

# Computer Networks

Sungchan Yi

February 2020

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	What is the Internet? . . . . .	3
1.2	Network Edge . . . . .	3
1.3	Network Core . . . . .	3
1.4	Internet Structure . . . . .	4
1.5	Performance Metrics . . . . .	5
1.6	Protocol Stack . . . . .	6
1.7	Network Security . . . . .	6
1.8	History of the Internet . . . . .	7
<b>2</b>	<b>Application Layer</b>	<b>8</b>
2.1	Principles of Application . . . . .	8
2.2	Web and HTTP . . . . .	10
2.3	Cookies and Web Caching . . . . .	12
2.4	SSL/TLS . . . . .	12
2.5	Electronic Mail . . . . .	13
2.6	Domain Name System . . . . .	14
2.7	Peer-to-Peer Application . . . . .	17
2.8	Video Streaming and CDNs . . . . .	17
<b>3</b>	<b>Transport Layer</b>	<b>19</b>
3.1	Transport Layer Services . . . . .	19
3.2	Multiplexing and Socket . . . . .	20
3.3	User Datagram Protocol . . . . .	20
3.4	Reliable Data Transfer Principles . . . . .	21
3.5	Transmission Control Protocol . . . . .	25
3.6	Congestion Control . . . . .	28
3.7	TCP Congestion Control Algorithm . . . . .	30
3.8	TCP vs UDP . . . . .	31

<b>4</b>	<b>Network Layer I</b>	<b>33</b>
4.1	Overview of Network Layer . . . . .	33
4.2	Inside of Router . . . . .	33
4.3	Internet Protocol Overview . . . . .	36
4.4	IP Addressing . . . . .	37
4.5	Datagram Forwarding . . . . .	38
4.6	Dynamic Host Configuration Protocol . . . . .	39
4.7	Network Address Translation . . . . .	40
4.8	IPv6 . . . . .	41
<b>5</b>	<b>Network Layer II</b>	<b>43</b>
5.1	Introduction to Routing . . . . .	43
5.2	Link State Routing . . . . .	44
5.3	Distance Vector Routing . . . . .	45
5.4	Intra-AS Routing: OSPF . . . . .	46
5.5	Inter-AS Routing: BGP . . . . .	47
5.6	Software Defined Networking . . . . .	48
5.7	OpenFlow . . . . .	48
5.8	Internet Control Message Protocol . . . . .	48

# 1 Introduction

## 1.1 What is the Internet?

- **Internet** = Inter- + net(work)
- Network of networks
- Various types of networks: mobile network, home network, global ISP, regional ISP

### 1.1.1 Components

- HW components
  - End hosts
  - Links: copper, fiber, radio, satellite
  - Interconnection devices: router, switch, repeater
- SW components
  - Operating software
  - Application programs
  - **Protocols**

### 1.1.2 Protocols

- A defined set of rules and regulations that determine how data is transmitted in telecommunications and computer networking
- All communication activity in Internet is governed by protocols
- Protocols define
  - Message format
  - Order of messages sent and received among network entities
  - Actions taken on message transmission and receipt

## 1.2 Network Edge

Network edge contain hosts, which are clients and servers.

## 1.3 Network Core

Network core is a mesh of interconnected routers or switches. They forward packets from one router (or switch) to the next along the path from the source to destination

### 1.3.1 Switching Mechanisms

- **Circuit Switching**

- End to end resources reserved for communication between source and destination
- Entire data flow through along the path like water
- Resources are dedicated to each connection - the circuit segment is idle if it is not being used
- Common in traditional telephone networks
- Channel allocation methods: frequency division multiplexing (FDM), time division multiplexing (TDM)

- **Packet Switching**

- Entire data is broken into small packets
- Each packet has its destination address
- Each packet is handled independently
- Packet is transmitted at full link capacity
- Takes  $L/R$  seconds to transmit  $L$ -bit packet into link at  $R$  bps
- *Store and forward*: entire packet must arrive at the router before it can be forwarded to the next
- End to end delay  $\approx 2L/R$
- If arrival rate of packets exceed the transmission rate of link, packets can be dropped (lost) if the packet queue inside the router is full

- Comparison

- Packet switching allows more users to use the network
- Circuit switching guarantees the quality of service for each connection

## 1.4 Internet Structure

Nobody or everybody is in charge of the Internet.

- End systems connect to the Internet via **access ISPs** (Internet Service Providers)
- Access ISPs in turn must be interconnected so that any two hosts can communicate with each other
- Resulting network of networks is very complex

### 1.4.1 Connecting Access ISPs

- How do we *interconnect* the access ISPs?
- Connecting each access ISP to every other one would not be scalable, since we need  $\mathcal{O}(n^2)$  connections
- Better: Keep a *global transit ISP* and connect each access ISP to it
- Competing global transit ISPs appear and they are also interconnected by peering link and Internet exchange point (IXP)
- Regional networks arise to connect access networks to global ISPs

## 1.5 Performance Metrics

### 1.5.1 Evaluation Metrics

- **Delay:** Packet delivering time from source to destination
- **Packet loss:** Ratio of lost packets to total sent packets<sup>1</sup>
  - If the queue is full, the arriving packets will be dropped
  - Lost packets may be re-transmitted by previous nodes, by source end system, or not at all
- **Throughput:** Amount of traffic delivered / unit time
  - Rate at which bits are transferred between source and destination
  - Can be measured instantaneously, or on average
  - *Bottleneck link:* The link on end-end path that constrains the throughput (usually the one with minimum capacity)

### 1.5.2 Four sources of delay

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

- Queueing Delay
  - Time waiting at output buffer for transmission
  - **Congestion** dependent
- Transmission Delay
  - $d_{trans} = L/R$  where  $L$ : packet length (bits),  $R$ : link bandwidth (bps)
- Processing Delay
  - Bit error checking
  - Decision of output link
  - Typically takes less than a few milliseconds (hardware acceleration)
- Propagation Delay
  - $d_{prop} = d/s$  where  $d$ : length of physical link,  $s$ : signal speed

Queueing and transmission delay take up most of the delay.

### 1.5.3 More on Queueing Delay

$$(\text{Traffic Intensity}) = \frac{La}{R}$$

where  $R$  is the link bandwidth (bps) or transmission rate,  $L$  is the average packet length (bits),  $a$  is the average packet arrival rate. As the traffic intensity  $\rightarrow 1$ , the average queueing delay will grow without bound.

---

<sup>1</sup>PDR: Packet delivery ratio

## 1.6 Protocol Stack

A communication protocol stack is composed of several **layers**. Each layer implements a service via its own internal actions, and by relying on services provided by the underlying layers.

Layering or **modularization** eases development, maintenance, and updating of the whole system. But this can be harmful in cases when a higher level layer needs information from the lower layers.<sup>2</sup>

### 1.6.1 Internet Protocol Stack

1. **Physical**: Bits on the wire
2. **Link**: Data transfer between neighboring network elements
3. **Network**: Routing of datagrams from source to destination
4. **Transport**: Process to process data transfer
5. **Session**: Synchronization, connection management, recovery of data exchange
6. **Presentation**: Allows applications to interpret the meaning of data
7. **Application**: Supporting network applications

## 1.7 Network Security

The field of network security arises from these questions:

- How can bad guys attack our computer networks?
- How do we defend our networks against those attacks?
- How do we design architectures that are immune to attacks?

### 1.7.1 Forms of Attacks

- Malware
  - virus: A self-replicating infection by receiving or executing an object
  - worm: A self-replicating infection by passively receiving object that gets itself executed
  - spyware
  - ransomware
- Packet Sniffing: Promiscuous network interface reads/records all packets passing by
- Denial of Service: Attackers make resources unavailable by sending a huge amount of fake traffic
- Fake Addresses: Send packets with fake source address
- Fake Wi-Fi AP: Steal user's credentials using fake AP

---

<sup>2</sup>Consider a navigation system, which uses "physical" information like actual traffic, when choosing the fastest route between two places. But this "physical" information wouldn't normally be visible to other layers.

## 1.8 History of the Internet

- Firstly developed as ARPAnet
- Internetworking architecture = autonomy + minimalism
- TCP/IP, WWW

## 2 Application Layer

Previously, we have seen that there are 5 layers in the Internet protocol. We will go through each layer, in a top-down approach. We will cover the application layer in this section.

### 2.1 Principles of Application

#### 2.1.1 Network Applications

- Types: email, web (server software, browser), P2P file sharing, SNS, messenger program, online games, streaming stored video (YouTube, Netflix)
- Run on (different) **end systems**: network core devices *do not* run user applications. Ex) Routers only transfer information
- Communication over network (between end systems)

#### 2.1.2 Application Architectures

There are two kinds of application structures: *client-server* model, and *peer-to-peer* model

- **Client-Server Model**

- The **server** is *always on*, has permanent IP address, since clients must be able to access the server anytime
- For scaling (to support a huge number of clients), a data center is typically built
- The **client** is a program that communicates with the server. It may be intermittently connected, and may have dynamic IP addresses. For communication, the client must send a request to the server first, using its IP address. Then the server will respond to that IP address.
- Clients do not communicate directly with each other

- **Peer-to-Peer Model (P2P)**

- There is no *always on* server. Each clients can function both as a client and a server. Arbitrary end systems will directly communicate with each other.
- Peers are intermittently connected and they can change IP addresses, so it is complex to manage.
- **Self scalability**: New peers bring new service capacity, as well as new service demands

#### 2.1.3 Application Layer Protocol

The application layer is on the top of the Internet protocol. Then, the applications on this layer will communicate with the application layer on another end device. The protocol defined for this communication is called the **application layer protocol**. There are two types of messages that network application protocols exchange - **requests** and **responses**.

- Message **Syntax**: What kinds of fields are there in the messages? How are the fields delineated?
- Message **Semantics**: What is the meaning of the information in the fields?
- Message **Pragmatics**: When and how do we process (send/respond) the messages?



### 2.1.4 Application Protocol

- Open Protocols
  - Standardized, open to public, for interoperability.<sup>1</sup>
  - Even if some non-authorized clients create a message that obeys the protocol, they can communicate with the server.
  - Ex) HTTP, SMTP
- Proprietary Protocols
  - A protocol specific for some program
  - Ex) Skype
  - We cannot communicate with Skype without using the Skype application

### 2.1.5 Requirements of Network Applications

Application	Data Loss	Throughput	Time Sensitive
File Transfer	×	elastic	×
Email	×	elastic	×
Web Documents	×	elastic	×
Real-time Audio/Video	loss-tolerant	Audio: 5kbps~1Mbps / Video: 10kbps~5Mbps	100ms
Stored Audio/Video	loss-tolerant	(same)	A few secs
Interactive Games	loss-tolerant	≥ Few kbps	100ms
Text Messaging	×	elastic	Yes and no

The applications that require correct communication of information have *elastic* throughput, since the correctness of the information is a lot more important than the speed of communication.

These requirements should be met by the *transport layer protocols*. The throughput and other requirements highly depend on the physical devices (routers, cables etc.) that hold up the network structure. The developers on the application layer cannot handle these requirements properly.

### 2.1.6 Transport Protocol Services

- **TCP Service** (Transmission Control Protocol)
  - **Error control**: In charge of reliable transport between sending and receiving processes
  - **Flow control**: The sender won't overwhelm the receiver (not too much data)
  - **Congestion control**: Throttle sender when network is overloaded
  - *Does not provide*: Timing, minimum throughput guarantee, security
  - *Connection oriented*: Setup is required between client and server processes
- **UDP Service** (User Datagram Protocol)

---

<sup>1</sup>**Interoperability** is a characteristic of a product or system, whose interfaces are completely understood, to work with other products or systems, at present or in the future, in either implementation or access, without any restrictions.

- Unreliable data transfer between sending and receiving processes
- Why do we use UDP? - Some programs may require UDP. For example, the error controlling in TCP lowers the throughput, but UDP will ignore this error, resulting in faster communication, which is suitable for multimedia programs

## 2.2 Web and HTTP

**Web pages** consist of **objects**. They can be an HTML file, JPEG image, Java applet, audio file, and more. Web page is described by **HTML-file(s)** which include several referenced objects. Each object is addressable by a **URL**(Uniform Resource Locator).

$$\underbrace{\text{www.someschool.edu}}_{\text{host name}} / \underbrace{\text{someDept/pic.gif}}_{\text{path name}}$$

### 2.2.1 HTTP Overview

- **HTTP** (HyperText Transfer Protocol)
  - Web's application layer protocol
  - **Hyperlink**: A reference to data the reader can directly follow by clicking
  - Uses client-server model
  - Client uses a browser that requests, receives, and displays the Web objects
  - The server is a web server that sends objects in response to request from clients
- Based on **TCP**
  1. Client initiates TCP connection (socket connection) to server
  2. Server accepts TCP connection from client
  3. HTTP messages are exchanged between browser and Web server
  4. TCP connection is closed
- *HTTP response time* (**RTT**: Round trip time)
  - 1 RTT to initiate TCP connection
  - 1 RTT for HTTP request and first few bytes of HTTP response to return
  - File transmission time
  - Total = 2 RTT + file transmission time

### 2.2.2 HTTP Version

- **HTTP/1** (1996)
  - *Non-persistent* HTTP: Single object per single TCP connection
  - Long latency
- **HTTP/1.1** (1999)
  - *Persistent* HTTP: Multiple objects over on TCP connection

- Decreased latency
- *Synchronous* order of response/request pairs over one TCP connection
- **HTTP/2 (2015)**
  - *Persistent* HTTP
  - *Asynchronous* (parallel) multiple response/request pairs over one TCP connection

### 2.2.3 HTTP Message

There are two types of messages: **requests** and **responses**.

- **Request Message:** Request line + Header lines + Body
- **Response Message:** Response line + Header lines + Body
- Request line: method, URL, version
- Response line: version, status code, status text
- Header lines: header field name, value
- Body: entity body

### 2.2.4 REST

- **REpresentational State Transfer**
  - HTTP should be **stateless**. If the state of client is kept in the server, it causes overhead for the server.
  - The server will not store each client's states. Instead, the server will store the information about the client in the HTTP message. The server should also store the method to interpret the message. With all these information, the client knows how to fetch the data.
  - We call a service **RESTful** if the service conforms to this architectural style
- **Architectural Constraints**
  - *Client-server architecture*: Separation of the user interface concerns from the data storage concerns
  - *Statelessness*: No client context should be stored on the server between requests
  - *Cacheability*: Server responses are cachable on client and intermediaries
  - *Layered system*: One should be unable to tell whether a client is directly connected to the end server or to an intermediary along the way
  - *Uniform interface*: Simplification and decoupling of the architecture (Use of standardized languages - HTML, XML, JSON - that is not restricted by some computer architecture)
  - *Code on demand* (optional): Should be able to transfer executable code such as Java applets and JavaScript

## 2.3 Cookies and Web Caching

### 2.3.1 Cookies

**Cookies** keep the states of clients.

1. Client has a cookie file
2. A usual HTTP request message is sent to the server (for the first time)
3. The server creates an ID for the user, and creates an entry in the database
4. The server responds with the ID, and tells the client to set the cookie with the given ID
5. Now the client can send HTTP requests using that ID inside the cookie file
6. The server (database) performs a cookie-specific action (distinguished by ID), and responds as usual
7. A week later, (if the cookie still exists) the cookie can be used again for communication with the server

### 2.3.2 Web Caching

- For some servers (sites) with lots of visitors, request and responses for the exact same site would cause a huge overhead.
- The server prepares a **proxy server** that caches this information
  - If the requested object is not in the cache, the proxy server requests the object from the origin server, and caches the data (and also responds to the client)
  - Otherwise, the proxy server will use the cached data to respond to the client request
- Effects of Web Caching
  - For clients, the response time is reduced
  - For servers, it can handle more users (reduced request overheads)
  - For local ISPs, the traffic to external server is reduced, so the *access link* can be used efficiently

## 2.4 SSL/TLS

### 2.4.1 Securing TCP

- We often use TCP and UDP for transport layer protocols. But when these protocols were developed back then, security considerations were not taken into account
- TCP and UDP have no encryption, even passwords traversed the Internet in clear text
- **SSL/TLS** provides encrypted TCP connection at the *application layer*
- Assures data integrity<sup>2</sup>, and end-point authentication
- **SSL** (Secured Socket Layer)
  - SSL v2.0 and v3.0 were released in 1995 and 1996, respectively

---

<sup>2</sup>Data doesn't change during transmission

- **TLS** (Transport Layer Security)
  - Improved version of SSL v3.0
  - More secure than SSL, but slower due to the two step communication processes

## 2.4.2 SSL/TLS Principle

1. When client connects to a server, the client asks for a secure SSL session
2. The server sends a certificate<sup>3</sup>, that contains the server's public key
3. The client sends a one time encryption key for the SSL session, encrypted with the server's public key
4. The server decrypts the session key using its private key and establishes a secure session
5. Now the client and the server can safely send/receive encrypted data

HTTPS is HTTP + SSL/TLS

## 2.5 Electronic Mail

### 2.5.1 Electronic Mail

- There are 3 components that consist electronic mail
  - *User agents* (clients): edit mail, read mail
  - *Mail servers*: Google, Daum, Naver etc.
  - *Protocols*: SMTP, POP3, IMAP
- Components of mail servers
  - A mailbox for incoming messages
  - A message queue for outgoing messages
  - A protocol for exchanging mail between mail servers

### 2.5.2 SMTP Protocol

- **SMTP** (Simple Mail Transfer Protocol) [RFC 2821] uses TCP as the transport layer protocol for reliable email delivery
- Three phases of transfer
  - Handshaking (Check sending server and receiving server)
  - Transfer of messages
  - Closure
- Command/Response interaction (like HTTP)
  - Commands are in ASCII text
  - Responses contain status codes and phrase

---

<sup>3</sup>The client must check that this certificate is valid, and this certificate has to be signed by someone that the client trusts.

### 2.5.3 Mail Access Protocol

- **POP3** (Post Office Protocol 3)
  - By default, deletes messages from the server after retrieving data
  - Disconnects from the server after download
  - Needs to reconnect to the server on each download
- **IMAP** (Internet Mail Access Protocol)
  - Keeps all messages at the server and allows user to organize message folders
  - Support synchronization across multiple devices
  - Stays connected until the mail client app is closed and downloads messages on demand
- **HTTP**
  - Web-based email
  - Used between browser and server
  - Hotmail in the mid 1990s
  - Google, Yahoo, etc.

## 2.6 Domain Name System

To connect on the Internet, the client must know the *address* of the server. This *address* is called an **IP address**, and it is represented by 4 numbers between 0 and 255. But since these 4 numbers are hard to remember, (imagine having to memorize addresses for each website you use!) people use the *name* of servers, instead of IP addresses.

### 2.6.1 Domain Name System

- DNS **translates** an hostname to an IP address
- Procedure (Simplified)
  1. The client asks the DNS for the IP address of some website
  2. The DNS responds with the IP address of that website
  3. The client uses that IP address to connect to the website
- Achieves *load distribution* - many IP addresses correspond to one hostname (replicated Web servers), thus requests to the same hostname can be distributed between multiple servers
- **Distributed database system**
- De-centralized system
  - Not scalable (for large number of clients)
  - Single point of failure: The whole system breaks down when the centralized system fails
  - Traffic volume: Bottleneck
  - Distant centralized database: Longer delay for connections from far places from the system

- Therefore uses **distributed & hierarchical** database
  - Root DNS server on the top
  - Top-level domains (TLD) on the next level (com, org, edu)
  - Authoritative domains on the next level
- Database Query Procedure (Simplified)
  1. Client wants IP for `www.amazon.com`
  2. Client queries root DNS server to find com DNS server
  3. Client queries com DNS server to get `amazon.com` DNS server
  4. Client queries `amazon.com` DNS server to get the IP address for `www.amazon.com`

### 2.6.2 DNS Hierarchy

- **Root Name Servers**
  - Directly answers requests for records in the root zone
  - Answers requests by returning a list of the authoritative name servers for the appropriate top-level domain
  - 13 of them worldwide
- **Top Level Domains (TLD)**
  - Responsible for com, org, net, edu ...
  - All top-level country domains like uk, fr, ca, kr ...
- **Authoritative Servers**
  - An organization's own DNS server(s), providing authoritative hostname to IP mappings for organization's named hosts
  - Maintained by organization itself or service provider
- **Local DNS Servers**
  - Does not strictly belong to DNS hierarchy
  - Each ISP has on also called “default name server”
  - When the client makes a DNS query, it is sent to its local DNS server
  - The local DNS server caches the results of recent queries (name to address translation pairs)
  - Can also act as a proxy, forwards query into hierarchy

### 2.6.3 DNS Name Resolution

- Situation: Host at `cis.poly.edu` wants IP for `gaia.cs.umass.edu`
- **Iterated Query**
  - Contacted server replies with the name of server to contact
  - *I don't know this name, but ask this server*

- In the procedure below, we assume that no results are cached. If there is a cached result during the process, that result can be used. Then some steps can be skipped

1. Client sends a request to local DNS server
2. The local DNS server asks the root DNS server for the edu TLD DNS server
3. The root DNS server replies
4. The local DNS server asks the edu TLD DNS server for the authoritative DNS server
5. The local DNS server asks the authoritative DNS server for the IP of `gaia.cs.umass.edu`
6. The local DNS server replies to the client with the IP address

- **Recursive Query**

- Contacted server requests another server
- Puts the burden of name resolution on the contacted name server
- Not recommended due to heavy load at upper levels of hierarchy and security issues
- Systems on upper levels must wait for the result of each query - may result in a DoS attack

1. Client sends a request to local DNS server
2. The local DNS server asks the root DNS server
3. The root DNS server asks the edu TLD DNS server
4. The edu TLD DNS server asks the authoritative DNS server
5. The authoritative DNS server replies with the IP address
6. The edu TLD DNS server responds to the root DNS server with the IP address
7. The root DNS server responds to the local DNS server with the IP address
8. The local DNS server replies to the client with the IP address

#### 2.6.4 Attacking DNS

- DDoS Attacks

- Attacking root DNS servers with traffic - doesn't work well since local DNS servers cache IPs of TLD servers (doesn't connect to the root server)
- Attacking TLD servers can be more dangerous

- Amplification Attacks

- Tricks the DNS server by putting the victim's IP inside the query
- The DNS server will send the result of the query to the victim's IP

- Pharming Attacks

- Private data + Farming
- Domain hijacking, DNS poisoning
- The attacker breaks into the local DNS and modifies some mapping to the attacker's fake website
- The user may connect to that fake website and may enter private information



## 2.7 Peer-to-Peer Application

### 2.7.1 P2P Architecture

- No always-on server
- Arbitrary end systems communicate directly
- Peers are intermittently connected and change IP addresses
- P2P is scalable, compared to client-server model

### 2.7.2 BitTorrent

- File is divided into 256 KB chunks
- Peers in torrent send/receive file chunks
- Torrent: group of peers exchanging chunks of a file
- Tracker: tracks peers participating in torrent
- A new client arrives, obtains list of peers from tracker, and begins exchanging file chunks with peers
- Chunk Receiving
  - At any given time, different peers have different subsets of file chunks
  - Periodically, client asks each peer for list of chunks that they have
  - Client requests missing chunks from peers, *rarest first*
  - Then the rarity of chunks will decrease, and be more available to other peers
- Chunk Sending: *tit-for-tat*
  - To prevent free-riders (people who download but do not provide content)
  - The sender sends chunks to 4 other peers that are currently sending chunks to itself at the highest rate
  - Other peers will not receive chunks from this sender
  - The top 4 peers are evaluated every 10 seconds
  - For newly connected peers to receive chunks, the sender selects a random peer every 30 seconds, and starts sending chunks to it
  - Now the newly connected peer may be included in the top 4 peers
  - More upload results in better peers, resulting in faster download of data

## 2.8 Video Streaming and CDNs

### 2.8.1 Content Distribution Networks

- Video traffic is the major consumer of Internet bandwidth
- Challenge: **Scalability** - If we only use a single video server ...
  - Single point of failure

- Point of network congestion
  - Longer path to distant clients
  - Multiple copies of same video sent over outgoing link
- Solution: *multiple copies of videos at multiple geographically distributed sites*
- CDN servers contain multiple copies of the same content, and lets the client stream the content data from the nearest/fastest CDN server

## 3 Transport Layer

### 3.1 Transport Layer Services

Some terminology:

- **Program**: An executable file containing a set of instructions written to perform a specific job (usually stored on a disk)
- **Process**: An executing *instance* of a program that resides on the primary memory. Several processes can be related to the same program at the same time
- **Thread**<sup>1</sup>: The smallest executable unit of a process

#### 3.1.1 Transport Layer Function

- Provides **logical communication between processes**
- The layer relies on and enhances services from the network layer<sup>2</sup>
- Sending Side
  - Applies **fragmentation** to application messages
  - Passes segments<sup>3</sup> to network layer
- Receiving Side
  - **Reassembles** segments into messages
  - Passes the assembled message to the application layer

#### 3.1.2 TCP and UDP

- **Transmission Control Protocol (TCP)**<sup>4</sup>
  - **Reliable, in-order**<sup>5</sup> delivery
  - **Connection oriented** service: connection setup, error control, flow control, congestion control
- **User Datagram Protocol (UDP)**
  - Unreliable, unordered delivery
  - Connection-less service: faster than TCP

---

<sup>1</sup>**Thread** of execution is the smallest sequence of programmed instructions that can be managed independently by a scheduler, which is typically a part of the operating system.

<sup>2</sup>The network layer provides logical communication between **hosts**

<sup>3</sup>We see terms like *frame*, *datagram*, *packet*, *segment* when we study computer networks. They are similar, but we use layer-specific terms to represent the information unit in that layer. In the transport layer, we use the term **segment**.

<sup>4</sup>Refer to 2.1.6.

<sup>5</sup>This doesn't mean that the order of sent messages is always preserved when receiving. (physically) In the application layer's point of view, it just *seems* like it's received in the same order.

## 3.2 Multiplexing and Socket

### 3.2.1 Multiplexing and Demultiplexing

- This is the most basic role of the transport layer
- **Multiplexing** at sender: the sender sends data from its own multiple applications through network
- Data from multiple services are sent through a single shared channel
- **Demultiplexing** at receiver: the receiver delivers data packets to their appropriate receivers among its own multiple applications
- **Port numbers** are used for demultiplexing
  - Different applications are assigned to different port numbers
  - Transport layer segments have fields for source/destination port numbers in common
  - Used to differentiate segments
- Note that for connection oriented protocols (like TCP), the source IP and port are also used to differentiate each connection<sup>6</sup>
- But how does the sender know the destination port on the receiver?

### 3.2.2 Sockets

- API between application layer and transport layer
- Processes send/receive messages to/from its **socket**
- Analogous to door
  - Sender passes message through the door
  - Sender relies on transport infrastructure on other side of the door to deliver message to socket at the receiving process
- The *socket* is provided as the form of APIs by the operating system

## 3.3 User Datagram Protocol

- Only provides the basic functions (multiplexing)
- **Connection-less** service
  - Each UDP segment is handled independently of others
  - *Unreliable*: UDP segments may be lost or delivered out of order to app
- But since UDP is fast, it is used for streaming multimedia applications, DNS, and SNMP
- For reliable transfer over UDP, the application must add that function (such as application specific error recovery)
- **UDP segment header**: Source port (16 bits), Destination port (16 bits), length, checksum, payload

---

<sup>6</sup>Multiple applications can listen on the same port.

- **Advantages**

- No connection establishment (no delay)
- Simple: No connection state at sender/receiver
- Small header size<sup>7</sup>
- No congestion control: UDP can blast away as fast as desired

- **UDP Checksum**

- Detects transmission errors
- UDP doesn't have to provide reliable connections, but this checksum can be used to provide additional features elsewhere
- Sender creates a 16 bit integer checksum code of segment contents including the header
- The receiver will compute the checksum of the received message, and checks if the computed value is equal to the received value

### 3.3.1 Checksum Method

- Checksum is the 16-bit one's complement of the one's complement sum of a pseudo header of information from the IP header, the UDP header, and the data, padded with zero octets at the end (if necessary) to make a multiple of two octets.<sup>8</sup>

## 3.4 Reliable Data Transfer Principles

### 3.4.1 Principles of Reliable Data Transfer

- Service abstraction (provided to the upper-layer) through a reliable channel
- *Service model of TCP*: No corruption and no loss of data, delivered in the order in which they were sent
- The lower layers of the network (below transport layer) is unreliable, but the TCP protocol in the transport layer will pre-process any existing errors and pass them onto the receiver.

### 3.4.2 Error Types and Solution

- **Bit Error**

- Some of the bits are changed
- This can be checked by comparing the checksum in every segment
- If the receiver successfully received the packet, **ACK**(acknowledgment) message is sent to the sender

- **Packet Loss** (Data or ACK)

- Packet is gone, the receiver doesn't receive the packet
- *Timeout* of sender's timer - the sender sends the packet and waits for ACK, but if ACK doesn't arrive, the sender will re-send the packet

---

<sup>7</sup>Compare this to TCP headers.

<sup>8</sup>RFC 768

- But consider the case where the returning ACK is lost - the sender will re-send the packet anyways, but how does the receiver know that this packet is a duplicate or not?
- To solve this problem, *packet sequence number* is used
- Suppose the receiver received packet  $k$ , and sent an ACK message. If packet  $k$  arrives again, the receiver will know that this packet was re-sent
- Also, for transmitting large data, the data will be segmented and labeled with a packet sequence number. Then when the receiver receives the data, it will be possible to re-assemble the original data
- Packet sequence numbers allow *ordered delivery* and data duplication prevention

### 3.4.3 Automatic Repeat Request

- For communication error recovery, we need a packet re-transmission method
- We use **ARQ**(Automatic Repeat reQuest)
- **Stop and wait**: sending and checking one segment at a time
- **Pipelining method**: sending and checking multiple segments at a time (*go back N*, *selective repeat* method)

#### Stop and Wait

- Sender sends a packet and waits for the receiver's response with ACK
- After receiving the ACK message, the sender will send the next packet
- If the sender doesn't receive the ACK message, the sender will re-send the previous packet
- The receiver responds with ACK if the calculated checksum matches the checksum value in the segment
- The length of *timeout*  $t$  is the main problem!
- If  $t$  is too long, the sender has to *wait* for that amount of time which will slow down the transmission.
- If  $t$  is too short, a normal transmission may be thought of as a timeout (premature timeout). This may happen when the network is too busy. The ACK message couldn't arrive on time. This will cause the sender to re-send the same message again (often), which is a waste of network resource
- Thus when setting the timeout time, the *round trip time* between the sender and the receiver should be considered
- Performance Analysis
  - 1 Gbps link, 15 ms propagation delay, 1 kB packet

$$d_{trans} = \frac{L}{R} = \frac{8000 \text{ bits}}{10^9 \text{ bits/s}} = 8\mu s$$

- **Utilization**: fraction of time sender busy sending

$$U = \frac{d_{trans}}{RTT + d_{trans}} = 0.00027$$

- If  $RTT^9$  is 30 ms, 1 kB packet is sent every 30 ms, so *33 kb/s throughput over 1 Gbps link*<sup>10</sup>
- Usually  $d_{trans}$  is very small (fast link), compared to RTT. Thus  $U \approx 0$ , which means that the time spent in sending the packet is nearly 0%
- Note that the size of ACK message is very small, thus it is ignored here

### Pipelined Protocols

- **Pipelining**: multiple packets can be sent and received
- Range of sequence numbers must be increased
- Buffering is necessary at sender and receiver
- Suppose we send  $n$  packets in the above example. ( $n$ -packet pipelining) Then the utilization<sup>11</sup> would be

$$U = \frac{n \cdot d_{trans}}{RTT + d_{trans}} = 0.00027n$$

which is a lot better. Now we try to increase the value of  $n$

### Go Back $N$

- The sender can have up to  $N$  unacknowledged packets in pipeline at once
- The receiver only sends **cumulative ACK** - if the receiver finds an error in packet  $k$ , then the sender has to re-send packets  $k, \dots, n$
- The cumulative ACK will mean: “transmit successful, up to these packets”
- The sender has a timer for *oldest* unacknowledged packet, and when timeout occurs, the sender will re-send all unacknowledged packets

### Selective Repeat

- The sender can have up to  $N$  unacknowledged packets in pipeline at once
- The receiver sends **individual ACK** for each packet - if the receiver finds an error in packet  $k$ , only that packet needs to be re-sent
- The sender has a timer for each unacknowledged packet, and when timeout for packet  $k$  occurs, only packet  $k$  needs to be re-sent

#### 3.4.4 Go Back $N$ - In Detail

##### Sender's Perspective

- Packet sequence number is contained in each packet header
- A **window** is defined - the number of consecutive unacknowledged packets allowed

---

<sup>9</sup>Round trip time

<sup>10</sup>What a waste of resources!

<sup>11</sup>The denominator doesn't change because we measure the time from [the moment that the first packet was sent], to [the time when ACK for the first packet arrived].

1. Let the window size be  $n$ , and suppose we have sent up to  $k$  ( $k \leq n$ ) packets.
  2. The receiver will receive the data, and send a cumulative ACK for packet  $i$  ( $i \leq k$ )
  3. Since the cumulative ACK for packet  $i$  has arrived, packets  $1, \dots, i$  are acknowledged
  4. The timer will be moved to packet  $i + 1$  and wait for the next ACK. If the ACK does not arrive in time, packets  $i + 1, \dots, k$  will be re-sent
- Note that  $k$  (number of sent packets) can keep changing in the process above, just keep in mind that the acknowledged packet number  $i$  should be less than  $k$
  - Furthermore, the window size  $n$  is not fixed. Hence the term **sliding window**
  - The larger the window size, higher the throughput, but the number of packets to re-send on an error will also increase
  - The window size should be controlled so that it does not cause network congestion or receiver buffer overflow

### Receiver's Perspective

- When packet  $k$  arrives,
  - If: Packet  $k + 1$  arrives shortly<sup>12</sup> after packet  $k$  has arrived, wait for the next packet
  - Else if: Packet  $k + 1$  doesn't arrive (due to congestion or other reasons), respond with *cumulative* ACK for packet  $k$  and wait for the next packet
  - Else: A packet other than  $k + 1$  arrives (ex. packet  $k + 1$  has been lost), discard the packet<sup>13</sup> and re-send the last sent ACK
- The *Else* case is the only problem. The sender would be expecting ACK for packet  $k + 1$ , but the sender will get the ACK for the last successful transmission<sup>14</sup>
- Then packet  $k + 1$  will cause timeout and the sender will resend the packets from  $k + 1, \dots$

### 3.4.5 Selective Repeat - In Detail

- The receiver will buffer the packets, in case some packet is lost - the receiver waits for that packet to be re-sent, and when the receiver gets it, it can pass that information to the application layer<sup>15</sup>
- **Maximum packet sequence number  $\geq 2 \times$  window size**
- For example, with 4 sequence numbers 0, 1, 2, 3 and window size 3, there could be a case where
  1. The sender sent packets 0, 1, 2
  2. *The receiver got those packets and now expects packets 3, 0, 1*
  3. Unfortunately, the ACK for packets 0, 1, 2, do not arrive to the sender
  4. The sender re-sends packets 0, 1, 2 - which are the same packets from step 1

<sup>12</sup>Depends on the receiver's settings

<sup>13</sup>It's going to be re-sent from the sender anyways

<sup>14</sup>This duplicate ACK will be ignored

<sup>15</sup>Recall that transport layers had to re-assemble the packets before passing them to the application layer, so that it would look like the packets were received in order



5. The receiver has no idea that this packet is a re-sent version!
  6. The receiver will accept the re-sent packet 0 as the next packet,
- The receiver doesn't know what's happening on the senders side. It can only distinguish the packets by the packet sequence number. If the packet sequence number is too small for the window size, data may be corrupted during transmission

## 3.5 Transmission Control Protocol

### 3.5.1 Overview

- *Point-to-point*: One sender, one receiver
- **Connection Oriented Service**<sup>16</sup>
  - Reliable transfer, delivery in order
  - Handshaking initializes sender and receiver state before data exchange
- **Pipelined Transmission**
  - Window size is set by TCP congestion and flow control
  - *Congestion control*: transmission rate is controlled to avoid congestion in network
  - *Flow control*: sender will not overwhelm the receiver
  - *MSS*: maximum segment size
- **Full-Duplex Connection**
  - In the same connection, data can flow bidirectionally

### 3.5.2 TCP Segment Structure

- Source/Destination port in the front
- Followed by sequence number, acknowledgment number, data offset<sup>17</sup>, flags, window size, header/data checksum, urgent pointer
- At least 20 bytes, optional fields can follow

### Sequence and Acknowledgment Number

- **Sequence number** is the *byte stream number* of the first byte in a segment's data
- Large data is fragmented into segments (with size of MSS) when they are sent through the network
- The receiver must know the original order of these segments to re-assemble the original data<sup>18</sup>
- If the first segment has sequence number  $n$ , the next sequence number will be  $n + MSS$ .<sup>19</sup>

<sup>16</sup>How many times has this been repeated...?

<sup>17</sup>Similar to header length - gives the offset of actual data in the segment

<sup>18</sup>Recall that the segments don't actually arrive to the receiver in order, it is the transport layer's job to actually *order* them

<sup>19</sup>The first sequence number isn't necessarily 0, it depends.

- **Acknowledgment number** is the sequence number of the next segment expected by the receiver, and functions as a *cumulative ACK*
- After the receiver gets the first segment, it will put  $n + MSS$  (next sequence number) as the acknowledgment number and send it to the sender
- The sender will know that the receiver has correctly received data up to that acknowledgment number
- There's a **ACK flag** in the TCP header. If the flag bit is 1, it means that the information in the acknowledgment number is actually meaningful and works as a cumulative ACK
- In *go back N* or *selective repeat*, the data packet and the ACK packet are the same in TCP
- TCP is *full duplex*. The receiver can send its own data to the sender, and also send an ACK number with a single TCP segment. The receiver can make use of the ACK number field in the header to contain the ACK number for the data received from the sender

### 3.5.3 TCP Reliable Data Transfer

- TCP provides reliable data transfer service on top of IP's *unreliable service*
- TCP provides **pipelined segments** (for throughput/performance), **cumulative ACKs**, and **single re-transmission timer**
- Using cumulative ACKs and a single transmission timer may seem similar to *go back N*, but the difference here is that TCP doesn't throw away packets. TCP will hold on to the received packet in its *receive buffer*<sup>20</sup>
- Also this *storage of packets* may seem similar to *selective repeat*<sup>21</sup>
- **TCP Sender**
  - When receiving data from application, create a segment with the sequence number and send it, start timer if not already running
  - When timer expires, re-transmit the segment that caused the timeout and restart timer
  - When receiving ACK, update acknowledged packet list and start timer if there are still unacknowledged segments
- **TCP Receiver** (When receiving data from sender)
  - If it discovers a bit error through checksum, drop the packet
  - Check the sequence number so that there would be no gaps between sequence numbers. If there isn't any problem, send cumulative ACK
  - If the packet is duplicated, drop the packet
- Note that *handling out of order segments* depends on implementation

<sup>20</sup>Its size is limited. But this is possible because the size of TCP's receiving buffer is larger than the congestion window. So the receiver can store the data, provided that it arrives correctly.

<sup>21</sup>Kind of a fusion of both methods.

## TCP Timeout Value

- Must be longer than RTT, but this is unknown before connection, and also may change depending on the current network congestion status
- In TCP, a **time out interval** is set by  $estimated\ RTT + (\text{safety margin})$
- The hosts send data and measure the *sample RTT* - the time from segment transmission until ACK arrival
- Using the *sample RTT* values, *estimated RTT* is *iteratively* calculated by moving averages, with some weight  $\alpha$

$$t_n = (1 - \alpha)t_{n-1} + \alpha s_n$$

where  $s_n$  is the  $n$ -th sample RTT and  $t_n$  is the estimated RTT after sampling  $n$  times

### 3.5.4 Fast Re-transmit

- From the above calculations, it turns out that the timeout period is relatively long, so we get a long delay before re-sending the lost packet
- Senders often send many segments back-to-back. But if a segment is lost, there will likely be many duplicate ACKs.
- For example, suppose we send packets 1, 2, 3, 4, 5, but packet 2 was lost and every other packets were received. Then the receiver would send the same ACK for packets 3 ~ 5
- *If sender receives 3 duplicate ACKs (except for the first normal ACK) with same data, the sender will re-send the unacknowledged segment with the smallest sequence number*
- TCP re-transmissions are triggered by: timeout events and duplicated ACKs

### 3.5.5 TCP Flow Control

- The receiver controls the sender, so that the sender won't overflow receiver's buffer by transmitting too much data too fast. This is called **flow control**
- The remaining size of the receiver's buffer will be transmitted in the **receive window** field in the TCP header
- The sender will limit the size of transmission to the value inside the *receive window* field

### 3.5.6 Connection Management

- 9 flags in the TCP header - 3 are not almost never used
- Remaining 6: URG(*urgent*), ACK(*acknowledge*), PSH(*push*), RST(*reset*), SYN(*synchronize*), FIN(*finale*)
- URG, PSH are also not used frequently
- **RST, SYN, FIN** flags contribute to TCP connection management

### Three Way Handshake

1. The client requests for a connection and sends a segment with SYN flag set, and a sequence number  $x$  to start with
2. The server will receive the segment and respond with ACK flag set, and ACK number will be  $x + 1$ . Moreover, since TCP is *full duplex*, the segment will have SYN flag set and the sequence number will be  $y$
3. The client responds with ACK flag set and ACK number  $y + 1$

### Closing Connection

1. The client will send a segment with FIN flag set
2. The server will respond with ACK flag set
3. Connection from client to server is closed<sup>22</sup>
4. The server will send a segment with FIN flag set
5. The client will respond with ACK flag set
6. Connection from server to client is closed

### RST Flag

- The client will ask for a port in the server to connect with
- If there are no processes listening to that port, the server will respond with a segment with RST flag set
- Used as a flag to reset the connection

## 3.6 Congestion Control

### 3.6.1 Network Congestion

- Congestion was not an issue before the 1970s, when the Internet was first developed. There weren't many nodes, and there wasn't much traffic
- **Congestion collapse**
  - When the offered load is smaller than the network capacity, the traffic is transmitted to the destination. (offered load  $\approx$  effective load)
  - But as the offered load increases and exceeds the network capacity, the effective load decreases and eventually goes to 0.<sup>23</sup>
- **Congestion:** Too many sources sending too much data too fast for the network to handle
- Difference between *flow control*
  - Flow control is a point-to-point issue, between a sender and a receiver

---

<sup>22</sup>Note that at this point, the connection from the server to the client still lives. The server can still send data.

<sup>23</sup>Network will have more data to process, then the delay for each packet will increase, causing timeouts a lot more often, and causes re-transmission. Eventually, the effective load will decrease.

- Flow control problem could be resolved only by knowing the receiver buffer size, which is in the receiver window field in the TCP header
- Congestion control is a *global* issue
- All the nodes that share the whole network have to yield to each other - only then the traffic will decrease
- Manifestations: lost packets (buffer overflow at routers), long delays (queueing in router buffers)

### 3.6.2 Congestion Control Approaches

#### Network Assisted

- Data from multiple sources go through routers
- The routers *notify* the end systems, and provide feedback that congestion is happening, and also provides the explicit rate for the sender to send data at
- Single bit indicating congestion (SNA, DECbit, TCP/IP<sup>24</sup>, ECN, ATM)
- But this is quite complicated and causes problems - TCP uses end-to-end

#### End to End

- No explicit feedback from the network
- Congestion is inferred from the observed loss and delay by the end system
- Approach taken by TCP

### 3.6.3 TCP Congestion Control Overview

- **AIMD** approach: sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs
  - *Additive increase*: increase cwnd by 1 MSS<sup>25</sup> every RTT until loss is detected
  - *Multiplicative decrease*: Cut cwnd in half after loss
- **Congestion window** (cwnd): Maximum amount of un-ACKed packets in transit, without causing congestion
- (Size of data between *last byte sent* and *last byte ACKed*)  $\leq \min(\text{cwnd}, \text{rwnd})$   
rwnd: Receive window (receiver buffer)
- cwnd is dynamic, changes with respect to network congestion
- How do you decide the size of cwnd when connection has just began?

#### Slow Start

- When connection begins, increase rate exponentially, until the first loss event occurs
  - Initially set cwnd as 1 MSS
  - Double cwnd every RTT, if the sent packet is ACKed
- Note that the *initial rate* is slow, but the size increases exponentially fast
- When should this increase stop?

<sup>24</sup>There is a flag in the TCP header for *explicit congestion notification*. But this is rarely used.

<sup>25</sup>Maximum segment size

## Congestion Avoidance

- If all nodes increase its `cwnd` exponentially, this would cause congestion very fast
- A variable `ssthresh` (slow start threshold) is used
- On a loss event, `ssthresh` is set to half of [`cwnd` value of the most recent successful transfer]<sup>26</sup>
- If the `cwnd` size exceeds `ssthresh`, `cwnd` increases by a single MSS every RTT
- Exponential increase if less than `ssthresh`, linear increase after exceeding `ssthresh`

## 3.7 TCP Congestion Control Algorithm

- The following TCP variants were or (currently) is the standard

### 3.7.1 TCP Tahoe (1988)

- Slow start, congestion avoidance
- When loss is detected, set `cwnd` to 1
- **Fast re-transmit**: re-transmission upon 3 duplicate ACKs
- Slow start begins on timeout and fast re-transmit

### 3.7.2 TCP Reno (1990)

- Limitations of Tahoe
  - Doesn't differentiate between losses indicated by timeout and duplicate ACKs
  - Always sets `cwnd` to 1 MSS (throughput decrease)
- Two (different) responses to packet losses by Reno
- Loss indicated by *timeout*
  - Set `cwnd` to 1 MSS
  - *Slow start* to `ssthresh`, then congestion avoidance begins
- Loss indicated by *fast re-transmit*
  - `cwnd` is cut in half
  - **Fast recovery**, then congestion avoidance

---

<sup>26</sup>`ssthresh` starts at some value (may be high) when the end system connects to the network for the first time. As time passes and the system transmits data through the network, this `ssthresh` value will be maintained

## Fast Recovery

- Half of `cwnd` size is maintained until a non-duplicate ACK is received
- Amount of packets in transit is kept from decreasing more than expected
- Operation
  1. Upon the event of 3 duplicate ACKs, reduce `cwnd` to  $1/2$
  2. Re-transmit the oldest un-ACKed packet
  3. Inflate the congestion window by the number of duplicated packets
  4. A non-duplicate ACK starts congestion avoidance

### 3.7.3 TCP NewReno (1999)

- Vulnerability of Reno's fast recovery - in case of multiple packet losses<sup>27</sup>, `cwnd` may be reduced too much (exponential way)
- Reno stops fast recovery on non-duplicated ACK
- But NewReno restricts the exit from recovery until all data packets from the initial congestion window are ACKed (resistant to multiple packet losses)
- Exits from fast recovery only when a new data ACK is received
- Keeps the sequence number of the last data packet sent before entering fast recovery

### 3.7.4 TCP SACK (1996)

- Limited information of cumulative ACKs
  - ACK of only the last in-order packet
  - Reno's fast recovery assumes loss of only 1 data packet
  - A duration of the recovery is directly proportional to the number of packet losses
- **SACK** (Selective ACK): receiver provides information about several packet losses in a single ACK message by reporting blocks of successfully delivered data packets
- Client sends packets 1, 2, 3, 4 and if packet 2 is lost, ACK message from the server will contain: ACK for packet 1, SACK for 3, 4
- Then the client only has to re-transmit packet 2

## 3.8 TCP vs UDP

- Fairness between connections
  - TCP has AIMD - difference of throughput between connections is reduced by half, when loss occurs, and eventually the difference tends to 0 (fairness)

---

<sup>27</sup>And usually, multiple packets are lost, not just single one

- UDP sends the amount the sender wants to send (unfair)
- Note that multimedia applications do not want the rate throttled by congestion control, so they use UDP and tolerate packet loss

<b>TCP</b>	<b>UDP</b>
<ul style="list-style-type: none"><li>- Slower but reliable transfers</li><li>- Email, web browsing</li><li>- Unicast</li></ul>	<ul style="list-style-type: none"><li>- Fast but non-guaranteed transfers</li><li>- VoIP, music streaming</li><li>- Unicast, multicast, broadcast</li></ul>



## 4 Network Layer I

### 4.1 Overview of Network Layer

#### 4.1.1 Network Layer

- Transport segment from sending host to receiving host
- Sending side encapsulates segments (from transport layer) into **datagrams**<sup>1</sup> and sends them
- Receiving side receives datagrams and delivers the segments to transport layer
- Network layer protocols exist in **every** host and router
- Router examines header fields in all IP datagrams passing through it - the destination is determined by **IP address**

#### 4.1.2 Two Key Functions of Network Layer

##### Routing

- Determines the route taken by packets from source to destination
- Routing algorithms help create the forwarding table

##### Forwarding

- Moves packets from router's input to appropriate router output
- Packet delivery to the next node

##### Traditional vs SDN Network

- In traditional IP networks, routing and forwarding is done in the same system (for each router and switch)
- In *software-defined network* (SDN), routing function and forwarding function are separated - The central server determines the optimal route, and the switches do the forwarding according to the route given by the central server

### 4.2 Inside of Router

- router input ports → high-speed switching fabric + routing processor → router output ports

#### 4.2.1 Input Port Functions

- physical layer → data link layer → decentralized switching
- **Physical layer** determines if the input is 0 or 1 (bit level)
- **Data link layer** - Explained later in Chapter 5. Checks errors, etc.
- IP datagram is created after these two layers
- IP datagrams may not be processed in real-time, because switching fabric may be busy - datagrams are stored in a buffer called **input queue**

---

<sup>1</sup>'Message' or 'packet' that is on the network layer level

- **Decentralized switching**

- Using header field values, lookup output port using forwarding table in input port memory
- The goal is to complete input port processing at 'line speed'
- *Queueing* is done if datagrams arrive faster than forwarding rate into switch fabric

### Input Port Queueing

- Fabric may be slower than input ports combined. Then queueing occurs at input queues
- Queueing delay and loss may occur due to input buffer overflow
- **Head-of-the-line blocking:** queued datagram at front of the queue may prevent other datagrams in queue from moving forward

### 4.2.2 Switching Fabrics

- Transferring packet from input buffer to appropriate output buffer
- *Switching rate:* rate at which packets can be transferred from inputs to outputs
- Switching rate is often measured as a multiple of input/output line rate
- $N$  inputs: switching rate of  $N$  times line rate is desirable
- Three types of switching fabrics: memory, bus, interconnection network

### Switching via Memory

- First generation routers
- Traditional computers with switching under direct control of CPU
- Packets are copied to system's (main) memory, and the IP in the header is read and sent to the destination
- Speed is limited by memory bandwidth (2 bus crossings per datagram)

### Switching via Bus

- Datagram from input port memory to output port memory via a shared bus (direct copying - faster)
- Any input port can be sent to any output port (decided by CPU)
- But because a shared bus is used, *bus contention* occurs: switching speed is limited by bus bandwidth
- When sending a datagram, an input port and an output port will be used exclusively (limitation of speed, because of a shared bus)

## Switching via Interconnection Network

- Multiple datagrams can be sent at the same time, if the datagrams do not share ports - if the output port has to be shared, collision is inevitable
- A lot faster than memory/bus method
- Banyan networks, crossbar, other interconnection nets are used for multiprocessing
- Banyan networks use less switches compared to crossbars so it may be cheaper, but more susceptible to HOL blocking or collisions

### 4.2.3 Output Port Functions

- **Buffering** is required when datagrams arrive from fabric faster than the transmission rate
- **Scheduling** discipline chooses among queued datagrams for transmission
- Scheduling: choosing the next packet to send on link
- **FIFO**<sup>2</sup>: send in order of arrival to queue
- Discard policy: if packets arrive to a full queue, which packet should be drop?
  - *Tail drop*: Drop the arriving packet
  - *Priority*: Drop the packet on priority basis
  - *Random*: Drop a random packet

### Priority Scheduling

- Sends highest priority queued packet
- Keep two queues, for high/low priority packets, respectively
- Multiple classes, with different priorities
- Priority class may depend on marking or other header info such as IP source/destination, port numbers etc.
- Packets in the low priority queue will be transmitted if the high priority queue is empty - causes problems

### Round Robin (RR) Scheduling

- Keep multiple classes and queues
- Scan all queues round and around. For each cycle, send a *single* packet if the queue is not empty
- Each queue is given a chance to send its packets so this method is a bit different compared to priority scheduling
- But if there are less (kinds of) packets corresponding to some class, the packets are likely to be processed a lot faster (less delay)

---

<sup>2</sup>First in first out, or sometimes FCFS, for first come first served

## Weighted Fair Queueing (WFQ) Scheduling

- Generalized round robin: RR + weights
- Each class gets weighted amount of service in each cycle

## 4.3 Internet Protocol Overview

- Network layer functions of host and router
- **Routing protocols** - path selection, RIP, OSPF, BGP<sup>3</sup>
- **IP protocol** - addressing conventions, datagram format, packet handling conventions
- **ICMP<sup>4</sup> protocol** - error reporting, router signaling

### 4.3.1 IP Datagram Format

- 20 bytes of header + options (if necessary)
  - *Version*: 4 or 6
  - *Header length*: offset of data / length of the header + options
  - *Type of service*: multimedia data, real-time data... rarely used
  - *Datagram length*: length of the datagram
  - *16-bit identifier, flags, fragment offset*: For fragmentation and reassembly of segments in routers
  - *Time to live(TTL)*: maximum number of remaining hops, decremented at each router
  - *Upper layer*: Upper layer protocol to deliver payload to (TCP or UDP)
  - *Header checksum*: cycle redundancy check (CRC)
  - *Source/Destination IP address*: 32 bit

### 4.3.2 IP Fragmentation & Reassembly

- Network links have *maximum transfer unit* (MTU)
- Each network uses different link types, so they have different MTUs
- Large IP datagrams should be divided (fragmented) within the network
- Reassembly will occur at the final destination
- IP header bits are used to identify and order related fragments
- *Identifier*: Datagrams fragmented from large datagrams contain the same ID
- *Fragment offset*: Offset of this datagram's data before fragmentation. 13 bits are used for offset, so the value is divided by 8 (flag uses 3 bits)<sup>5</sup>
- *Flag*: 3 bits. First bit is reserved, second bit is for *don't fragment flag* (DF), last bit is for *more fragment flag* (MF)
- If MF flag of a packet is 0, it is the last datagrams among fragmented datagrams.

---

<sup>3</sup>Routing information protocol, open shortest path first, border gateway protocol

<sup>4</sup>Internet control message protocol

<sup>5</sup>MTU - (header size) is divided by 8 to calculate the offset

### 4.3.3 TTL

- TTL field (8 bytes)
- Included to ensure that datagrams do not circulate forever (ex. a long-lived routing loop) in the network
- Decrement by one each time the datagram is processed by a router
- If TTL reaches 0, a router must drop that datagram

## 4.4 IP Addressing

- **IP address:** 32 bit identifier for host, router interface
- *Interface:* connection between host/router and physical link
  - Routers typically have multiple interfaces
  - Hosts typically have one or two interfaces
- IP address associated with each interface
- IP address is *hierarchically addressed*

### 4.4.1 IP Address Class

- **ICANN**(Internet Corporation for Assigned Names and Numbers): allocates addresses, manages DNS, assigns domain names, resolves disputes
- IP: *network ID + host ID*
  - *Network ID:* ID representing a network group, or a group of hosts
  - *Host ID:* ID of each host in a network group
- Three classes: A, B, C
  - Class A: starts with bit 0, network 7 bits + host 24 bits (First byte 0 ~ 127)
  - Class B: starts with bit 10, network 14 bits + host 16 bits (First byte 128 ~ 191)
  - Class C: starts with bit 110, network 21 bits + host 8 bits (First byte 192 ~ 223)

### 4.4.2 Subnets

- **Subnet:** a logical sub-division of an IP network
- Datagram forwarding is performed by routers
- Hosts in the same network can reach each other without intervening router
- Too many hosts in a network may increase maintenance overhead
- *Divide and conquer*
- Division of IP address: Host ID = Subnet number + Host number
- Deciding the bit size of subnet number

- *Subnet mask*: Indicates the bits that will be used as the network number
- If the bit at some position is 1, the bit will be used as the network number
- Ex. 255.255.255.0 → 24 bits will be used as the network number
- The actual bit size of the subnet number will be: number of 1s in subnet mask - bit size of network ID (depending on IP address class)

#### 4.4.3 Classless Inter-Domain Routing (CIDR)

- *Supernetting*
- Limitations of IP address classes
  - Class A: 128 possible network IDs, too small
  - Class B: approximately 16,000 network IDs → small nowadays
  - Class C: Only 256 possible hosts, too small
- CIDR address format: A.B.C.D/x where x is the number of bits in the network portion of address
- Multiple class C addresses can be grouped together to function as a single network group
- The hosts in each class C address can be assigned to individual hosts inside the network

#### 4.5 Datagram Forwarding

- IP packets are often called datagrams, hence IP network is often called datagram network
- IP network has no call setup (recall that TCP had handshakes)
- Router looks at the destination address in the datagram and decides the next route for the datagram
- Routers don't have any state about end-to-end connections
- 2 methods of datagram forwarding
- *Destination based forwarding*: Based only on destination IP address (traditional)
  - What is the next most appropriate node, considering the destination address?
- *Generalized forwarding*: Based on any set of header field values (SDN)
  - Generalized, in the sense that factors other than the destination address can be taken into account when the datagram is forwarded
  - For instance, load balancing, congestion etc. can be taken into account
  - This is made possible by software defined networks
  - More about this later

### 4.5.1 Destination Based Forwarding

- When a datagram arrives, the router will look at the destination address
- Then with the local forwarding table, the router will decide the output link
- But since there are  $2^{32}$  IP addresses, it is inefficient to store all of them inside the table
- Instead, it will only contain the *network ID*
- When the datagram is forwarded to the destination network, the destination network will forward the datagram to the host
- Can we do better? There are  $2^{21}$  class C addresses...
- The router will send a range of destination addresses to a single output link
- **Longest prefix matching**: when looking for a forwarding table entry for given destination address, use the *longest* address prefix that matches the destination address
  - The longest matching prefix will have more information about the destination
  - So the longest matching prefix will be the most appropriate output link

## 4.6 Dynamic Host Configuration Protocol

- Host dynamically obtains its IP address from the network server when it joins the network
  - Allows reuse of addresses (only hold the address while connected)
  - Can renew its lease on addresses in use
  - Support for mobile users who want to join the network
  - Dynamic allocation of IP addresses for each connected user

### 4.6.1 DHCP Overview

#### 1. DHCP Discover: host broadcast

- Broadcast: Is there a DHCP server out there?
- Destination address: 255.255.255.255 (broadcast to all devices on network)

#### 2. DHCP Offer: DHCP server responds

- Broadcast: I'm a DHCP server, here's an IP address you can use!
- yiaddr contains the IP address to use
- The client checks transaction ID to realize that this is a response to its request<sup>6</sup>

#### 3. DHCP Request: Host requests IP address

- Broadcast: OK, I'll take that IP address!

#### 4. DHCP ACK: DHCP server sends IP address

---

<sup>6</sup>DHCP offer is also broadcast to all devices on network

- Broadcast: OK, you can use that IP address!
- IP obtained from DHCP has a *lifetime*. Should be renewed before expiration
- Why take 4 steps? Can't we use the address given from the offer?
  - There may be multiple DHCP servers on the network
  - Client will receive multiple offers from multiple servers
  - The client decides on which offer to take, and requests it to the server
  - *Broadcast* - the decision will be broadcast, to let other devices on the network know about this decision
- DHCP server will give client's IP address, IP address of the first-hop router for client, and name/address of DNS server - as these information are necessary for the client to connect to the network

#### 4.6.2 Advantages of DHCP

- Efficient use of IP addresses
  - IPv4 address exhaustion
  - If every device used static IPv4 addresses, each device would have to hold on to that address even when it is not connected to the network
  - IP is also considered as a network resource, should be used efficiently
- Allows the "plug & play" of a computer system
  - In the past, we had to manually type in static IP addresses, subnet masks, default gateway
  - Dynamic IP, DNS server address allocation provides abstraction for people without network knowledge
  - Plug in the network cable and everything is set up automatically

### 4.7 Network Address Translation

- In the past, each household used a single IP for its desktop computer. Nowadays, each household has more network devices such as smartphones etc. How to we assign IP addresses to all the devices?
- *Motivation*: local network uses just one IP address as far as outside world is concerned...
  - Range of IP addresses is not needed from ISP - just one IP address for all devices
  - But this IP address gets translated inside the local network, without notifying the outside world
  - Can change ISP without changing addresses of devices in local network
  - Devices inside local network is not explicitly addressable or visible by the outside world (security)

#### 4.7.1 Network Address Translation

- **Private IP**: Meaningful only in the local/home network
  - A class: 10.x.x.x
  - B class: 172.16.x.x ~ 172.31.x.x
  - C class: 192.168.0.x ~ 192.168.255.x



- All datagrams leaving the local network have the same single source NAT IP address (**public IP**), but different source port numbers
  - Datagrams with source or destination in the network use private IP for source and destination
1. Host in local network sends datagram to destination, using *private IP* as source
  2. The router uses **NAT translation table** to change datagram source address from private IP to *public IP* (update table if necessary)
  3. Reply arrives with destination address using the *public IP*
  4. The router uses NAT translation table again to change the destination address of datagram from public IP to *private IP*
  5. The host in local network receives the datagram

## 4.8 IPv6

- Motivation
  - DHCP, NAT partially solved the IPv4 address exhaustion issue, but it isn't the fundamental solution
  - The number of addresses should be increased
- Additional Motivation
  - Network technology has advanced since the use of IPv4
  - Why not change the header format to speed up processing and forwarding?
  - Header changes to facilitate QoS

### 4.8.1 IPv6 Datagram Format

- Uses fixed-length 40 byte header
- Comparison with IPv4 datagram
  - IPv4 used 20 bytes, with 8 bytes for source and destination. Other fields consist the remaining 12 bytes
  - IPv6 address is 16 bytes, using 32 bytes for source and destination address - this is inevitable, and the remaining 8 bytes are used to contain other values, which is a *decrease*
- No fragmentation allowed (fragmentation offset, flag, id is gone)
- *version*: either v4 or v6
- *priority*: identify priority among datagrams in flow
  - Router does priority scheduling
  - On buffer overflow, throw away packets with lower priority
- *flow label*: identify datagrams in same 'flow'
  - IPv4 used to consider each datagrams independently

- Define a flow label between source/destination processes, and the datagrams between these processes will have the same flow label
- Treats subsequent IP datagrams as a whole
- Important in software defined networks (SDN)
- *next header*: identify upper layer protocol for data (UDP/TCP)
- *hop limit*: same as TTL in IPv4
- Checksum has been removed entirely to reduce processing time at each hop
  - We use optical fibers now - reduced error rate
  - TTL value decreases every hop, thus checksum has to be recalculated every hop
  - Checking and recalculating checksum is an overhead

#### 4.8.2 Migration from IPv4 to IPv6

- Not all routers can be upgraded simultaneously
- How to operate with mixed IP versions?
- **Tunneling**: IPv6 datagram is carried as payload in IPv4 datagram among IPv4 routers
  - Pad IPv4 header onto the IPv6 datagram
  - The IPv6 routers should provide dual protocol
- Long time for deployment!
- Necessary for the future - preparation for the IoT era, more objects will need IP addresses

## 5 Network Layer II

### 5.1 Introduction to Routing

- **Routing**: determine the route taken by packets from source to destination
- **Forwarding**: move packets from router's input to appropriate router output
- **Control plane**: routing algorithm determines end to end path through network
- **Data plane**: forwarding table determines local forwarding at this router
- *Per-router control plane*: individual routing algorithm components in each and every router interact with each other in control plane to compute forwarding tables
- *Logically centralized control plane*: a distinct (typically remote) controller interacts with local control agents (CA) in routers to compute forwarding tables
  - A remote controller controls the forwarding process
  - Forwarding path decided with more information, thus efficient
  - Often used in software defined networks

#### 5.1.1 Routing Protocols

- **Goal of Routing**: determine 'good' paths from source to destination through network of routers
  - *Path*: a sequence of routers that packets will traverse, in going from initial source to final destination
  - Paths that take the least cost, that are fastest, that are the least congested
- Abstraction as a graph  $G = (N, E)$  (can be a directed graph)
  - $N$  is the set of vertices - *routers* in this case
  - $E$  is the set of edges - *links* in this case
  - Cost function  $c : N \times N \rightarrow \mathbb{R}$  where  $c(v, w)$  is the cost of using the link between router  $v$  and  $w$

#### 5.1.2 Routing Algorithm Classification

##### Static vs Dynamic

- Static
  - Routes are changed manually by the administrator
- Dynamic
  - Routes are changed more quickly
  - Periodic update in response to link cost changes

## Global vs Decentralized

- Global
  - All routers have complete topology of the network and all link cost information
  - *Link state algorithms*
- Decentralized
  - Router knows directly connected neighbors
  - Router knows the path costs to other nodes, and the next router that should be taken to reach some other node
  - *Distance vector algorithm*

## 5.2 Link State Routing

- Need to know the complete topology of the network - **link state advertisement**(LSA) message
  - Contains neighbor node information
  - Contains link information to others (capacity, delay, etc.)

### Operation

1. A router constructs LSA
  - When a neighbor changes, link goes up or down (activated/deactivated)
  - Periodically
2. The LSA is broadcast to **all routers** in the network
3. Based on LSA from all other nodes, each router constructs the complete topology
4. Least cost paths between two nodes are computed by **Dijkstra's algorithm**

### 5.2.1 Dijkstra's Algorithm

- Single source shortest path algorithm
- Computes least cost paths from the source to all other nodes
- $k$  iterations will give the least cost path to  $k$  destinations

### Steps

1. Let  $S = \{u\}$  (source)
2. For all  $v \in N(S)$  (neighbor nodes) set  $d(v) = c(u, v)$ . If  $v \notin N(S)$ , set  $d(v) = \infty$
3. while  $S \neq N$ 
  - (a) Find  $w \notin S$  such that  $d(w)$  is minimum, and  $S \leftarrow S \cup \{w\}$
  - (b) For all  $x \in N(w)$  and  $x \notin S$ , update  $d(x) = \min\{d(x), d(w) + c(w, x)\}$
4.  $d(x)$  contains the cost of shortest path from  $u$  to  $x$

## 5.3 Distance Vector Routing

### 5.3.1 Bellman-Ford Equation

- Let  $d_x(y)$  be the cost of the least-cost path from  $x$  to  $y$
- Then,

$$d_x(y) = \min_{v \in N(x)} \{d_v(y) + c(x, v)\}$$

### 5.3.2 Distance Vector Algorithm

- Node  $x$  initially
  - Knows cost to each neighbor  $v$ :  $c(x, v)$
  - Maintains its neighbors' distance vectors. For  $v \in N(x)$ , node  $x$  maintains

$$d_v = \{d_v(y) : y \in N\}$$

- From time to time, each node sends its own distance vector estimate to neighbors
- When  $x$  receives new distance vector estimate from neighbor, it updates its own distance vector using Bellman-Ford equation

$$d_x(y) \leftarrow \min_{v \in N(x)} \{d_v(y) + c(x, v)\} \text{ for each node } y \in N$$

- Each node
  - *Waits* for changes in local link cost or messages from neighbor
  - *Recomputes* estimates
  - If distance vector to any destination has changed, *notifies* neighbors

### Good News Travels Fast

- When link cost changes,
  - Nodes detect local link cost change
  - Nodes update routing info, and recalculates distance vector
  - If distance vector changes, notify neighbors
- Link cost decrease travels fast to other nodes
- Link cost increase travels slowly to other nodes - *count to infinity* problem
  - Takes a lot of iterations before algorithm stabilizes

### 5.3.3 Link State vs Distance Vector

Link State	Distance Vector
Incremental updates	Entire routing table is sent as an update
Fast convergence	Slow convergence
Updates are broadcast to entire network	Updates are sent to directly connected neighbor only
Routers know the entire network topology of that area	Routers don't know the entire network topology
No routing loops	Count to infinity problem
Complex configuration (for large network)	Simple configuration
Example: <i>OSPF</i> (Open Shortest Path First)	Example: <i>RIP</i> (Routing Information Protocol)

## 5.4 Intra-AS Routing: OSPF

### 5.4.1 Scalability of Routing

- So far, we have assumed that all routers are identical and the network is 'flat'
- In practice, we have billions of destinations. We can't store all destinations in routing tables, and routing table exchange would swamp links
- There are **scalability issues** in routing
- Administrative autonomy
  - Internet is a network of networks
  - Each network admin may want to control routing in its own network

### 5.4.2 Autonomous System

- Aggregate routers into regions known as **autonomous systems** (AS)
- *Autonomous system* is a collection of connected Internet protocol routing prefixes under the control of one or more network operators on behalf of a *single administrative entity*

### Intra-AS Routing

- Routing among hosts, routers in same AS (network)
- All routers in AS must run the same intra-domain protocol
- Routers in different AS can run different intra-domain routing protocol
- Gateway router (border router) at the edge of its own AS has links to routers in other AS
- RIP (distance vector), OSPF (link state)

### Inter-AS Routing

- Routing among multiple AS
- Gateway routers perform inter-domain routing (as well as intra-domain routing)
- BGP (path vector)

## Interconnected AS

- Forwarding table is configured by both intra-AS and inter-AS routing algorithm
  - Intra-AS routing algorithm determines entries for destinations within the AS
  - Inter-AS and intra-AS algorithm determine entries for external destinations

### 5.4.3 Routing Protocols

#### Open Shortest Path First (OSPF)

- Most common intra-AS routing protocol
- 'open': publicly available open protocol (documentation available)
- Uses link state algorithm
  - LSAs to all other routers in entire AS - carried directly over IP
  - Topology map is known at each node
  - Route computation using Dijkstra's algorithm
- Security issue: all OSPF messages should be authenticated, to prevent malicious intrusion
  - OSPF packet contains authentication type and authentication
  - If authentication fails, the message is ignored
  - *Black hole attack*: sends false information about shortest paths, which may cause congestion

#### Hierarchical OSPF

- Hierarchical OSPF in large domains
- Minimizes routing update traffic
- Localizes the impact of topology changes to an area
- *Two level hierarchy*: local area, backbone
  - LSAs are only shared in the same area
  - Each node has detailed area topology, only knows direction (shortest path) to nets in other areas
- Area border routers: summarize distances to nets in own area, advertise to other area border routers
- Backbone routers: run OSPF routing limited to backbone area, and connect to other AS

## 5.5 Inter-AS Routing: BGP

- 5.6 Software Defined Networking**
- 5.7 OpenFlow**
- 5.8 Internet Control Message Protocol**