrity, authentication.

### Transport Layer Protocols

Two main protocols for the Internet.

#### • Transmission Control Protocol (TCP)

- Reliable data transfer, Connection-oriented service: A handshake required.
- Flow control, Congestion control: Throttle sender when network overloaded
- Security: Can be enhanced at the app layer with Secure Sockets Layer
- Does not provide: Timing and throughput guarantee

#### • User Datagram Protocol (UDP)

- Unreliable data transfer
- Connectionless: No handshake

- No flow control, no congestion control
- Does not provide: Timing and throughput guarantee, security.

### Application-Layer Protocols

An application-layer protocol defines:

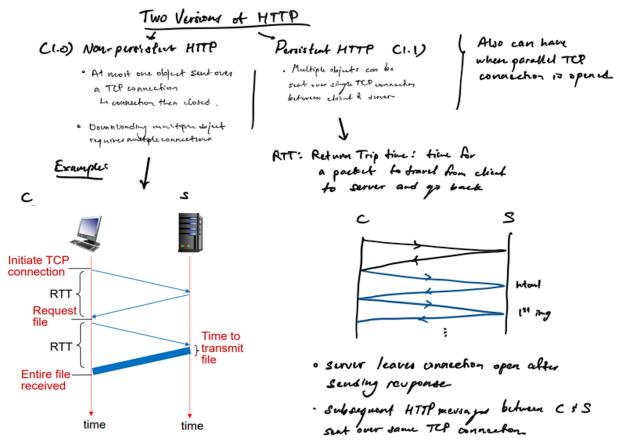
- Types of messages exchanged, e.g. request and response messages.
- Syntax of message types, e.g. fields and how they are delineated.
- Semantics of the fields, i.e. what the information means.
- Rules for when and how to send a message and respond to messages.

### Web & HTTP

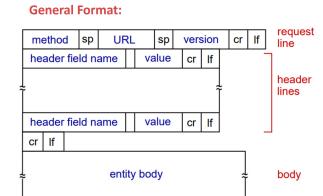
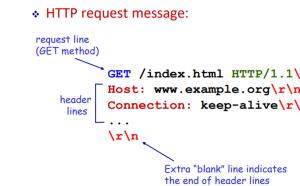
- **Webpage:** Consists of base HTML files and referenced objects. Addressable by URL.
- URL made up of hostname as well as path name. (E.g. <http://www.comp.nus.edu.sg/cs2105/img/doge.jpg>)
- **HyperText Transfer Protocol** is Web's app-layer protocol.
- **Client-server model:** Client requests, receives and displays Web objects, server is Web server that sends objects in response.
- **Stateless:** server maintains no information about clients, and **Reliable:** Over TCP.
- **Three-way Handshake:** Client sends small TCP segment to ask for connection, server acknowledges and responds, client acknowledges and sends it back with request message.

### HTTP versions

- **RTT:** Round trip time, time taken for packet to travel from server and back to client, does not include transmission delay.
- **Non-persistent HTTP:** Response time =  $2 * RTT + \text{file transmission time (per object)}$
- **Persistent HTTP:** Server leaves connection open after sending response, subsequent HTTP sent over same TCP connection. Also uses pipelining, send requests back to back.



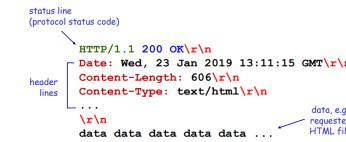
### HTTP Request



- **HTTP 1.0 Methods:** GET (gets object), POST (posts form data), HEAD (gets header without body).
- **HTTP 1.1 Methods:** GET, POST, HEAD, PUT (uploads file to path specified), DELETE

### HTTP Response

#### Example HTTP Response Message

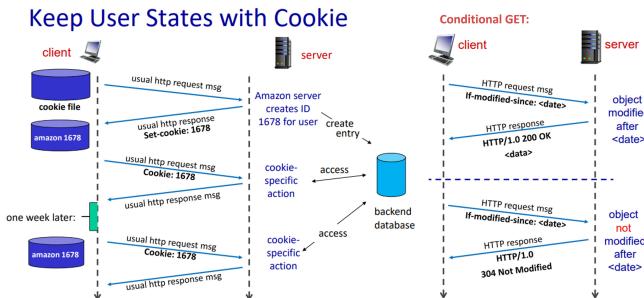


- Status code appears in 1<sup>st</sup> line in server-to-client response message.
- Some sample codes:
  - 200 OK**
    - request succeeded, requested object later in this msg
  - 301 Moved Permanently**
    - request succeeded, new location specified later in this msg (Location:)
  - 403 Forbidden**
    - server declines to show the requested webpage
  - 404 Not Found**
    - requested document not found on this server

## Cookies

- HTTP is designed to be “stateless”, server maintains no information about past client requests.
- Good to maintain states over multiple transactions, e.g. shopping carts
- **Cookie:** http messages carry “state”:
  - 1) cookie header field of HTTP req/res messages
  - 2) cookie file kept on user host, managed by browser
  - 3) back-end database at Web site
- **Conditional GET:** In cache, specify date of cached copy in HTTP request. Server response contains no object if cached copy is up to date.

### Keep User States with Cookie



## Domain Name System

DNS translates between hostname and IP addresses. Client must carry out a DNS query to determine the IP address corresponding to the server name.

### • DNS: Resource Records (RR)

- Mapping between host names and IP addresses (and others) are stored as resource records (RR).

RR format: `(name, value, type, ttl)`

#### type = A

- **name** is hostname
- **value** is IP address

#### type = NS

- **name** is domain (e.g., **nus.edu.sg**)
- **value** is hostname of authoritative name server for this domain

#### type = CNAME

- **name** is alias name (e.g. **www.nus.edu.sg**) for some “canonical” (the real) name
- **value** is canonical name (e.g. **mgnzsqc.x.incapdns.net**)

#### type = MX

- **value** is name of mail server associated with **name**

- **Distributed, Hierarchical Database:** DNS servers form hierarchy to distribute mappings. Contains root servers (for Top Level Domain TLD servers), TLD servers (e.g. uk, sg), Authoritative servers.

- Local DNS servers have local cache, acts as proxy, forward query into hierarchy if answer not found locally.
- **DNS Caching:** Cache mapping, which expires after some time (TTL: time to live).
- DNS runs over **UDP**.

## DNS Name Resolution



♦ This is known as **iterative query**.

♦ This is known as **recursive query**.

• rarely used in practice

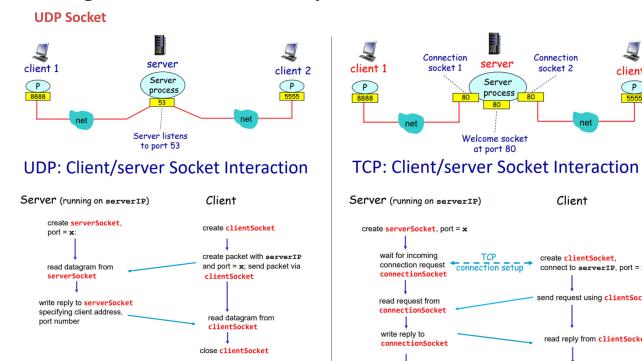
## Socket Programming

Applications treat the Internet as black box, send/receive message through sockets.

- Two types of sockets
- **TCP:** reliable, byte stream-oriented socket
- **UDP:** unreliable datagram socket

### UDP vs. TCP Socket

- **UDP Socket:** Sender attach des IP address + port no. to each packet. (OS inserts add. info source IP and port). Receiver extracts sender IP + port number from packet.
- **TCP Socket:** Attempts to establish TCP connection to server first. Server TCP contacted creates new socket to communicate with client, allows server to talk with multiple clients individually.



### • TCP vs. UDP Differences

- In TCP, two processes communicate as if pipe between them. The pipe remains in place until one of two processes closes it. Sending process doesn't need to attach a

destination IP / port number to the bytes in sending attempt as the logical pipe has been established

- In UDP, programmers need to form UDP datagram packets explicitly and attach destination IP address / port number to every packet.

# 3. Transport

## Transport Layer Services

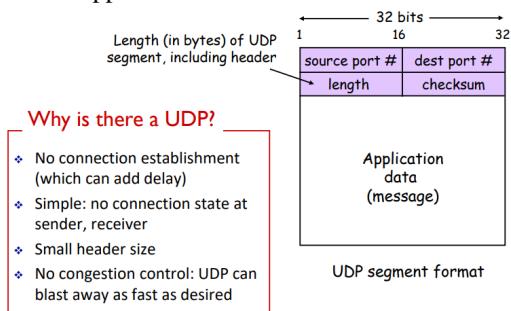
- Transport layer protocols run in hosts.
- Sender side: breaks app message into segments (as needed), passes them to network layer (aka IP layer).
- Receiver side: reassembles segments into message, passes it to app layer.
- Packet switches (routers) in between: only check destination IP address to decide routing. running on different hosts

## Transport and Network Layer

- **Transport** layer takes care of logical communication between **processes**.
- **Network** layer takes care of logical communication between **hosts**. (best-effort, unreliable)
- **IP Datagram**: Contains source and dest IP addresses, carries one transport layer segment that contains source and dest port numbers.

## UDP: Connectionless Transport

- UDP adds very little service on top of IP.
- **Multiplexing at sender**: UDP gathers data, forms packets, passes to IP.
- **De-multiplexing at receiver**: UDP receives packets from lower layer, checks dest port, and dispatches them to right processes.
- **Unreliable**: UDP transmission used by loss tolerant and rate sensitive apps.



## UDP Checksum

- Allows for error detection, but not correction. May be bit errors when segments are stored and passed in router

memory.

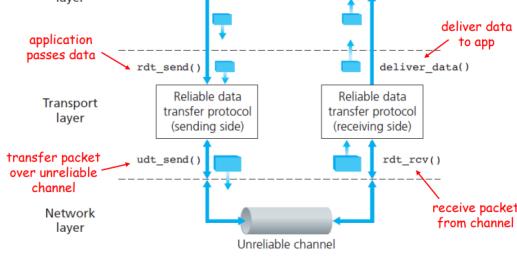
- Treat UDP segment as a sequence of 16-bit integers.
- Apply binary addition on every 16-bit integer (checksum field currently 0).
- If carry from MSB, add 1 to result (wrap).
- Compute 1's complement to get UDP checksum.

## Principles of Reliable Data Transfer (rdt)

We need to build a reliable transport layer protocol on top of unreliable communication.

- Factors: **Packet corruption, Packet loss, Packet reordering, Packet (Long) Delay.**
- Finite State Machines to describe protocol.

## Reliable Data Transfer: Service Model



## rdt 1.0 (Perfectly Reliable)

- Assumption: Underlying channel perfectly reliable.
- Sender creates packet and sends, Receiver extracts and deliver data to application.

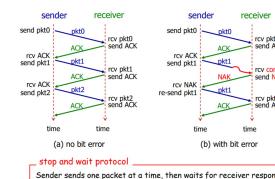
## rdt 2.0 (Corruptable Data)

- Assumption: Underlying channel may flip bits. Use **stop and wait** (for receiver response) protocol.
- Receiver uses checksum to detect bit errors, sends NAK if corrupted. Sender resends if NAK received.
- **Problem: If ACK or NACK corrupted**, no guaranteed way to recover. If packet resent, the receiver will not know it's a duplicate.

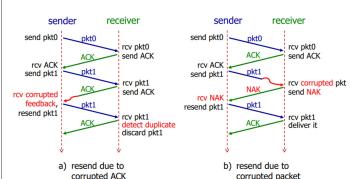
## rdt 2.1

- Add **sequence number to packet**, alternate 1 & 0. Sequence number detects duplicates.
- Same as rdt2.0, but receiver knows if it is duplicate.

## rdt 2.0 In Action



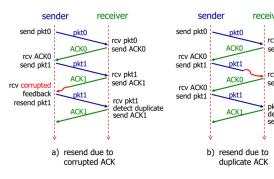
## rdt 2.1 In Action



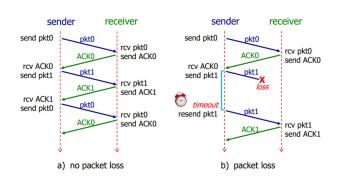
## rdt 2.2

- Use **ACK of last packet sequence number for NAK**.
- Receiver explicitly include seq. no, duplicate ACKs results in retransmit current pkt.

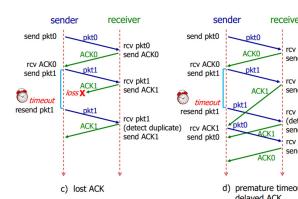
## rdt 2.2 In Action



## rdt 3.0 In Action



## rdt 3.0 In Action



## rdt 3.0 (Corruptable, Lossy, Delay)

- Assume corruption, packet loss/delay, no re-order.
- To detect packet loss, use **sender timeout**. Sender retransmits if no ACK received till timeout.
- If packet/ACK just delayed, retransmission may generate duplicates but receiver can use seq. no. to detect.
- **rdt 3.0 performance**: Utilisation rate of sender low. For RTT 30ms,  $L = 8000b$ , link 1Gbps,  $d_{trans} = 0.008ms$ , send 8000 bits per 30.008ms. Utilisation 0.027%.
- Stop and Wait limits use of physical resources.

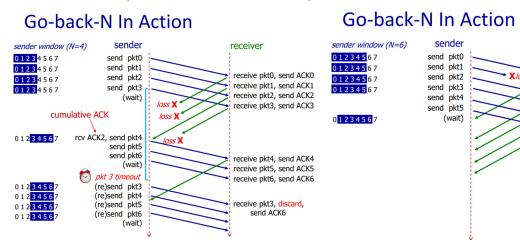
## Pipelining

- **Pipelined Protocols**: sender allows multiple, “in-flight”, yet to-be-acknowledged packets.
  - range of sequence numbers must be increased
  - buffering at sender and/or receiver
- Number of packets sent at once is called **window size**.

- Benchmarked Pipelined Protocols:** Go-Back-N (GBN), Selective Repeat (SR).
- Assumption of corruption, packet loss / delay.

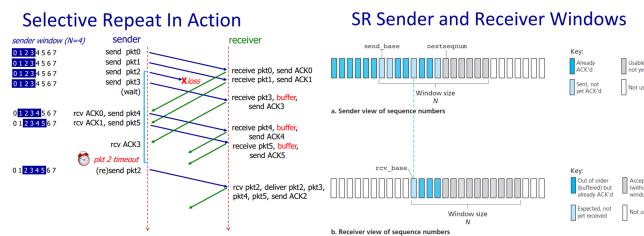
## Go-back-N

- Sliding window, slides forward only when ACK received for the leftmost packet in window.
- Requires  $k$  bits in packet header for  $2^k$  sequence number.
- Sender keeps only 1 timer for oldest unACKed packet.
- Receiver only accepts ACK packets that arrive in order, discards out-of-order packets. ACK last in-order sequence number. (cumulative ACK).



## Selective Repeat

- Receiver individually acknowledges all correctly received packets.
- Buffers out-of-order packets, for eventual in-order delivery to upper layer.
- Sender maintains timer for each unACKed packet. When timer expires, retransmit only unACKed packet.



## TCP: Connection-oriented Transport

- Connection oriented:** handshake before sending data.
- Point to point:** one sender, one receiver. The connection is duplex (bidirectional data flow). **Reliable in-order.**
- TCP socket is fully identified by four-tuple: (source IP addr, source port no., dest IP addr, dest port no.).
- Multiplexing:** TCP gathers data from processes, form transport-layer segments including app data and 4-tuple

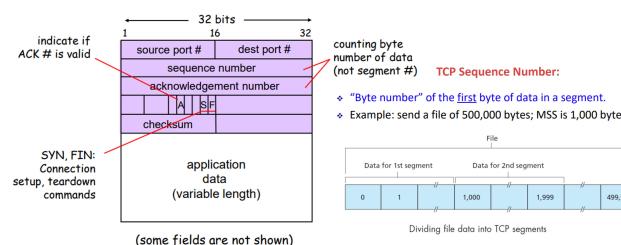
pass to network layer.

- Demultiplexing:** Connection socket created, server already noted 4-tuple. Subsequent packets directed, or demultiplexed, to the appropriate socket using those 4 values.
- TCP creates **buffers** after handshaking.

## TCP Segment / Header

- The maximum segment size (MSS) depends on maximum transmission unit (MTU).
- Generally MSS is 1460 bytes, (MTU is 1500 bytes for Ethernet and PPP link-layer protocols.) 40 bytes split half for TCP and IP header.
- TCP Seq. no is "byte no.", first b of data in segment.
- TCP Ack. no is "seq no." of next b expected by receiver.
- Checksum computation uses 1s complement (UDP same)

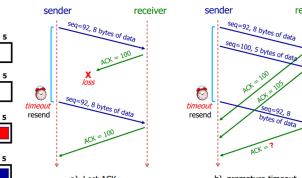
### TCP Header



### TCP ACK Generation [RFC 2581]

Event at TCP receiver	TCP receiver action
Arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	Delayed ACK: wait up to 200ms for next segment. If no next segment, send ACK
Arrival of in-order segment with expected seq #. One other segment has ACK pending	Immediately send single cumulative ACK, ACKing both in-order segments
Arrival of out-of-order segment higher-than-expected seq. (gap detected)	Immediately send duplicate ACK, indicating seq. # of next expected byte
Arrival of segment that partially or completely fills gap	Immediately send ACK, provided that segment starts at lower end of gap

### TCP Timeout / Retransmission



- Random initial sequence number: Minimise probability of some segment from previous connection mistaken as from current connection

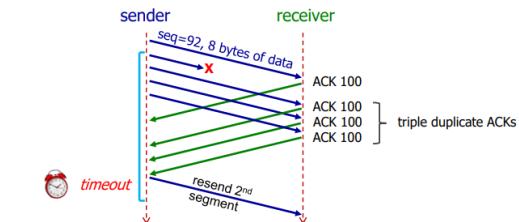
## TCP Timeout Value

- Determining TCP appropriate timeout value:
- too short timeout: premature timeout and unnecessary retransmissions.
- too long timeout: slow reaction to segment loss. Timeout interval must be longer than RTT – but RTT varies!

- TCP computes (and keeps updating) timeout interval based on estimated RTT. (TimeoutInterval = EstimatedRTT + 4\*DevRTT)

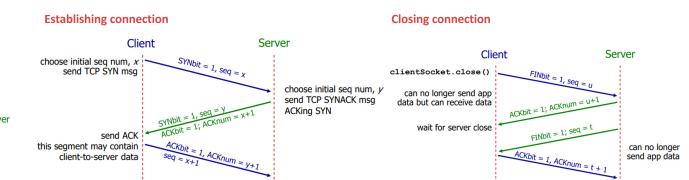
## TCP Fast Retransmission

- Timeout period is often relatively long. long delay before resending lost packet.
- Fast retransmission:** If sender receives 4 ACKs for same segment, suppose segment is lost, resend segment (even before timer expires).



## TCP Handshake / Closing

- Before exchanging app data, TCP sender and receiver "shake hands", agree on connection and exchange connection parameters.
- Closing: Client, server each close their side of connection, send TCP segment with FIN bit = 1



## 4. Network

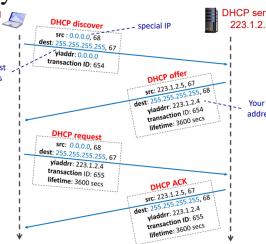
### Network Layer Services

- Network layer delivers packets to receiving hosts.
- Routers examine header fields of IP datagrams passing it.
- Forwarding:** Moving of incoming packet to appropriate output link.
- Routing:** Calculation of path taken by packets from sender to receiver.
- IP Address** used to identify host / (router), 32-bit integer expressed in binary/decimal.

- Host gets an IP address either through manual configuration by sys admin, or auto assigned by a **DHCP** server.

### Some Special IP Addresses

Special Addresses	Present Use
0.0.0.0/8	Non-routable meta-address for special use
127.0.0.0/8	Loopback address. A datagram sent to an address within this block loops back inside the host. This is ordinarily implemented using only 127.0.0.1/32.
10.0.0.0/8 172.16.0.0/12 192.168.0.0/16	Private addresses, can be used without any coordination with IANA or an Internet registry.
255.255.255.255/32	Broadcast address. All hosts on the same subnet receive a datagram with such a destination address.



## DHCP: Dynamic Host Configuration Protocol

- DHCP** allows a host to dynamically obtain its IP address from DHCP server when it joins network.
- IP address is renewable, allow reuse of addresses (only hold address while connected), support mobile users to join network.
- DHCP: 4-step process: Host broadcasts “DHCP discover” message, server responds with “DHCP offer” message, Host requests IP address: “DHCP request” message, DHCP server sends address: “DHCP ACK” message
- DHCP may provide host additional network information, e.g. IP of first-hop router, local DNS server, network mask.
- DHCP runs over **UDP**. DHCP server port 67, client port 68.

## IP Address & Network Interface

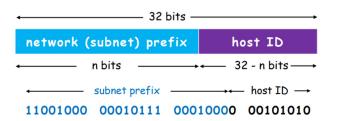
- IP address is associated with a network interface.
- Host usually has one or two network interfaces (e.g. wired Ethernet and WiFi), A router typically has multiple interfaces.
- IP Addr** comprises network/subnet prefix and host ID.

## Subnet

- Subnet** is a network formed by a group of “directly” interconnected hosts.
- Hosts in same subnet have same network prefix of IP addr, can physically reach each other without intervening router. They connect to the outside world through a router
- Classless Inter-domain Routing (CIDR)**: Internet’s IP address assignment strategy.
- Subnet prefix of IP addr of arbitrary length, Address format:  $a.b.c.d/x$ , where  $x$  is the no. of bits in subnet prefix of IP addr.

- Subnet mask** is used to determine which subnet an IP address belongs to.
- made by setting all subnet prefix bits to “1”s and host ID bits to “0”s.

- An IP address logically comprises two parts:

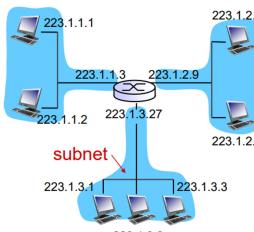


this subnet contains  $2^n$  IP addresses  
subnet prefix: 200.23.16.42/23

/23 indicates the no. of bits of subnet prefix

IP address in binary: 11001000 00010111 00010000 00101010

Subnet mask in decimal: 11111111 11111111 11111110 00000000



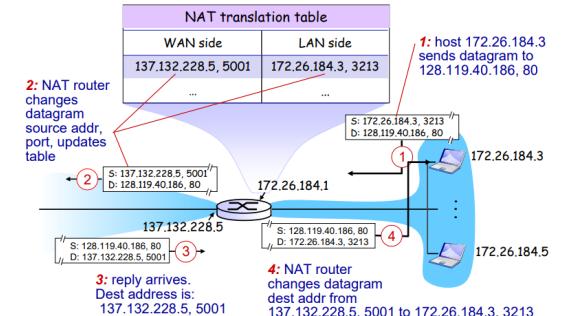
A network consisting of 3 subnets  
(first 24 bits of IP addr. are network prefix)

## Network Address Translation (NAT)

- Map IP addr space by modifying network addr info in packets IP header through traffic routing device. NAT Routers must:
- Replace** (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #), **Remember** (in NAT translation table) the mapping and **Replace** destination fields of every incoming datagram with that stored in NAT translation table.

- Benefits:** Single public IP for NAT router allows multiple private IP address. Change (private IP) addr of hosts in local network without notifying outside world.
- ISP change w/o changing local host addresses in local network.
- Hosts inside local network not explicitly addressable or visible by outside world (security plus).

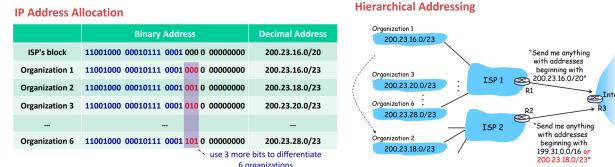
## NAT: Illustration



## IP Address Allocation

- Organization can buy from registry / rent from ISP’s addr space to obtain block of IP addr.

- Hierarchical Addressing:** Allows efficient way of routing.



- Longest Prefix Match:** Choose one with longer match. If IP addr matches all 23 bits for org, packet forwarded, else forwarded to latter.

## Routing Algorithms

- The Internet as “**network-of-networks**”, hierarchy of **Autonomous Systems** (AS), e.g., ISPs, each owns routers and links.
- Due to size and decentralized administration of Internet, routing is done hierarchically.
- Intra-AS routing:** Finds good path btwn two routers within AS. Commonly used protocols: RIP, OSPF
- Inter-AS routing (not covered)** Handles the interfaces between ASs, standard protocol: BGP.

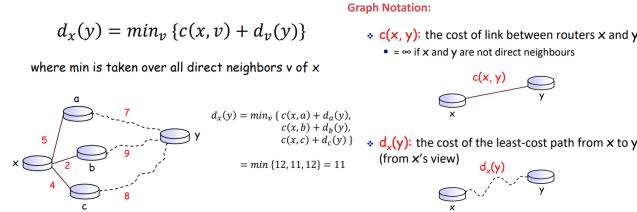
## Intra-AS Routing

- Abstractly view a network of routers as a graph, vertices are routers, edges are physical links between routers.

- Associate cost to each link. (cost = 1, or inversely related to bandwidth, or related to congestion)
- Routing:** find least cost path btwn two vertices in graph.
- Link state Algorithmnns:** Centralised routing algo, all routers have complete knowledge of network topology and link costs. Routers periodically broadcast link costs to each other.
- Use Dijkstra algorithm compute least cost path locally! (using global map).
- Distance vector Algorithms:** Decentralised routing algo, Routers know physically-connected neighbors and link costs to neighbors.
- Routers exchange “local views” with neighbors, update own “local views”. Iterative computation: Swap local view with direct neighbours, Update own view, Repeat till no further change.

## Distance Vector Algo (Bellman Ford)

### Bellman-Ford Equation



- $d_x(y) = \min_v \{c(x, v) + d_v(y)\}$
- To find **least cost path**, x needs to know cost from each of its direct neighbour to y. Each neighbour v sends its distance vector (y, k) to x, telling x that the cost from v to y is k.
- Every router, x, y, z, sends its distance vectors to directly connected neighbors. When x finds y is advertising cheaper path to z than known, x update distance vector to z accordingly and note down all packets for z should be sent to y. Info used to create forwarding table of x.
- After every router exchanged several rounds of updates with direct neighbors, all routers will know least-cost paths to all other routers.

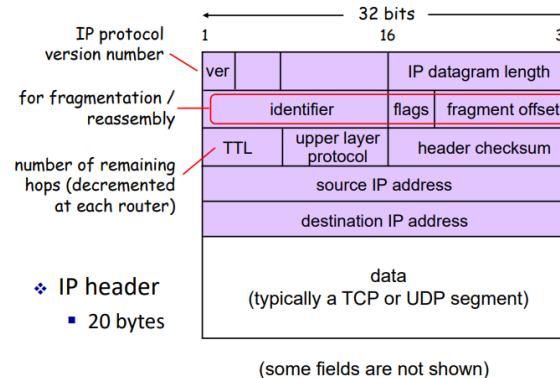
## RIP (Routing Information Protocol)

- RIP implements the DV algorithm. Uses hop count as the cost metric (insensitive to network congestion).

- Exchange routing table every 30 seconds over UDP port 520.
- “Self-repair”: if no update from a neighbour router for 3 minutes, assume neighbour failed.

## Internet Protocol (IP): IPv4

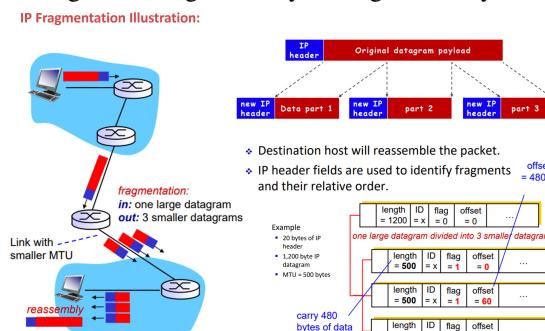
### IPv4 Datagram Format



- Length: of IP datagram + 20b header. Identifier, flags, fragment offset: support fragmentation & reassembly
- TTL: Prevent infinite circulation. Upper layer protocol: Only used at final dest, determine if UDP/TCP (for Internet). Checksum also uses 1s complement.
- IPv6:** 40b header with 128b IP addr.

### IP Fragmentation & Reassembly

- Different links, different MTU (Max Transfer Unit, max amt of data link-level frame can carry).
- “Too large” IP datagrams may be fragmented by routers.



- Flag(frag flag)** is set to 1 if next fragment from same segment, 0 if this is the last fragment.

- Offset** is expressed in unit of 8-bytes

## Internet Control Message Protocol

- ICMP:** used by hosts & routers to communicate network-level information.
- Error reporting:** unreachable host / network / port / protocol.
- Echo request/reply (used by ping).
- ICMP messages carried in IP datagrams, ICMP header starts after IP header.

### ICMP Type and Code

- ICMP header: Type + Code + Checksum + others.

Type	Code	Description
8	0	echo request (ping)
0	0	echo reply (ping)
3	1	dest host unreachable
3	3	dest port unreachable
11	0	TTL expired
12	0	bad IP header

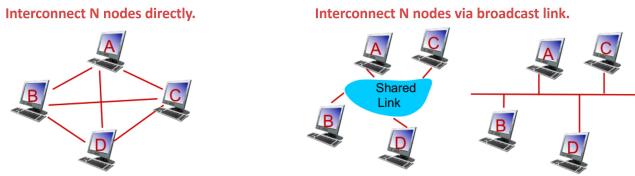
Selected ICMP Type and subtype (Code)

## 5. Link

- **Link Layer:** Concerned about electronics of sending and receiving binary data over a communication channel.
- **Communication Channel:** transmission medium of data signals (e.g. copper wire, satellite, optical fiber)
- **Node:** Devices exchanging data. (E.g. hosts, routers).
- **Link:** Comm. channels that connect adjacent nodes.

### Connecting N nodes via cable

- **Aim:** Send data between N nodes via cable.
- Interconnect N nodes directly: Each link needs to be addressed, N-1 Links needed, does not scale.
- Interconnect via broadcast link: Each link needs to be addressed, need to define protocol, need to handle errors.
- Protocol: Framing, Link Access Control. Errors: Detection, Reliability.



### Link Layer

- **Link layer** sends datagram between adjacent nodes (hosts or routers) over a single link. Responsible for transfer of datagram from one node to physically adjacent node over link.
- **Frames (layer 2 packet):** IP datagrams are encapsulated in link-layer **frames** for transmission.
- **Protocols:** Different link-layer protocols may be used on different links, each protocol may provide a different set of services.

### Link Layer Services

- **Framing:** Encapsulate datagram into frame, adding header and trailer.
- **Link access control:** When multiple nodes share a single link, need to coordinate which nodes can send frames at a certain point of time.
- **Error detection:** Errors are usually caused by signal attenuation or noise. Receiver detects presence of errors, and may signal sender for retransmission or simply drops frame.

- **Error correction:** Receiver identifies and corrects bit error(s) without resorting to retransmission.
- **Reliable delivery:** Seldom used on low bit-error link (e.g., fiber) but often used on error-prone links (e.g., wireless link).

### Network Adapter

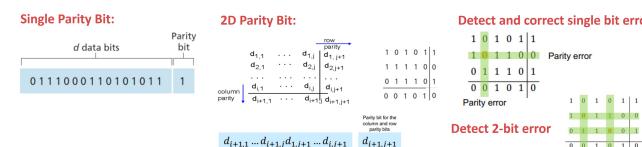
- **network adapter**, aka network interface card (NIC), is a single, special-purpose chip that implements the link-layer services above. (e.g. Ethernet card, Wifi Adapter).
- Semi-autonomous, implementing both link & physical layers. Many services implemented in hardware.

### Error Detection and Correction Techniques

- **EDC:** Error Detection and Correction Bits.
- **D:** Data protected by error checking, may include header.
- Larger EDC fields added to link layer frame yields better detection (and correction), but larger overhead.
- **Common error detection schemes:** Checksum (used in TCP/UDP/IP), Parity checking, CRC (link layer).
- **Checksum review:** treat segment contents as 16 bit int sequences, get 1s complement of sum of segment contents.

### Parity Checking

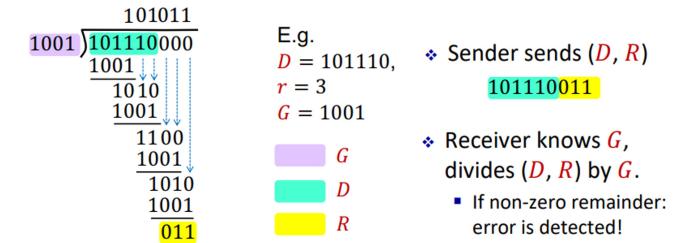
- **Single Bit Parity:** 1 parity bit for the data.
  - **Even Parity Scheme:** Choose the bit to make the total number of 1s even.
  - **Odd Parity Scheme:** Same but we make the number odd.
- Single bit parity can detect odd number of single bit errors, but cannot detect even single bit errors.
- Errors often clustered together in “bursts”, probability of undetected errors in a frame can approach 50%.
- **Two-Dimensional Parity:** Divide the data into rows and columns, and repeat the above but have parity bits for each column and each row.
  - Can detect and correct single bit errors in data.
  - Can detect two-bit errors.



### Cyclic Redundancy Check

- “Long division” but with division being replaced by a bitwise XOR operation.
- Generally done by hardware, so very fast.
- **D: d-bit data**, which is also the dividend.
- **G: Generator of r + 1 bits**, which is also the divisor.
- **R: r-bit CRC**, which is also the remainder.
- The resultant  $d + r$  bits is “divisible” by G, so the receiver can check for a zero remainder.
- Can detect all odd number of single bit errors.
- CRC of r bits can detect all burst errors of less than  $r + 1$  bits, burst errors greater than  $r$  bits with probability  $1 - 0.5^r$ . Aka polynomial code.

#### Cyclic Redundancy Check (CRC)



### Multiple Access Links and Protocols

- **Multiple Access Protocols:** Categorisable into three broad classes: Random Access, “Taking Turns”, Channel Partitioning.
- **Ideal MAP:** Collision free, Efficient, Fairness, Fully Decentralized.
- Additionally, coordination about channel sharing must use channel itself (no out-of-channel signalling).

### Types of Network Links (2)

- **Point-to-point link:** Sender and receiver connected by a dedicated link. No need for multiple access control.
- **Broadcast link:** Multiple nodes connected to same shared broadcast channel. When any one node transmits a frame, all other nodes in the channel receives a copy. We need a Multiple Access Protocol to prevent frame collisions.

## Multiple Access Protocols

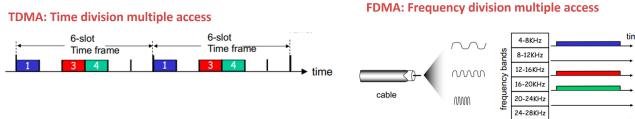
### Channel Partitioning Protocols

#### TDMA: Time Division Multiple Access

- Each node gets fixed length time slot (time frame), where length = frame transmission time. This repeats in rounds. Unused slots go idle.

#### FDMA: Frequency Division Multiple Access

- Channel spectrum is divided into frequency bands, and each node is assigned one band. Bandwidth has thus decreased, thus transmission is slower. Unused transmission time in frequency bands go idle.
- Both TDMA and FDMA: Collision Free, Inefficient, Perfectly Fair and fully Decentralized.



### Taking Turns Protocols

#### Polling

- A master node invites each of the other nodes (slaves) to transmit in turns. Minor polling overhead. A single point of failure, which is the master node.

#### Token Passing (Token Ring) / Round Robin

- Control token is passed from one node to the next sequentially. There is overhead for the token and single point of failure as well, which is the token.
- Even if only a few of the nodes have data to send, it can still be quite efficient.

### Random Access Protocols

Generally these protocols specify how to detect and recover from collisions (When two or more transmitting nodes). Thus, no need for centralised coordination, thus no single point of failure.

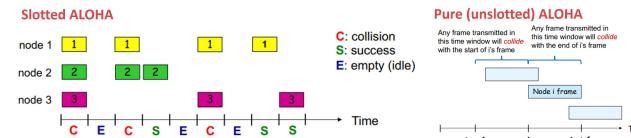
### Slotted ALOHA

- Assume all frames are of the same size, and split the time into slots of equal length, where length = time to transmit 1 frame =  $L(\text{bits})/R(\text{rate})$ .
- A node will only transmit at the start of a slot.

- Each node listens to the channel while transmitting. If a collision occurs, it retransmits in each subsequent slot with probability  $p$  until success.
- $p$  depends on network congestion
- Effectiveness:** Not collision free, Efficiency high when only one node is active, but maximum efficiency falls to 37% when many active nodes (collision & empty slots). Perfectly fair and decentralized.

### Pure (Unslotted) ALOHA

- Like ALOHA but no slots nor synchronisation, just transmit when there's a fresh frame.
- Chance of collision increases, as now it can collide with frames both in front and behind ( $t_0 - 1, t_0 + 1$ ).
- Effectiveness:** Not collision free, Efficiency high when only one node is active, but maximum efficiency falls to 18% when many active nodes (collision & empty slots). Perfectly fair and decentralized.



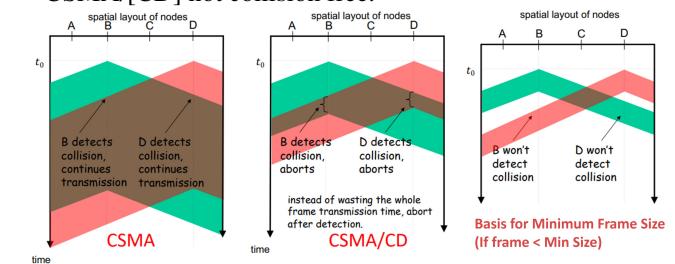
### Carrier Sense Multiple Access (CSMA)

- Sense if the channel is idle, if so, transmit frame.
- Still collides when two nodes sense that a channel is idle.
- Collision: propagation delay, nodes far apart, two nodes may not hear each other's transmission immediately.
- CSMA transm. does not stop despite collision detection.

### Carrier Sense MA/ Collision Detection (CSMA/CD)

- Backoff Algorithm:** Same as CSMA except the moment a collision is detected, the node stops transmission.
- The node then retransmits after some random amount of time. (probability  $p$  for each sub. frame until success).
- Binary Exponential Backoff:** More collisions implies heavier load, adapt retransmission attempt to estimated current load. Longer back-off interval with more collisions.
- After the  $m^{\text{th}}$  collision, we choose K at random from  $0, 1, \dots, 2m - 1$ , then wait  $K * 512$  bit transmission times before retransmitting.

- Effective:** Both Efficient, Fair, Decentralized, but both CSMA/[CD] not collision free.



### Minimum Frame Size

- For CSMA and CSMA/CD above, need minimum frame size so that collisions can always be detected.
- Ethernet has a minimum size of 64 bytes.

### CSMA / Collision Avoidance (CSMA/CA)

- Collision detection can be hard for wireless LANs, as energy levels drop too quickly.
- Hidden Node Problem: When two nodes cannot detect each other but a node in-between encounters a collision.
- As such, an ACK is required from the receiver as well.

### Bit Times

- Common unit used in multiple access protocols is bit times.
- Bit transmission time, is time taken to transmit 1 bit.
- Often, used as such: propagation delay is equals to 800 bit times. This means the propagation delay is equals to the time it takes to transmit 800 bits onto the link.

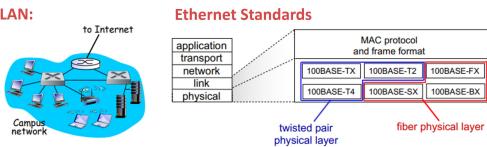
# Switched Local Area Networks

## MAC Address

- Every adapter (NIC) has a unique MAC address (aka physical or LAN address).
- **MAC:** Media Access Control.
- Used to send and receive link layer frames, when adapter receives a frame, it checks if the destination MAC address of the frame matches its own MAC address.
- If yes, adapter extracts enclosed datagram and passes it to the protocol stack.
- If no, adapter simply discards the frame without interrupting host.
  - **48 bits:** Burned in NIC ROM (read-only memory), e.g. 5C-F9-DD-E8-E3-D2
  - **IEEE:** Administers the MAC address allocation. The first three bytes of the MAC identifies the vendor of the adapter.
  - If somehow MAC is manually configured to be not unique, and the NICs are on the same subnet, transmission will be severely affected.

## Local Area Network (LAN)

- **LAN** is a computer network that interconnects computers within a geographical area such as office building or university campus.
- Example LAN technologies: IBM Token Ring: IEEE 802.5 standard, Ethernet: IEEE 802.3 standard, Wi-Fi: IEEE 802.11 standard



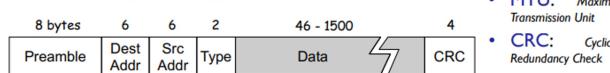
## Ethernet

- **Ethernet:** “Dominant” wired LAN technology, developed in mid 1970s, Standardized by Xerox, DEC, and Intel in 1978.
- Simpler and cheaper than token ring and ATM (asynch transf mode).
- MAC protocol, frame format remain unchanged over years.

## Ethernet Frame Structure

- Sending NIC (adapter) encapsulates IP datagram in Ethernet frame.

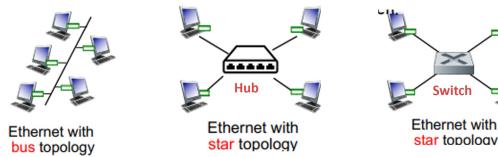
Ethernet Frame Structure



- **Preamble:** 7 bytes with pattern 10101010 ( $AA_{Hex}$ ), Followed by 1 byte with pattern 10101011 ( $AB_{Hex}$ ). Also called “start of frame”. Allows sender and receiver to synchronise clock rates, as alternating bits form a square wave. Lets receiver know how long 1 bit is.
- **Source and dest MAC address:** If NIC receives frame with matching destination address or with broadcast address, passes data in frame to network layer protocol, otherwise, NIC discards frame.
- **Data / Payload:** The maximum size is 1500 bytes, which is link MTU mentioned in IP fragmentation. Minimum size is 46 bytes, to ensure that collision will always be detected.
- **CRC:** For corruption detection.
- **Type:** Higher level protocol used, usually IP. (others e.g. ARP, AppleTalk), permits Ethernet to multiplex network-layer protocols.

## Ethernet Topology

- **Bus Topology:** (broadcast LAN) All nodes are connected and can collide with each other.
- **Star Topology:** Switch / Hub in the centre and nodes are connected to that switch. Do not collide with each other.



## Ethernet Delivery

- **Connectionless:** No handshaking btwn. sender & receiver.
- **Unreliable:** NIC does not send ACK/NAK. Data in dropped frames recovered only if initial sender uses higher layer rdt (e.g. TCP).
- **Ethernet's multiple access protocol:** CSMA/CD with binary (exponential) backoff.

## Ethernet CSMA/CD Algorithm

1. NIC receives datagram from network layer, creates frame.
2. If NIC senses that the channel is idle, start frame transmission. Else, wait until idle.
3. If NIC transmits the entire frame without detecting another transmission, NIC is done.
4. If another transmission is detected, NIC aborts and sends a **jam signal**, tells all other nodes that a collision has been detected and NIC will be retransmitting.
5. After aborting, NIC enters exponential (binary) backoff, and repeat from step 2

## Ethernet Hub

- **Hub:** Physical-layer, acting on indiv. bits. (not frames)
- When bit arrives from one interface, hub re-creates the bit, boosts energy strength, and transmits the bit onto all the other interfaces.
- **Adv:** Cheap, easy to maintain (modular design of network).
- **Disadv:** Slow, not ideal for larger networks (collisions).

## Ethernet Switch

- **Switch** is link-layer device used in LAN that also stores and forwards Ethernet frames.
- **Layer 2 device:** Unlike routers, which is a layer 3 device (i.e. it goes up to the network layer), switches only have 2 layers, i.e. up to link layer. **Switches act on frames**.
- **No IP address:** For the reason above, it has no IP address.
- **Transparent to hosts:** Hosts unaware of switch presence.
- **Collision-free:** Each host has dedicated connection to the switch, (separate collision domains). Connection has two channels, i.e. fully duplex and frames sent two-way simultaneously.
- **Store and Forward packet switch:** Switch buffers frames, e.g. if currently forwarding another frame to the outgoing link.

## Switch Forwarding Tables

- A switch has multiple interfaces, needs to know which nodes are reachable via which interface.
- Done via **switch forwarding tables**, which have entry format:

<MAC address of host, interface to reach host, TTL>

- **Self-Learning:** Whenever the switch receives a frame from host A, it will record that interface for A for future frames.
- **Broadcast:** If destination host not found in the switch forwarding table, the switch will broadcast the frame to all outgoing links.

## Routers vs. Switches

- **Routers:** Check IP address, Store-and-forward, Compute routes to destination.
- **Switches:** Check MAC address, Store-and-forward, Forward frame to outgoing link or broadcast

## Address Resolution Protocol (ARP)

How to know MAC address of receiving host, knowing its IP address? Use ARP [RFC 826].

- **ARP** provides query mechanism to learn MAC address.
- All nodes have an ARP table containing mappings of IP addresses and MAC addresses of other neighbouring nodes in the same subnet.
- **Plug & Play:** nodes create ARP tables without intervention from network admin.
- **Entry format:** (TTL typically a few minutes)  
 $\langle \text{IP address}; \text{MAC address}; \text{TTL} \rangle$

### Sending Frame in Same / Another Subnet:

1. If A and B in same subnet, A knows B's MAC address from its ARP table:
  - (a) A just creates frame with B's MAC address and send.
  - (b) Only B will process frame, all other hosts ignore.
2. If A and B in same subnet, A does not know B's address:
  - (a) A broadcasts an ARP query packet, containing B's IP address. The destination MAC address is set to FF-FF-FF-FF-FF-FF.
  - (b) All other nodes in the same subnet will receive this ARP query packet, but only B will reply it.
  - (c) A caches B's IP-to-MAC address mapping in its ARP table (until TTL expires).
3. If A and B in different subnets (assume router R directly connecting the two subnets, and A and B):
  - (a) A will need to send a frame with R's MAC address but B's IP address as destination.
  - (b) R will realise it needs to forward this frame as the IP doesn't match when MAC matches.
  - (c) R will forward the datagram to an outgoing link and construct a new frame with B's MAC address.

## IP Addresses vs. MAC Addresses

- **IP address:** 32 bits in length, network-layer address used to move datagrams from source to dest.
- Dynamically assigned, hierarchical (to facilitate routing), Analogy: postal address.
- **MAC address:** 48 bits in length, link-layer address used to move frames over every single link.
- Permanent, to identify the hardware (adapter), Analogy: NRIC number.
- **ARP** resolves mapping from network layer (IP) address to link layer (MAC) address.