

Prepared by
G .Kiran Kumar
Assistant Professor,
ECE

INTRODUCTION TO PHYSICAL LAYER


Introduction: Networks, network types, internet history, standards and administration; Network models: Protocol layering, TCP/IP protocol suite, the OSI model; Introduction to physical layer: Data and signals, transmission impairment, data rate limits, performance; Transmission media: Introduction, guided media, unguided media; Switching: Introduction, circuit switched networks, packet switching.

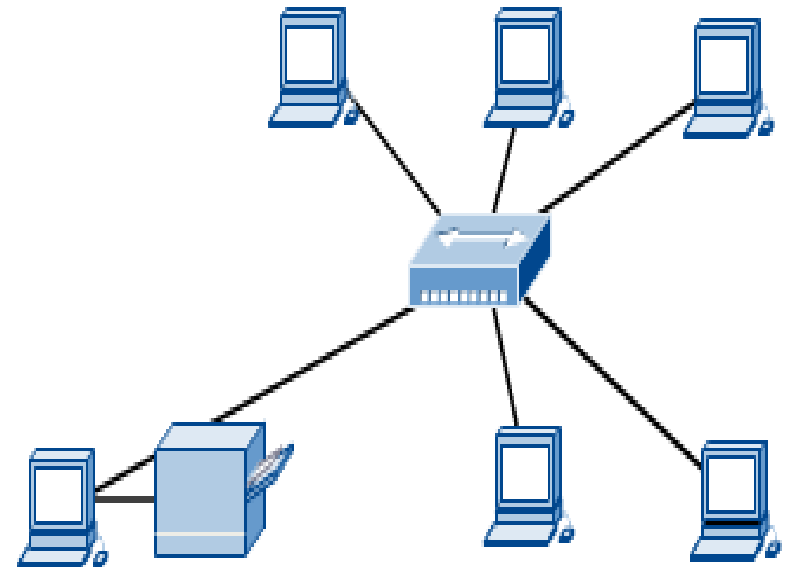


INTRODUCTION TO COMPUTER NETWORKS

Computer Networks

 **Computer network connects two or more autonomous computers.**

 **The computers can be geographically located anywhere.**



LAN, MAN & WAN



- ❏ Network in small geographical Area (Room, Building or a Campus) is called LAN (Local Area Network)
- ❏ Network in a City is call MAN (Metropolitan Area Network)
- ❏ Network spread geographically (Country or across Globe) is called WAN (Wide Area Network)

Applications of Networks

Resource Sharing

-  Hardware (computing resources, disks, printers)
-  Software (application software)

Information Sharing

-  Easy accessibility from anywhere (files, databases)
-  Search Capability (WWW)


Communication

-  Email
-  Message broadcast

Remote computing

Distributed processing (GRID Computing)

Network Topology

 The network topology defines the way in which computers, printers, and other devices are connected. A network topology describes the layout of the wire and devices as well as the paths used by data transmissions.



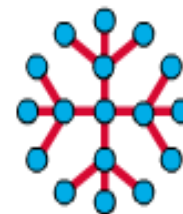
Bus Topology



Ring Topology



Star Topology



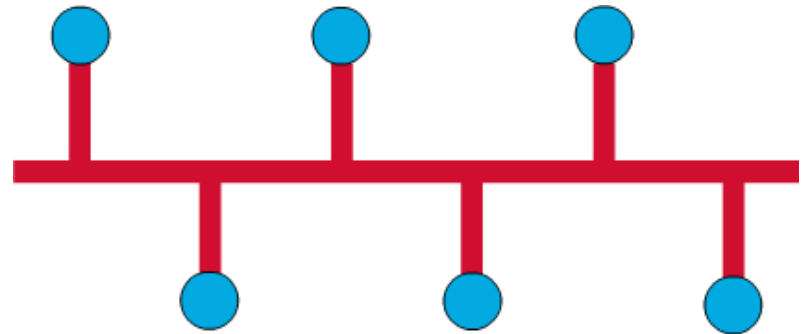
Extended Star Topology



Mesh Topology

Bus Topology


Commonly referred to as a linear bus, all the devices on a bus topology are connected by one single cable.

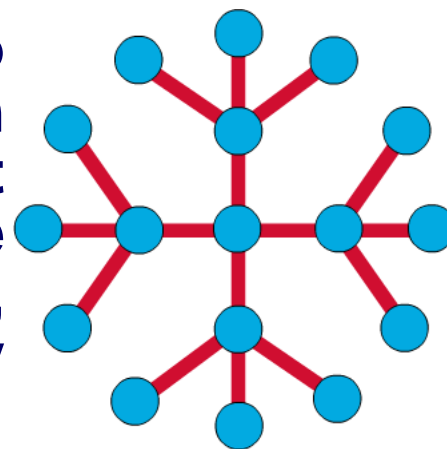
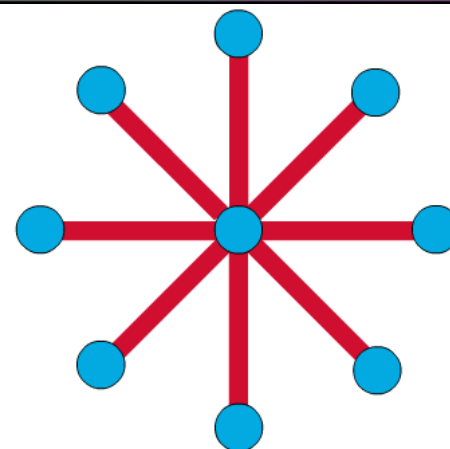


Star & Tree Topology


 The star topology is the most commonly used architecture in Ethernet LANs.


 When installed, the star topology resembles spokes in a bicycle wheel.

 Larger networks use the extended star topology also called tree topology. When used with network devices that filter frames or packets, like bridges, switches, and routers, this topology significantly reduces the traffic on the wires by sending packets only to the wires of the destination host.




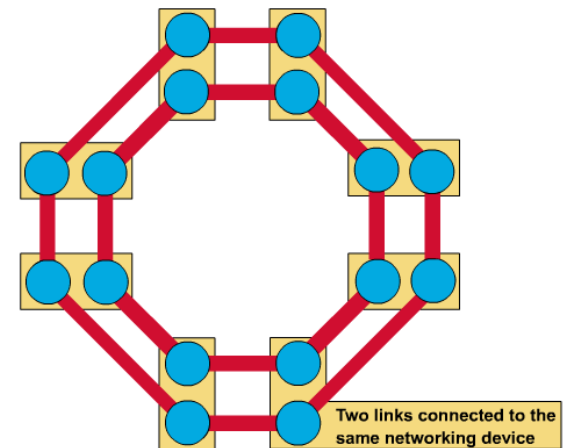
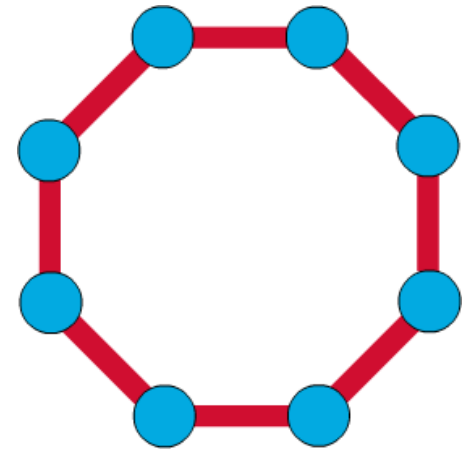
Ring Topology

 A frame travels around the ring, stopping at each node. If a node wants to transmit data, it adds the data as well as the destination address to the frame.

 The frame then continues around the ring until it finds the destination node, which takes the data out of the frame.

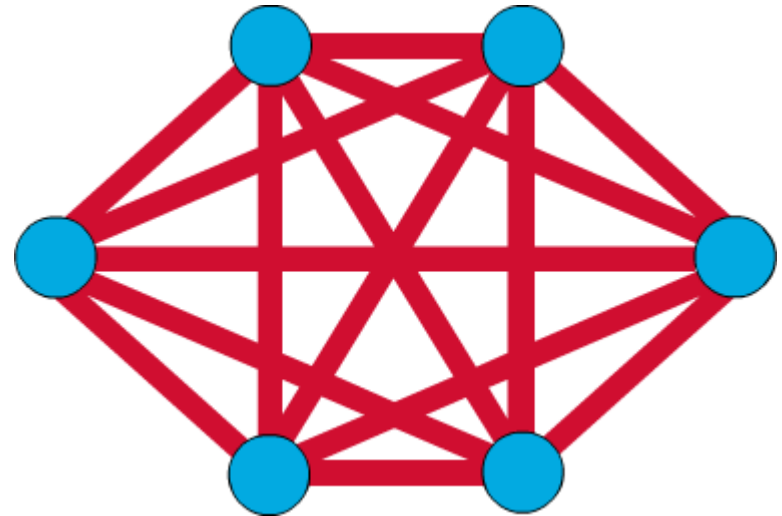
 Single ring – All the devices on the network share a single cable

 Dual ring – The dual ring topology allows data to be sent in both directions.



Mesh Topology

- ❑ The mesh topology connects all devices (nodes) to each other for redundancy and fault tolerance.
- ❑ It is used in WANs to interconnect LANs and for mission critical networks like those used by banks and financial institutions.
- ❑ Implementing the mesh topology is expensive and difficult.




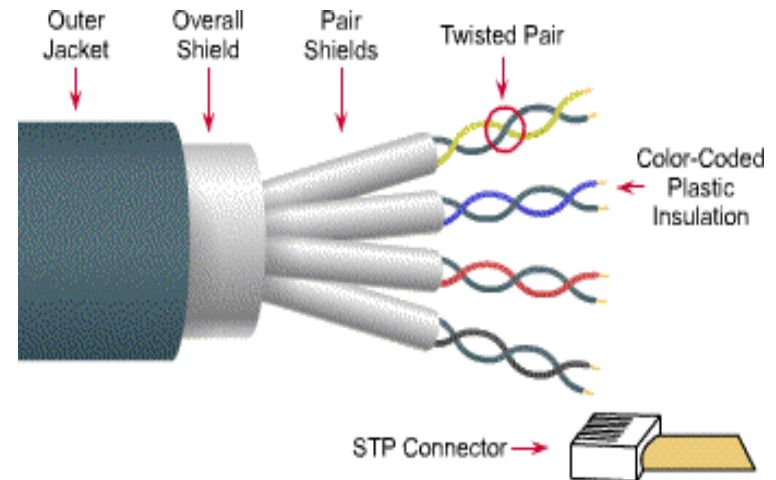
Network Components



-  **Physical Media**
-  **Interconnecting Devices**
-  **Computers**
-  **Networking Software**
-  **Applications**

Networking Media

 Networking media can be defined simply as the means by which signals (data) are sent from one computer to another (either by cable or wireless means).




- Speed and throughput: 10-100 Mbps
- Cost per node: Moderately expensive
- Media and connector size: Medium to Large
- Maximum cable length: 100m (short)


Networking Devices

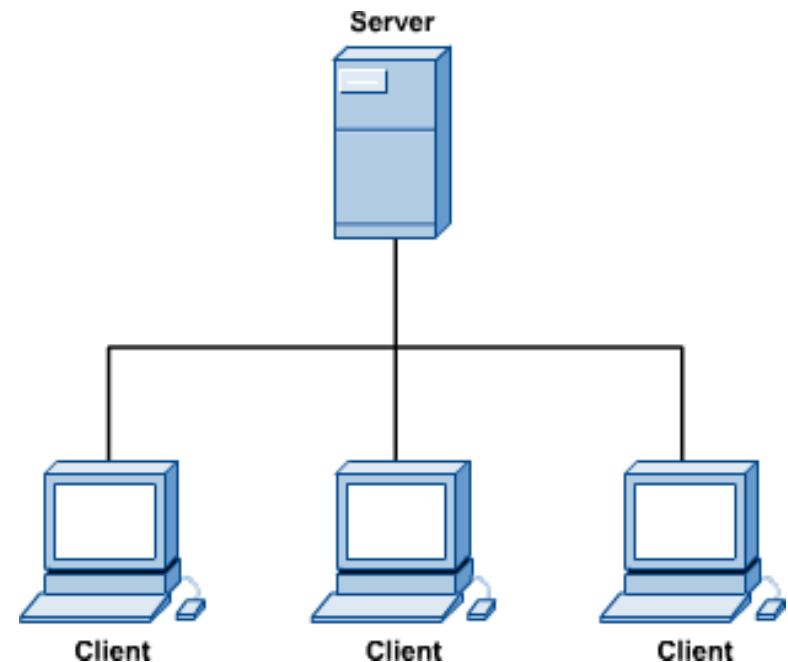
■ HUB, Switches, Routers, Wireless Access Points, Modems etc.



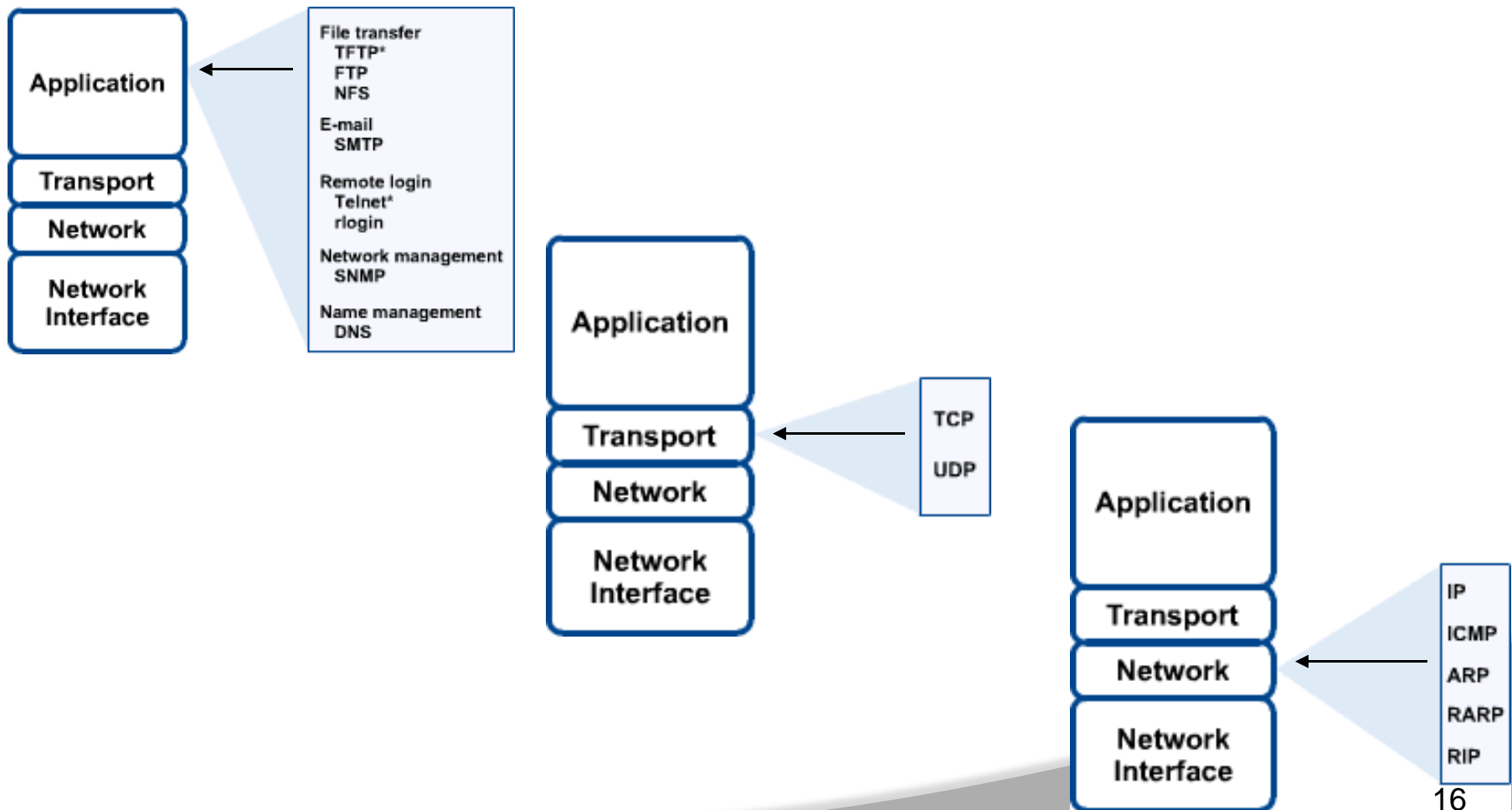
Computers: Clients and Servers

 In a client/server network arrangement, network services are located in a dedicated computer whose only function is to respond to the requests of clients.

 The server contains the file, print, application, security, and other services in a central computer that is continuously available to respond to client requests.



Networking Protocol: TCP/IP



Applications

- ❑ E-mail
- ❑ Searchable Data (Web Sites)
- ❑ E-Commerce
- ❑ News Groups
- ❑ Internet Telephony (VoIP)
- ❑ Video Conferencing
- ❑ Chat Groups
- ❑ Instant Messengers
- ❑ Internet Radio



Networking

Computer network

A collection of computing devices connected in order to communicate and share resources

Connections between computing devices can be physical using wires or cables or wireless using radio waves or infrared signals

Can you name some of the devices in a computer network?

Networking

Node (host)

Any device on a network

Data transfer rate (bandwidth)

The speed with which data is moved from one place to another on a network

Why is bandwidth so key?

Networking

Computer networks have opened up an entire frontier in the world of computing called the **client/server model**

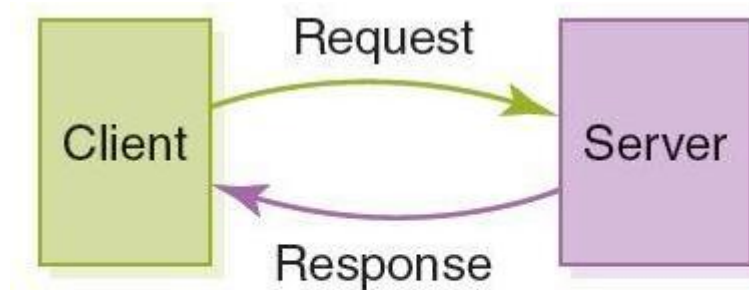


FIGURE 15.1 Client/server interaction

Networking

Protocol

A set of rules that defines how data is formatted and processed on a network

File server

A computer dedicated to storing and managing files for network users

Web server

A computer dedicated to responding to requests for web pages

P2P model

A decentralized approach that shares resources and responsibilities among many “peer” computers

Types of Networks

Local-area network (LAN)

A network that connects a relatively small number of machines in a relatively close geographical area

Ring topology connects all nodes in a closed loop on which messages travel in one direction

Star topology centers around one node to which all others are connected and through which all messages are sent

Bus topology nodes are connected to a single communication line that carries messages in both directions

Types of Networks

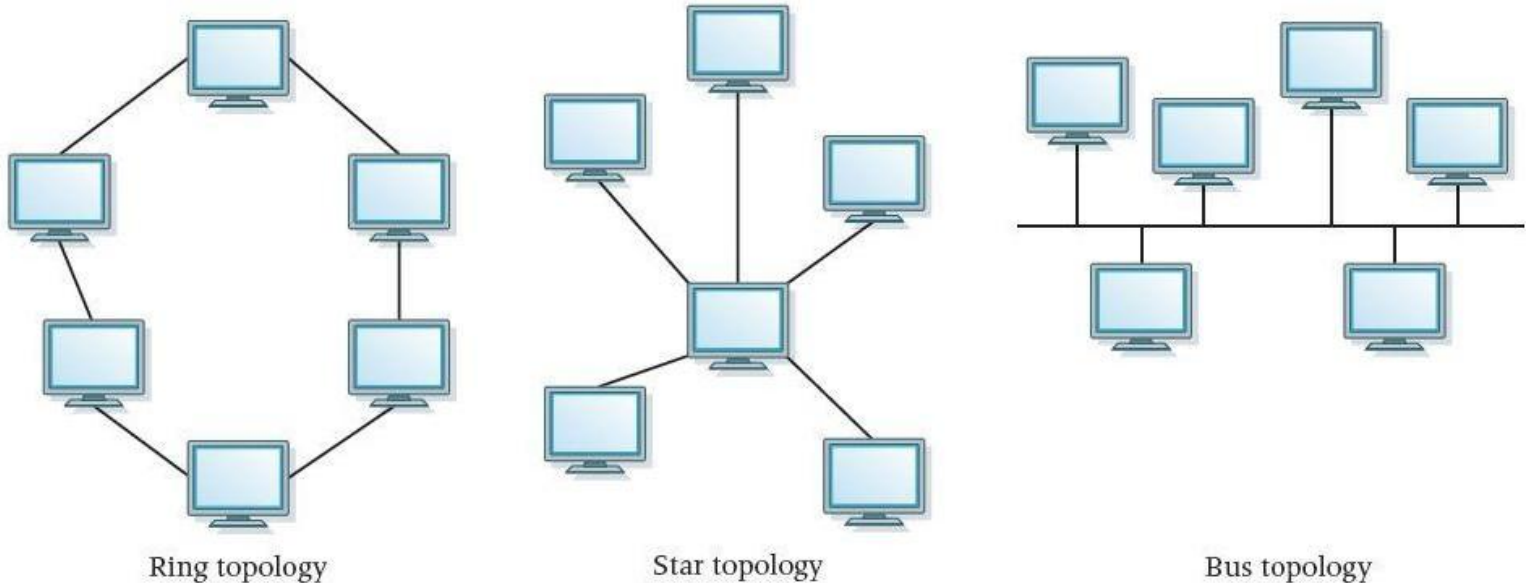


FIGURE 15.2 Network topologies

Ethernet

The industry standard bus technology for local-area networks

Types of Networks

Wide-area network (WAN)

A network that connects local-area networks over a potentially large geographic distance

Metropolitan-area network (MAN)

The communication infrastructures that have been developed in and around large cities

Gateway

One particular set up to handle all communication going between that LAN and other networks

Types of Networks

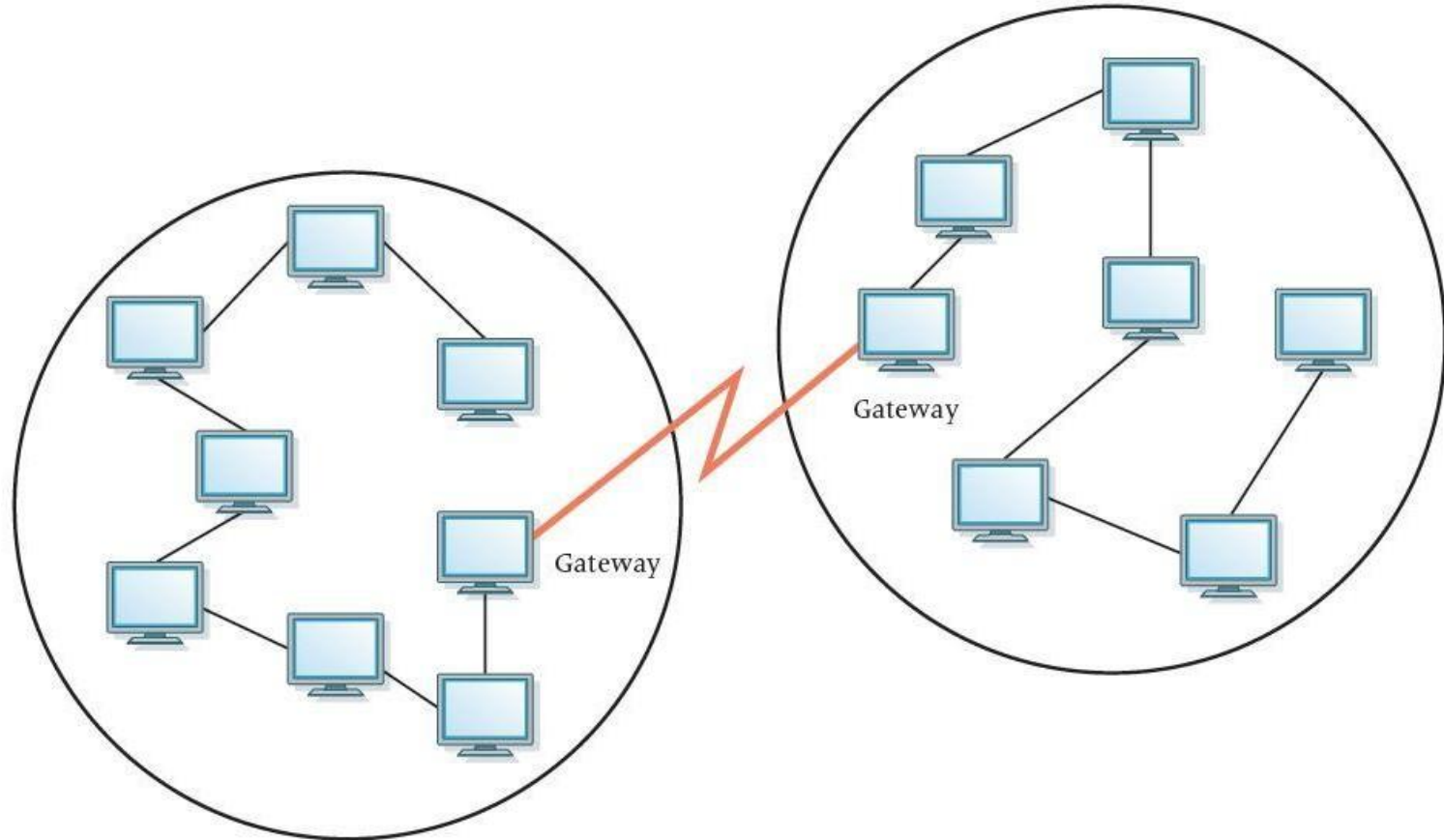


FIGURE 15.3 Local-area networks connected across a distance to create a wide-area network

Types of Networks

Internet

A wide area network that spans the planet

So, who owns the Internet?

Internet Connections

Wireless network

A network in which devices communicate with other nodes through a wireless access point

Bluetooth

A technology used for wireless communication over short distances

Internet Connections

Internet backbone

A set of high-speed networks that carry Internet traffic, provided by companies such as AT&T, Verizon, GTE, British Telecom, and IBM

Internet service provider (ISP)

An organization providing access to the Internet

Internet Connections

Various technologies available to connect a home computer to the Internet

Phone modem converts computer data into an analog audio signal for transfer over a telephone line, and then a modem at the destination converts it back again into data

Digital subscriber line (DSL) uses regular copper phone lines to transfer digital data to and from the phone company's central office

Cable modem uses the same line that your cable TV signals come in on to transfer the data back and forth

Internet Connections

Broadband

A connection in which transfer speeds are faster than 768 kilobits per second

- DSL connections and cable modems are broadband connections
- The speed for **downloads** (getting data from the Internet to your home computer) may not be the same as **uploads** (sending data from your home computer to the Internet)

Packet Switching

Packet

A unit of data sent across a network

Router

A network device that directs a packet between networks toward its final destination

Packet switching

Messages are divided into fixed-sized, numbered packets; packets are individually routed to their destination, then reassembled

Packet Switching

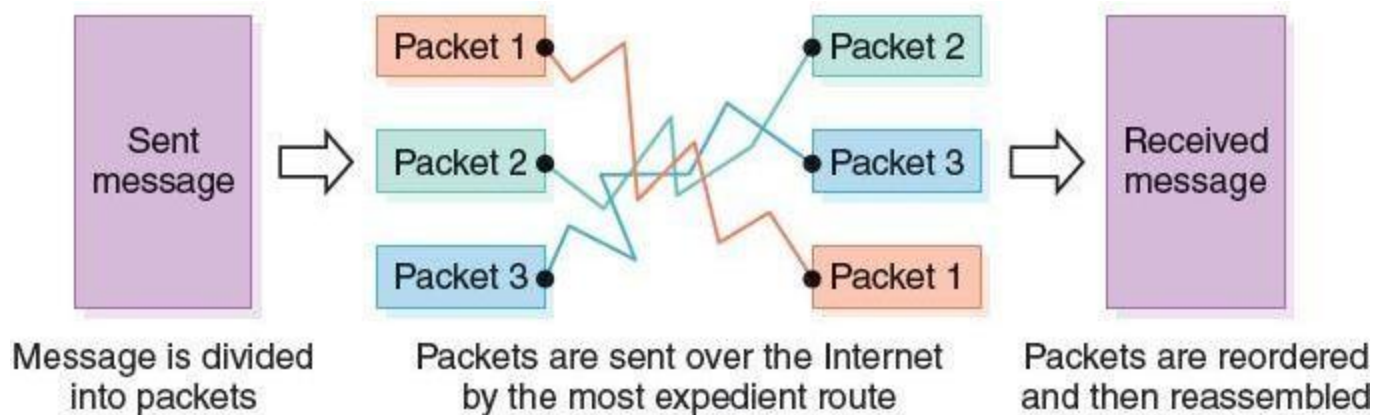


FIGURE 15.4 Messages sent by packet switching

Take a message, break it into three packets, and simulate this process

Open Systems

A logical progression...

Proprietary system

A system that uses technologies kept private by a particular commercial vendor

Interoperability

The ability of software and hardware on multiple machines and from multiple commercial vendors to communicate

Open systems

Systems based on a common model of network architecture and a suite of protocols used in its implementation

Open Systems

Number	Layer
7	Application layer
6	Presentation layer
5	Session layer
4	Transport layer
3	Network layer
2	Data Link layer
1	Physical layer

FIGURE 15.5 The layers of the OSI Reference Model

Open Systems Interconnection Reference Model

A seven-layer logical break down of network interaction to facilitate communication standards

Each layer deals with a particular aspect of network communication

Network Protocols

- Network protocols are layered such that each one relies on the protocols that underlie it
- Sometimes referred to as a **protocol stack**

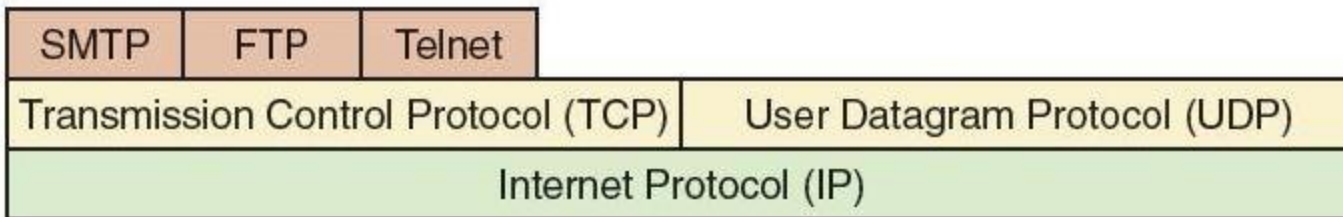


FIGURE 15.6 Layering of key network protocols

TCP/IP

Transmission Control Protocol (TCP)

Software that breaks messages into packets, hands them off to the IP software for delivery, and then orders and reassembles the packets at their destination

Internet Protocol (IP)

Software that deals with the routing of packets through the maze of interconnected networks to their final destination

TCP/IP

User Datagram Protocol (UDP)

An alternative to TCP that is faster but less reliable

Ping

A program used to test whether a particular network computer is active and reachable

Traceroute

A program that shows the route a packet takes across the Internet

Traceroute in Action

```

C:\WINDOWS\System32\cmd.exe
Microsoft Windows XP [Version 5.1.2600]
(C) Copyright 1985-2001 Microsoft Corp.

C:\Documents and Settings\UserName>tracert google.com

Tracing route to google.com [64.233.187.99]
over a maximum of 30 hops:

  1    1 ms    <1 ms    <1 ms    192.168.1.1
  2    1 ms    <1 ms    <1 ms    GATEWAY1.ORLANDO.dimenoc.com [66.193.174.11]
  3    1 ms    1 ms    1 ms    POS4-1.GW5.ORTL.ALTER.NET [63.122.161.105]
  4    1 ms    1 ms    1 ms    500.at-1-1-0.CL2.ORTL.ALTER.NET [152.63.80.102]

  5    15 ms   13 ms   13 ms    0.so-7-0-0.XL2.ATL4.ALTER.NET [152.63.86.109]
  6    13 ms   13 ms   13 ms    0.so-7-0-0.BR1.ATL4.ALTER.NET [152.63.86.173]
  7    15 ms   15 ms   14 ms    so-1-1-0.gar2.Atlanta1.Level3.net [4.68.127.177]

  8    16 ms   15 ms   15 ms    ae-21-52.car1.Atlanta1.Level3.net [4.68.103.34]

  9    14 ms   16 ms   15 ms    4.78.208.2
 10    15 ms   16 ms   15 ms    66.249.95.125
 11    16 ms   16 ms   19 ms    216.239.49.226
 12    16 ms   16 ms   16 ms    64.233.187.99

Trace complete.

C:\Documents and Settings\UserName>

```

FIGURE 15.7 The traceroute utility

Used with permission from Microsoft

High-Level Protocols

Other protocols build on TCP/IP protocol suite

Simple Mail Transfer Protocol (SMTP) used to specify transfer of electronic mail

File Transfer Protocol (FTP) allows a user to transfer files to and from another computer

Telnet used to log onto one computer from another

Hyper Text Transfer Protocol (http) allows exchange of Web documents

Which of these have you used?

High-Level Protocols

Protocol	Port
Echo	7
File Transfer Protocol (FTP)	21
Telnet	23
Simple Mail Transfer Protocol (SMTP)	25
Domain Name Service (DNS)	53
Gopher	70
Finger	79
Hypertext Transfer Protocol (HTTP)	80
Post Office Protocol (POP3)	110
Network News Transfer Protocol (NNTP)	119
Internet Relay Chat (IRC)	6667

Port

A numeric designation that corresponds to a particular high-level protocol

FIGURE 15.8 Some protocols and the ports they use

MIME Types

MIME type

A standard for defining the format of files that are included as email attachments or on websites

What does MIME stand for?

Multipurpose Internet Mail Extension

Firewalls

Firewall

A gateway machine and its software that protects a network by filtering the traffic it allows

Access control policy

A set of rules established by an organization that specifies what types of network communication are permitted and denied

*Have your messages ever been
returned undelivered, blocked by a firewall?*

Firewalls

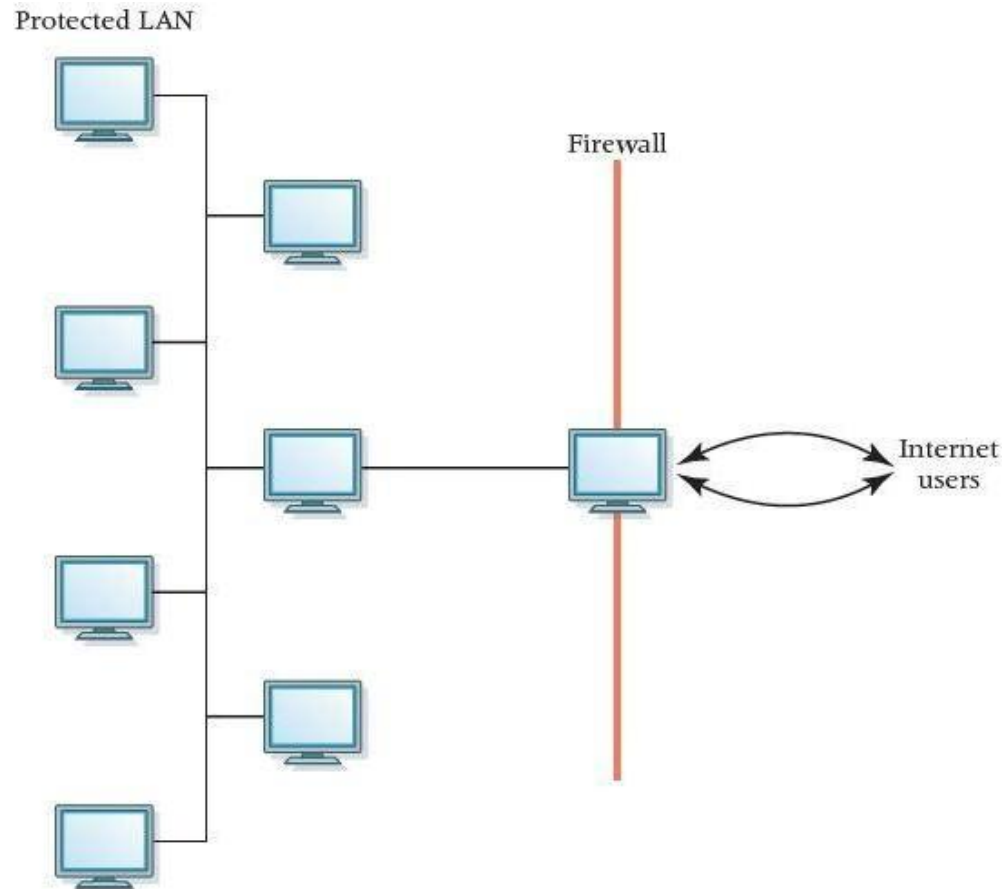


FIGURE 15.9 A firewall protecting a LAN.

Network Addresses

Hostname

A name made up of words separated by dots that uniquely identifies a computer on the Internet

IP address

An address made up of four one-byte numeric values separated by dots that uniquely identifies a computer on the Internet

Is there a correspondence between the parts of a hostname and an IP address?

Network Addresses

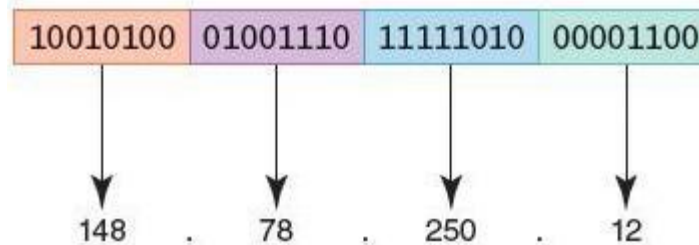


FIGURE 15.10 An IP address stored in four bytes

What is wrong with the IP4 strategy?

How did smartphones contribute to the problem?

Network Addresses

IPv4

The last block was assigned in 2011

IPv6

32 bits organized into 4 groups of 8

FE80:0000:0000:0000:0202:B3FF:FE1E:8329

They work in parallel

Domain Name System

Host number

The part of the IP address that specifies a particular host (machine) on the network *Yes, but what is it?*

Domain name

The part of a hostname that specifies a specific organization or group

Top-level domain (TLD)

The last section of a domain name that specifies the type of organization or its country of origin

Domain Name System

Domain name system (DNS)

A distributed system for managing hostname resolution

Domain name server

A computer that attempts to translate a hostname into an IP address

Domain Squatting

Ransoming domain names

Should the tables containing hostname/IP mappings be sorted or unsorted? Why?

Domain Name System

Top-Level Domain	General Purpose
.aero	Aerospace industry
.biz	Business
.com*	U.S. commercial (unrestricted)
.coop	Cooperative
.edu*	U.S. educational
.gov*	U.S. government
.info	Information (unrestricted)
.int*	International organizations
.jobs	Employment
.mil*	U.S. military
.museum	Museums
.name	Individuals and families
.net*	Network (unrestricted)
.org*	Nonprofit organization (unrestricted)
.pro	Certain professions

FIGURE 15.11 Some top-level domains and their general purpose (* indicates an original TLD)

Domain Name System

Organizations based in countries other than the United States use a top-level domain that corresponds to their two-letter country codes

Country Code TLD	Country
.au	Australia
.br	Brazil
.ca	Canada
.gr	Greece
.in	India
.ru	Russian Federation
.uk	United Kingdom

FIGURE 15.12 Some of the top-level domain names based on country codes.

*Have you
emailed
someone
in another
country?*

Domain Name System

social	furniture	dental	paris	media
career	town	rocks	cooking	rodeo
nyc	trade	webcam	vote	actor
vacations	industries	wiki	productions	flights
rentals	catering	dating	bargains	cool
pics	guitars	tax	dance	email
farm	education	ninja	coffee	shoes
menu	kitchen	land	support	associates
institute	camp	center	directory	florist

A very small, random selection of new TLDs that are available as of mid-2014

Who Controls the Internet?

Control of IP addresses and domain names

- Internet began as ARPANET, a project of the US Dept. of Defense
- Control subcontracted to ICANN in 1998
- US gov't to further reduce role as early as 2015

FCC proposal

- Would allow ISPs to provide “premium” access to certain customers, perhaps by deliberately slowing down data transfer for others
- **Net neutrality** - The principle that ISPs should deliver data to everyone equally, as fast as the technology allows

Cloud Computing

- **Public clouds** are accessible by any subscriber
- **Private clouds** are established for a specific group or organization
- **Community clouds** are shared among two or more organizations with the same needs
- **Hybrid clouds** are some combination of the others

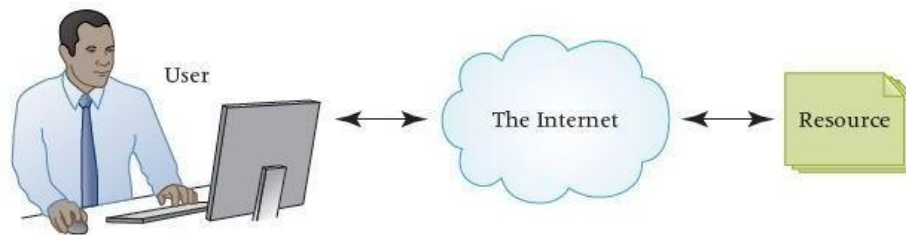


FIGURE 15.13 Internet communication depicted using a cloud

Ethical Issues

Effects of Social Networking

What are some examples of popular social networking sites?

Who uses social networking?

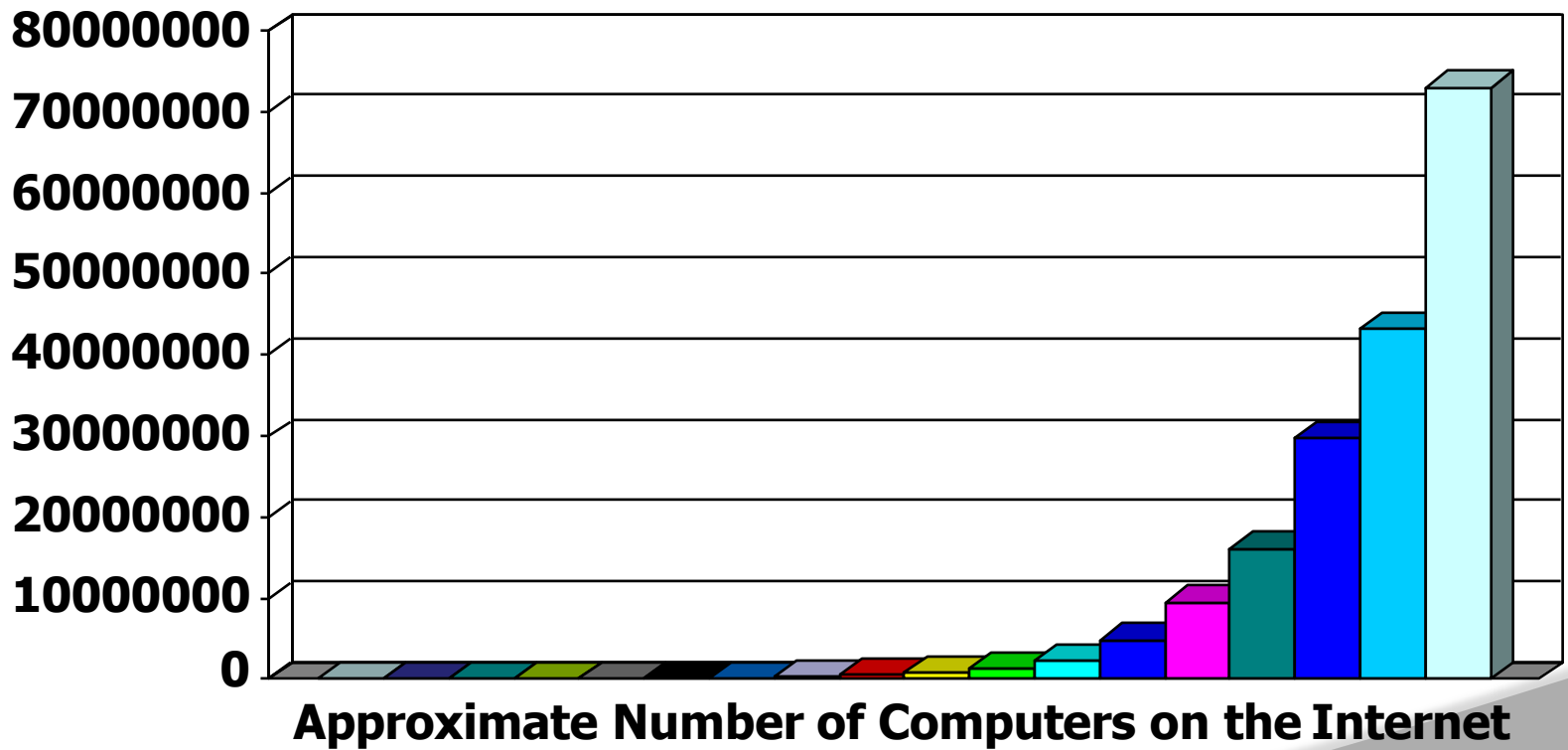
*What are the **benefits** and the **disadvantages** of using these social networking sites?*

Do the benefits of social networking out weigh the potential costs?

Internet Impact

- Check weather
- Buy goods
- Play music
- Find the shortest route
- Give a lecture...

The Incredible Growth



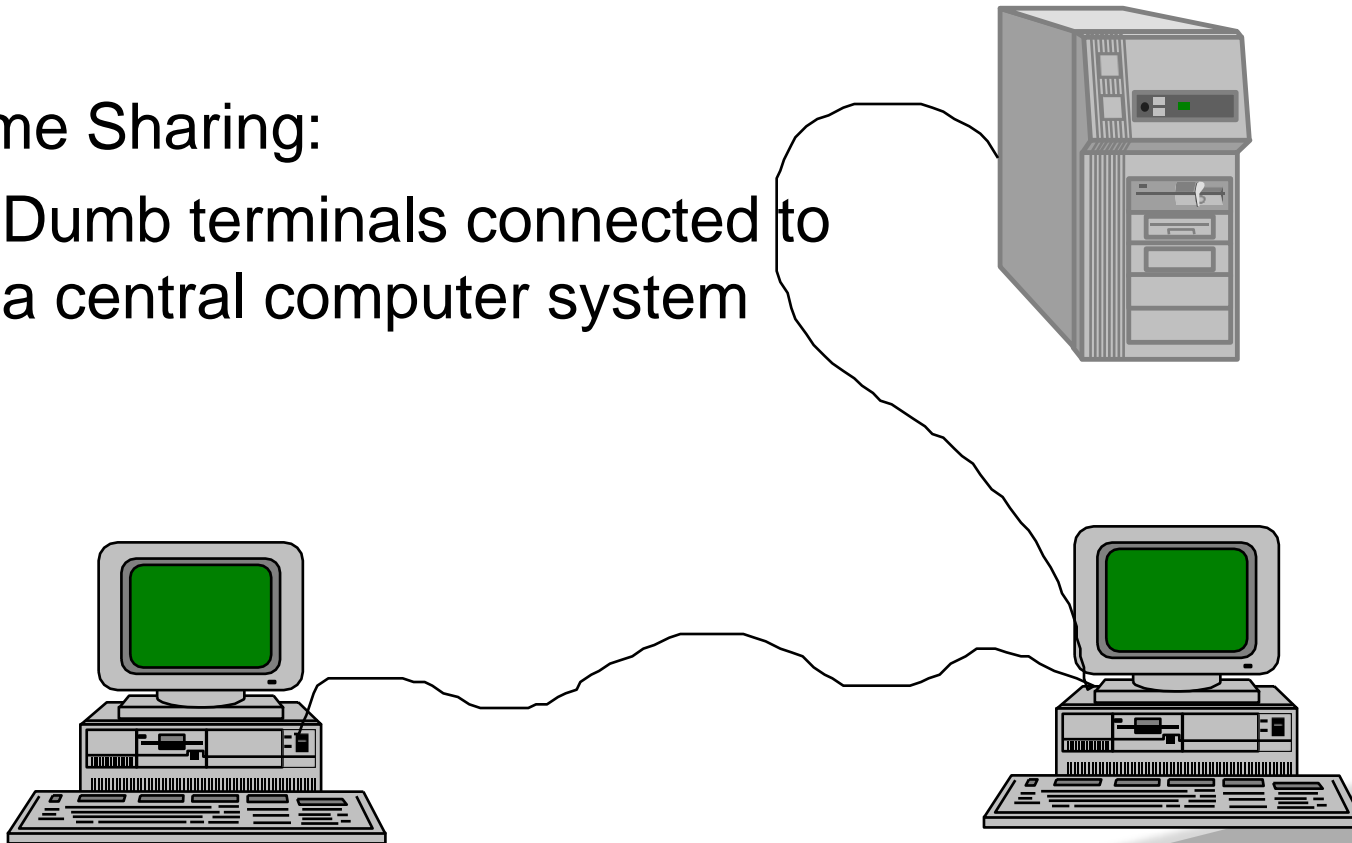
Brief Internet History

- Batch Environment - 1950s
 - No direct interaction between users and their programs during execution
- Time Sharing - 1960s
 - Users were able to interact with the computer and could share its information processing resources
 - Marked the beginning of computer communications

Brief Internet History (cont.)

Time Sharing:

Dumb terminals connected to
a central computer system



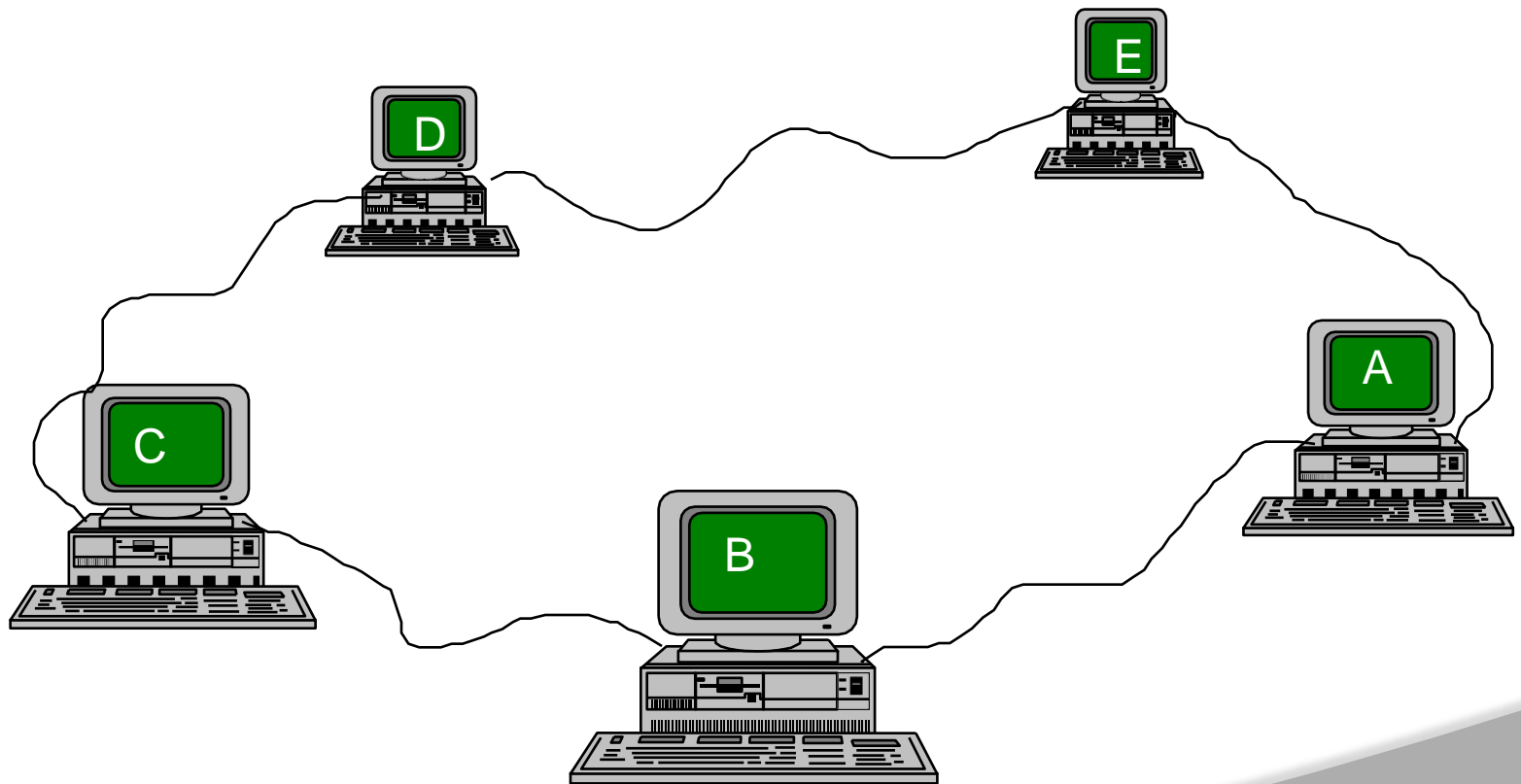
Brief Internet History (cont.)

- Late 1960s: ARPANET
 - ARPA (Advanced Research Projects Agency) commissioned an experimental computer network
- 1970s:
 - Distributed Processing: minicomputers;
 - Communication between neighbor processors and applications via networks
 - Growth of ARPANET and Invention of Email

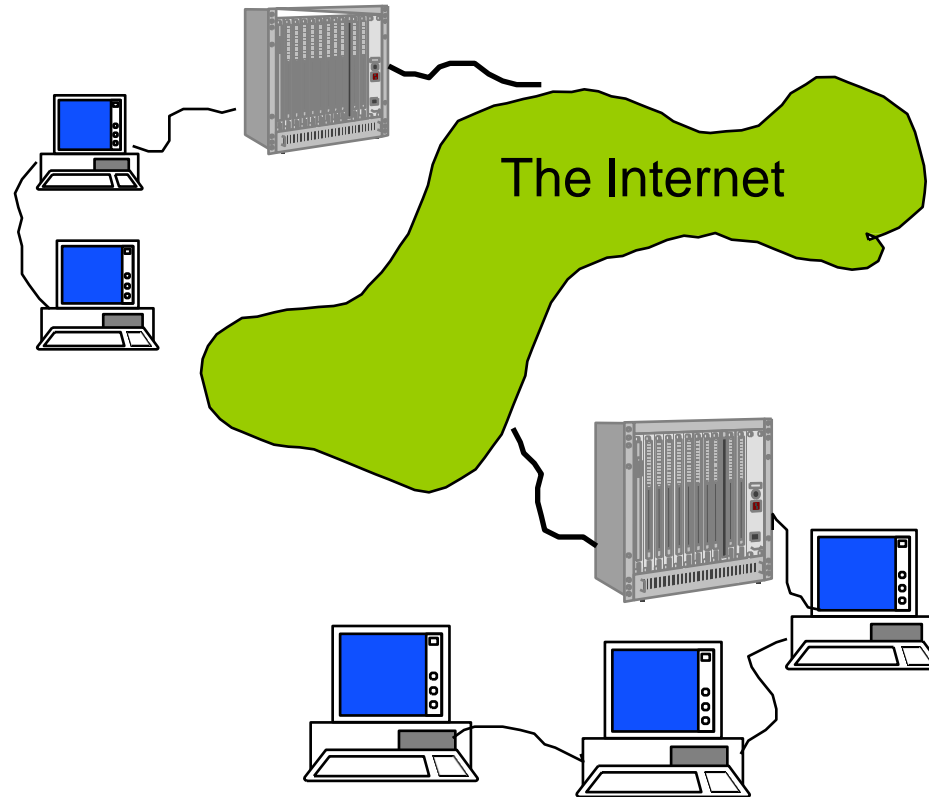
Brief Internet History (cont.)

- 1980s:
 - WAN and LAN
 - Prototype Internet
 - TCP/IP: Allows different networks to interconnect

A LAN Example



Internet: a network of networks



Brief Internet History (cont.)

- 1990s: WWW
 - HTTP and HTML
 - Marc Andreessen: Mosaic (all-in-one solution)
 - Commercial traffic allowed ECommerce

Networking Questions

- Over what medium?
- At what speed?
- How to address computers?
- Which path?
- How to handle (detect & recover) errors?
- What services?
- How to address documents?
- What data format?

Outline

- Introduction
- OSI Model
- TCP/IP Model
- IPv4 vs. IPv6

What is a Protocol?

- A standard that allows entities (i.e. application programs) from different systems to communicate
- Shared conventions for communicating information
- Includes syntax, semantics, and timing

Standardized Protocol Architectures

- Vendors like standards because they make their products more marketable
- Customers like standards because they enable products from different vendors to interoperate
- Two protocol standards are well-known:
 - TCP/IP: widely implemented
 - OSI: less used, still useful for modeling/conceptualizing

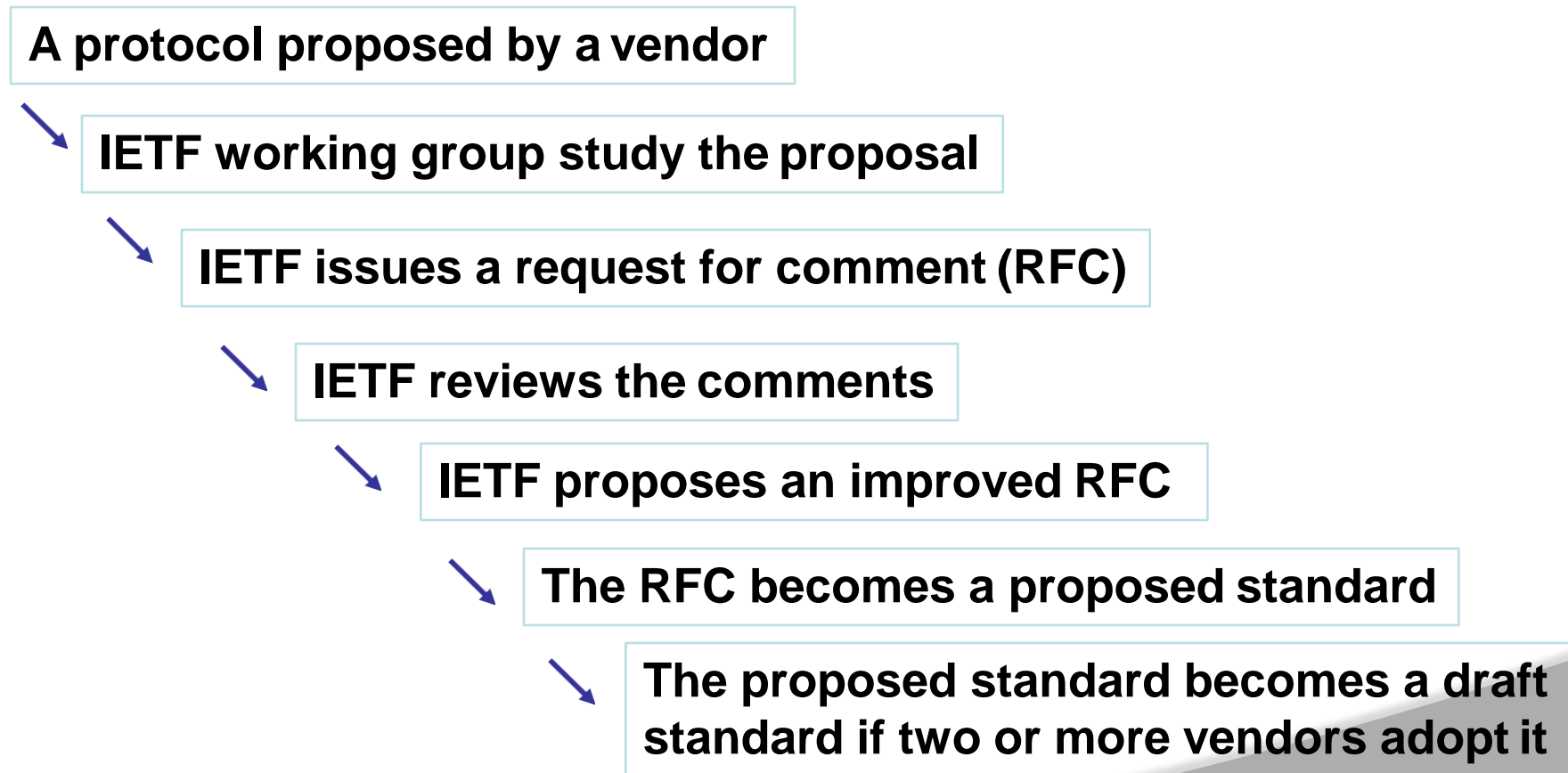
Internet Standards

- Email related standards
 - IMAP, POP, X.400, SMTP, CMC, MIME, binhex, uuencode
- Web related standards
 - http, CGI, html/xml/vrml/sgml
- Internet directory standards
 - X.500, LDAP
- Application standards
 - http, FTP, telnet, gopher, wais
- Videoconferencing standards
 - H.320, H.323, Mpeg-1, Mpeg-2

Telecommunication Standards Organizations

- International Telecommunications Union - Telecommunication Standardization Sector (**ITU-TSS**). Formerly called the Consultative Committee on International Telegraph and Telephone (**CCITT**)
- International Organization for Standards (**ISO**). Member of the ITU, makes technical recommendations about data communications interfaces.
- American National Standards Institute (ANSI)
- Institute of Electrical and Electronics Engineers (**IEEE**)
- Internet Engineering Task Force (**IETF**)
- Electronic Industries Association (EIA)
- National Institute of Standards and Technology (NIST)
- National Exchange Carriers Association (NECA)
- Corporation for Open Systems (COS)
- Electronic Data Interchange -(EDI) of Electronic Data Interchange for Administration Commerce and Transport (EDIFACT).

***Internet Engineering Task Force**



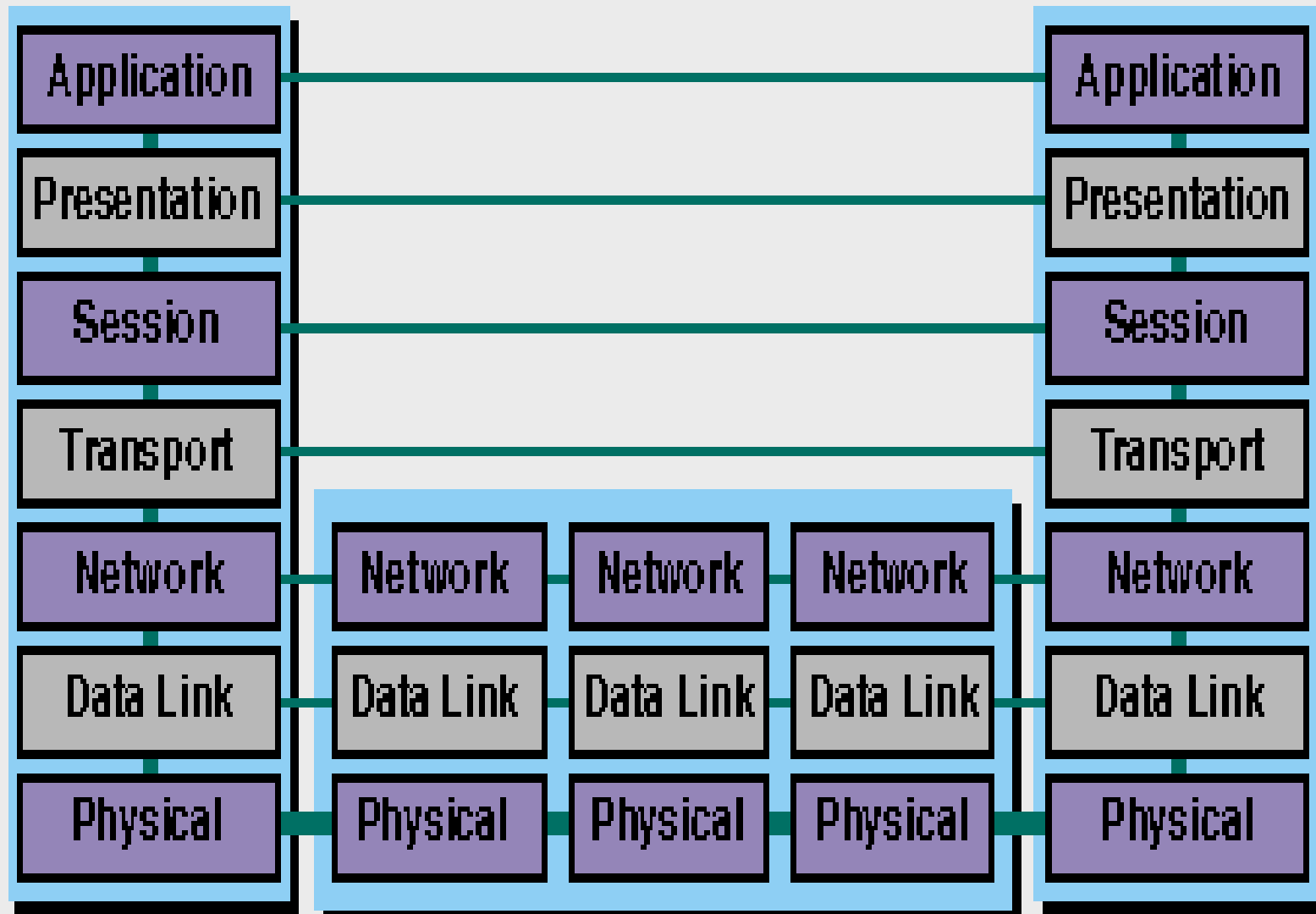
What is OSI?

- Developed by the International Organization for Standardization (ISO) in 1984
- The primary architectural model for intercomputer communications.
- A conceptual model composed of seven layers, each specifying particular network functions.
- Describes how information from a software application in one computer moves through a network medium to a software application in another computer.

Why Study OSI?

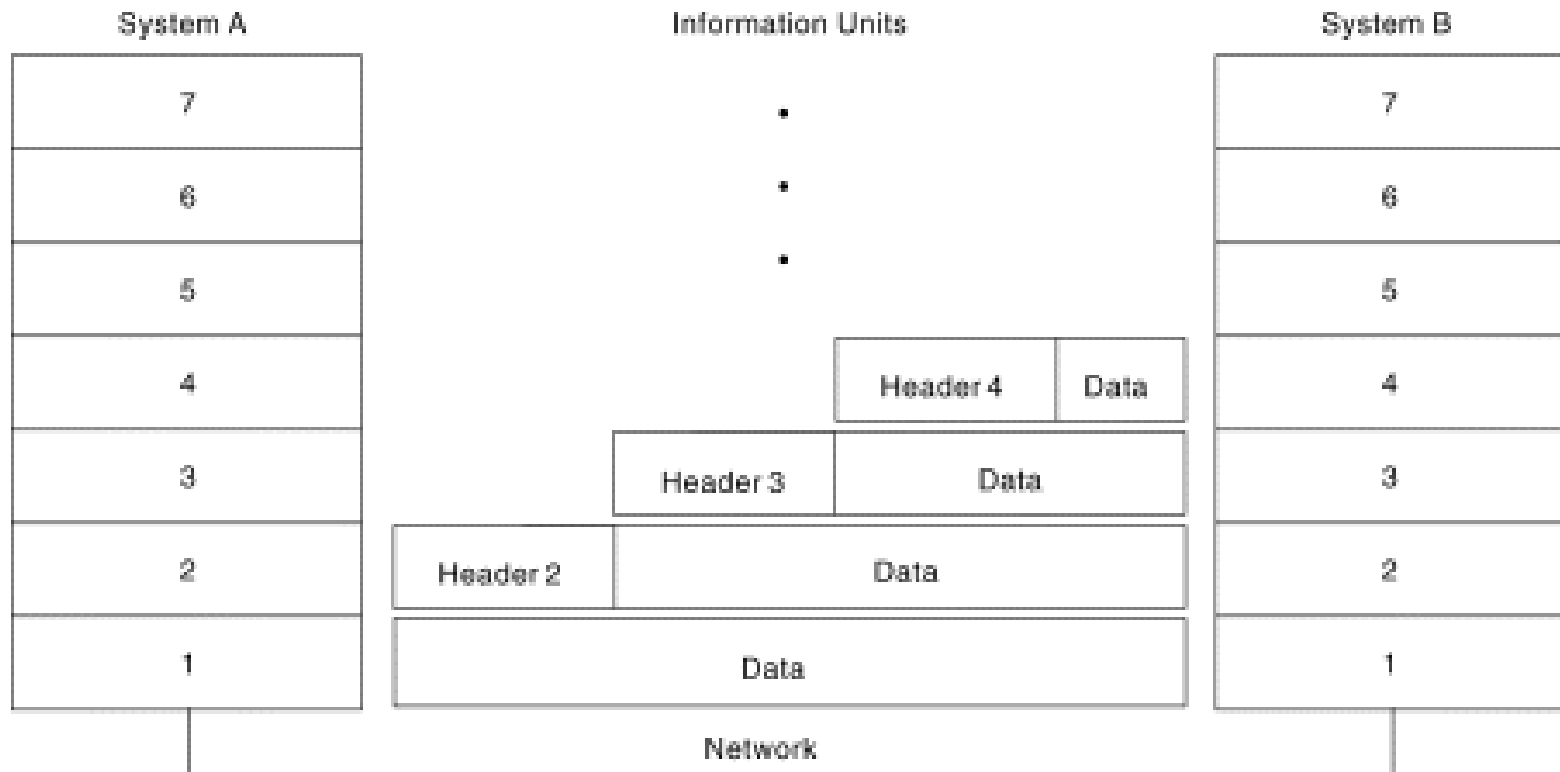
- Still an excellent model for conceptualizing and understanding protocol architectures
- Key points:
 - Modular
 - Hierarchical
 - Boundaries between layers=interfaces

OSI Stack



Intermediate nodes in a network
service (example AT&T APS)

Headers and Data



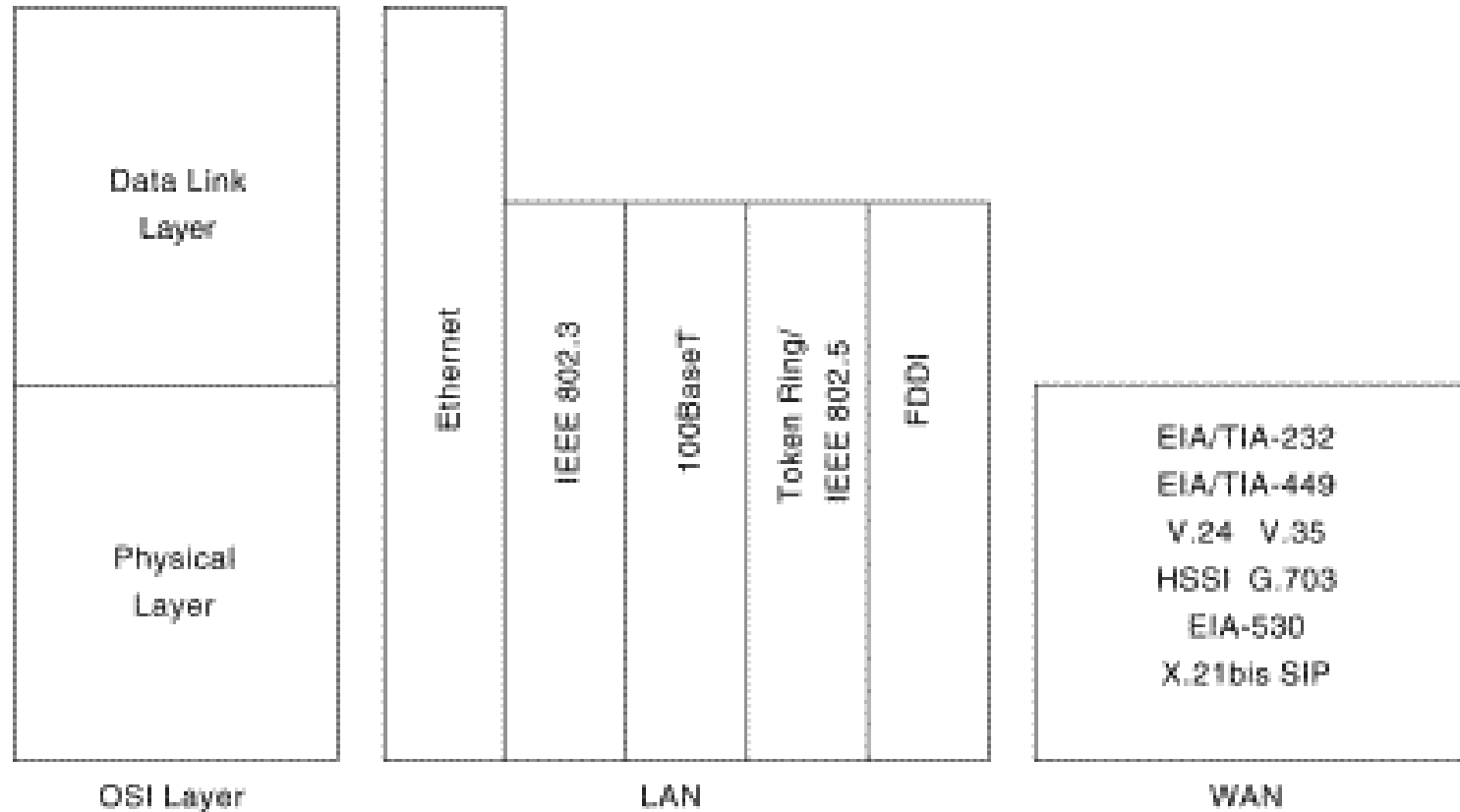
OSI Lower Layers

- Physical – Layer 1
- Data Link – Layer 2
- Network – Layer 3

OSI Physical Layer

- Responsible for transmission of bits
- Always implemented through hardware
- Encompasses mechanical, electrical, and functional interfaces
- e.g. RS-232

*Physical-layer Implementation

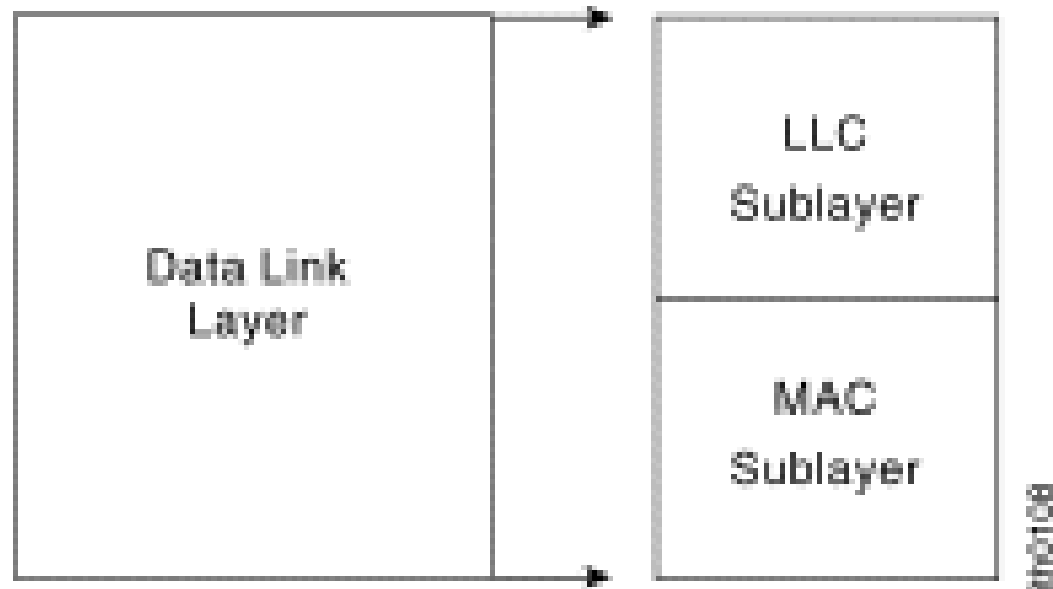


Physical Layer Implementations

OSI Data Link Layer

- Responsible for error-free, reliable transmission of data
- Flow control, error correction
- e.g. HDLC

OSI Data Link Layer



IEEE has subdivided data link layer into two sub-layers.

OSI Network Layer

- Responsible for routing of messages through network
- Concerned with type of switching used (circuit v. packet)
- Handles routing between networks, as well as through packet-switching networks

Network Access Layer

- Concerned with exchange of data between computer and network
- Includes addressing, routing, prioritizing, etc
- Different networks require different software at this layer
- Example: X.25 standard for network access procedures on packet-switching networks

OSI Upper Layers

- Transport
- Session
- Presentation
- Application

OSI Transport Layer

- Isolates messages from lower and upper layers
- Breaks down message size
- Monitors quality of communications channel
- Selects most efficient communication service necessary for a given transmission

Transport Layer

- Concerned with reliable transfer of information between applications
- Independent of the nature of the application
- Includes aspects like flow control and error checking

OSI Session Layer

- Establishes logical connections between systems
- Manages log-ons, password exchange, log-offs
- Terminates connection at end of session

OSI Presentation Layer

- Provides format and code conversion services
- Examples
 - File conversion from ASCII to EBDIC
 - Invoking character sequences to generate bold, italics, etc on a printer

OSI Application Layer

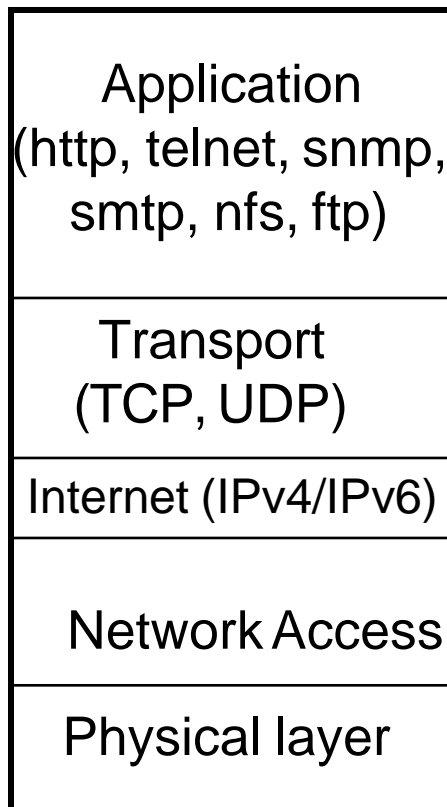
- Provides access to network for end-user
- User's capabilities are determined by what items are available on this layer
- Logic needed to support various applications
- Each type of application (file transfer, remote access) requires different software on this layer

TCP/IP

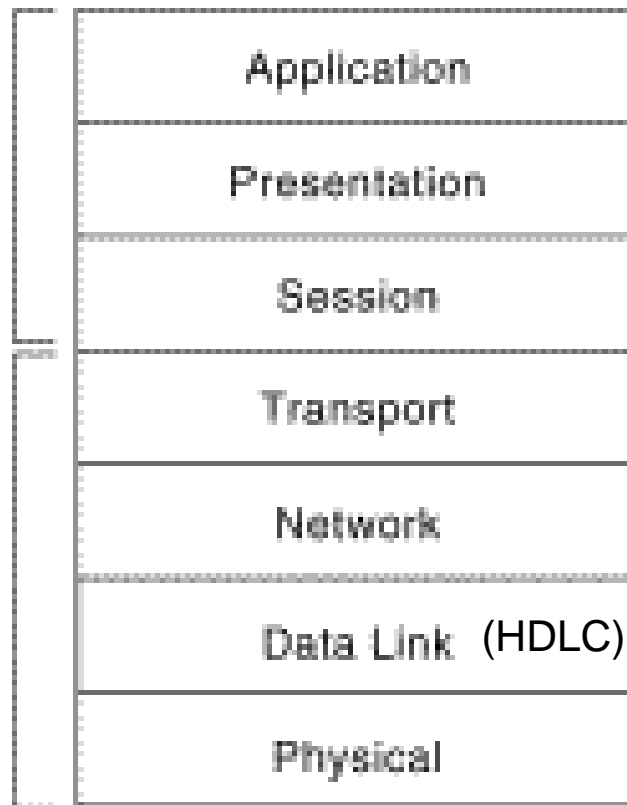
- Transmission control Protocol/Internet Protocol
- Developed by DARPA
- No official protocol standard
- Can identify five layers
 - Application
 - Host-to-Host (transport)
 - Internet
 - Network Access
 - Physical

An OSI View of TCP/IP

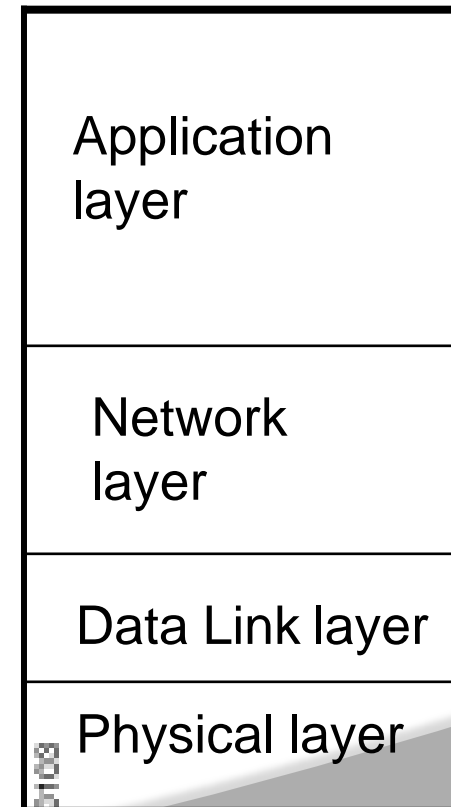
Internet Model



OSI Model

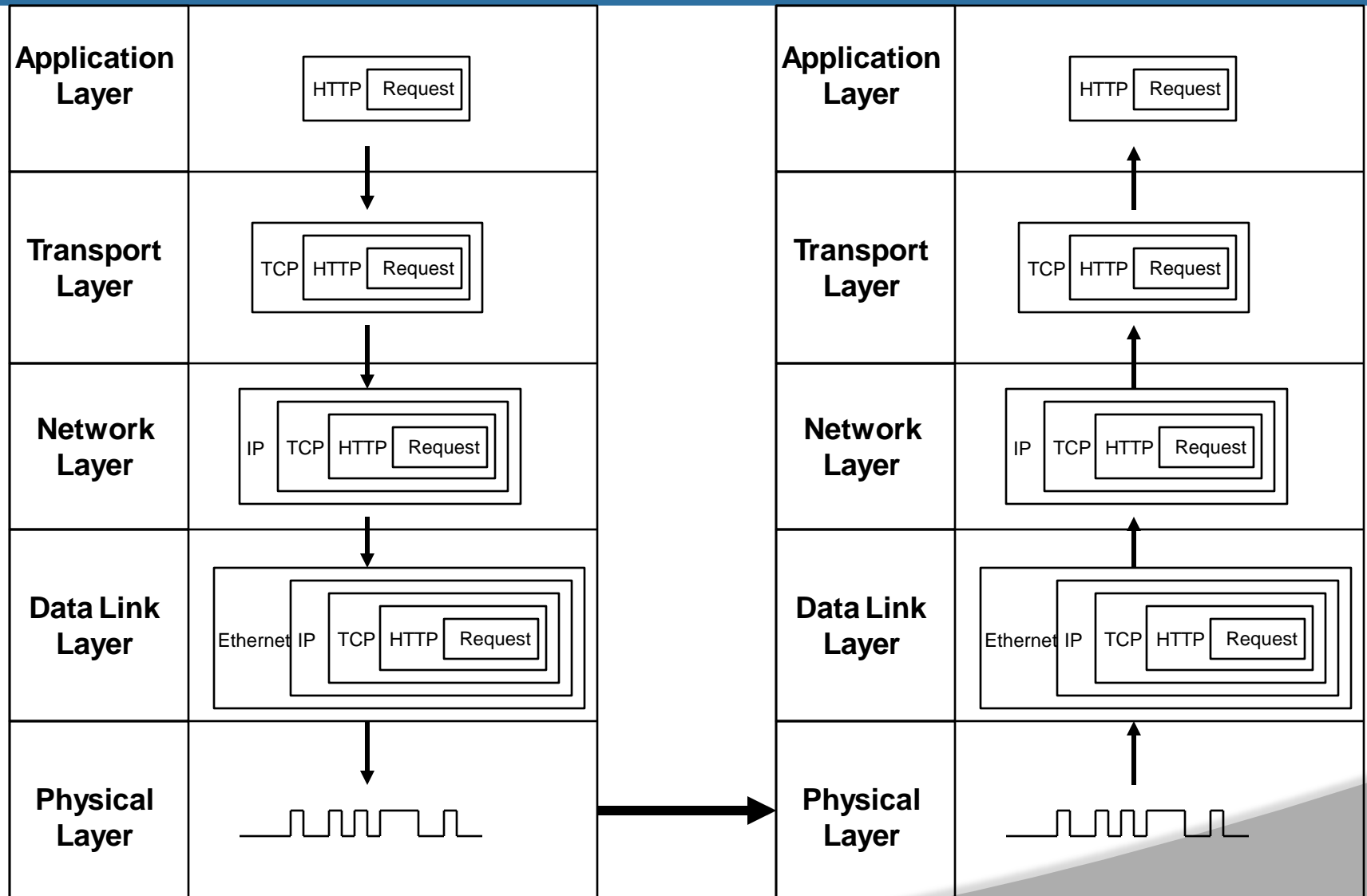


F-D's Model



Sender

Receiver



TCP/IP Network Access Layer

- Exchange of data between end system and network
- Address of host and destination
- Prioritization of transmission
- Software at this layer depends on network (e.g. X.25 vs. Ethernet)
- Segregation means that no other software needs to be concerned about net specifics

TCP/IP Internet Layer

- An Internet is an interconnection of two or more networks
- Internet layer handles tasks similar to network access layer, but between networks rather than between nodes on a network
- Uses IP for addressing and routing across networks
- Implemented in workstations *and* routers

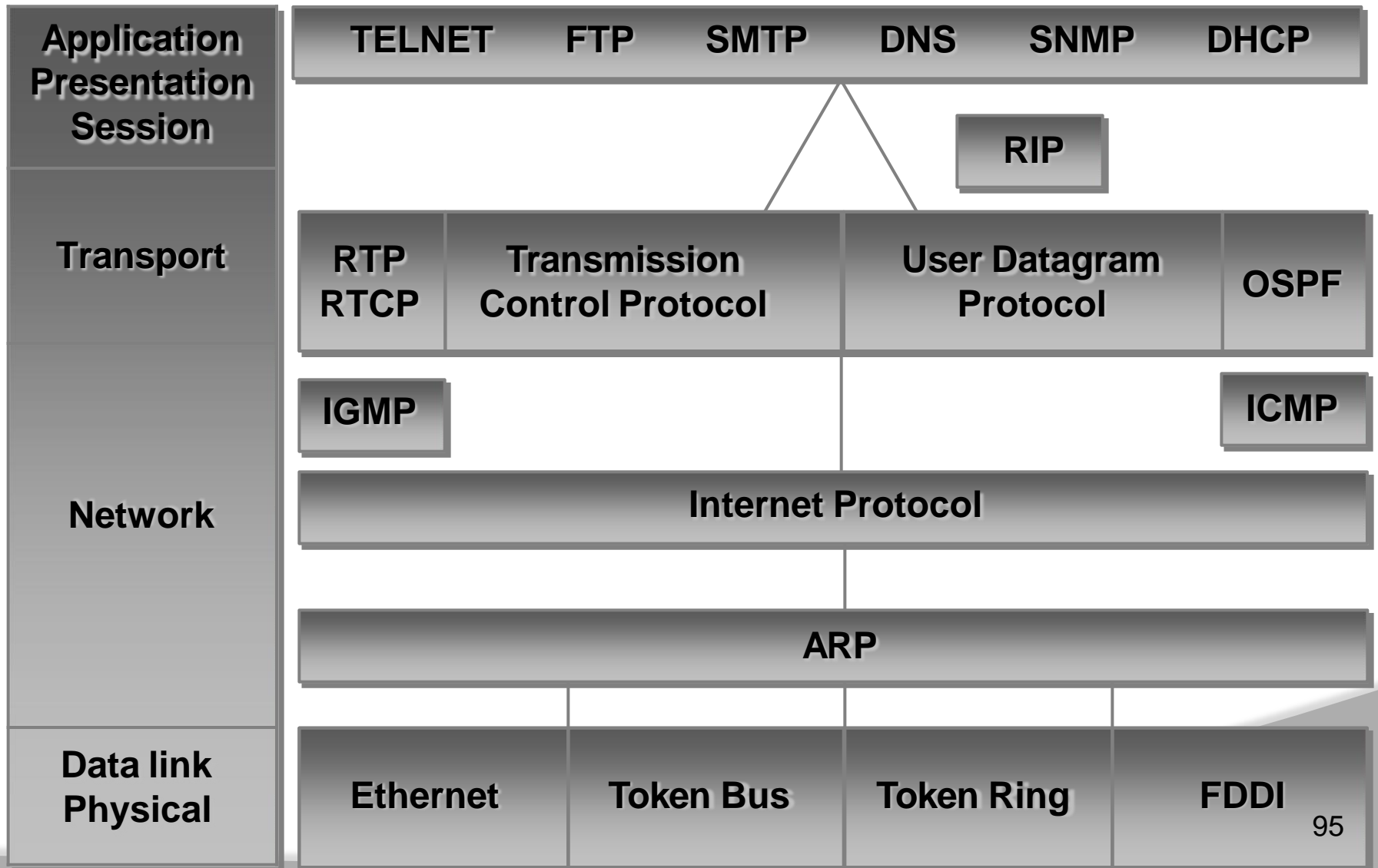
TCP/IP Transport Layer

- Also called host-to-host layer
- Reliable exchange of data between applications
- Uses TCP protocols for transmission

TCP/IP Application Layer

- Logic needed to support variety of applications
- Separate module supports each type of application (e.g. file transfer)
 - FTP
 - HTTP
 - Telnet
 - News
 - SMTP

*TCP/IP



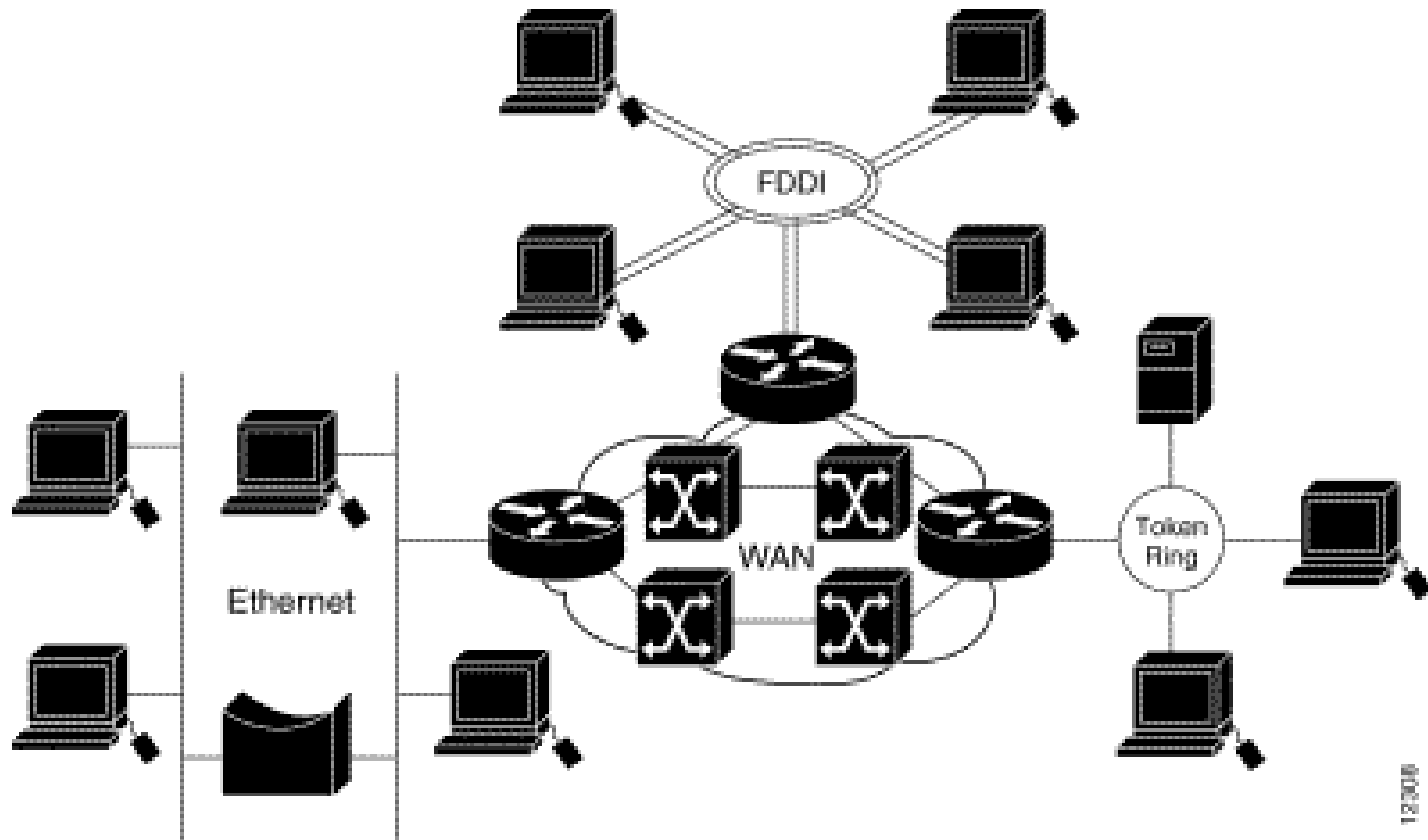
TCP & UDP

- Most TCP/IP applications use TCP for transport layer
- TCP provides a connection (logical association) between two entities to regulate flow check errors
- UDP (User Datagram Protocol) does not maintain a connection, and therefore does not guarantee delivery, preserve sequences, or protect against duplication

Internetworking

- Interconnected networks, usually implies TCP/IP
- Can appear to users as a single large network
- The global Internet is the largest example, but intranets and extranets are also examples

Internetworking



00021

TCP Segment (TCP PDU)

- Source port (16 bits)
- Destination port (16 bits)
- Sequence number (32 bits)
- Acknowledgment number (32 bits)
- Data Offset (4 bits)
- Reserved (6 bits)
- Flags (6 bits) : URG, ACK, PSH, RST, SYN, FIN
- Window (16 bits)
- Checksum (16 bits)
- Urgent Pointer (16 bits)
- Options (variable)

The size of TCP header is 192 bits = 24 bytes.

IPv4 and IPv6

- IP (IPv4) provides for 32-bit source and destination addresses, using a 192-bit header
- IPv6 (1996 standard) provides for 128-bit addresses, using a 320-bit header.
- Migration to IPv6 will be a very slow process

Attenuation

- Means loss of energy -> weaker signal
- When a signal travels through a medium it loses energy overcoming the resistance of the medium
- Amplifiers are used to compensate for this loss of energy by amplifying the signal.

Measurement of Attenuation

- To show the loss or gain of energy the unit “decibel” is used.

$$\text{dB} = 10\log_{10} P_2/P_1$$

P_1 - input signal

P_2 - output signal

Distortion

- Means that the signal changes its form or shape
- Distortion occurs in **composite** signals
- Each frequency component has its own **propagation speed** traveling through a medium.
- The different components therefore arrive with **different delays** at the receiver.
- That means that the signals have **different phases** at the receiver than they did at the source.

Noise

- There are different types of noise
 - **Thermal** - random noise of electrons in the wire creates an extra signal
 - **Induced** - from motors and appliances, devices act as transmitter antenna and medium as receiving antenna.
 - **Crosstalk** - same as above but between two wires.
 - **Impulse** - Spikes that result from power lines, lightning, etc.

Signal to Noise Ratio (SNR)

- To measure the quality of a system the SNR is often used. It indicates the strength of the signal wrt the noise power in the system.
- It is the ratio between two powers.
- It is usually given in dB and referred to as SNR_{dB} .



Example

3.31

The power of a signal is 10 mW and the power of the noise is 1 μ W; what are the values of SNR and SNR_{dB} ?

Solution

The values of SNR and SNR_{dB} can be calculated as follows:

$$\text{SNR} = \frac{10,000 \mu\text{W}}{1 \text{ mW}} = 10,000$$

$$\text{SNR}_{\text{dB}} = 10 \log_{10} 10,000 = 10 \log_{10} 10^4 = 40$$



Example

3.32

The values of SNR and SNR_{dB} for a noiseless channel are

$$\text{SNR} = \frac{\text{signal power}}{0} = \infty$$
$$\text{SNR}_{\text{dB}} = 10 \log_{10} \infty = \infty$$

We can never achieve this ratio in real life; it is an ideal.

Transmission Media

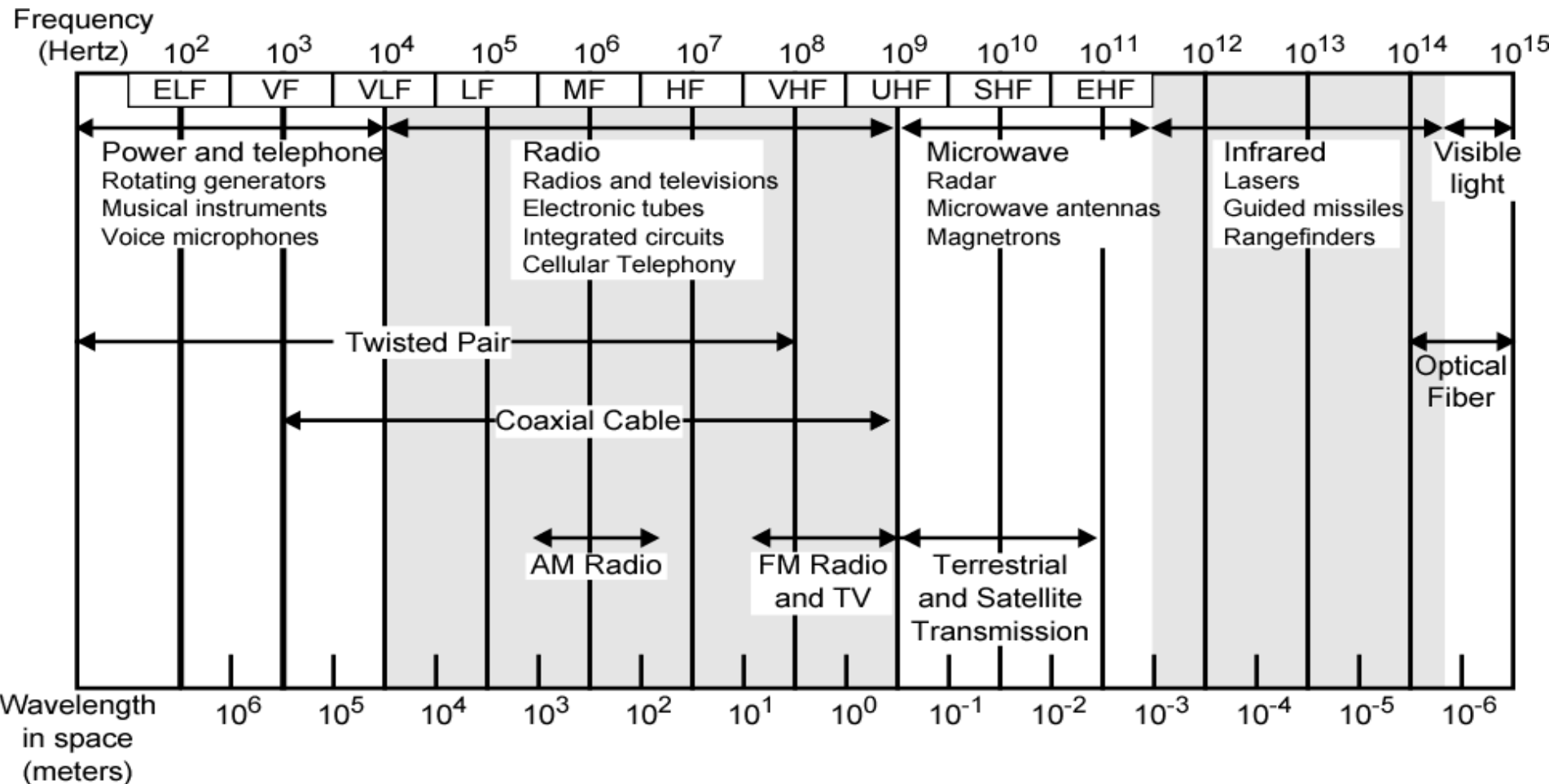
Overview

- Guided - wire
- Unguided - wireless
- Characteristics and quality determined by medium and signal
- For guided, the medium is more important
- For unguided, the bandwidth produced by the antenna is more important
- Key concerns are data rate and distance

Design Factors

- Bandwidth
 - Higher bandwidth gives higher data rate
- Transmission impairments
 - Attenuation
- Interference
- Number of receivers
 - In guided media
 - More receivers (multi-point) introduce more attenuation

Electromagnetic Spectrum



ELF = Extremely low frequency

VF = Voice frequency

VLF = Very low frequency

LF = Low frequency

MF = Medium frequency

HF = High frequency

VHF = Very high frequency

UHF = Ultrahigh frequency

SHF = Superhigh frequency

EHF = Extremely high frequency

Guided Transmission Media

- Twisted Pair
- Coaxial cable
- Optical fiber

Twisted Pair

- Separately insulated
- Twisted together
- Often "bundled" into cables
- Usually installed in building during construction



(a) Twisted pair

Twisted Pair - Applications

- Most common medium
- Telephone network
 - Between house and local exchange (subscriber loop)
- Within buildings
 - To private branch exchange (PBX)
- For local area networks (LAN)
 - 10Mbps or 100Mbps

Twisted Pair - Pros and Cons

- Cheap
- Easy to work with
- Low data rate
- Short range

Near End Crosstalk

- Coupling of signal from one pair to another
- Coupling takes place when transmit signal entering the link couples back to receiving pair
- i.e. near transmitted signal is picked up by near receiving pair

Unshielded and Shielded TP

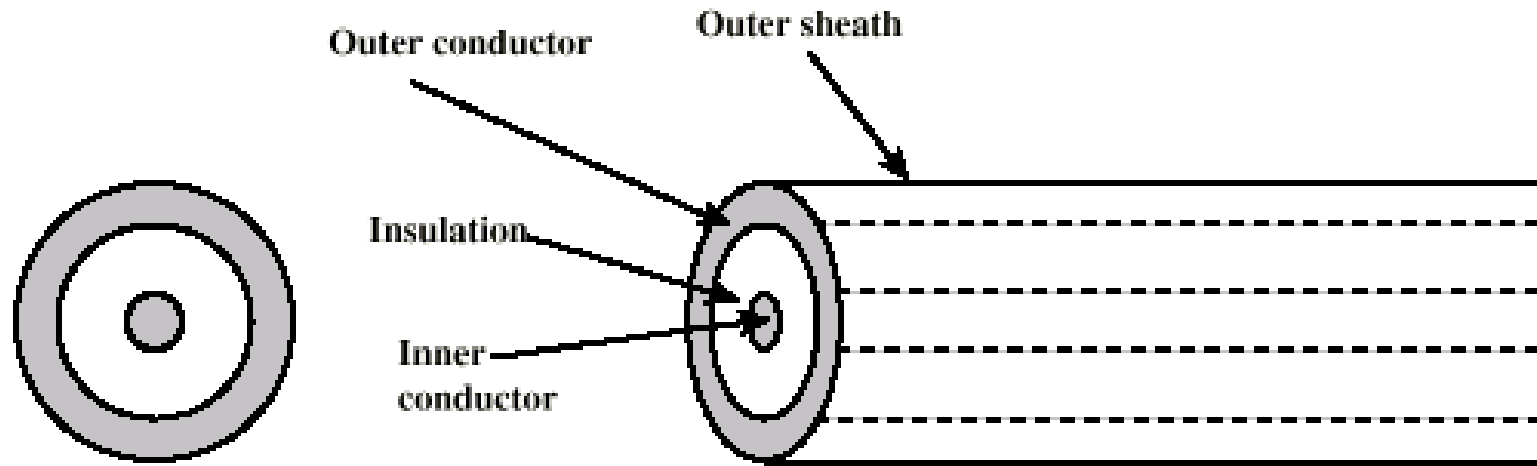
- Unshielded Twisted Pair (UTP)
 - Ordinary telephone wire
 - Cheapest
 - Easiest to install
 - Suffers from external EM interference
- Shielded Twisted Pair (STP)
 - Metal braid or sheathing that reduces interference
 - More expensive
 - Harder to handle (thick, heavy)

UTP Categories

- Cat 3
 - up to 16MHz
 - Voice grade found in most offices
 - Twist length of 7.5 cm to 10 cm
- Cat 4
 - up to 20 MHz
- Cat 5
 - up to 100MHz
 - Commonly pre-installed in new office buildings
 - Twist length 0.6 cm to 0.85 cm
- Cat 5E (Enhanced) –see tables
- Cat 6
- Cat 7

	Category 3 Class C	Category 5 Class D	Category 5E	Category 6 Class E	Category 7 Class F
Bandwidth	16 MHz	100 MHz	100 MHz	200 MHz	600 MHz
Cable Type	UTP	UTP/FTP	UTP/FTP	UTP/FTP	SSTP
Link Cost (Cat 5 =1)	0.7	1	1.2	1.5	2.2

Coaxial Cable



- Outer conductor is braided shield
- Inner conductor is solid metal
- Separated by insulating material
- Covered by padding

Coaxial Cable Applications

- Most versatile medium
- Television distribution
 - Ariel to TV
 - Cable TV
- Long distance telephone transmission
 - Can carry 10,000 voice calls simultaneously
 - Being replaced by fiber optic
- Short distance computer systems links
- Local area networks

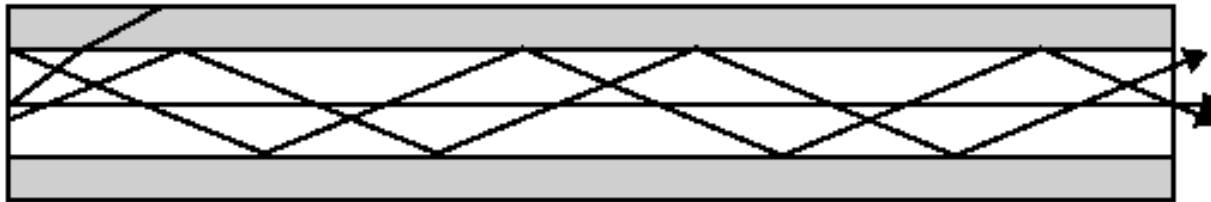
Optical Fiber - Benefits

- Greater capacity
 - Data rates of hundreds of Gbps
- Smaller size & weight
- Lower attenuation
- Electromagnetic isolation
- Greater repeater spacing
 - 10s of km at least

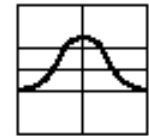
Optical Fiber - Applications

- Long-haul trunks
- Metropolitan trunks
- Rural exchange trunks
- Subscriber loops
- LANs

Input pulse

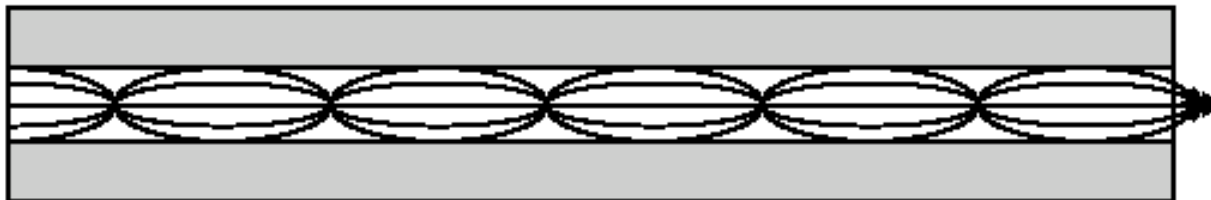


Output pulse

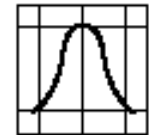


(a) Step-index multimode

Input pulse

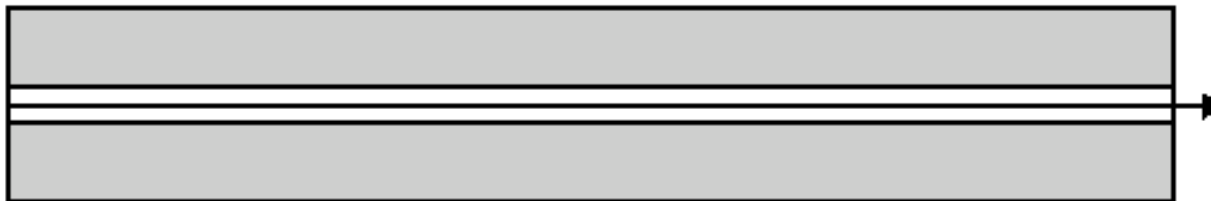
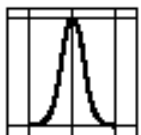


Output pulse

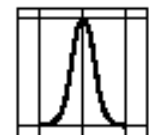


(b) Graded-index multimode

Input pulse



Output pulse

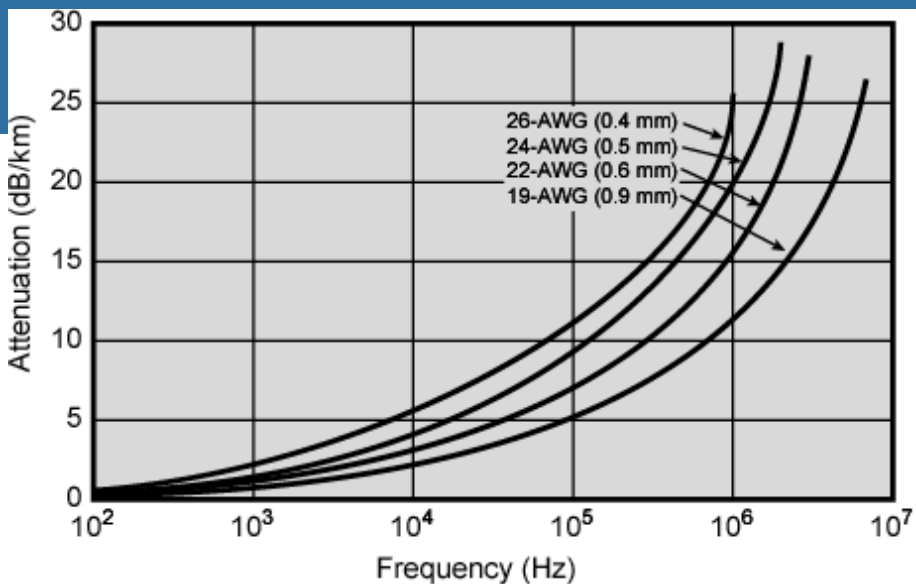


(c) Single mode

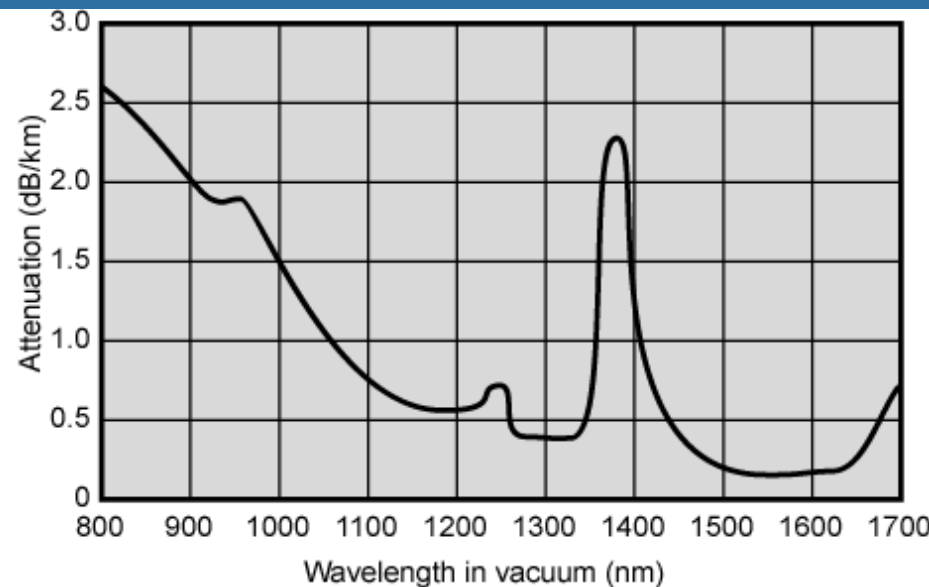
Frequency Utilization for Fiber

Applications

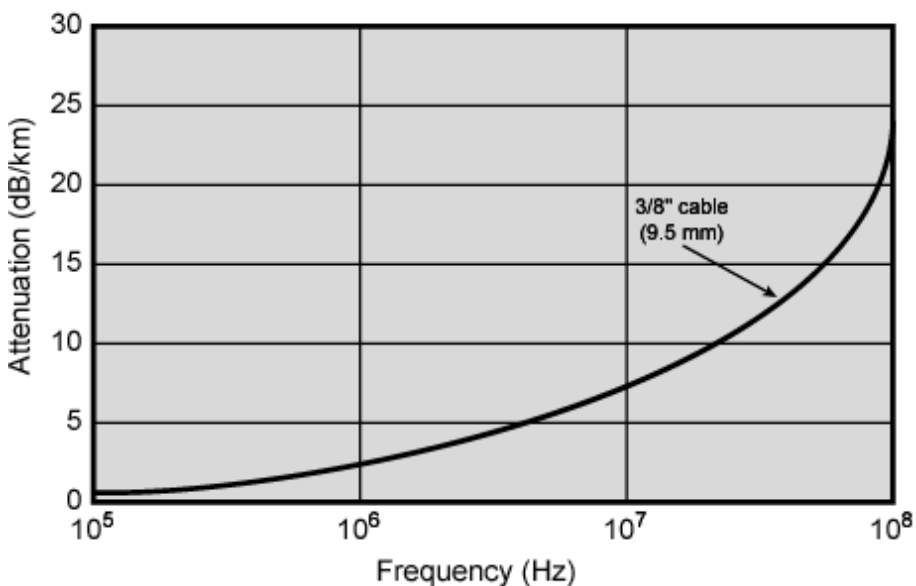
Wavelength (in vacuum) range (nm)	Frequency range (THz)	Band label	Fiber type	Application
820 to 900	366 to 333		Multimode	LAN
1280 to 1350	234 to 222	S	Single mode	Various
1528 to 1561	196 to 192	C	Single mode	WDM
1561 to 1620	185 to 192	L	Single mode	WDM



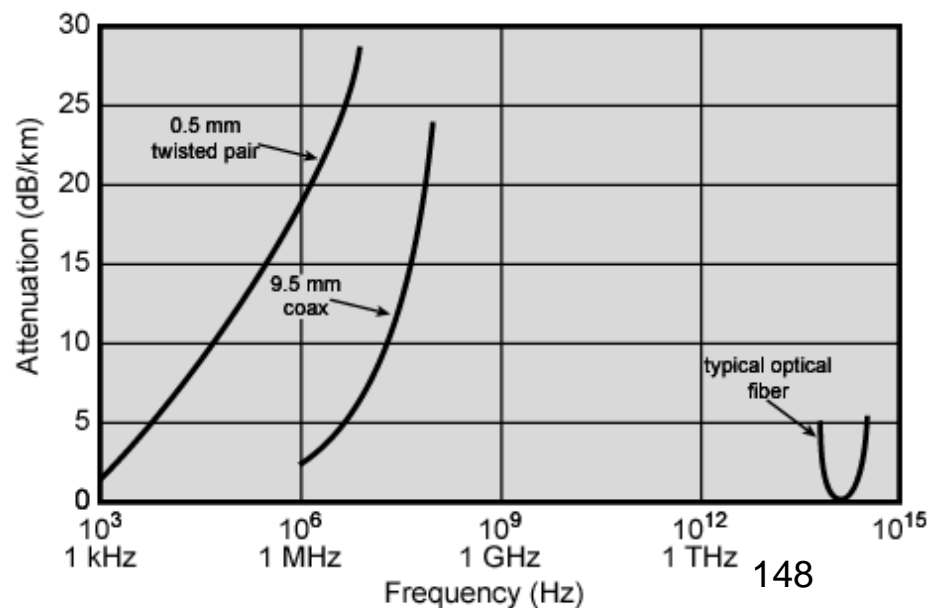
(a) Twisted pair (based on [REEV95])



(c) Optical fiber (based on [FREE02])



(b) Coaxial cable (based on [BELL90])



(d) Composite graph

Wireless Transmission Frequencies

- 2GHz to 40GHz
 - Microwave
 - Highly directional
 - Point to point
 - Satellite
- 30MHz to 1GHz
 - Omnidirectional
 - Broadcast radio
- 3×10^{11} to 2×10^{14}
 - Infrared
 - Local

UNIT II

INTRODUCTION TO DATA LINK LAYER

Link layer

our goals:

- understand principles behind link layer services:
 - error detection, correction
 - sharing a broadcast channel: multiple access
 - link layer addressing
 - local area networks: Ethernet, VLANs
- instantiation, implementation of various link layer technologies

Link layer, LANs: outline

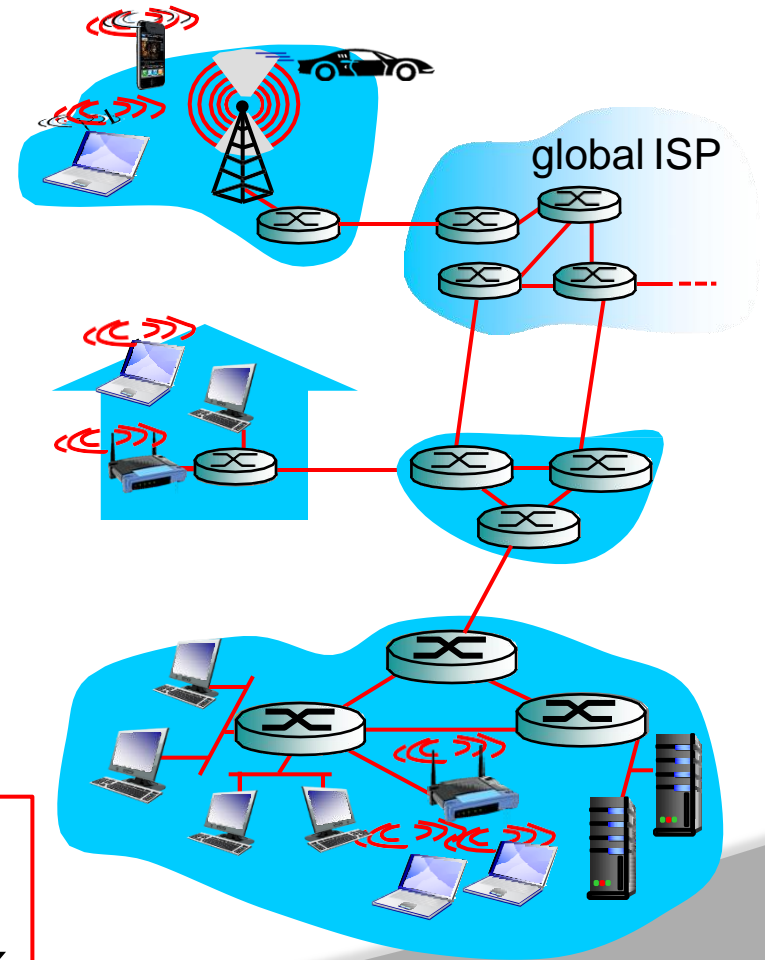
1. introduction, services
2. error detection, correction
3. multiple access protocols
4. link-layer addressing
5. Ethernet, LANs
6. LAN switches
7. a day in the life of a web request

Link layer: introduction

terminology:

- hosts and routers: **nodes**
- communication channels that connect adjacent nodes along communication path: **links**
 - wired links
 - wireless links
 - LANs
- layer-2 packet: **frame**, encapsulates datagram

data-link layer has responsibility of transferring datagram from one node to **physically adjacent** node over a link



Link layer: context

- datagram transferred by different link protocols over different links:
 - e.g., Ethernet on first link, frame relay on intermediate links, 802.11 on last link
- each link protocol provides different services
 - e.g., may or may not provide rdt over link

transportation analogy:

- trip from Princeton to Lausanne
 - limo: Princeton to JFK
 - plane: JFK to Geneva
 - train: Geneva to Lausanne
- tourist = **datagram**
- transport segment = **communication link**
- transportation mode = **link layer protocol**
- travel agent = **routing algorithm**

Link layer services

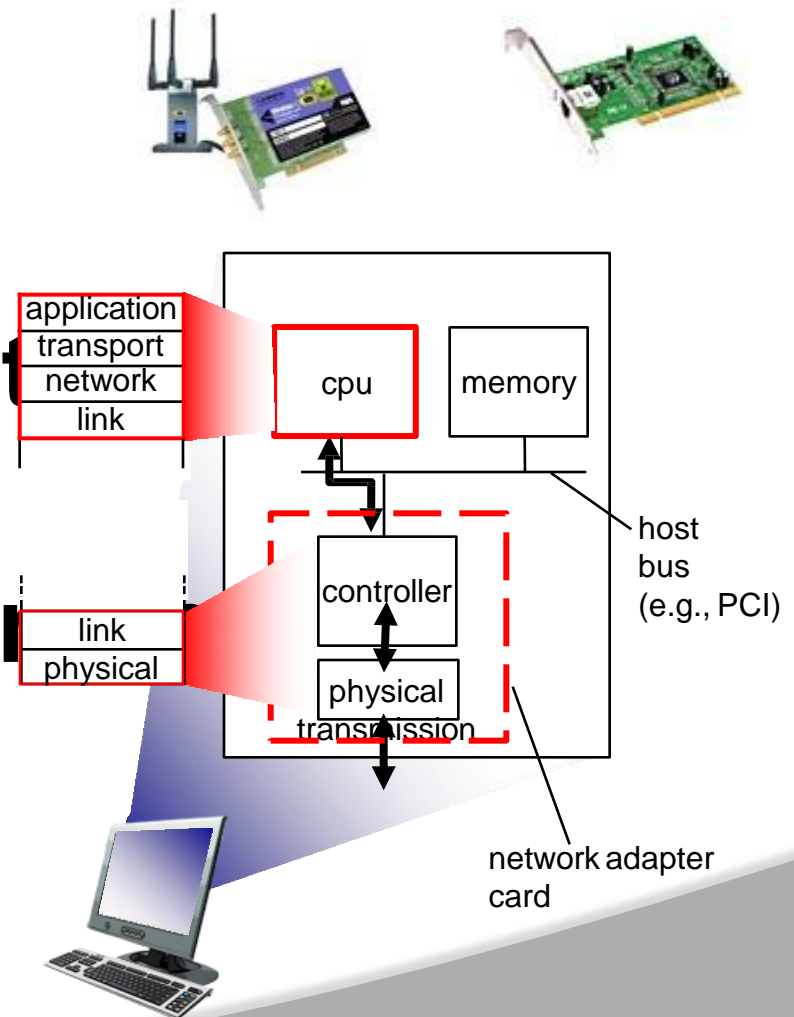
- *framing, link access:*
 - encapsulate datagram into frame, adding header, trailer
 - channel access if shared medium
 - “MAC” addresses used in frame headers to identify source, dest
 - different from IP address!
- *reliable delivery between adjacent nodes*
 - we learned how to do this already (chapter 3)!
 - seldom used on low bit-error link (fiber, some twisted pair)
 - wireless links: high error rates
 - *Q:* why both link-level and end-end reliability?

Link layer services (more)

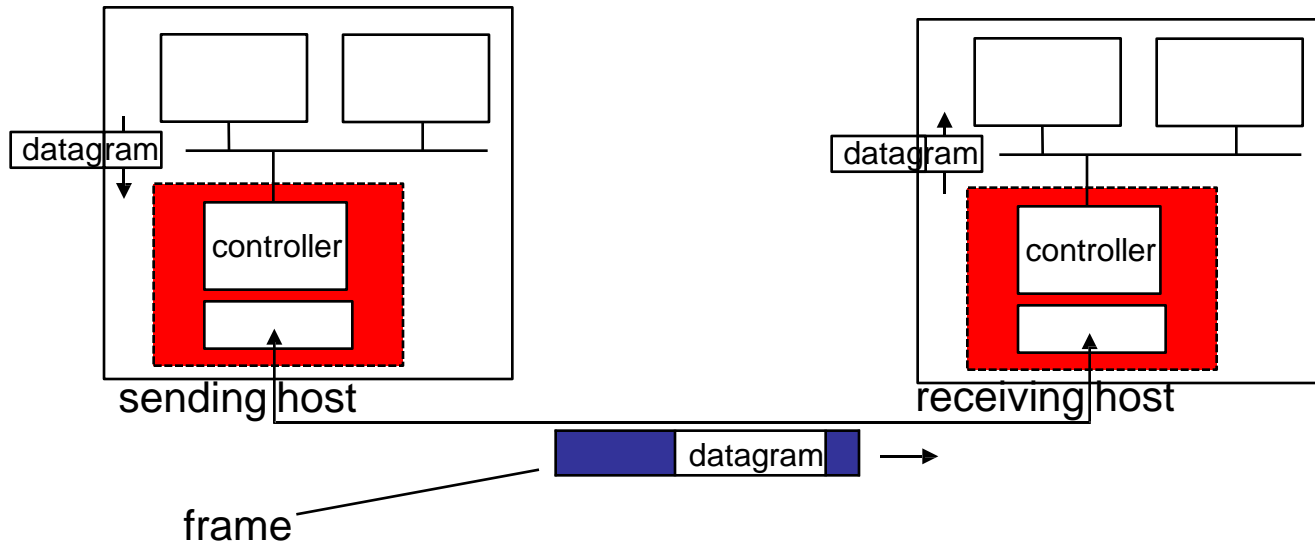
- *flow control:*
 - pacing between adjacent sending and receiving nodes
- *error detection:*
 - errors caused by signal attenuation, noise.
 - receiver detects presence of errors:
 - signals sender for retransmission or drops frame
- *error correction:*
 - receiver identifies *and corrects* bit error(s) without resorting to retransmission
- *half-duplex and full-duplex*
 - with half duplex, nodes at both ends of link can transmit, but not at same time

- in each and every host
- link layer implemented in “adaptor” (aka *network interface card* NIC)
 - Ethernet card 802.11 card
 - implements link, physical layer
- attaches into host's system buses
- combination of hardware, software, firmware

Where is the layer implemented



Adaptors communicating



- sending side:
 - encapsulates datagram in frame
 - adds error checking bits, rdt, flow control, etc.
- receiving side
 - looks for errors, rdt, flow control, etc
 - extracts datagram, passes to upper layer at receiving side

Link layer, LANs: outline

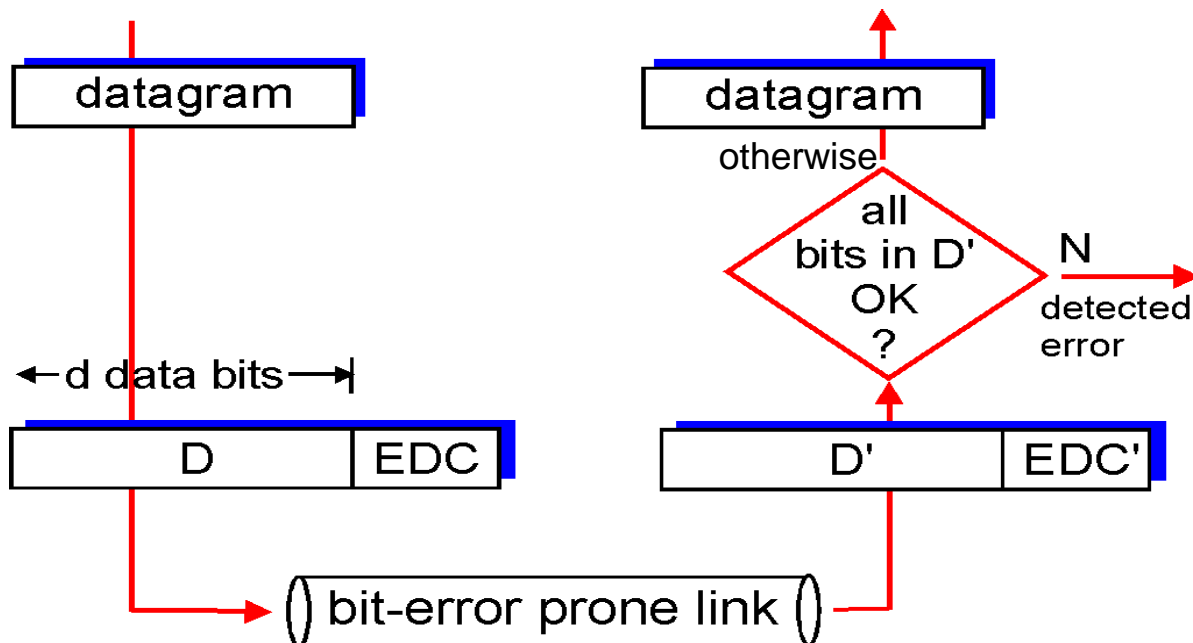
1. introduction, services
2. error detection, correction
3. multiple access protocols
4. link-layer addressing
5. Ethernet, LANs
6. LAN switches
7. a day in the life of a web request

Error detection

EDC= Error Detection and Correction bits (redundancy)

D = Data protected by error checking, may include header fields

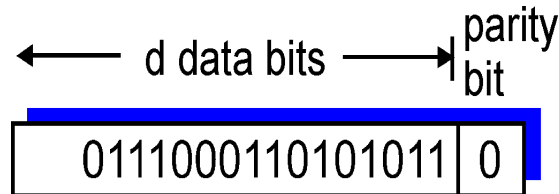
- Error detection not 100% reliable!
 - protocol may miss some errors, but rarely
 - larger EDC field yields better detection and correction



Parity checking

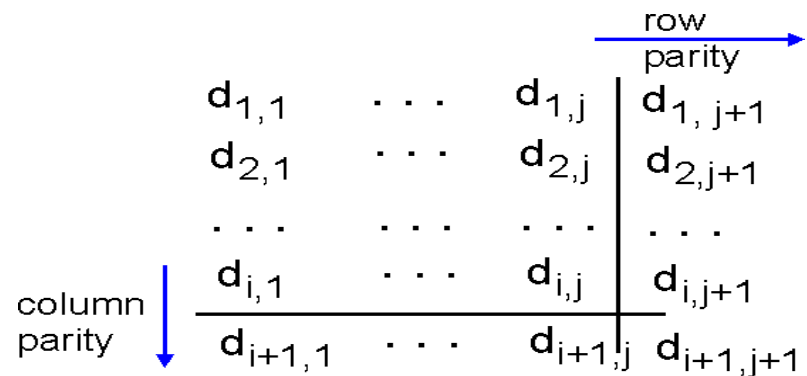
single bit parity:

- ❖ detect single bit errors



two-dimensional bit parity:

- ❖ detect and correct single bit errors



1	0	1	0	1	1
1	1	1	1	0	0
0	1	1	1	0	1
1	0	1	0	1	0

no errors

1	0	1	0	1	1
1	1	1	1	0	0
0	1	1	1	0	1
1	0	1	0	1	0

parity error

*correctable
single bit error*

Internet checksum (review)

goal: detect “errors” (e.g., flipped bits) in transmitted packet
(note: used at transport layer only)

sender:

- treat segment contents as sequence of 16-bit integers
- checksum: addition (1's complement sum) of segment contents
- sender puts checksum value into UDP checksum field

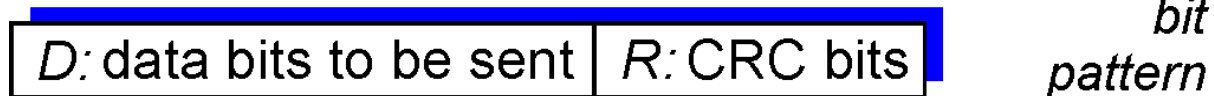
receiver:

- compute checksum of received segment
- check if computed checksum equals checksum field value:
 - NO - error detected
 - YES - no error detected. *But maybe errors nonetheless?*

Cyclic redundancy check

- more powerful error-detection coding
- view data bits, **D**, as a binary number
- choose $r+1$ bit pattern (generator), **G**
- goal: choose r CRC bits, **R**, such that
 - $\langle D, R \rangle$ exactly divisible by G (modulo 2)
 - receiver knows G , divides $\langle D, R \rangle$ by G . If non-zero remainder: error detected!
 - can detect all burst errors less than $r+1$ bits
- widely used in practice (Ethernet, 802.11 WiFi, ATM)

← d bits → ← r bits →



$$D * 2^r \text{ XOR } R$$

mathematical formula

CRC

example

want:

$$D \cdot 2^r \text{ XOR } R = nG$$

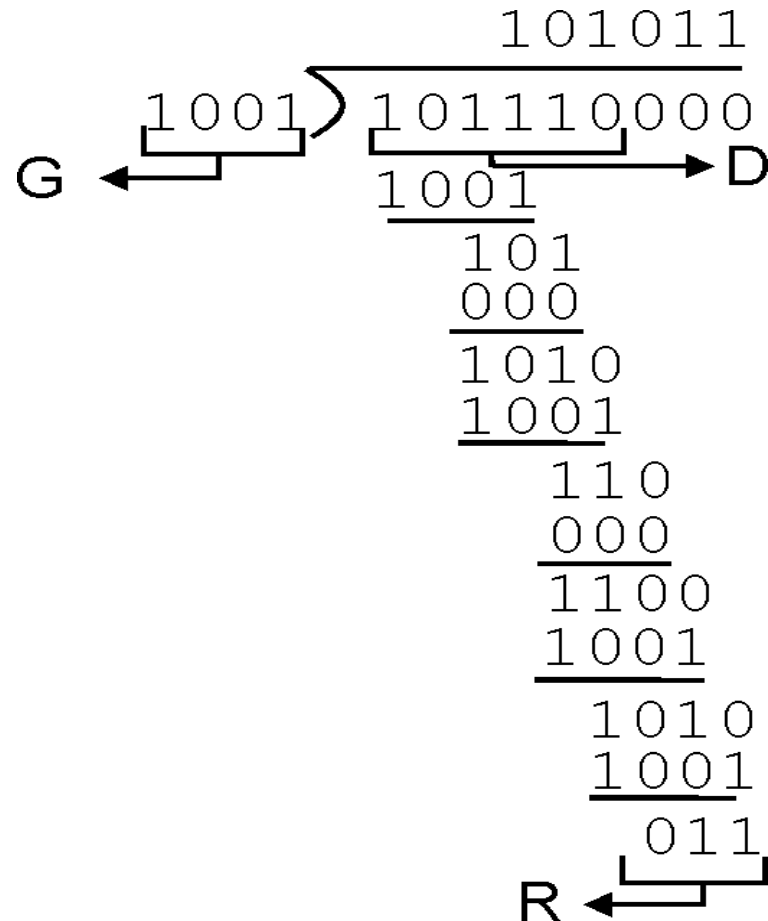
equivalently:

$$D \cdot 2^r = nG \text{ XOR } R$$

equivalently:

if we divide $D \cdot 2^r$
by G , want
remainder R to
satisfy:

$$R = \text{remainder}\left[\frac{D \cdot 2^r}{G}\right]$$



Link layer, LANs: outline

1. introduction, services
2. error detection, correction
3. multiple access protocols
4. link-layer addressing
5. Ethernet, LANs
6. LAN switches
7. a day in the life of a web request

Multiple access protocols

- single shared broadcast channel
- two or more simultaneous transmissions by nodes: interference
 - *collision* if node receives two or more signals at the same time

multiple access protocol

- distributed algorithm that determines how nodes share channel, i.e., determine when node can transmit
- communication about channel sharing must use channel itself!
 - no out-of-band channel for coordination

An ideal multiple access protocol

given: broadcast channel of rate R bps

desiderata:

1. when one node wants to transmit, it can send at rate R .
2. when M nodes want to transmit, each can send at average rate R/M
3. fully decentralized:
 - no special node to coordinate transmissions
 - no synchronization of clocks, slots
4. simple

MAC protocols: taxonomy

three broad classes:

- *channel partitioning*
 - divide channel into smaller “pieces” (time slots, frequency, code)
 - allocate piece to node for exclusive use
- *random access*
 - channel not divided, allow collisions
 - “recover” from collisions
- *“taking turns”*
 - nodes take turns, but nodes with more to send can take longer turns

Random access protocols

- when node has packet to send
 - transmit at full channel data rate R .
 - no *a priori* coordination among nodes
- two or more transmitting nodes →
“collision”,
- **random access MAC protocol** specifies:
 - how to detect collisions
 - how to recover from collisions (e.g., via delayed retransmissions)
- examples of random access MAC protocols:
 - slotted ALOHA
 - ALOHA
 - CSMA, CSMA/CD, CSMA/CA

Slotted ALOHA

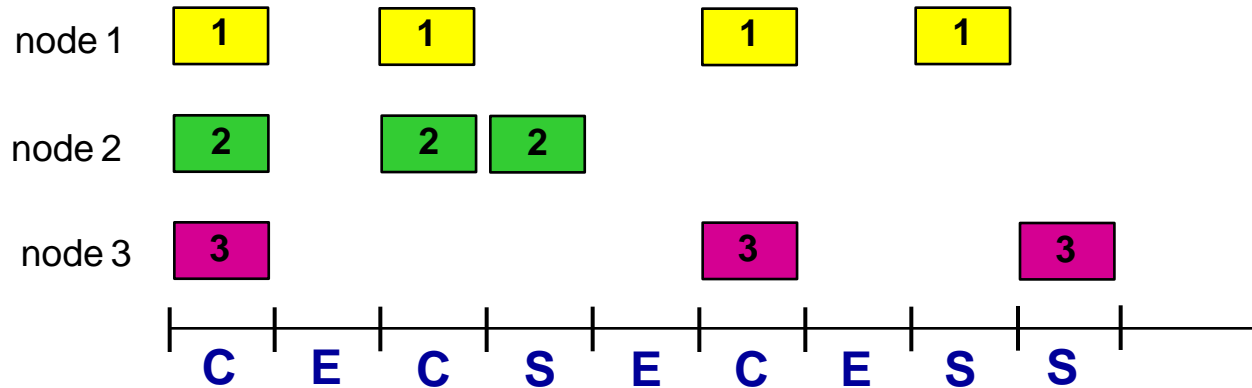
assumptions:

- all frames same size
- time divided into equal size slots (time to transmit 1 frame)
- nodes start to transmit only slot beginning
- nodes are synchronized
- if 2 or more nodes transmit in slot, all nodes detect collision

operation:

- when node obtains fresh frame, transmits in next slot
 - *if no collision:* node can send new frame in next slot
 - *if collision:* node retransmits frame in each subsequent slot with prob. p until success

Slotted ALOHA



Pros:

- single active node can continuously transmit at full rate of channel
- highly decentralized: only slots in nodes need to be in sync
- simple

Cons:

- collisions, wasting slots
- idle slots
- nodes may be able to detect collision in less than time to transmit packet
- clock synchronization

Slotted ALOHA: efficiency

efficiency: long-run fraction of successful slots (many nodes, all with many frames to send)

- *suppose:* N nodes with many frames to send, each transmits in slot with probability p
- prob that given node has success in a slot = $p(1-p)^{N-1}$
- prob that *any* node has a success = $Np(1-p)^{N-1}$

- max efficiency: find p^* that maximizes $Np(1-p)^{N-1}$
- for many nodes, take limit of $Np^*(1-p^*)^{N-1}$ as N goes to infinity, gives:

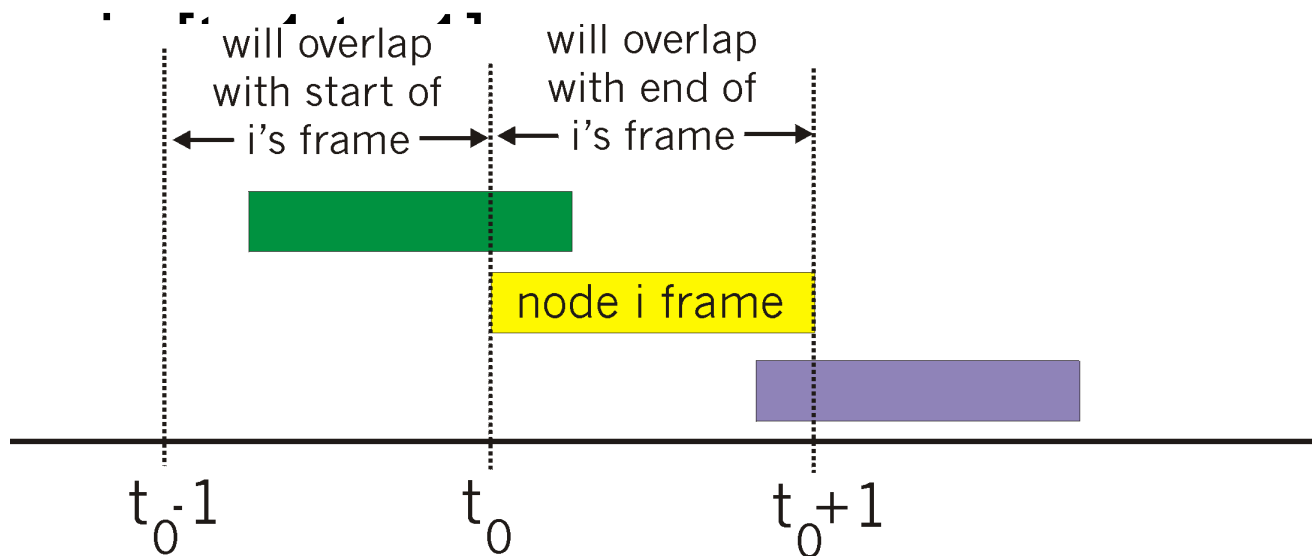
$$\text{max efficiency} = 1/e =$$

^{37%}
at best: channel used for useful transmissions 37% of time!



Pure (unslotted) ALOHA

- unslotted Aloha: simpler, no synchronization
- when frame first arrives
 - transmit immediately
- collision probability increases:
 - frame sent at t_0 collides with other frames sent



Pure ALOHA efficiency

$P(\text{success by given node}) = P(\text{node transmits}) \cdot$

$P(\text{no other node transmits in } [t_0-1, t_0]) \cdot$

$P(\text{no other node transmits in } [t_0-1, t_0])$

$$= p \cdot (1-p)^{N-1} \cdot (1-p)^{N-1}$$

$$= p \cdot (1-p)^{2(N-1)}$$

$\rightarrow \infty$

... choosing optimum p and then letting n

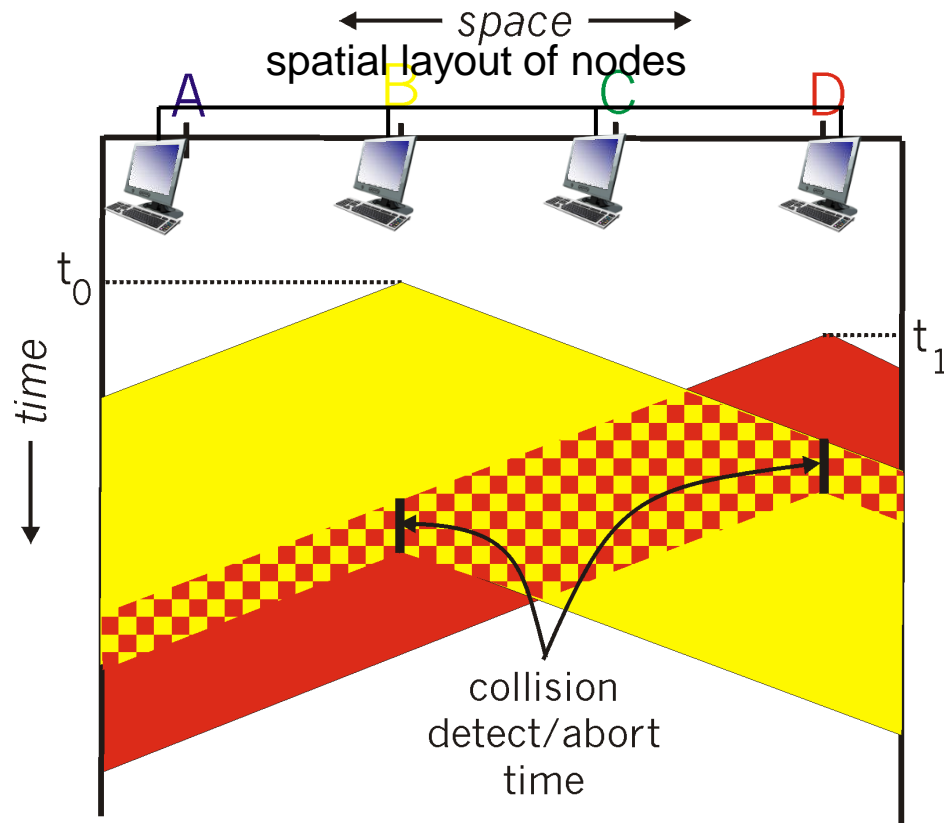
$$= 1/(2e) = .18$$

even worse than slotted Aloha!

CSMA/CD (collision detection)

- CSMA/CD*: carrier sensing, deferral as in CSMA
- collisions *detected* within short time
 - colliding transmissions aborted, reducing channel wastage
 - collision detection:
 - easy in wired LANs: measure signal strengths, compare transmitted, received signals
 - difficult in wireless LANs: received signal strength overwhelmed by local transmission strength
 - human analogy: the polite conversationalist

CSMA/CD (collision detection)



“Taking turns” MAC protocols

channel partitioning MAC protocols:

- share channel *efficiently* and *fairly* at high load
- inefficient at low load: delay in channel access, $1/N$ bandwidth allocated even if only 1 active node!

random access MAC protocols

- efficient at low load: single node can fully utilize channel
- high load: collision overhead

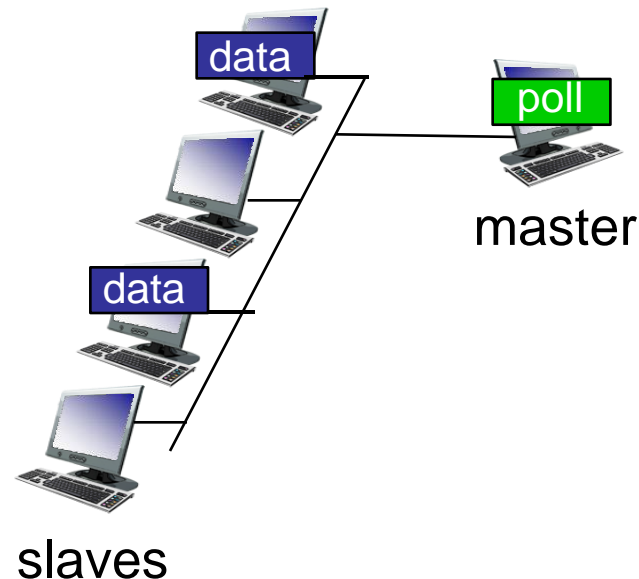
“taking turns” protocols

look for best of both worlds!

“Taking turns” MAC protocols

polling:

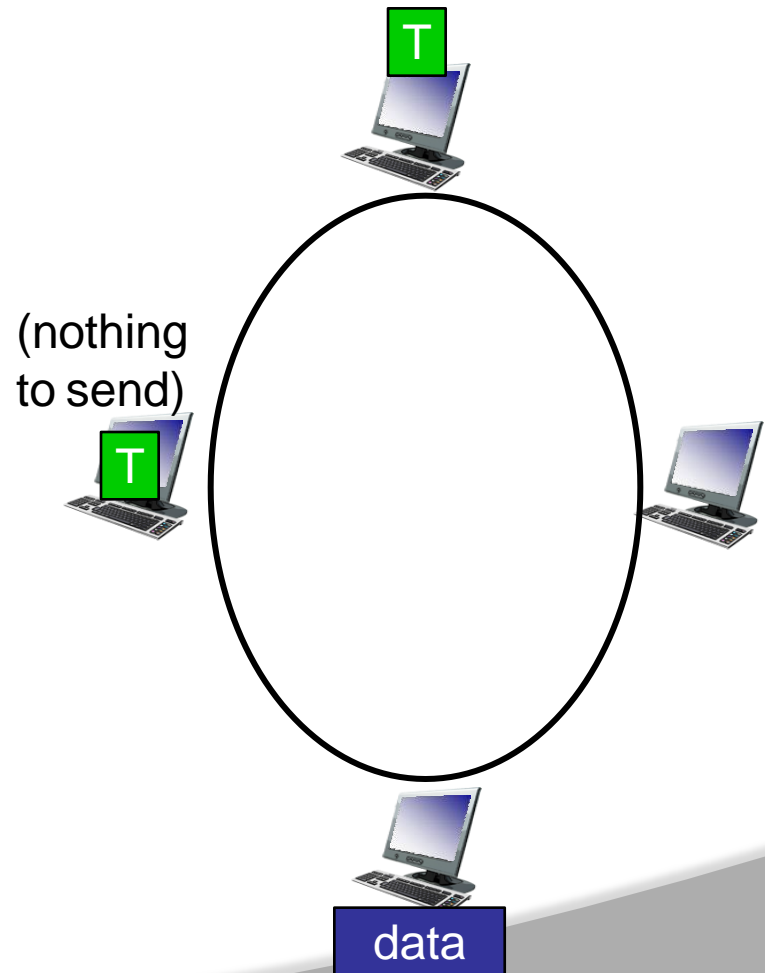
- master node “invites” slave nodes to transmit in turn
- typically used with “dumb” slave devices
- concerns:
 - polling overhead
 - latency
 - single point of failure (master)



“Taking turns” MAC protocols

token passing:

- ❖ control **token** passed from one node to next sequentially.
- ❖ token message
- ❖ concerns:
 - token overhead
 - latency
 - single point of failure (token)



Summary of MAC protocols

- *channel partitioning*, by time, frequency or code
 - Time Division, Frequency Division
- *random access* (dynamic),
 - ALOHA, S-ALOHA, CSMA, CSMA/CD
 - carrier sensing: easy in some technologies (wire), hard in others (wireless)
 - CSMA/CD used in Ethernet
 - CSMA/CA used in 802.11
- *taking turns*
 - polling from central site, token passing
 - bluetooth, FDDI, IBM token ring

Link layer, LANs: outline

1. introduction, services
2. error detection, correction
3. multiple access protocols
4. link-layer addressing
5. Ethernet, LANs
6. LAN switches
7. a day in the life of a web request

MAC addresses and ARP

- 32-bit IP address:
 - ❖ *network-layer* address
 - ❖ datagram to destination used to get IP subnet
- MAC (or LAN or physical or Ethernet) address:
 - ❖ function: *get frame from one interface to another physically-connected interface (same network, in IP-addressing sense)*
 - ❖ 48 bit MAC address (for most LANs) burned in NIC ROM, also sometimes software settable
 - ❖ e.g.: 1A-2F-BB-76-09-AD

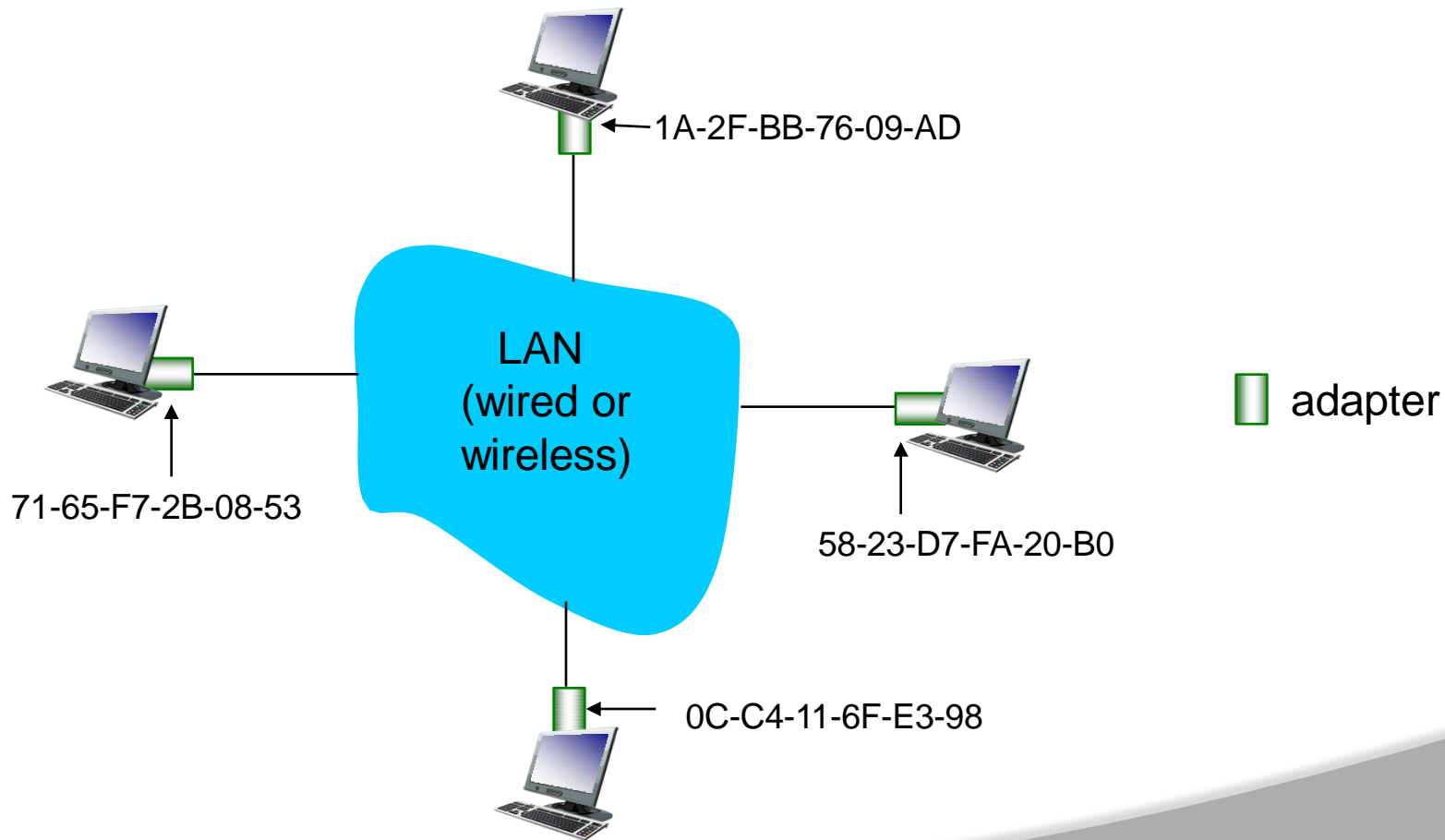
/

hexadecimal (base 16) notation

- Why two addresses for node ??
(each number represents 4 bits)

LAN addresses and ARP

each adapter on LAN has unique LAN address

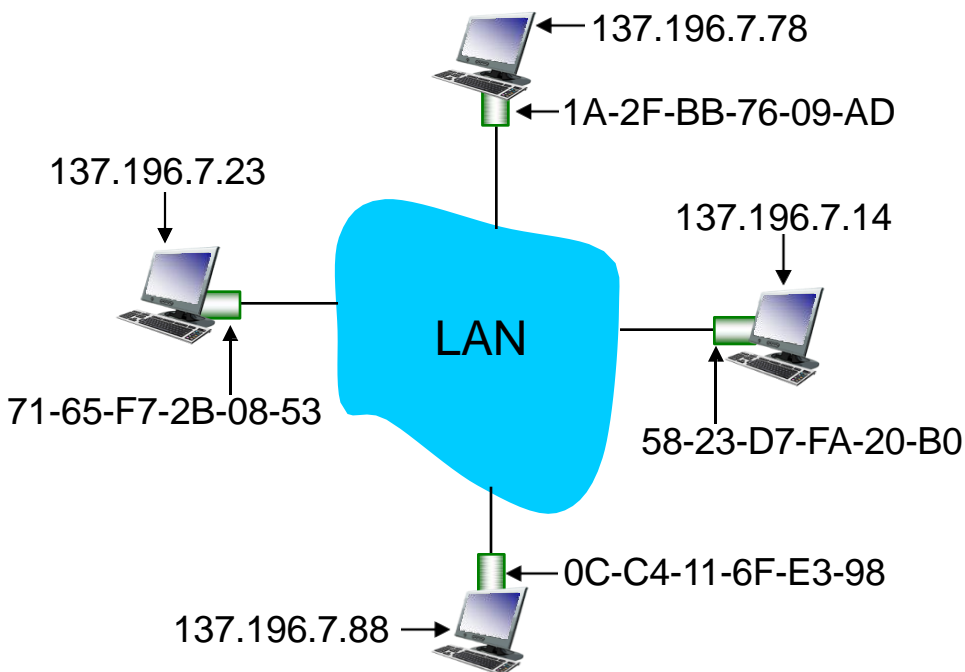


LAN addresses (more)

- MAC address allocation administered by IEEE
- manufacturer buys portion of MAC address space (to assure uniqueness)
- analogy:
 - MAC address: like Social Security Number
 - IP address: like postal address
- MAC flat address → portability
 - can move LAN card from one LAN to another
- IP hierarchical address *not* portable
 - address depends on IP subnet to which node is attached

ARP: address resolution protocol

Question: how to determine MAC address of B knowing B's IP address?



- each IP node (host, router) on LAN has **ARP** table
 - IP/MAC address mappings for some LAN nodes:
< IP address; MAC address; TTL >
 - TTL (Time To Live): time after which address mapping will be forgotten (typically 20 min)

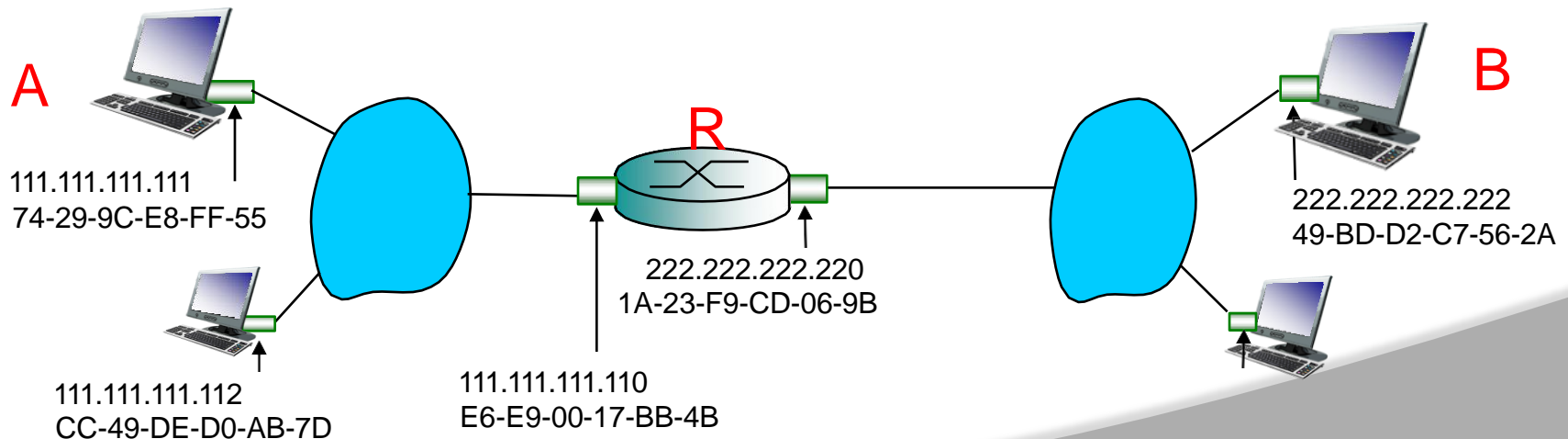
ARP protocol: same LAN

- A wants to send datagram to B
 - B's MAC address not in A's ARP table.
- A **broadcasts** ARP query packet, containing B's IP address
 - dest MAC address = FF-FF-FF-FF-FF-FF
 - all machines on LAN receive ARP query
- B receives ARP packet, replies to A with its (B's) MAC address
 - frame sent to A's MAC address (unicast)
- A caches (saves) IP-to- MAC address pair in its ARP table until information becomes old (times out)
 - soft state: information that times out (goes away) unless refreshed
- ARP is “plug-and-play”:
 - nodes create their ARP tables *without intervention from net*

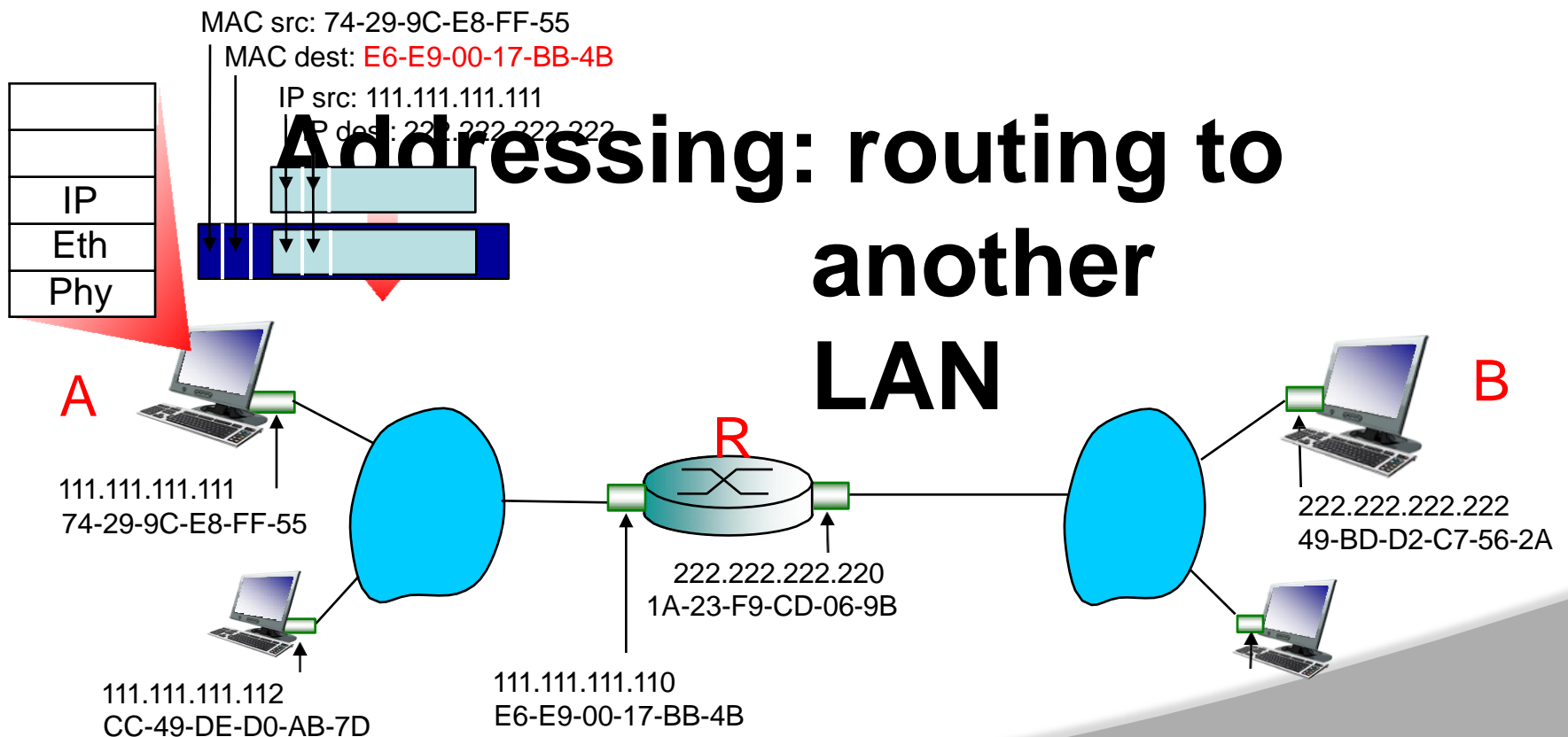
Addressing: routing to another LAN

walkthrough: **send datagram from A to B via R**

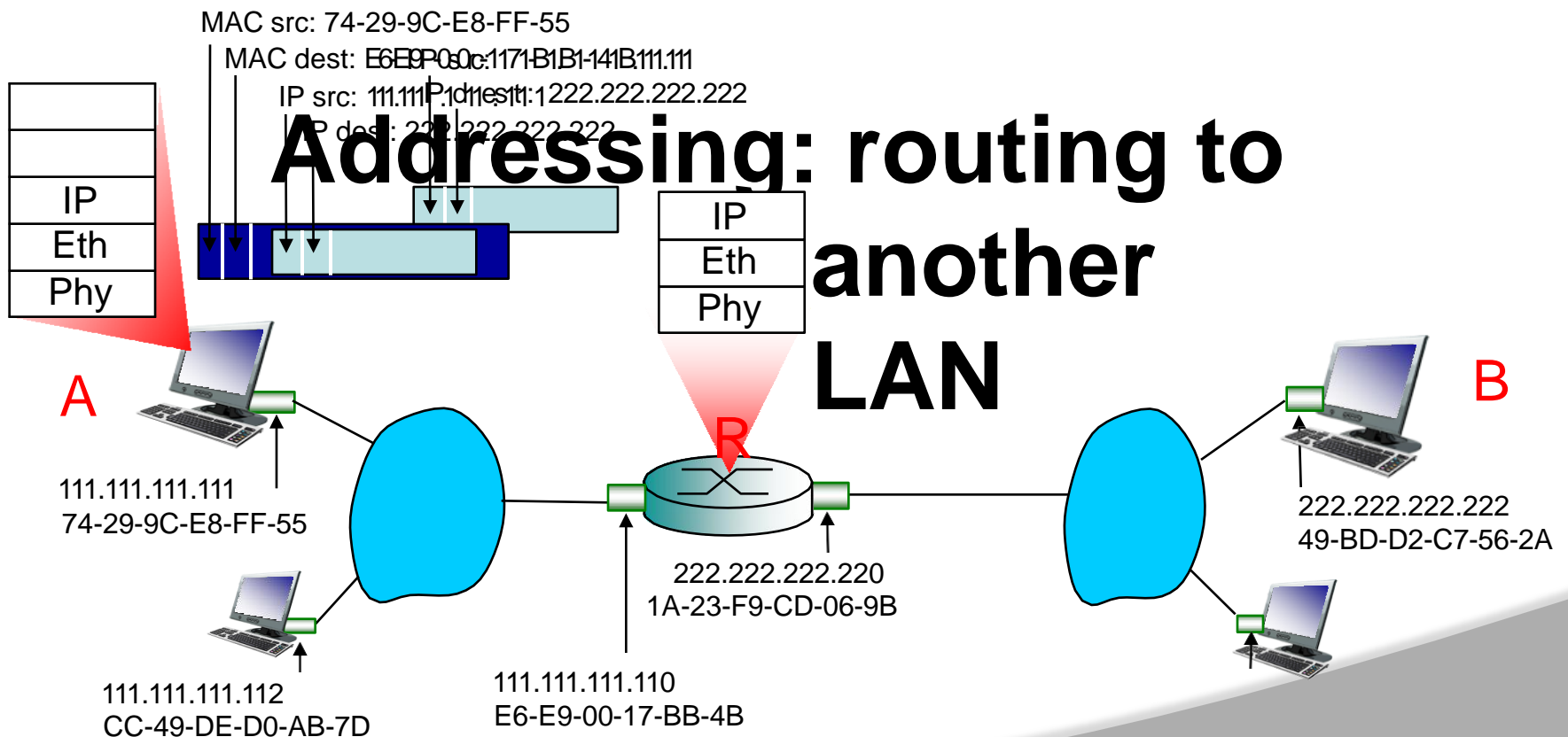
- focus on addressing - at both IP (datagram) and MAC layer (frame)
- assume A knows B's IP address
- assume A knows IP address of first hop router, R (how?)
- assume A knows R's MAC address (how?)



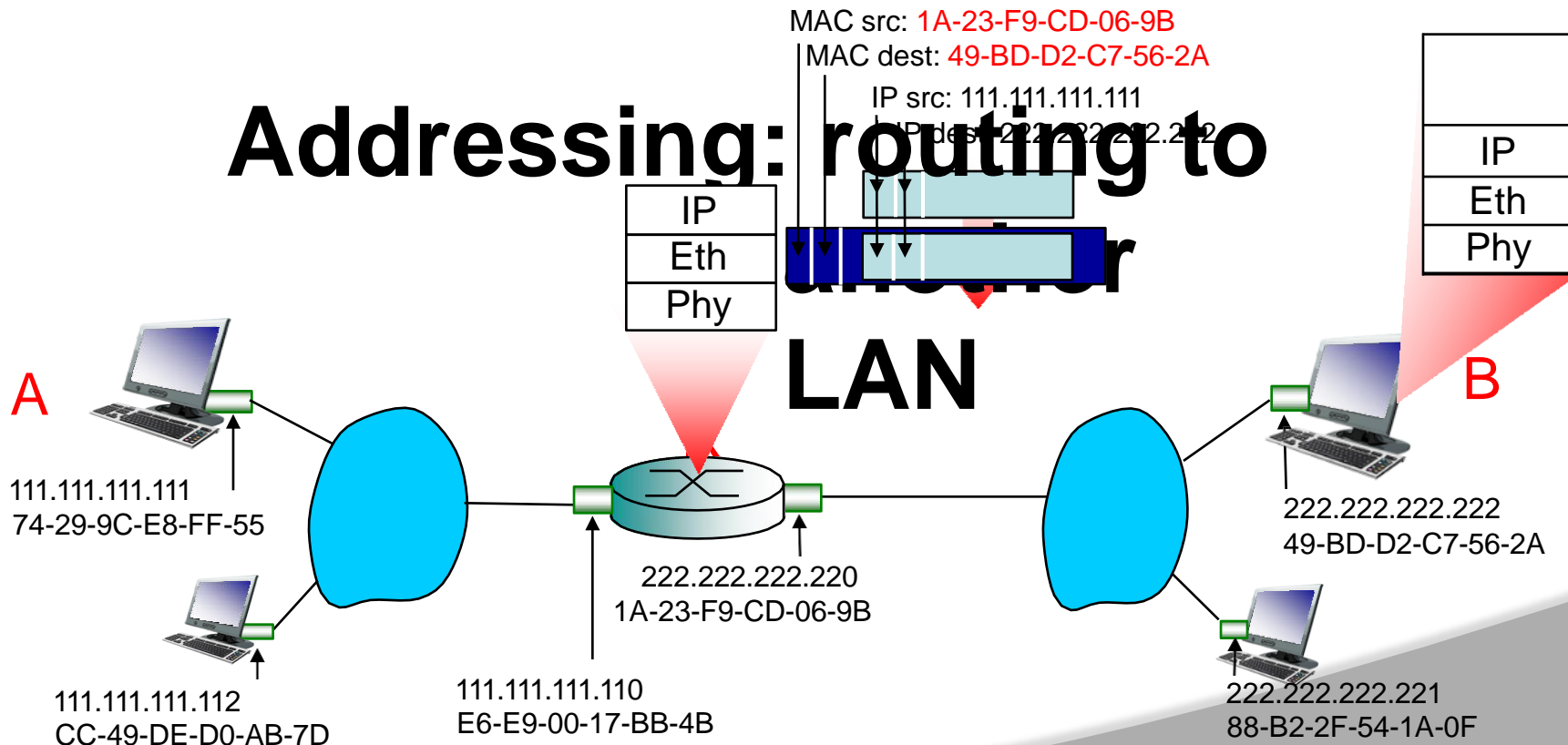
- ❖ A creates IP datagram with IP source A, destination B
- ❖ A creates link-layer frame with R's MAC address as dest, frame contains A-to-B IP datagram



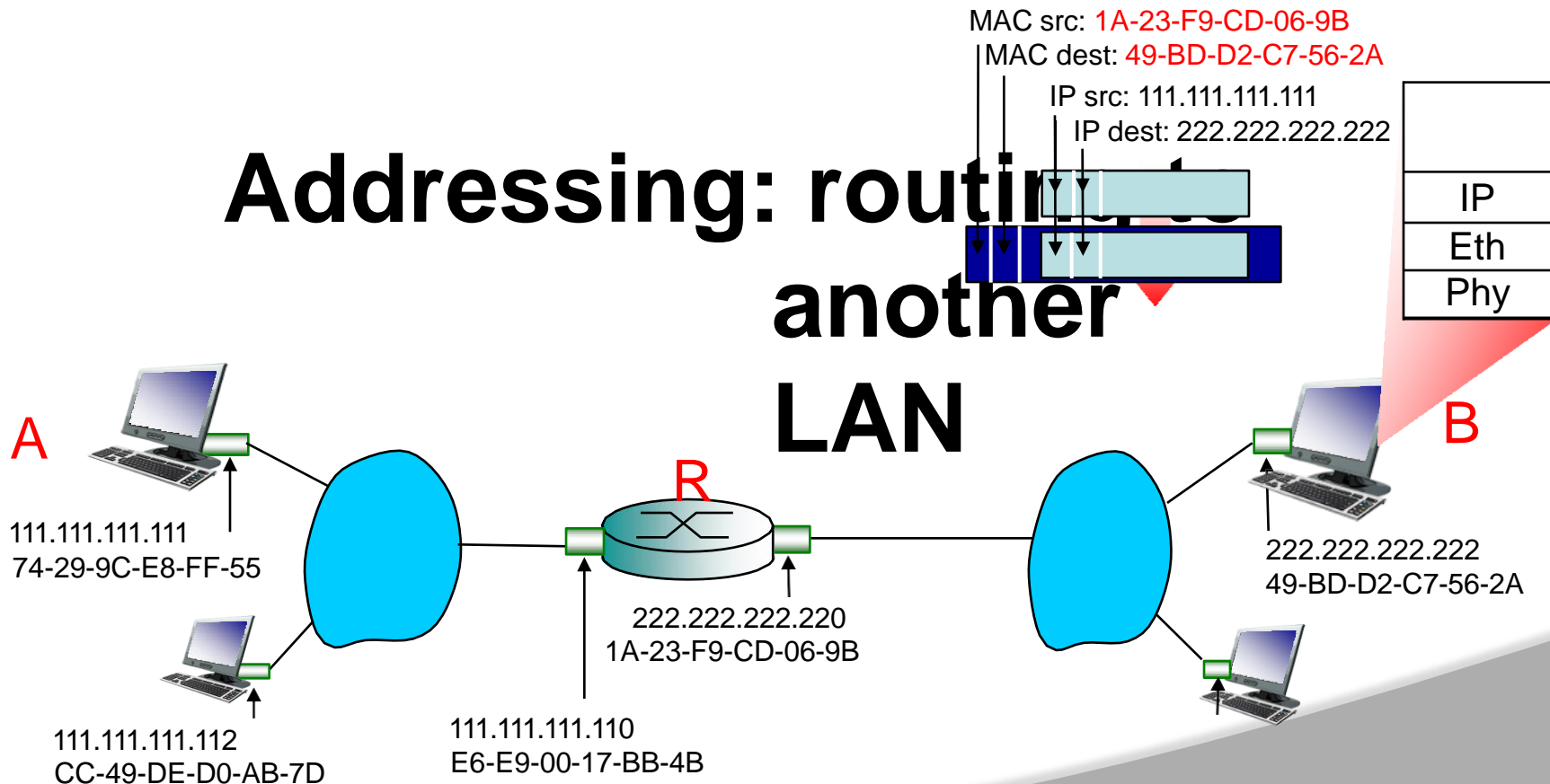
- ❖ frame sent from A to R
- ❖ frame received at R, datagram removed, passed up to IP



- ❖ R forwards datagram with IP source A, destination B
- ❖ R creates link-layer frame with B's MAC address as dest, frame contains A-to-B IP datagram



- ❖ R forwards datagram with IP source A, destination B
- ❖ R creates link-layer frame with B's MAC address as dest, frame contains A-to-B IP datagram



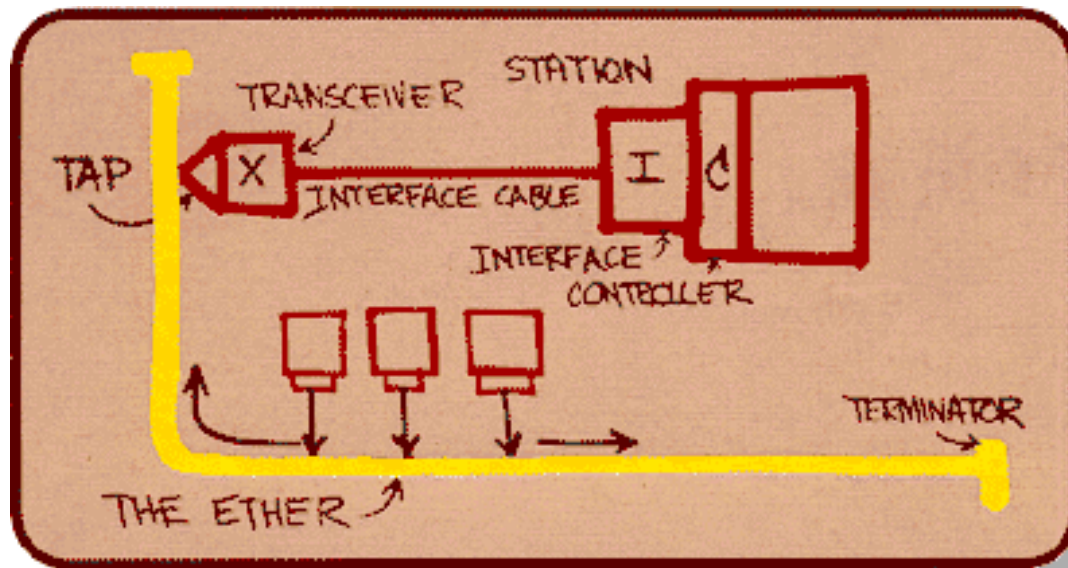
Link layer, LANs: outline

1. introduction, services
2. error detection, correction
3. multiple access protocols
4. link-layer addressing
5. Ethernet, LANs
6. LAN switches
7. a day in the life of a web request

Ethernet

“dominant” wired LAN technology:

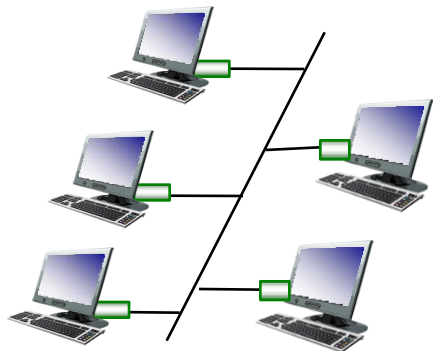
- cheap \$20 for NIC
- first widely used LAN technology
- Developed in the mid-1970s by researchers at the Xerox Palo Alto Research Centers (PARC)
- simpler, cheaper than token LANs and ATM
- kept up with speed race: 10 Mbps – 10 Gbps



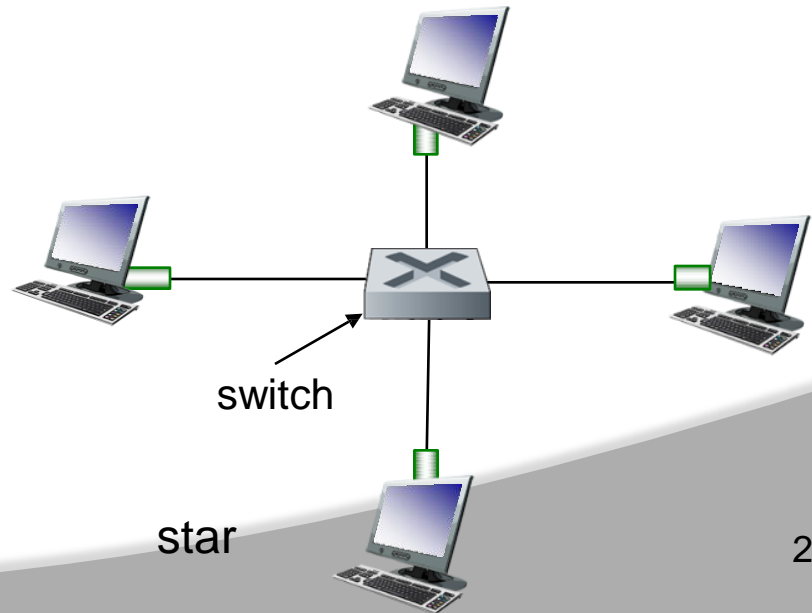
Metcalfe's Ethernet sketch

Star topology

- bus topology popular through mid 90s
 - all nodes in same collision domain (can collide with each other)
- today: star topology prevails
 - active **switch** in center
 - each “spoke” runs a (separate) Ethernet protocol (nodes do not collide with each other)

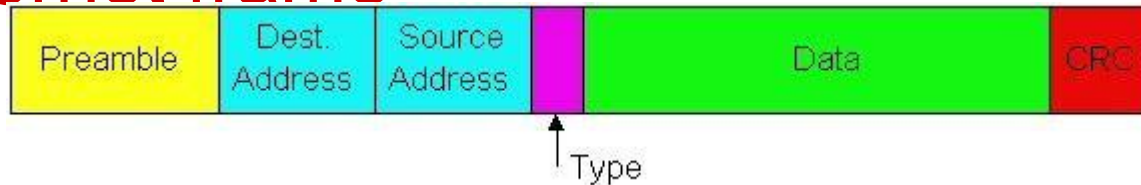


bus: coaxial cable



Ethernet frame structure

Sending adapter encapsulates IP datagram (or other network layer protocol packet) in **Ethernet frame**

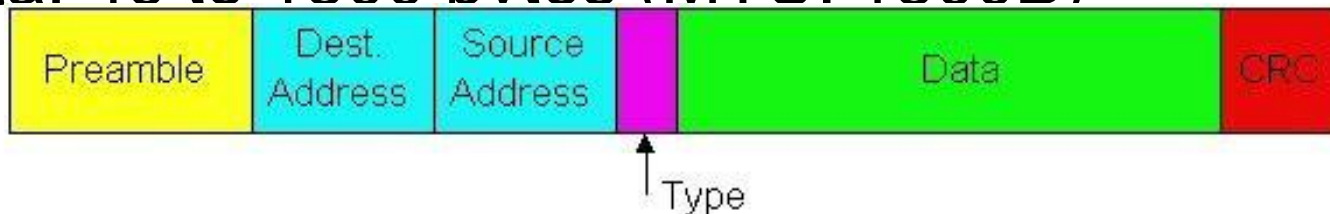


Preamble:

- 7 bytes with pattern 10101010 followed by one byte with pattern 10101011
- used to synchronize receiver, sender clock rates

Ethernet frame structure (more)

- **addresses:** 6 bytes
 - if adapter receives frame with matching destination address, or with broadcast address (e.g. ARP packet), it passes data in frame to network layer protocol
 - otherwise, adapter discards frame
- **type:** indicates higher layer protocol (mostly IP but others possible, e.g., Novell IPX, AppleTalk)
- **CRC:** checked at receiver, if error is detected, frame is dropped
- Data: 46 to 1500 bytes (MTU: 1500B)



Ethernet: unreliable, connectionless

- *connectionless*: No handshaking between sending and receiving NICs
- *unreliable*: receiving NIC doesn't send acks or nacks to sending NIC
 - stream of datagrams passed to network layer can have gaps (missing datagrams)
 - gaps will be filled if app is using TCP
 - otherwise, app will see gaps
- Ethernet's MAC protocol: unslotted *CSMA/CD*

Ethernet CSMA/CD algorithm

1. NIC receives datagram from network layer, creates frame
2. If NIC senses channel idle, starts frame transmission If NIC senses channel busy, waits until channel idle, then transmits
3. If NIC transmits entire frame without detecting another transmission, NIC is done with frame !
4. If NIC detects another transmission while transmitting, aborts and sends 48-bit jam signal
5. After aborting, NIC enters *exponential backoff*: after m th collision, NIC chooses K at random from $\{0, 1, 2, \dots, 2^m - 1\}$. NIC waits $K \cdot 512$ bit times, returns to Step 2

CSMA/CD efficiency

- T_{prop} = max prop delay between 2 nodes in LAN
- t_{trans} = time to transmit max-size frame
- efficiency goes to 1
 - as t_{prop} goes to 0
 - as t_{trans} goes to infinity
- better performance than ALOHA: and simple, cheap, decentralized!

$$\text{efficiency} = \frac{1}{1 + 5t_{\text{prop}}/t_{\text{trans}}}$$

Link layer, LANs: outline

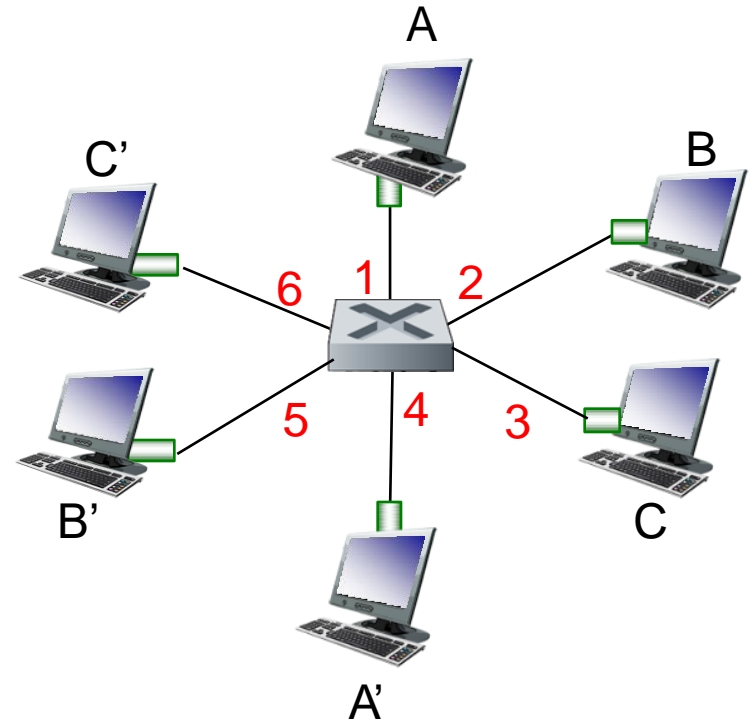
1. introduction, services
2. error detection, correction
3. multiple access protocols
4. link-layer addressing
5. Ethernet, LANs
6. LAN switches
7. a day in the life of a web request

Ethernet switch

- link-layer device: takes an *active* role
 - store, forward Ethernet frames
 - examine incoming frame's MAC address, *selectively* forward frame to one-or-more outgoing links when frame is to be forwarded on segment, uses CSMA/CD to access segment
- *transparent*
 - hosts are unaware of presence of switches
- *plug-and-play, self-learning*
 - switches do not need to be configured

Switch table

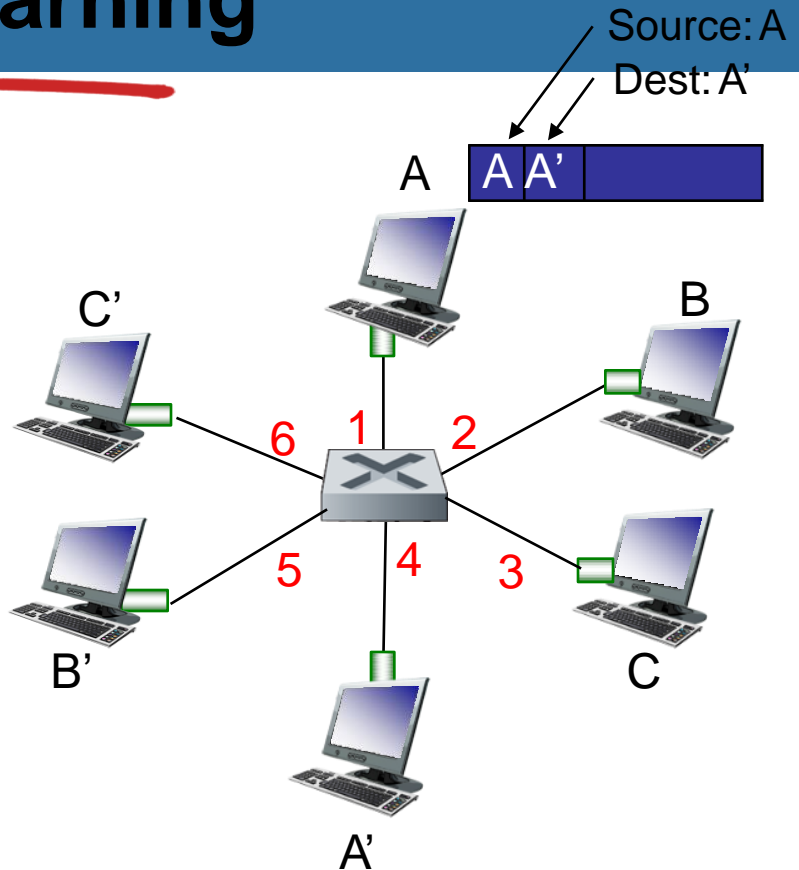
- Q: how does switch know that A' reachable via interface 4, B' reachable via interface 5?
- A: each switch has a **switch table**, each entry:
 - (MAC address of host, interface to reach host, time stamp)
- looks like a routing table!
- Q: how are entries created, maintained in switch table?
 - something like a routing protocol?



switch with six interfaces
(1,2,3,4,5,6)

Switch: self-learning

- switch *learns* which hosts can be reached through which interfaces
 - when frame received, switch “learns” location of sender: incoming LAN segment
 - records sender/location pair in switch table



MAC addr	interface	TTL
A	1	60

Switch table
(initially empty)

Switch: frame filtering/forwarding

When frame received:

1. record link associated with sending host
2. index switch table using MAC dest address
3. if entry found for destination
 then {
 if dest on segment from which frame arrived
 then drop the frame
 else forward the frame on interface indicated
 }
 else flood

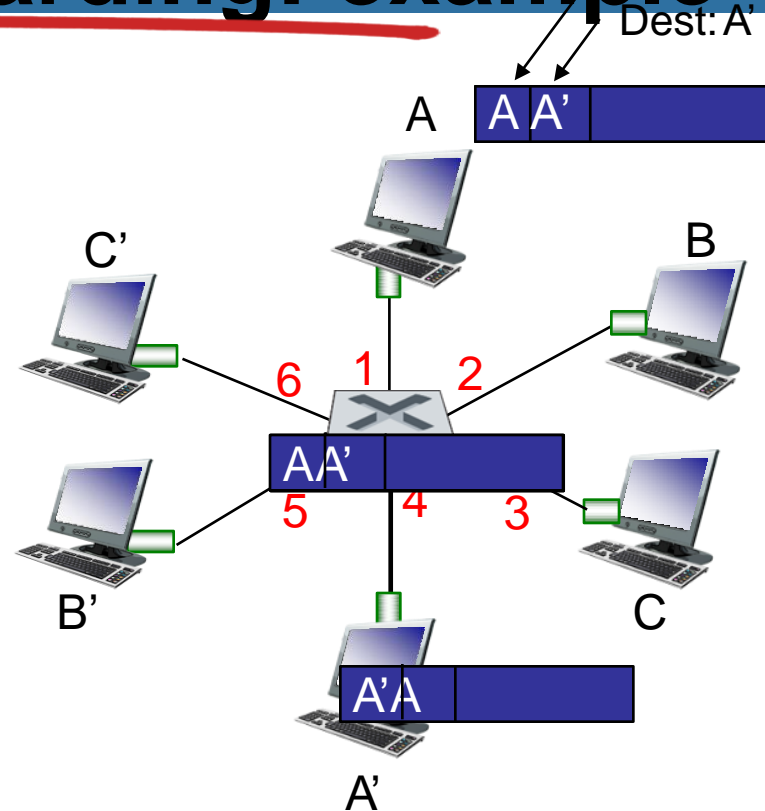


forward on all but the interface
on which the frame arrived

Self-learning, forwarding: example

Source: A

- frame destination unknown: flood
- ❖ destination A location known: selective send

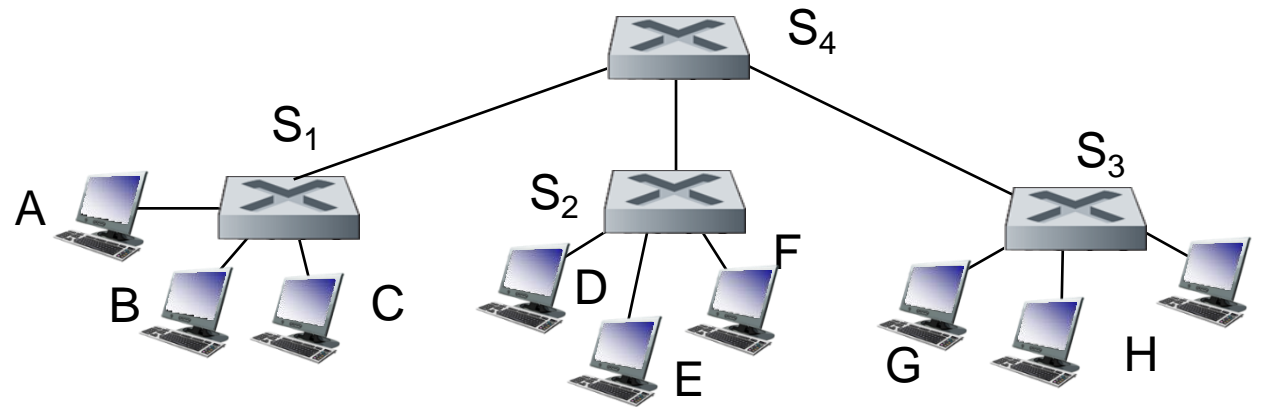


MAC addr	interface	TTL
A	1	60
A'	4	60

Switch table
(initially empty)

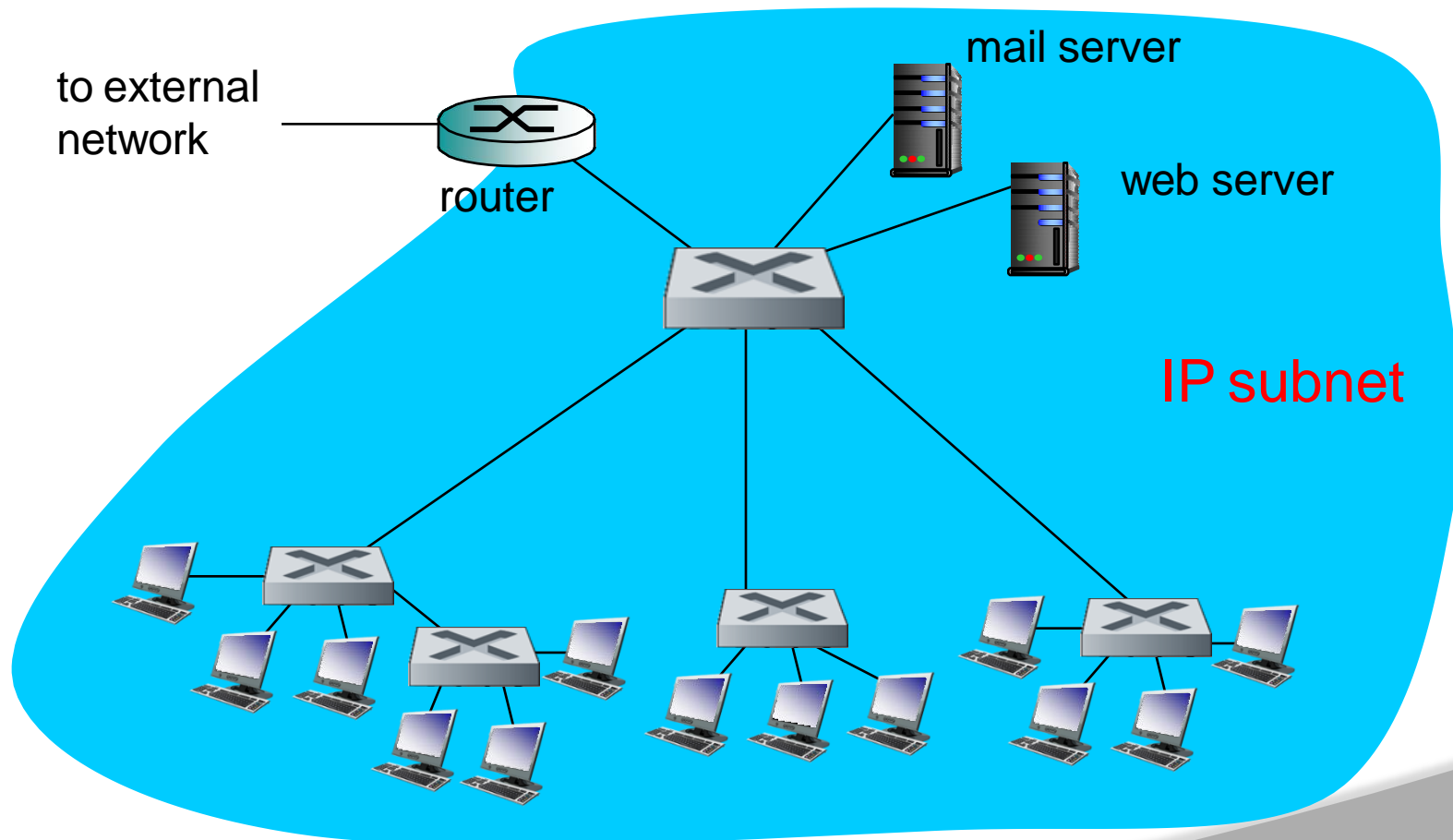
Interconnecting switches

- switches can be connected together



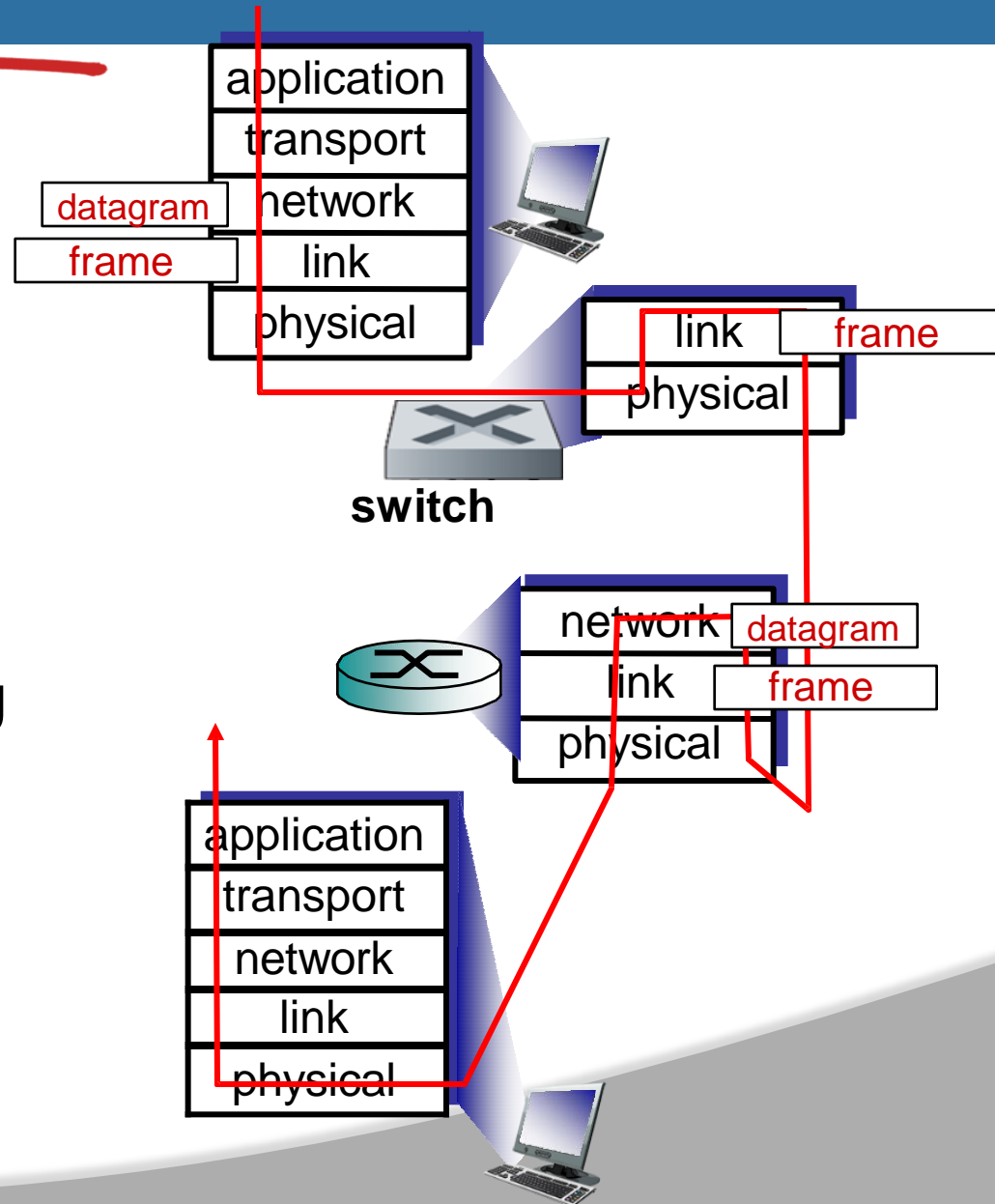
- ❖ Q: sending from A to G - how does S₁ know to forward frame destined to F via S₄ and S₃?
- ❖ A: self learning! (works exactly the same as in single-switch case!)

Institutional network



Switches vs. Routers

- both store-and-forward devices
 - routers: network-layer devices (examine network-layer headers)
 - switches are link-layer devices (examine link-layer headers)
- routers maintain routing tables, implement routing algorithms
- switches maintain switch tables, implement filtering, learning algorithms



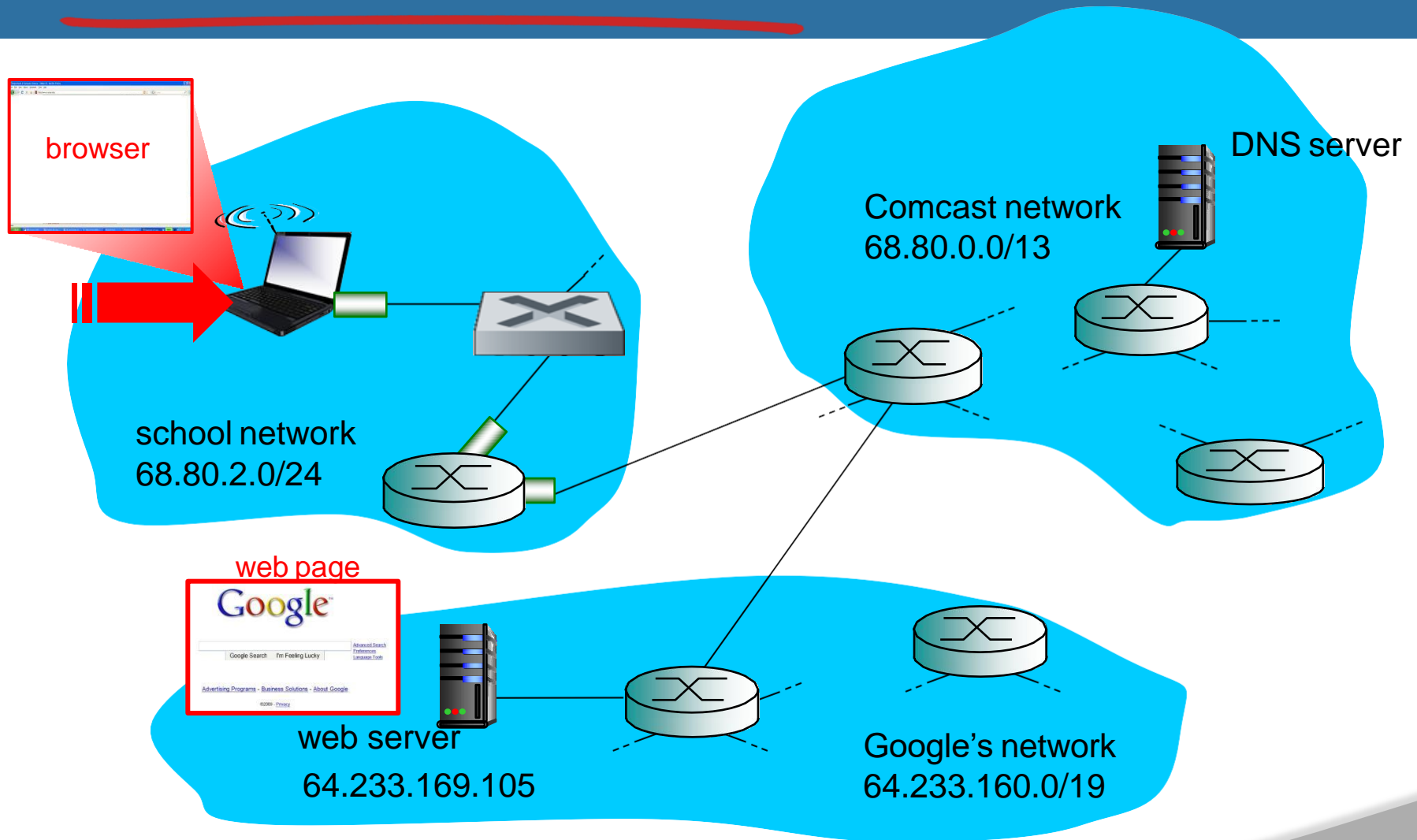
Link layer, LANs: outline

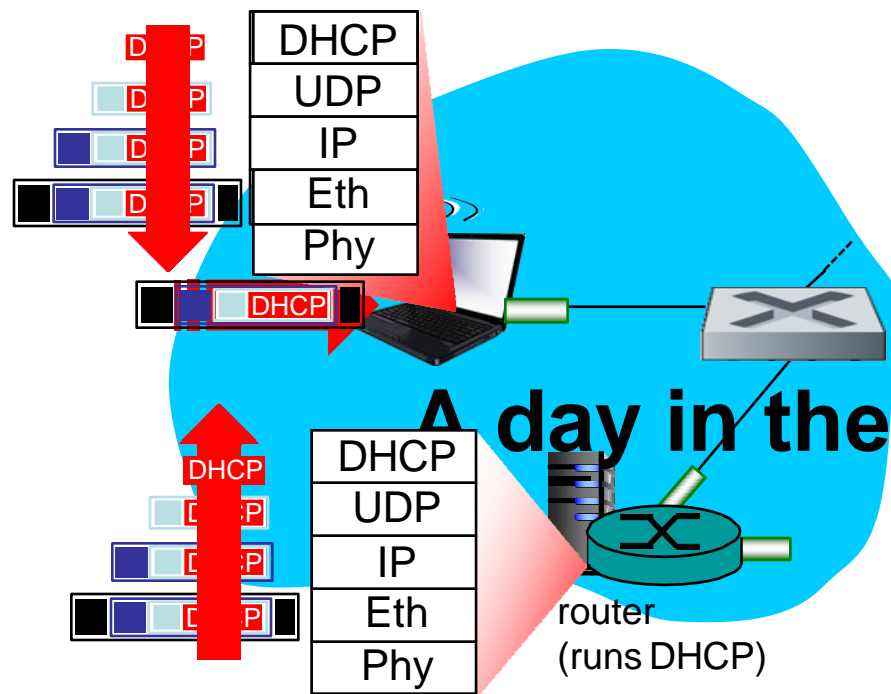
1. introduction, services
2. error detection, correction
3. multiple access protocols
4. link-layer addressing
5. Ethernet, LANs
6. LAN switches
7. a day in the life of a web request

Synthesis: a day in the life of a web request

- journey down protocol stack complete!
 - application, transport, network, link
- putting-it-all-together: synthesis!
 - *goal:* identify, review, understand protocols (at all layers) involved in seemingly simple scenario: requesting www page
 - *scenario:* student attaches laptop to campus network, requests/receives www.google.com

A day in the life: scenario

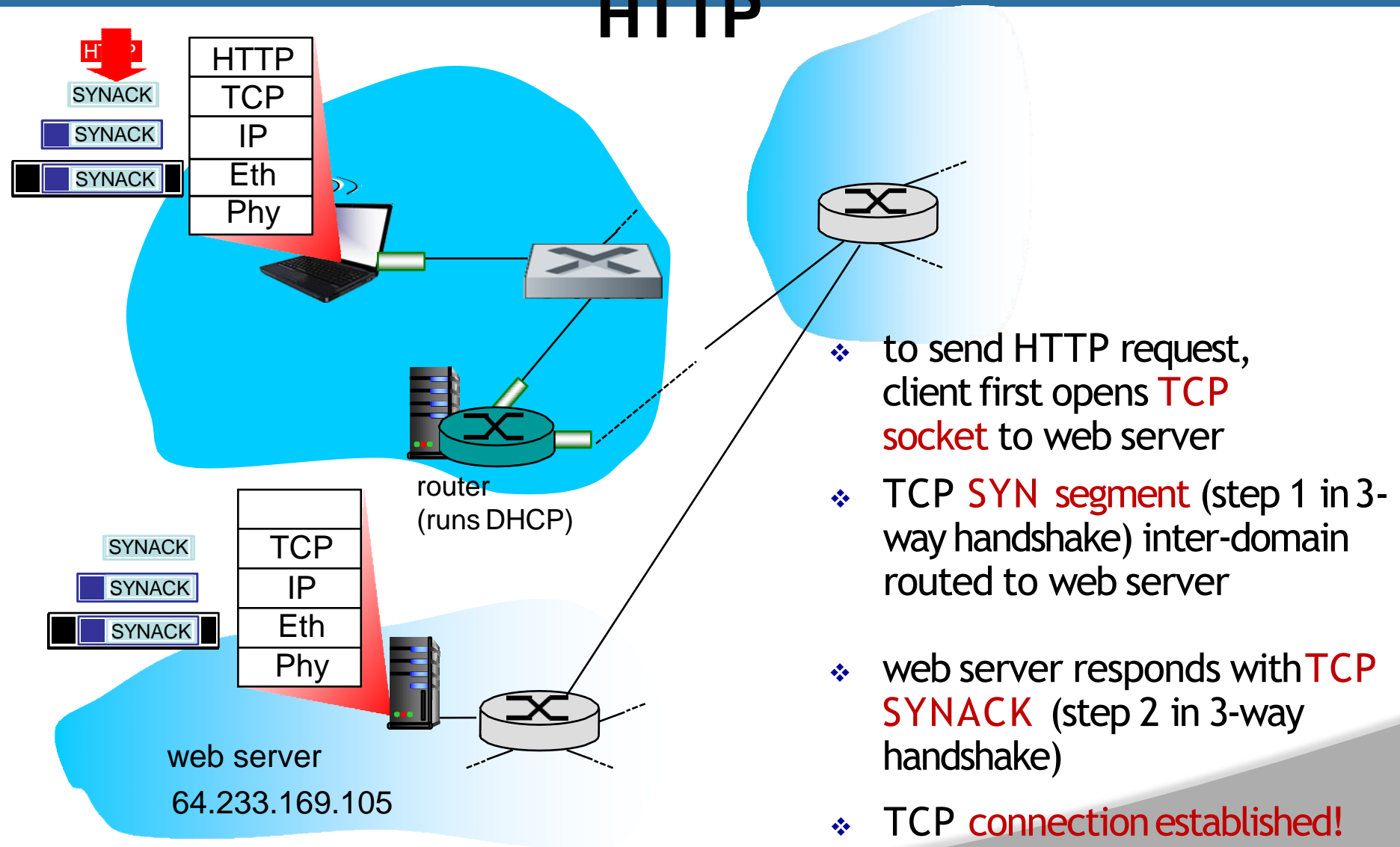




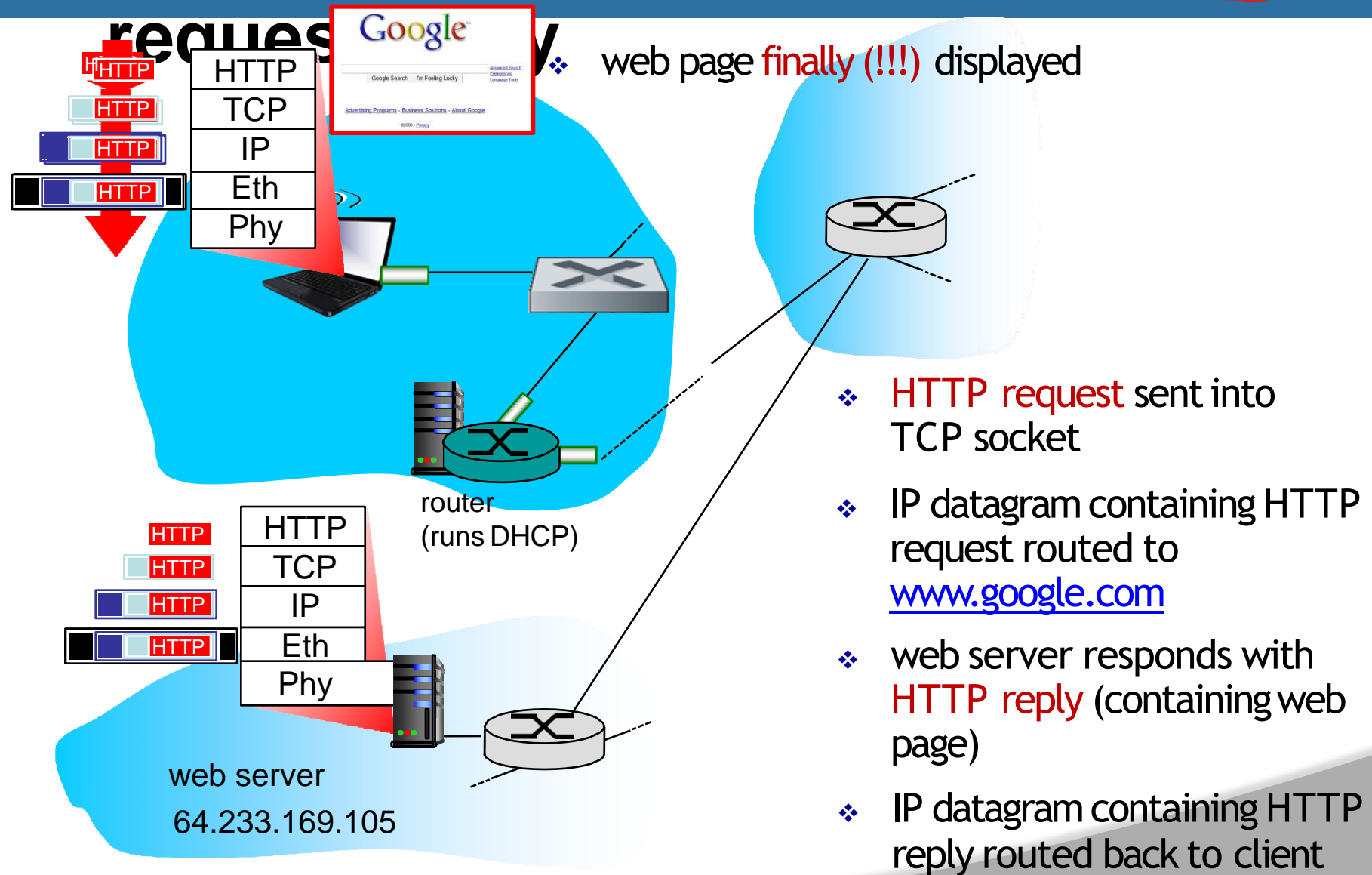
A day in the life... connecting to the Internet

- connecting laptop needs to get its own IP address, addr of first-hop router, addr of DNS server: use **DHCP**
- ❖ DHCP request **encapsulated** in **UDP**, encapsulated in **IP**, encapsulated in **802.3 Ethernet**
- ❖ Ethernet frame **broadcast** (dest: FFFFFFFFFFFFFFFF) on LAN, received at router running **DHCP** server
- ❖ Ethernet **demuxed** to IP demuxed, UDP demuxed to DHCP

A day in the life...TCP connection carrying HTTP



A day in the life... HTTP



Chapter 5: Summary

- principles behind data link layer services:
 - error detection, correction
 - sharing a broadcast channel: multiple access
 - link layer addressing
- instantiation and implementation of various link layer technologies
 - Ethernet
 - switched LANS
- synthesis: a day in the life of a web request

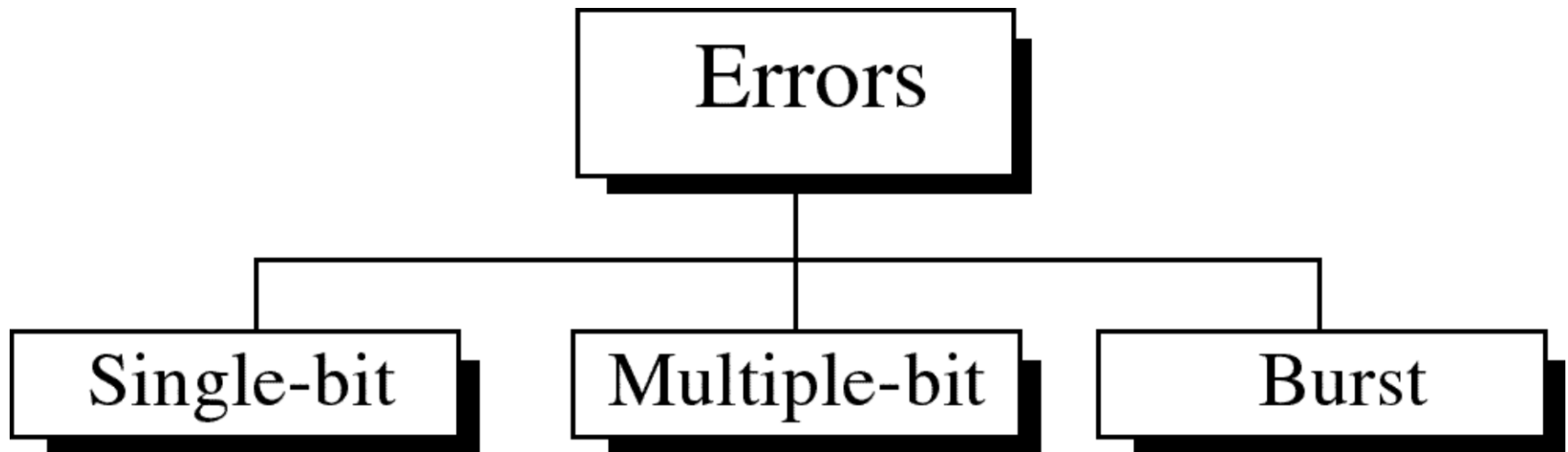
Chapter 5: let's take a breath

- journey down protocol stack *complete* (except PHY)
- solid understanding of networking principles, practice
- could stop here but *lots* of interesting topics!
 - wireless
 - multimedia
 - security
 - network management

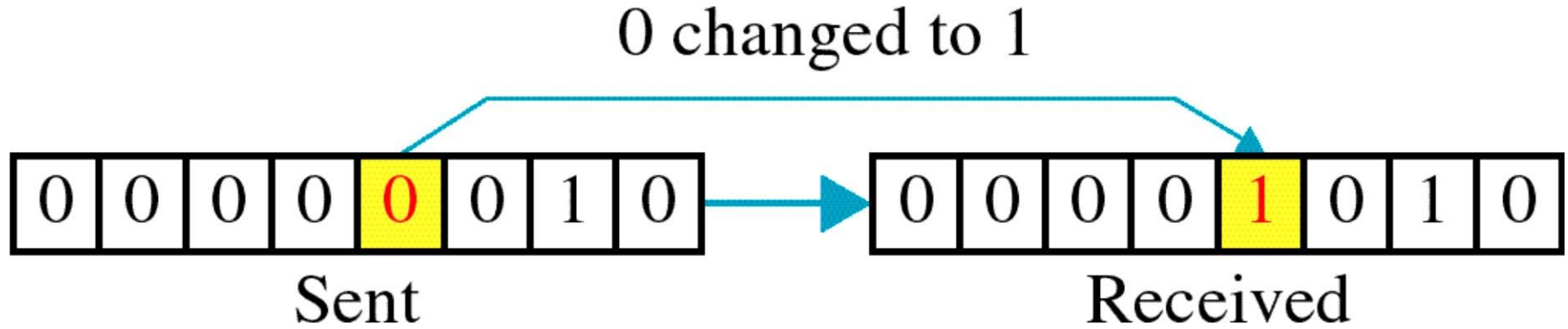
Basic concepts

- ★ Networks must be able to transfer data from one device to another with complete accuracy.
- ★ Data can be corrupted during transmission.
- ★ For reliable communication, errors must be detected and corrected.
- ★ **Error detection and correction** are implemented either at the **data link layer** or the **transport layer** of the OSI model.

Types of Errors



Single-bit error

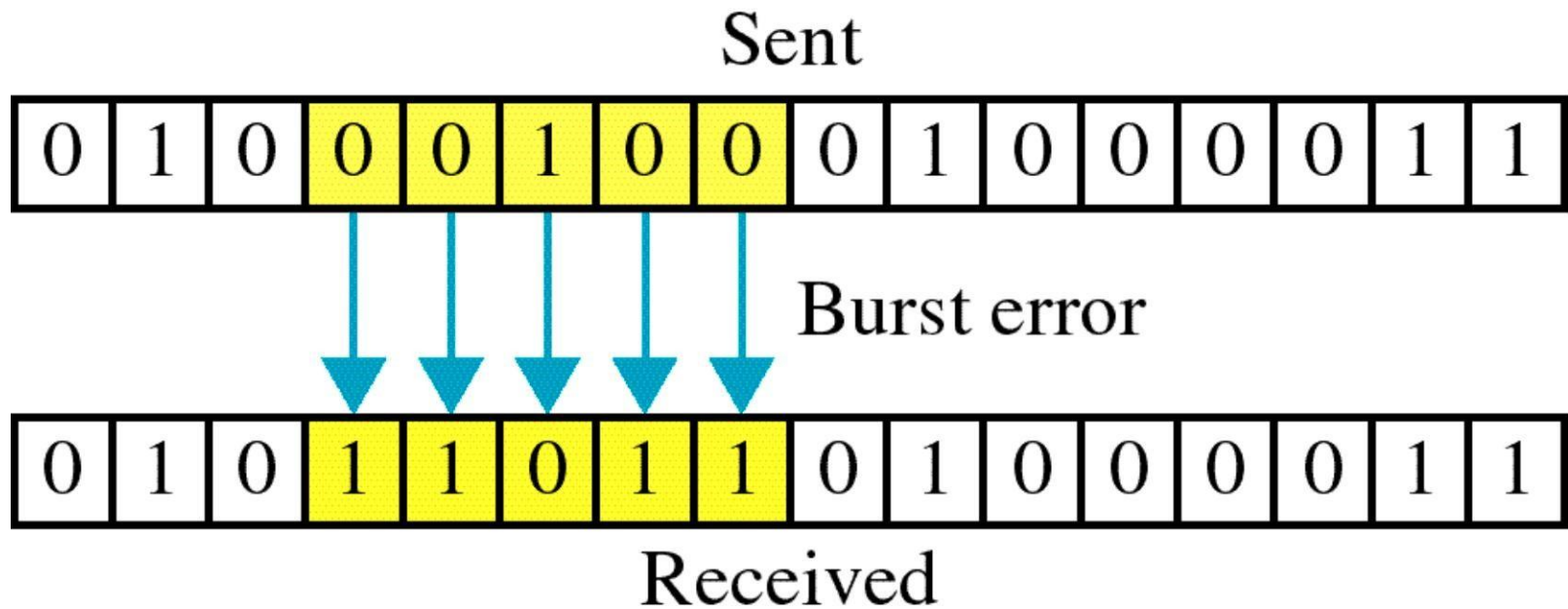


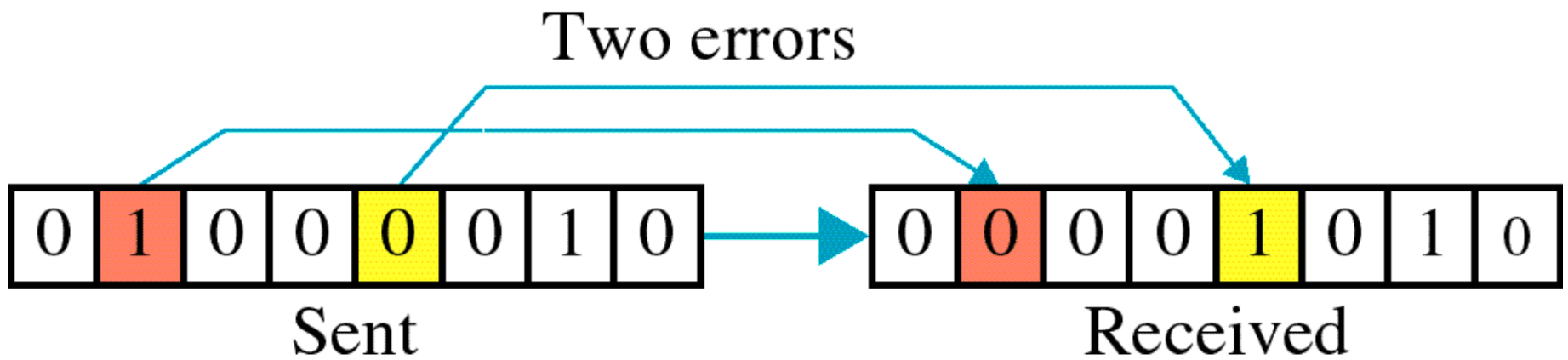
Single bit errors are the **least likely** type of errors in serial data transmission because the noise must have a very short duration which is very rare. However this kind of errors can happen in parallel transmission.

Example:

- ★ If data is sent at 1Mbps then each bit lasts only $1/1,000,000$ sec. or $1\ \mu\text{s}$.
- ★ For a single-bit error to occur, the noise must have a duration of only $1\ \mu\text{s}$, which is very rare.

Burst error





The term **burst error** means that two or more bits in the data unit have changed from 1 to 0 or from 0 to 1.

Burst errors does not necessarily mean that the errors occur in consecutive bits, the length of the burst is measured from the first corrupted bit to the last corrupted bit. Some bits in between may not have been corrupted.

★ **Burst error is most likely to happen in serial transmission** since the duration of noise is normally longer than the duration of a bit.

★ The number of bits affected depends on the data rate and duration of noise.

Example:

□ If data is sent at rate = 1Kbps then a noise of 1/100 sec can affect 10 bits. $(1/100 * 1000)$

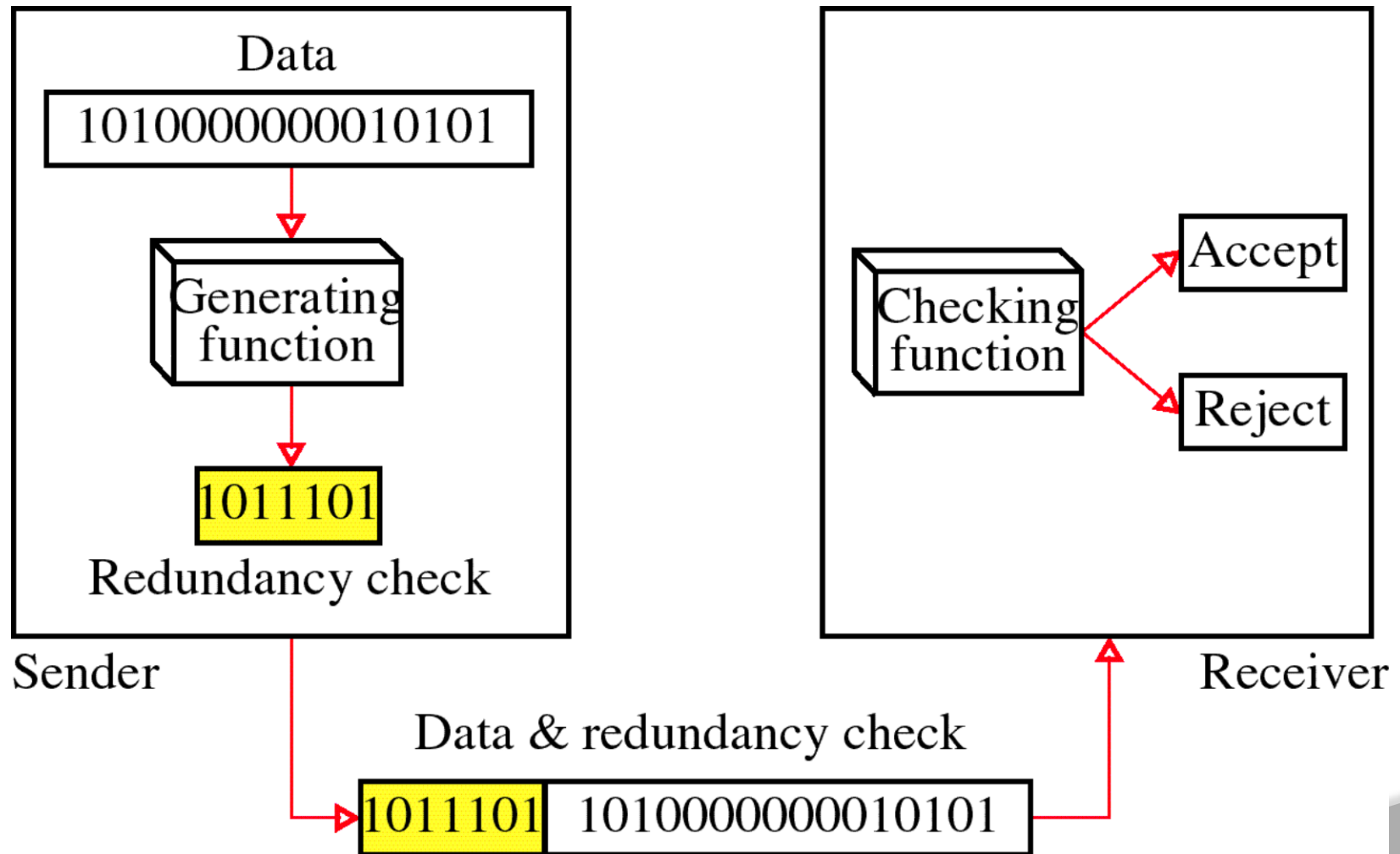
□ If same data is sent at rate = 1Mbps then a noise of 1/100 sec can affect 10,000 bits. $(1/100 * 10^6)$

Error detection

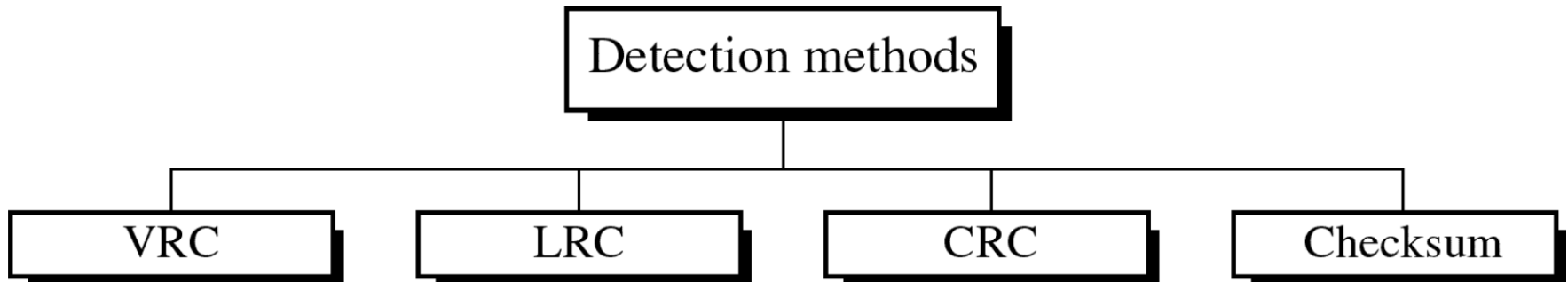
Error detection means to decide whether the received data is correct or not without having a copy of the original message.

Error detection **uses the concept of redundancy, which means** adding extra bits for detecting errors at the destination.

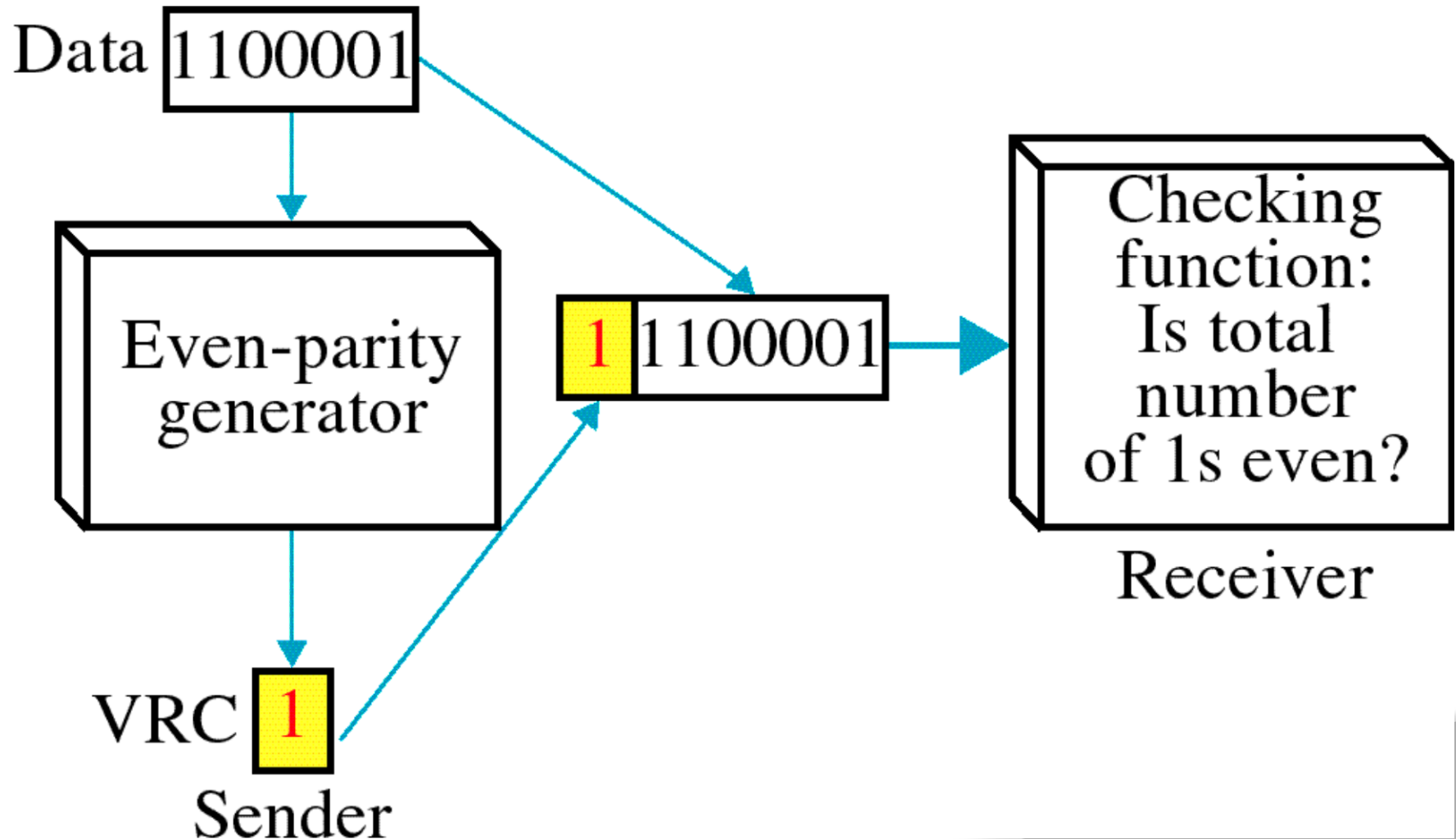
Redundancy



Four types of redundancy checks are used in data communications



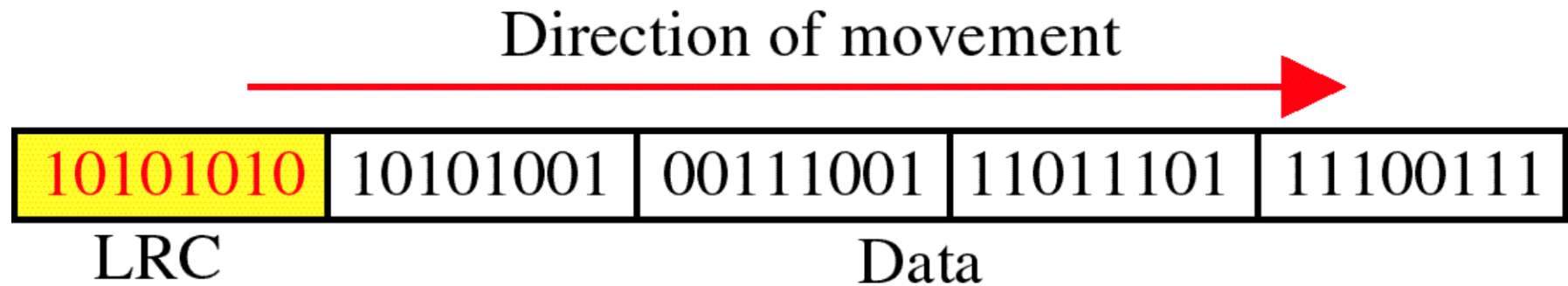
Vertical Redundancy Check VRC



Performance

- It can detect single bit error
- It can detect burst errors only if the total number of errors is odd.

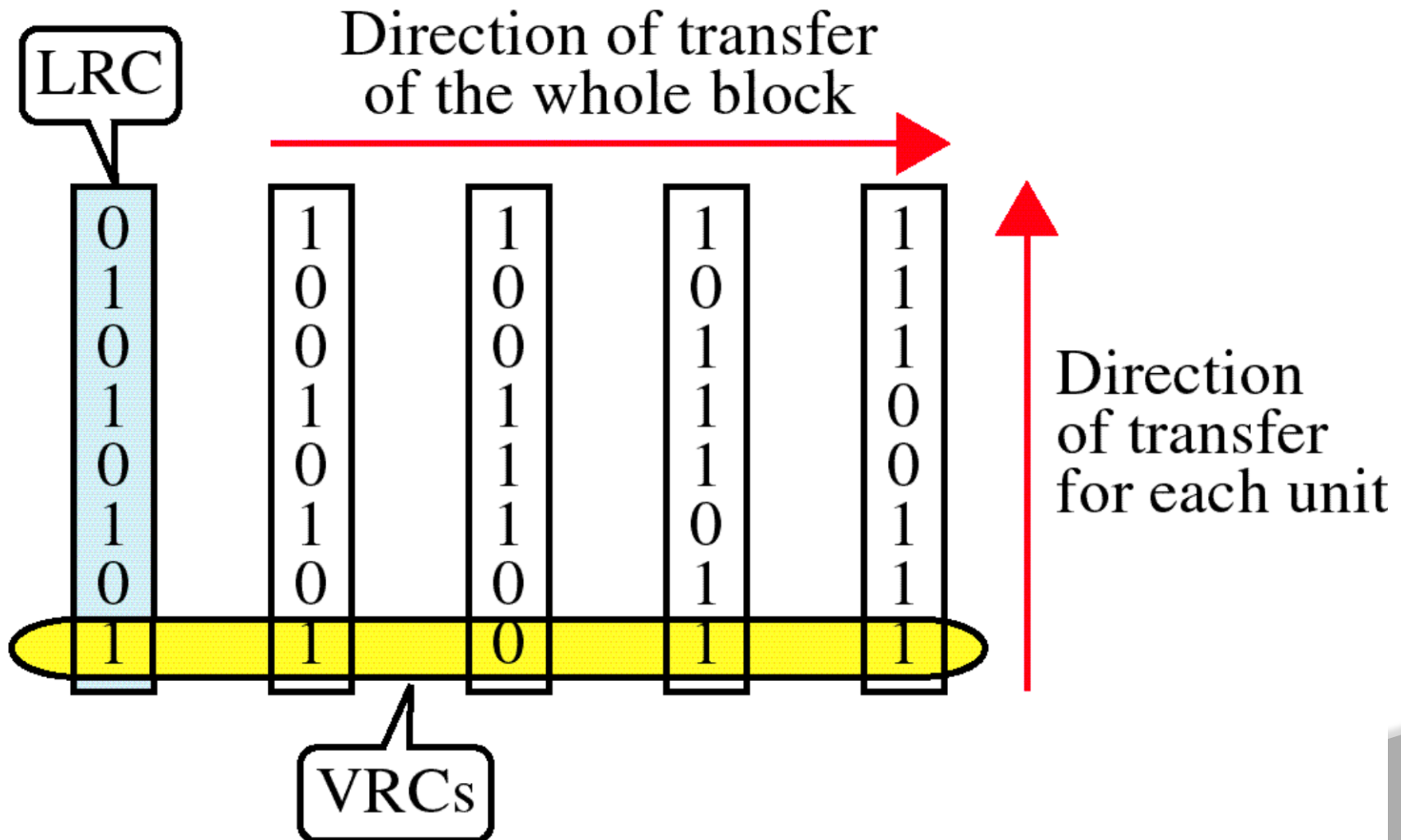
Longitudinal Redundancy Check LRC



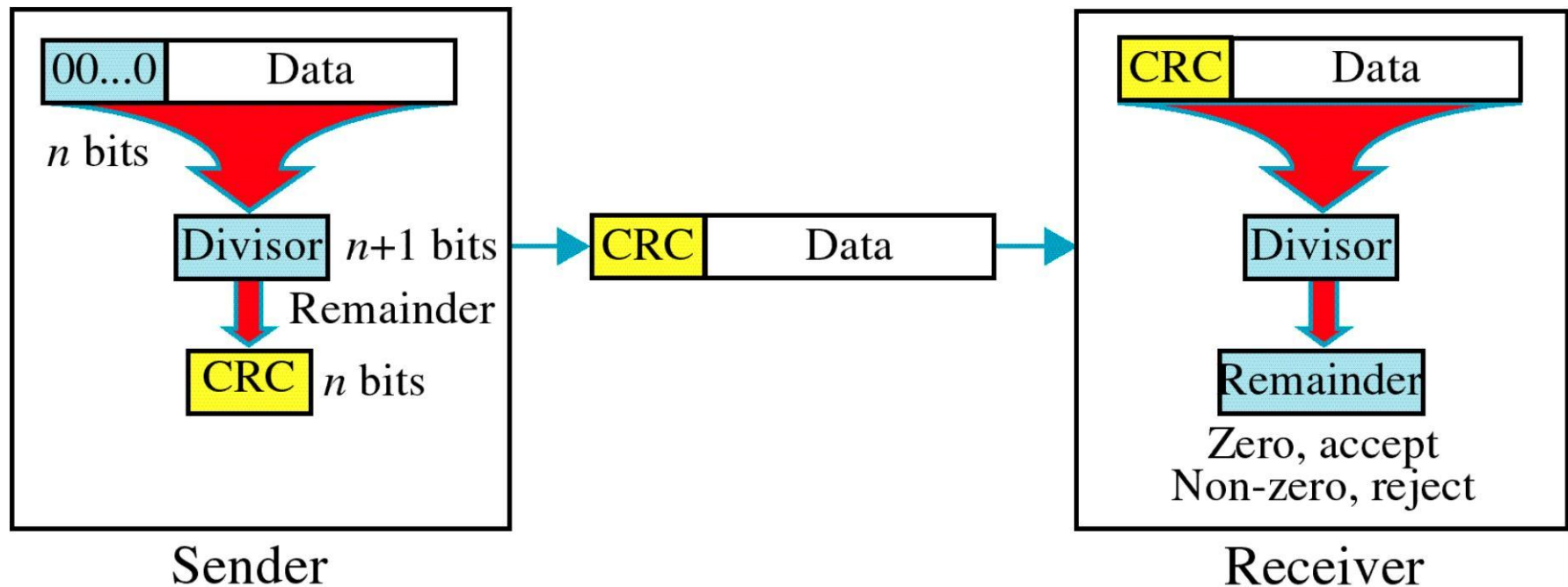
Performance

- L C R increases the likelihood of detecting burst errors.
- If two bits in one data units are damaged and two bits in exactly the same positions in another data unit are also damaged, the LRC checker will not detect an error.

VRC and LRC



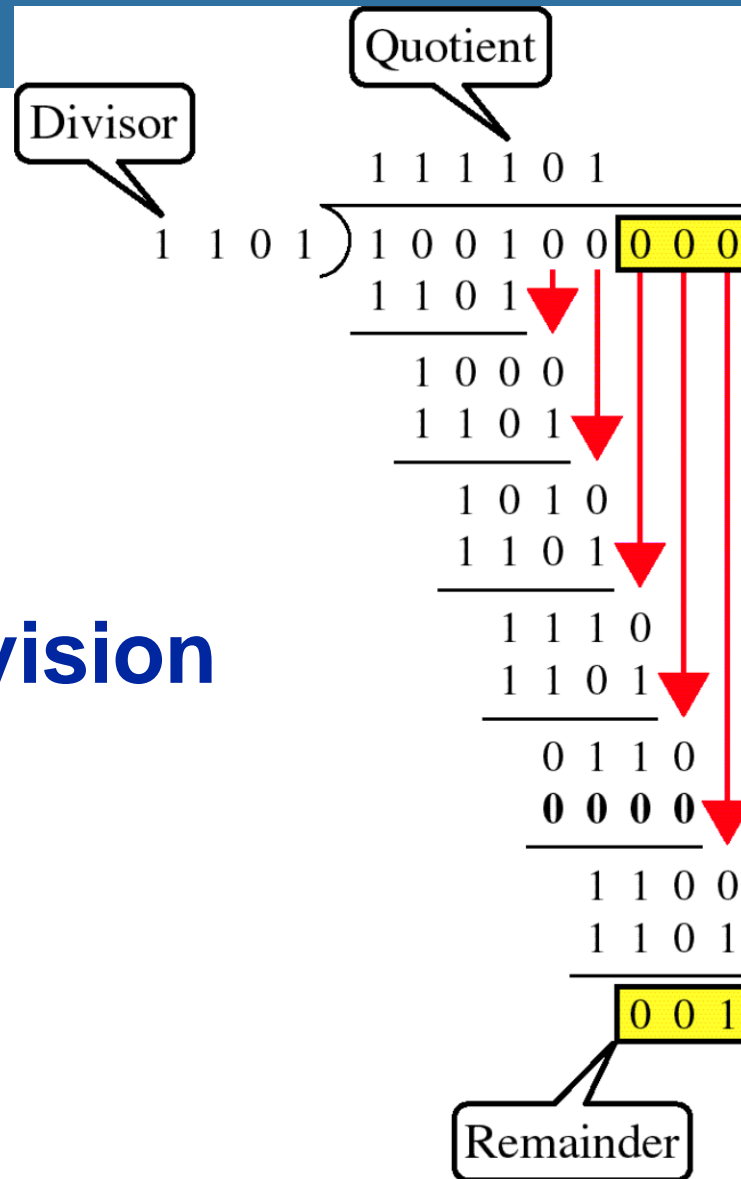
Cyclic Redundancy Check CRC



Cyclic Redundancy Check

- Given a k -bit frame or message, the transmitter generates an n -bit sequence, known as a *frame check sequence (FCS)*, so that the resulting frame, consisting of $(k+n)$ bits, is exactly divisible by some predetermined number.
- The receiver then divides the incoming frame by the same number and, if there is no remainder, assumes that there was no error.

Binary Division



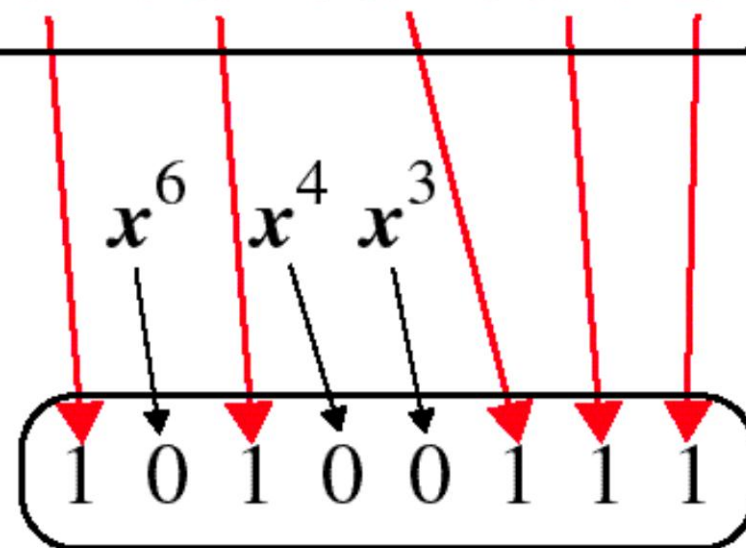
Polynomial

$$x^7 + x^5 + x^2 + x + 1$$

Polynomial and Divisor

Polynomial

$$x^7 + x^5 + x^2 + x + 1$$



Divisor

Standard Polynomials

CRC-12

$$x^{12} + x^{11} + x^3 + x + 1$$

CRC-16

$$x^{16} + x^{15} + x^2 + 1$$

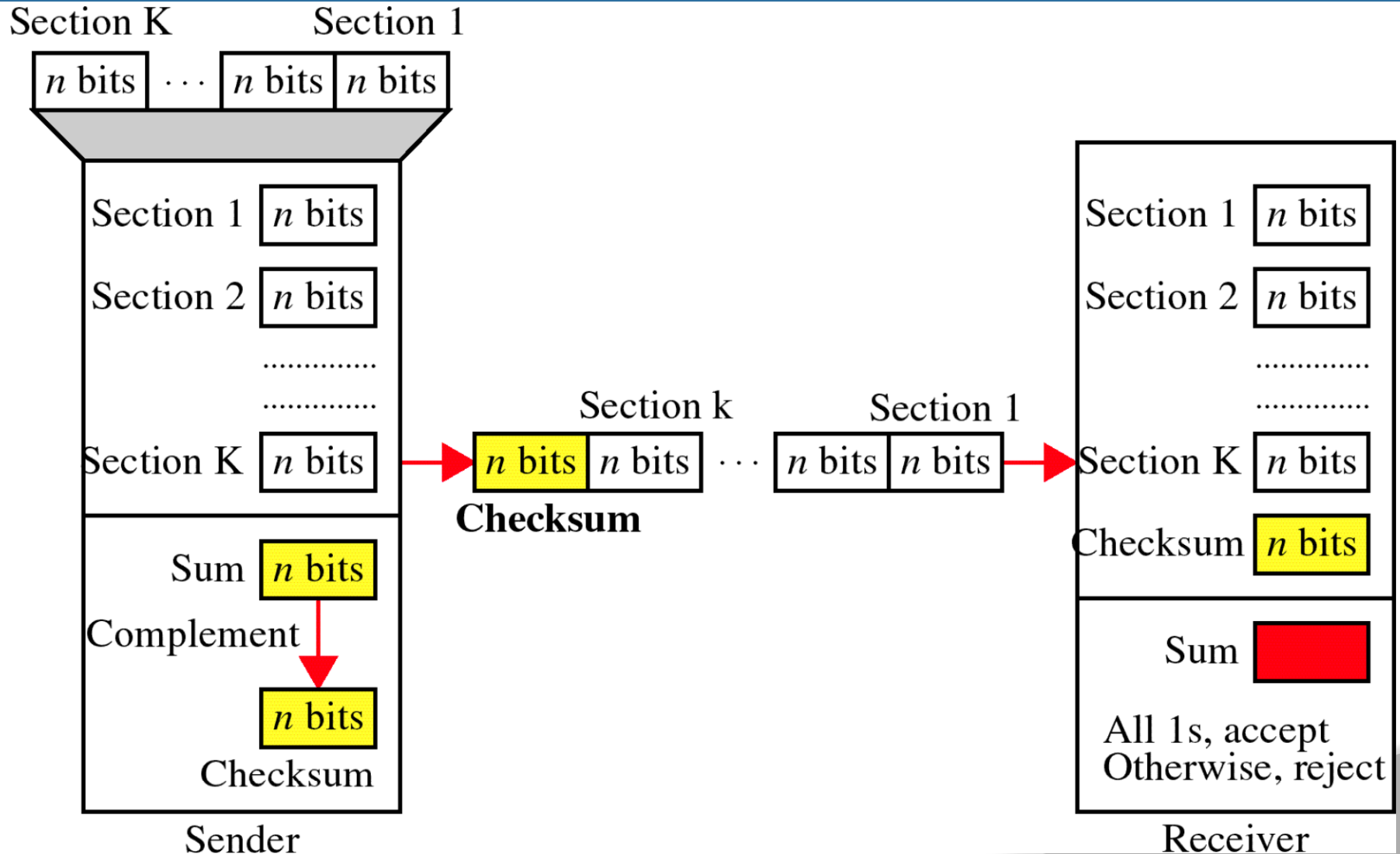
CRC-ITU

$$x^{16} + x^{12} + x^5 + 1$$

CRC-32

$$x^{32} + x^{26} + x^{23} + x^{22} + x^{16} + x^{12} + x^{11} + x^{10} + x^8 + x^7 + x^5 + x^4 + x^2 + x + 1$$

Checksum



At the sender

- ➡ The unit is divided into k sections, each of n bits.
- ➡ All sections are added together using one's complement to get the sum.
- ➡ The sum is complemented and becomes the checksum.
- ➡ The checksum is sent with the data

At the receiver

- ➡ The unit is divided into k sections, each of n bits.
- ➡ All sections are added together using one's complement to get the sum.
- ➡ The sum is complemented.
- ➡ If the result is zero, the data are accepted otherwise, they are rejected.

Performance

- ❓ The checksum detects all errors involving an odd number of bits.
- ❓ It detects most errors involving an even number of bits.
- ❓ If one or more bits of a segment are damaged and the corresponding bit or bits of opposite value in a second segment are also damaged, the sums of those columns will not change and the receiver will not

detect a problem.

Error Correction

It can be handled in two ways:

- 1) receiver can have the sender retransmit the entire data unit.
- 2) The receiver can use an error-correcting code, which automatically corrects certain errors.

Single-bit error correction

To correct an error, the receiver reverses the value of the altered bit. To do so, it must know which bit is in error.

Number of redundancy bits needed

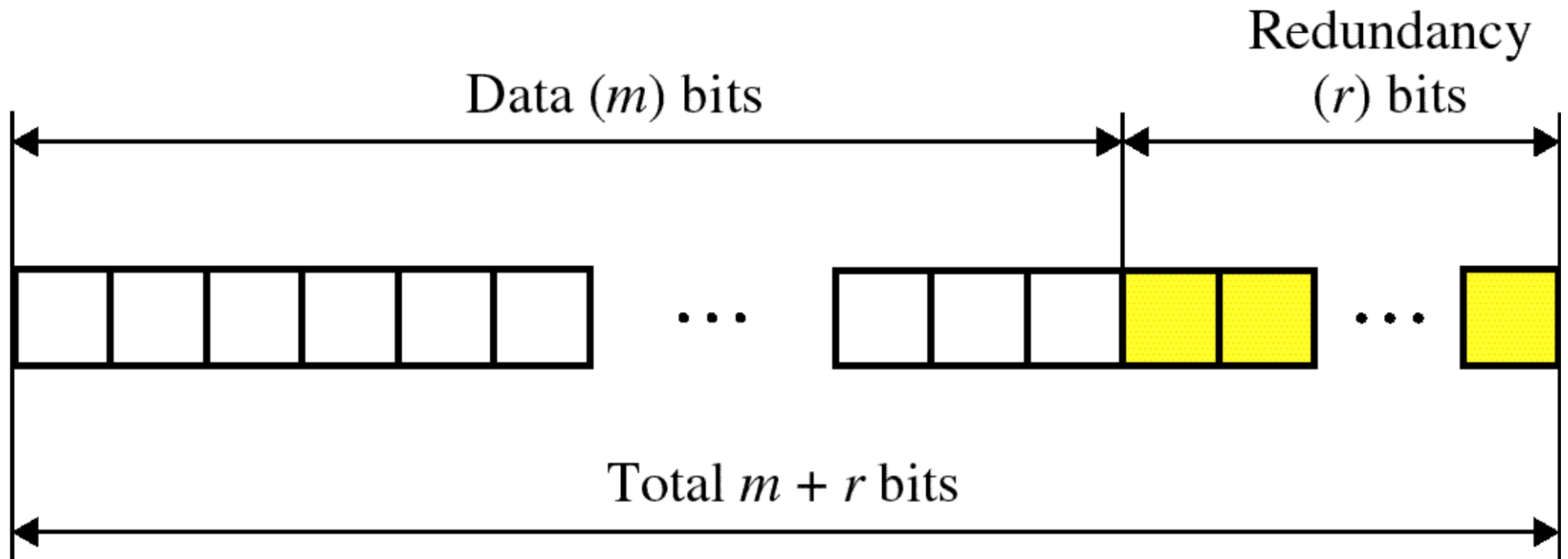
- Let data bits = m
- Redundancy bits = r

\therefore Total message sent = $m+r$

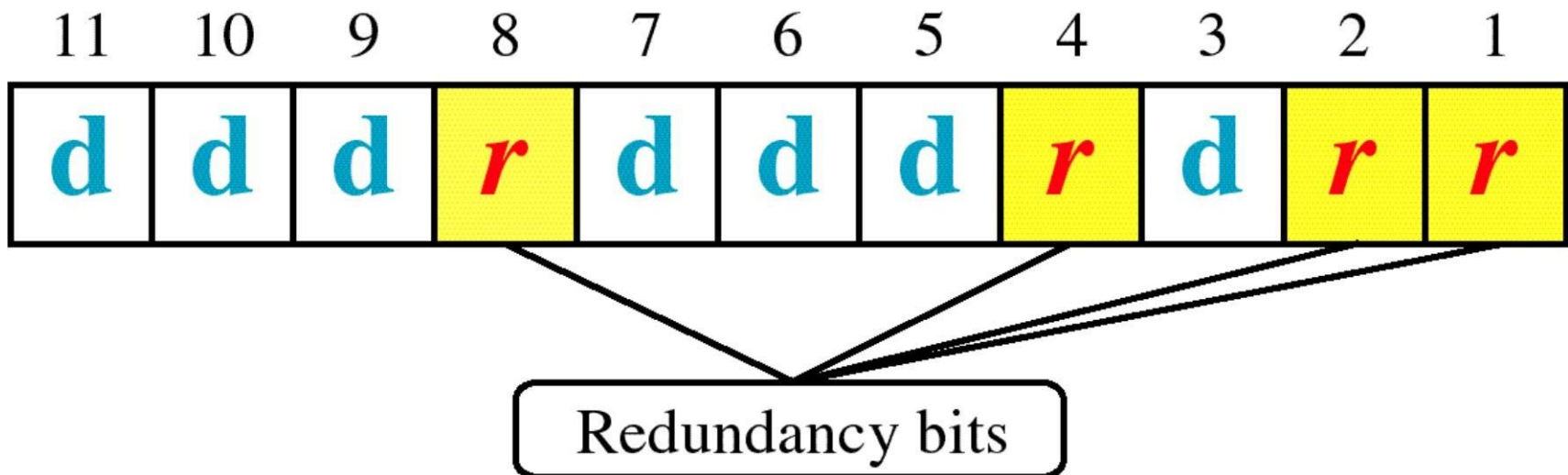
The value of r must satisfy the following relation:

$$2^r \geq m+r+1$$

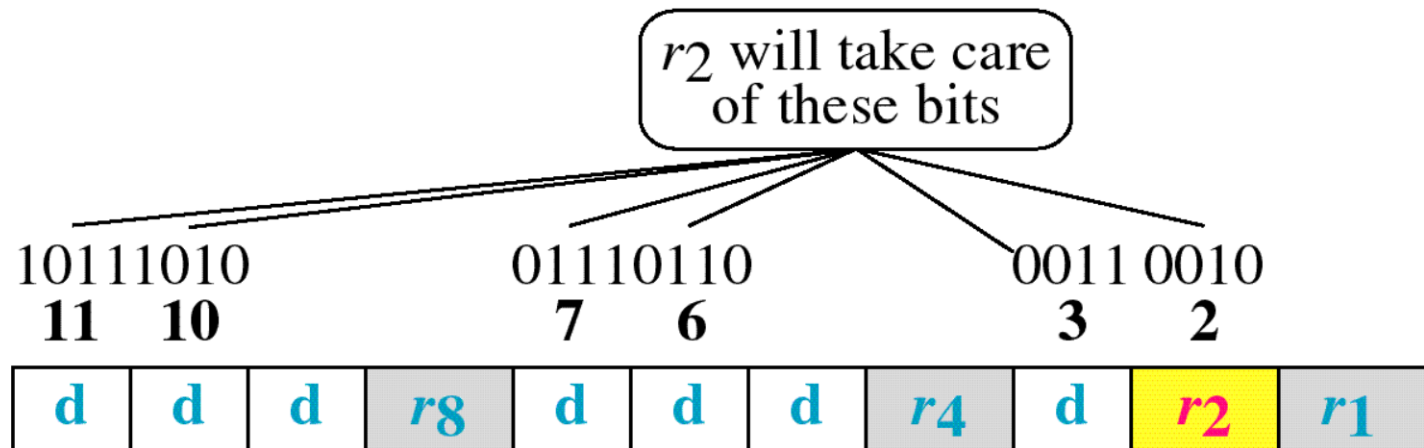
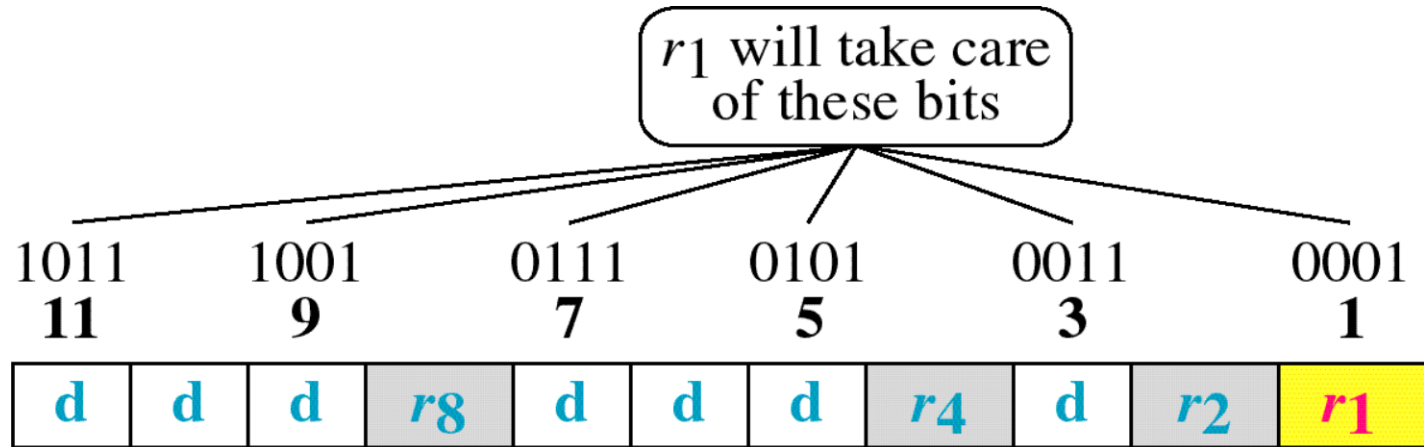
Error Correction



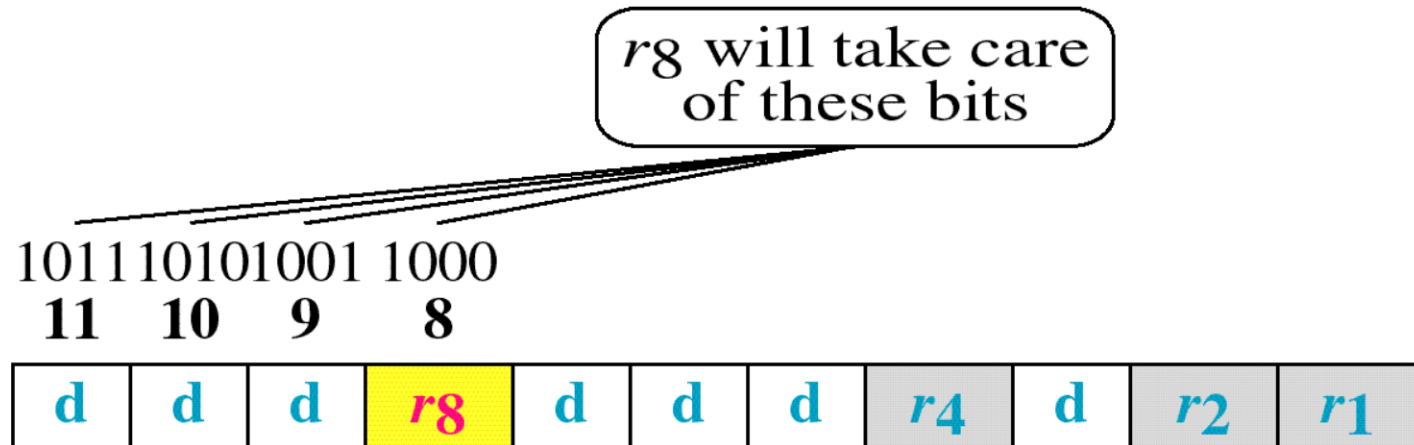
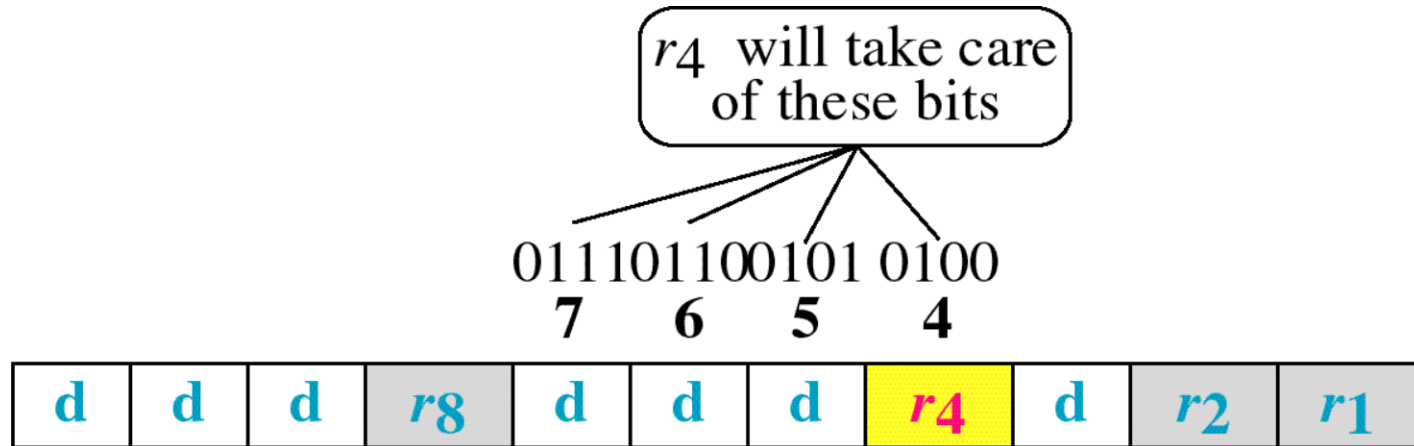
Hamming Code



Hamming Code



Hamming Code



Example of Hamming Code

Data: 1 0 0 1 1 0 1



Data

1	0	0		1	1	0		1		
---	---	---	--	---	---	---	--	---	--	--

Adding r_1

1	0	0		1	1	0		1		1
---	---	---	--	---	---	---	--	---	--	---

Adding r_2

1	0	0		1	1	0		1	0	1
---	---	---	--	---	---	---	--	---	---	---

Adding r_4

1	0	0		1	1	0	0	1	0	1
---	---	---	--	---	---	---	---	---	---	---

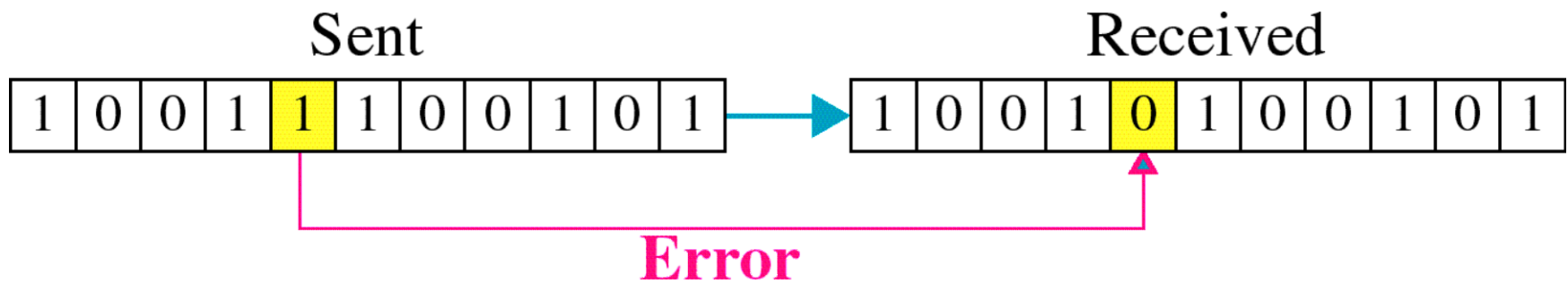
Adding r_8

1	0	0	1	1	1	0	0	1	0	1
---	---	---	---	---	---	---	---	---	---	---

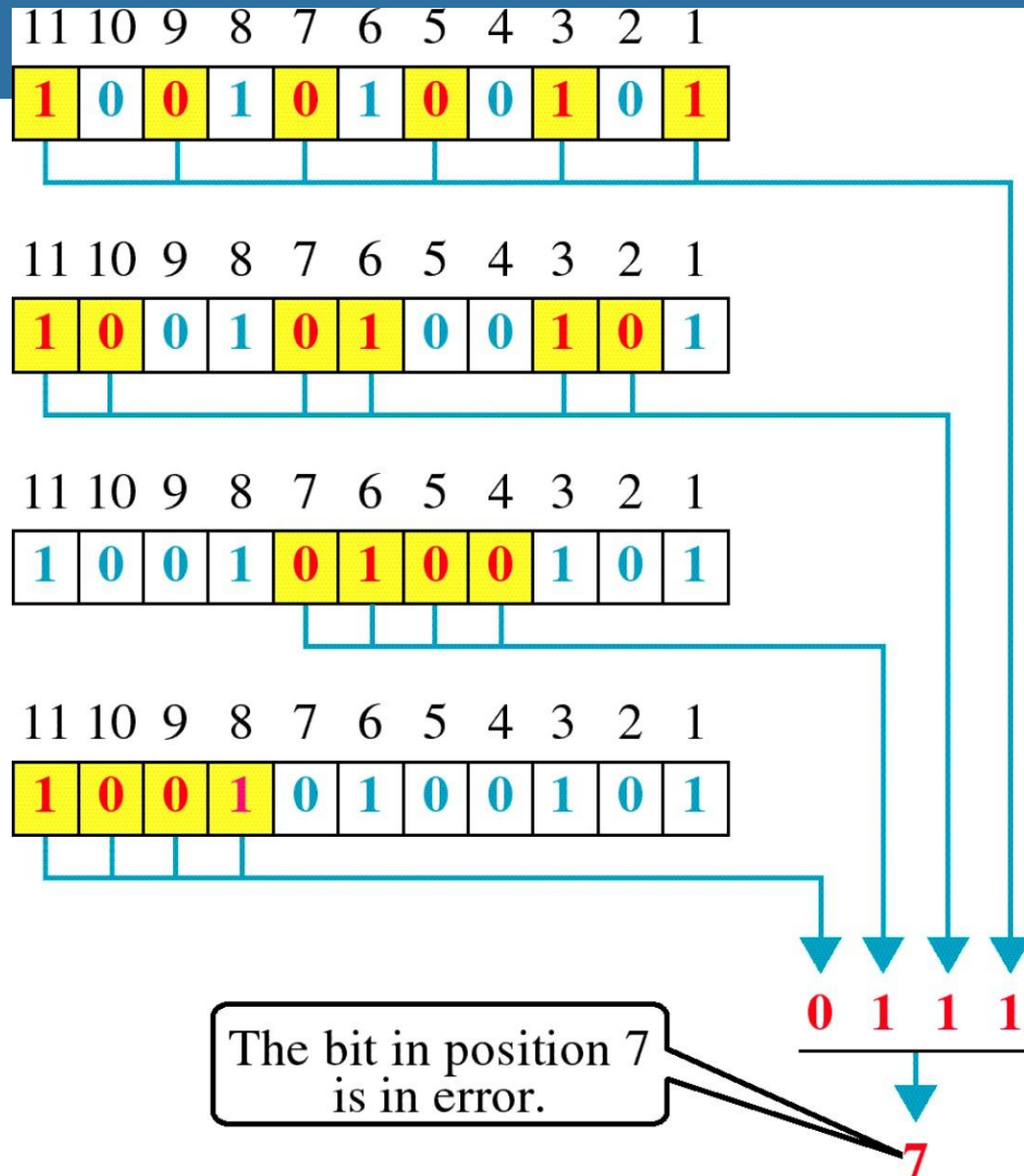


Code: 1 0 0 1 1 1 0 0 1 0 1

Single-bit error



Error Detection



Data Link Control

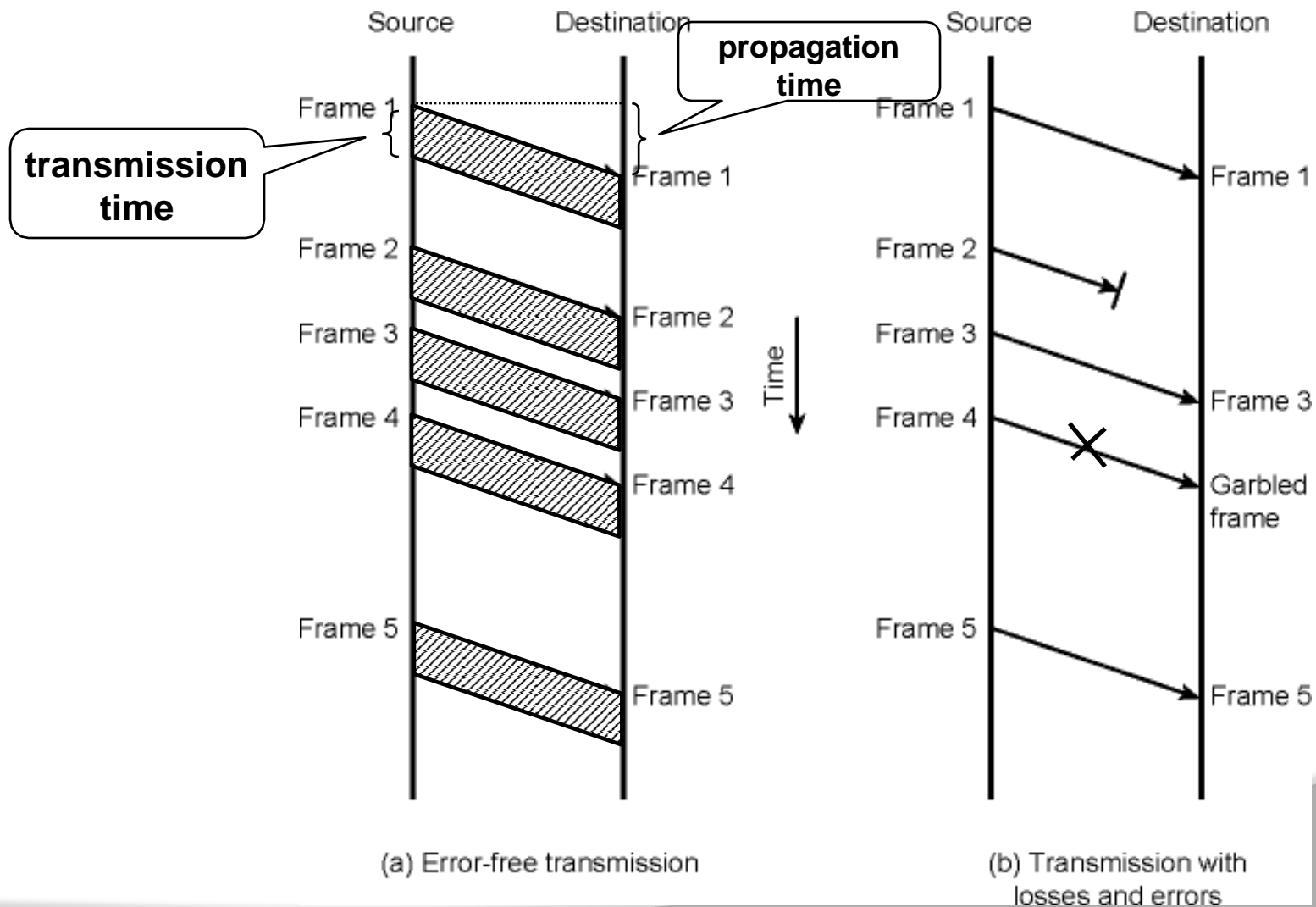
Announcements

- Midterm: November 28, Monday, 11:40 – 13:30
 - Places:
 - FENS G032 if (lastName[0] >= 'A' && lastName[0] <= 'D')
 - FASS G022 if (lastName[0] >= 'E' && lastName[0] <= 'Ö')
 - FASS G049 if (lastName[0] >= 'P' && lastName[0] <= 'Z')
- Exam will be closed book, closed notes
 - calculators are allowed
 - you are responsible all topics I covered in the class even if some of them are not in the book (I sometimes used other books) and not in the ppt files (I sometimes used board and showed applications on the computer)

Flow Control

- In Data Link Layer, we deal with issues related to point to point links
 - Flow control is one of these issues
- Flow control is needed since the sending entity should not overwhelm the receiving entity
 - Recipient needs some time to process incoming packets
 - If sender sends faster than recipient

Model of Frame Transmission

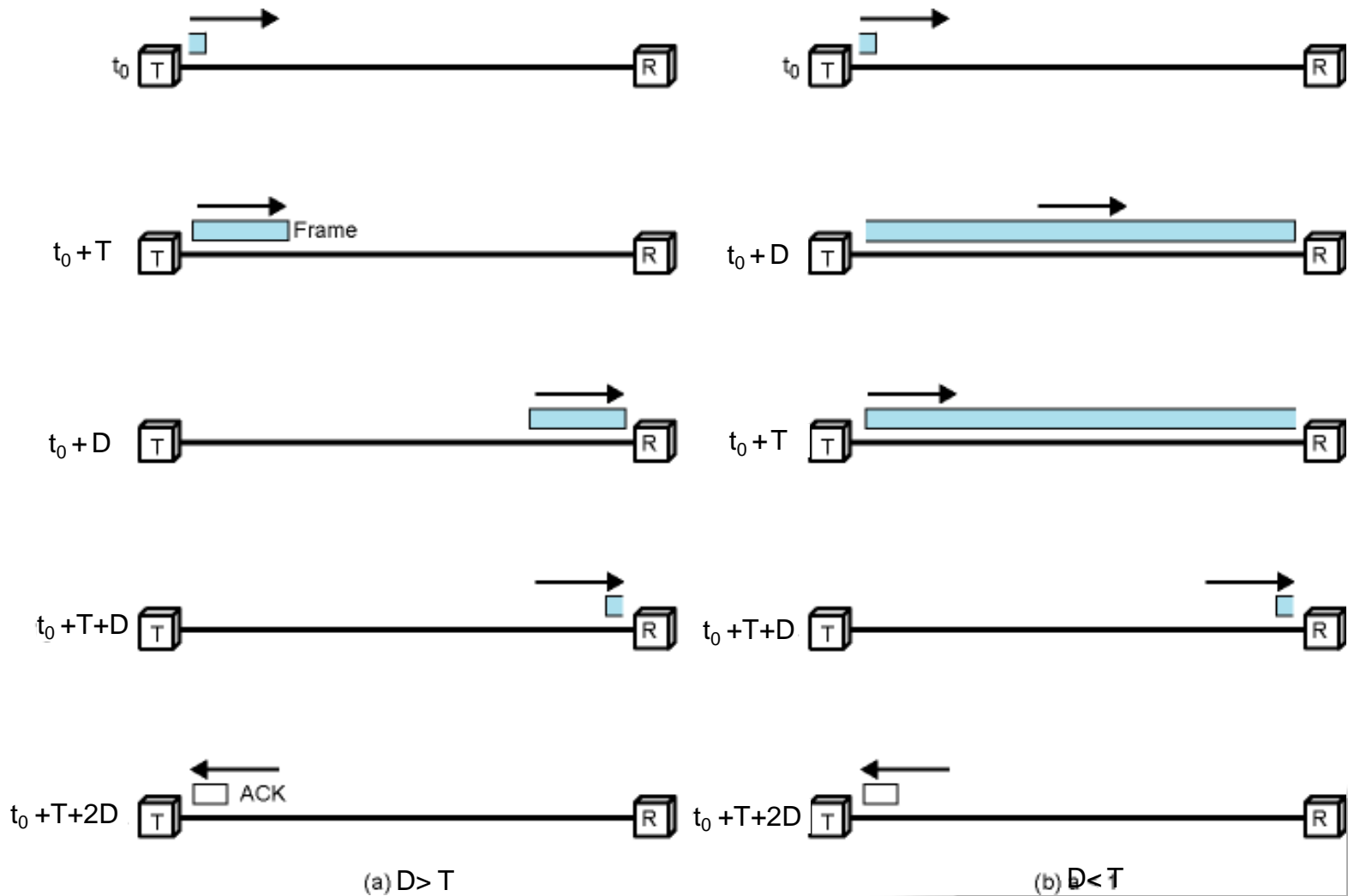


Stop and Wait Flow Control

- Source transmits frame
- Destination receives frame and replies with acknowledgement (ACK)
- Source waits for ACK before sending next frame
- Destination can stop flow by not sending ACK
- Works well for large frames
- Inefficient for smaller frames

Stop and Wait Flow Control

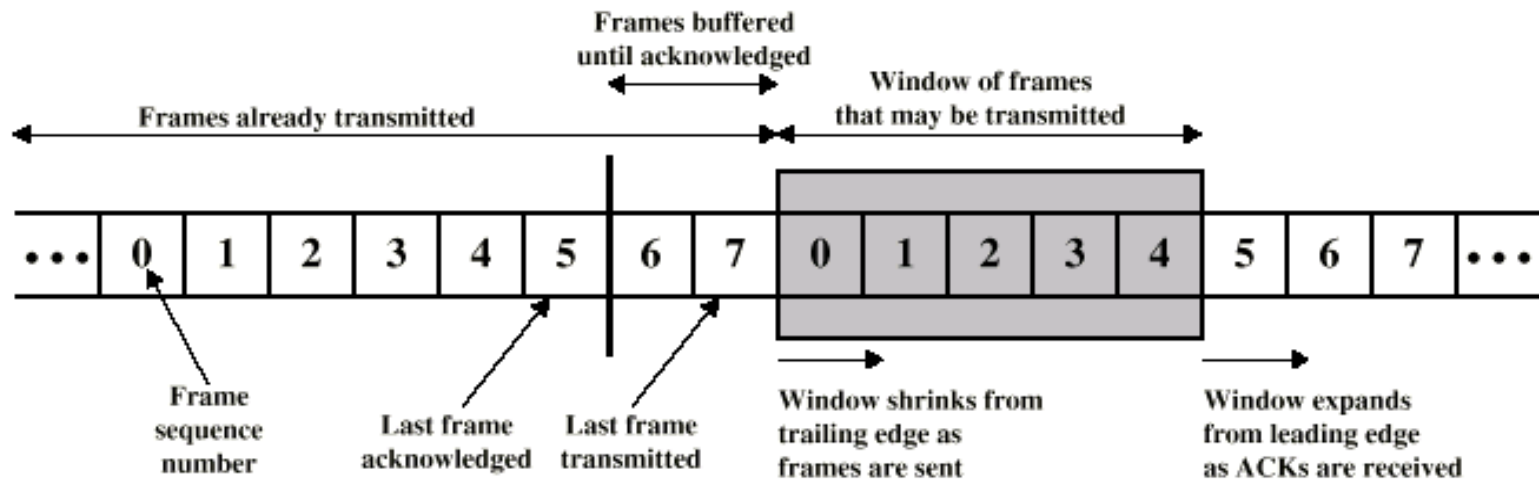
- However, generally large block of data split into small frames
 - Called “Fragmentation”
 - Advantages are
 - Limited buffer size at receiver
 - Errors detected sooner (when whole frame received)
 - On error, retransmission of smaller frames is needed
 - Prevents one station occupying medium for long periods
- Channel Utilization is higher when
 - the transmission time is longer than the
propagation time
 - frame length is larger than the bit length of the



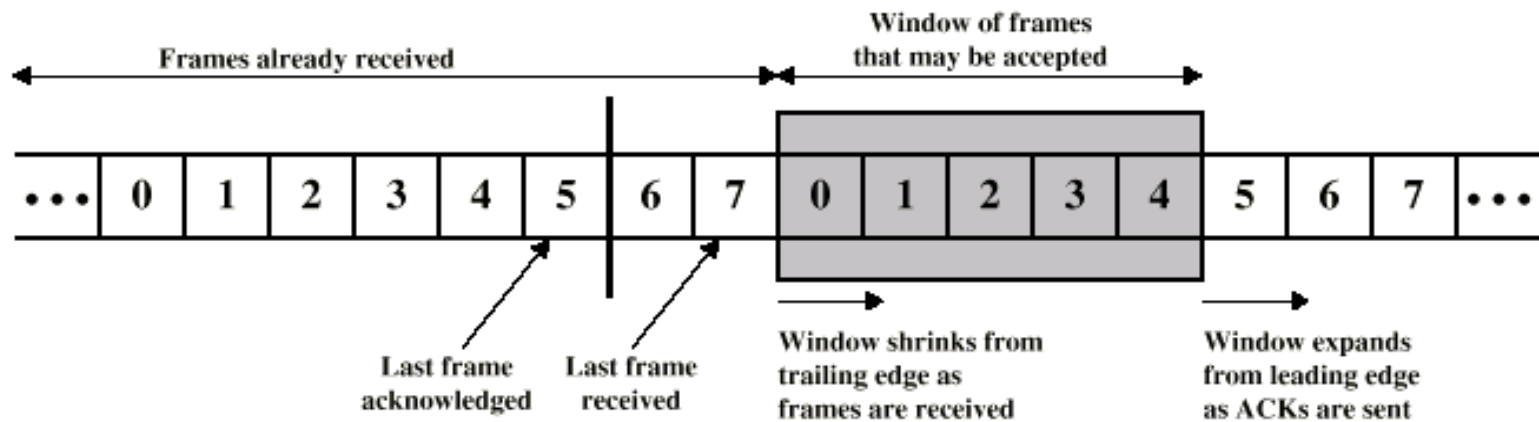
propagation time = D , transmission time = T

Sliding Window Flow Control

- The problem of “Stop and Wait” is not able to send multiple packets
- Sliding Window Protocol allows multiple frames to be in transit
- Receiver has buffer of W (called window size) frames
- Transmitter can send **up to W frames** without ACK
- Each frame is numbered
 - Sequence number bounded by size of the sequence number field (k bits)
 - thus frames are numbered modulo 2^k ($0 \dots 2^k-1$)
- ACK includes number of next frame expected



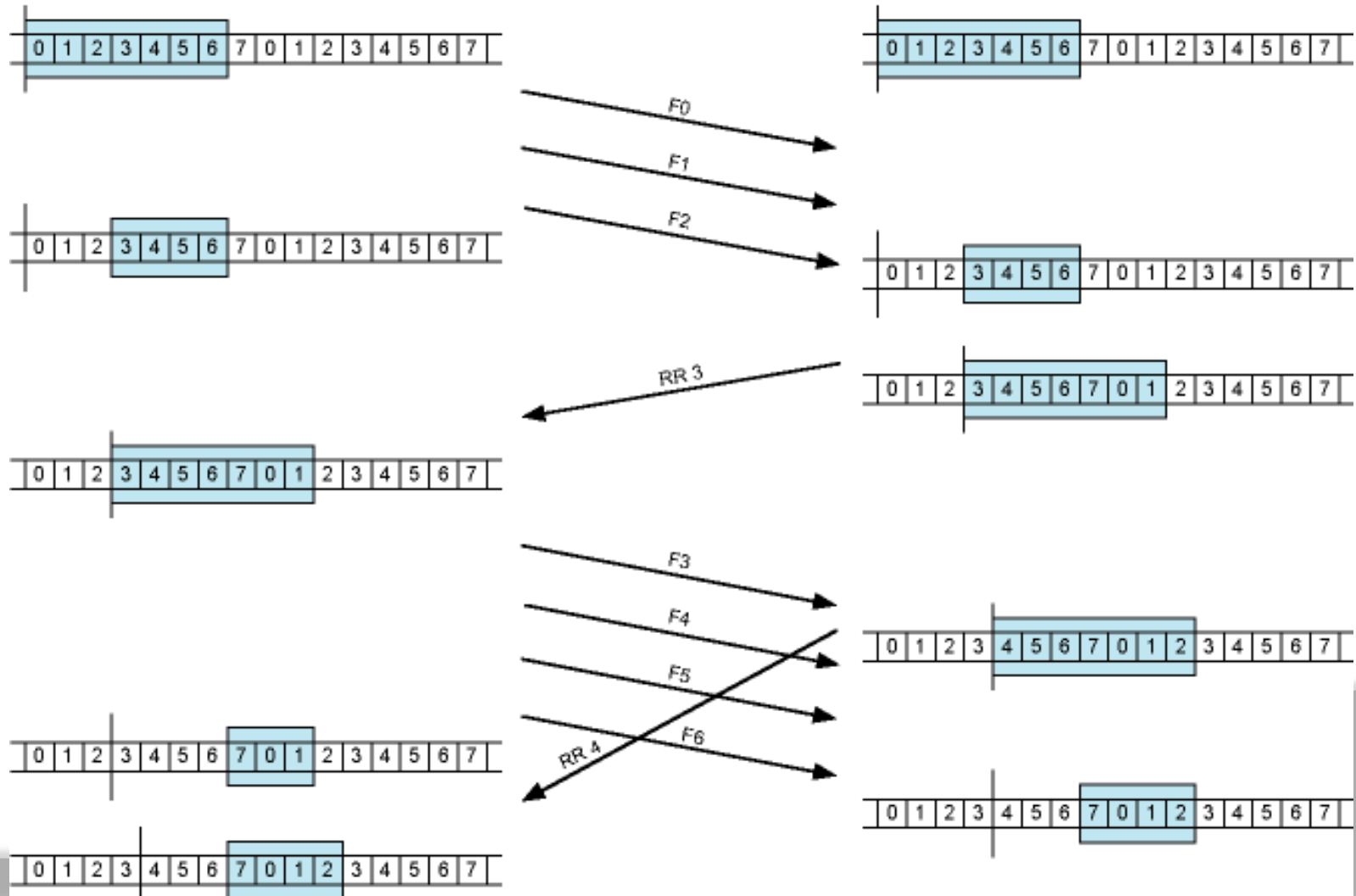
(a) Sender's perspective



(b) Receiver's perspective

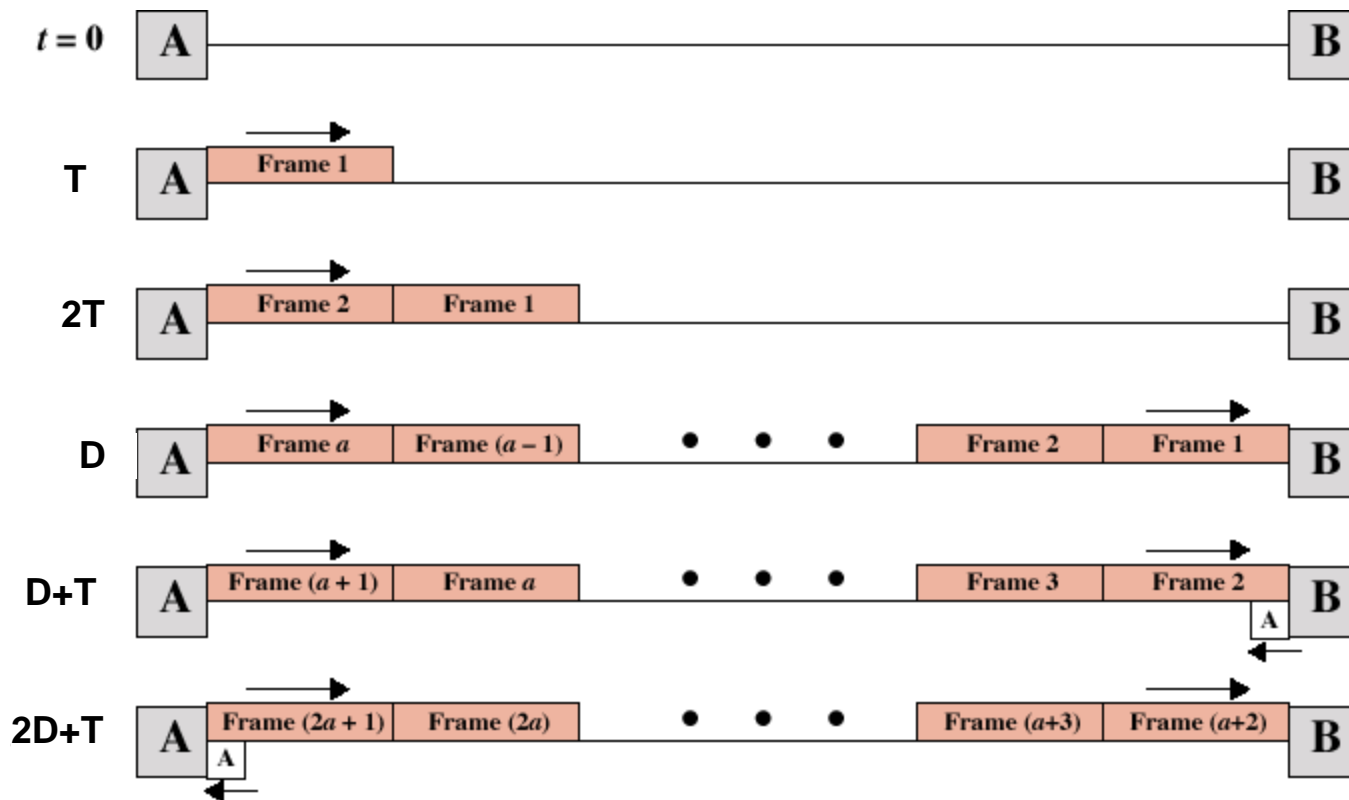
Source System A

Destination System B

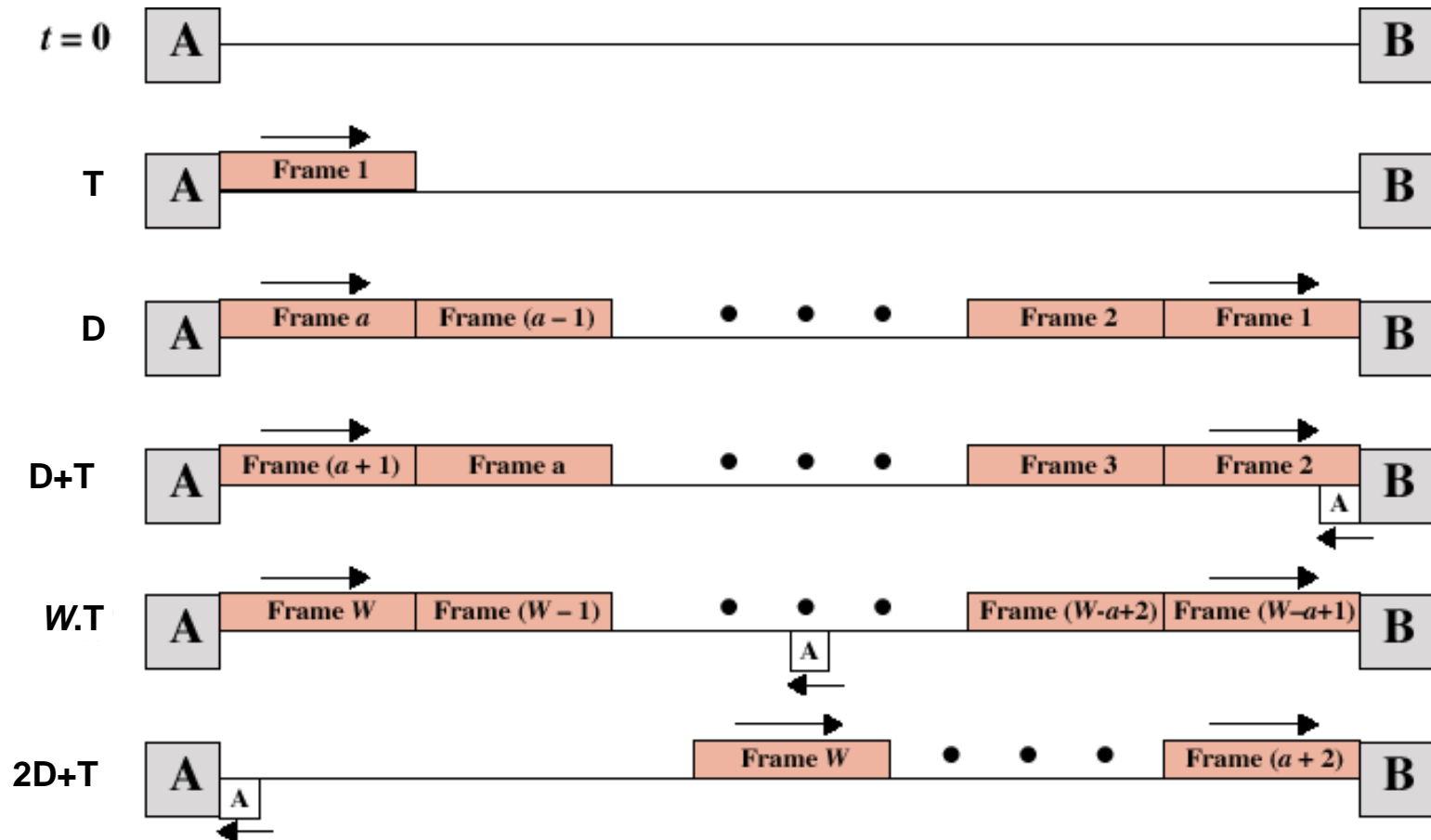


Sliding Windows Performance - 1

- two cases: $W \geq 2a+1$ and $W < 2a+1$, where $a=D/T$
- details are on board



(a) $W \geq 2a + 1$ ($W.T \geq 2D+T$)



Error Detection and Control

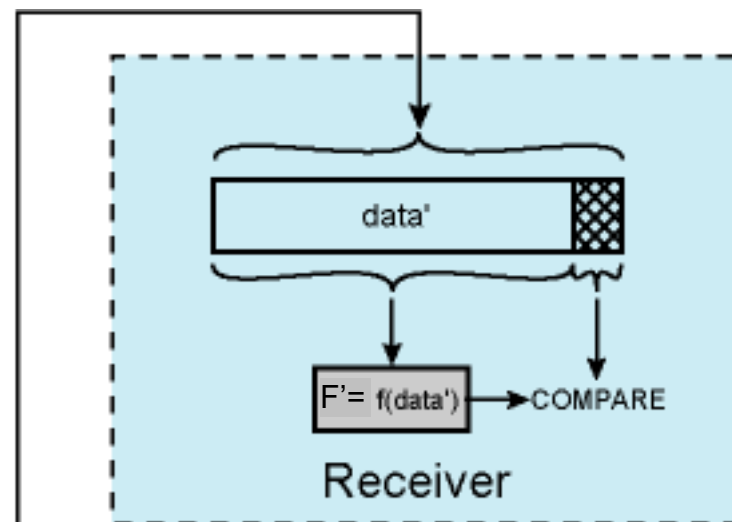
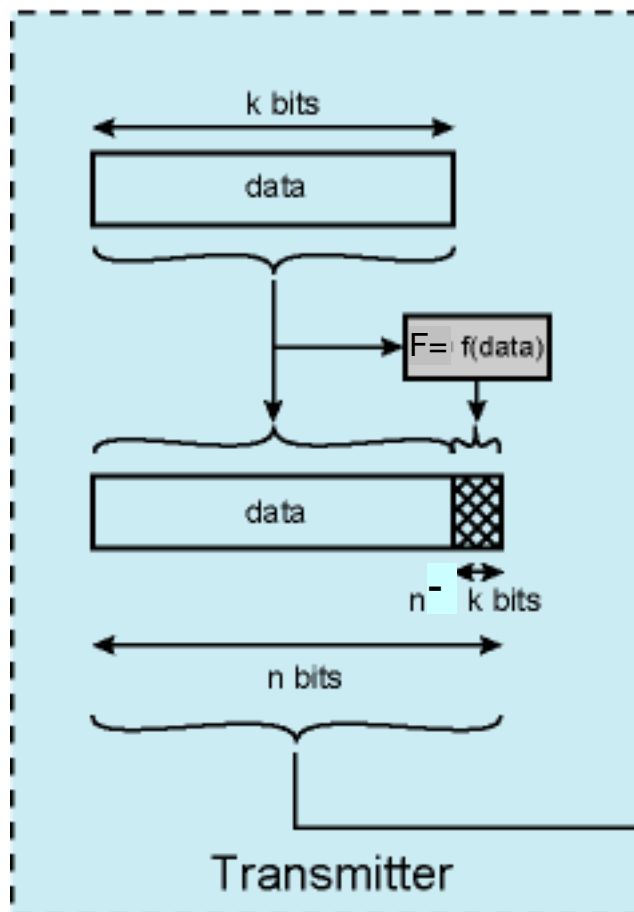
- So far we have seen flow control mechanisms where frames are transmitted without errors
 - in real life any transmission facility may introduce errors
- So we have to
 - detect errors
 - if possible, correct errors (not in the scope of CS 408)
 - adopt flow control algorithms such that erroneous frames are retransmitted

Types of Errors

- Single bit errors
 - isolated errors
 - affects (flips) one bit, nearby bits are not altered
 - not so common in real life
- Burst errors
 - a sequence of bits are affected
 - most common case
 - a burst error of length B is a contiguous sequence of B bits in which the first and the last and some intermediate bits are erroneously flipped.

Error Detection

- Additional bits added by transmitter as *error detection code*
 - receiver checks this code
- Parity
 - single bit added to the end of the data
 - Value of parity bit is such that data and parity have even (even parity) or odd (odd parity) number of ones
 - Even number of bit errors goes undetected



E, E' = error-detecting codes
 f = error-detecting code function

Error Control

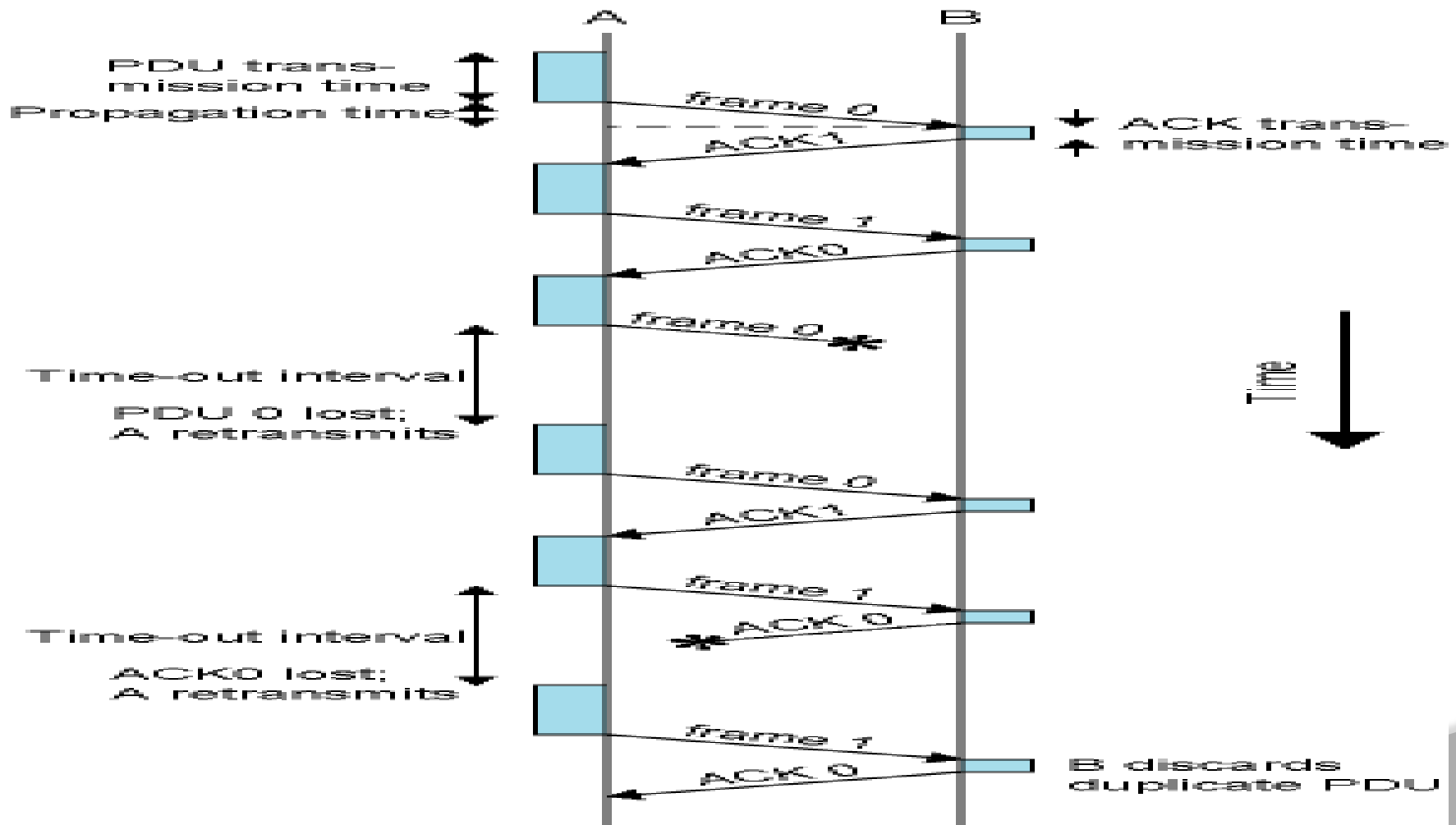
- Actions to be taken against
 - Lost frames
 - Damaged frames
- Automatic repeat request (ARQ) mechanism components
 - Error detection
 - Positive acknowledgment
 - Retransmission after timeout
 - Negative acknowledgement and retransmission

Stop and Wait ARQ

- Source transmits single frame
- Wait for ACK
- If received frame is damaged, discard it
 - If transmitter receives no ACK within timeout, retransmits
- If ACK damaged, transmitter will not recognize it
 - Transmitter will retransmit after timeout
 - Receiver gets two copies of frame, but disregards one of them
 - Use ACK_0 and ACK_1

Stop-and-Wait A Example

RQ –



Go-Back-N ARQ

- Based on sliding window
- If no error, ACK as usual with next frame expected
 - ACK_i means “I am ready to receive frame i ” and “I received all frames between i and my previous ack”
- Sender uses window to control the number of unacknowledged frames
- If error, reply with rejection (negative ack)
 - Discard that frame and all future frames until the frame in error is received correctly
 - Transmitter must go back and retransmit that

Go-Back-N ARQ- Damaged/Lost

Acknowledgment

- Receiver gets frame i and sends acknowledgment ($i+1$) which is lost
- Acknowledgments are cumulative, so next acknowledgement ($i+n$) may arrive before transmitter times out on frame i
==> NO PROBLEM
- If transmitter times out, it sends acknowledgment request with P bit set, as before

Selective Reject

- Also called *selective retransmission*
- Only rejected frames are retransmitted
- Subsequent frames are accepted by the receiver and buffered
- Minimizes retransmissions
- Receiver must maintain large enough buffer
- Complex implementation

The Network Layer

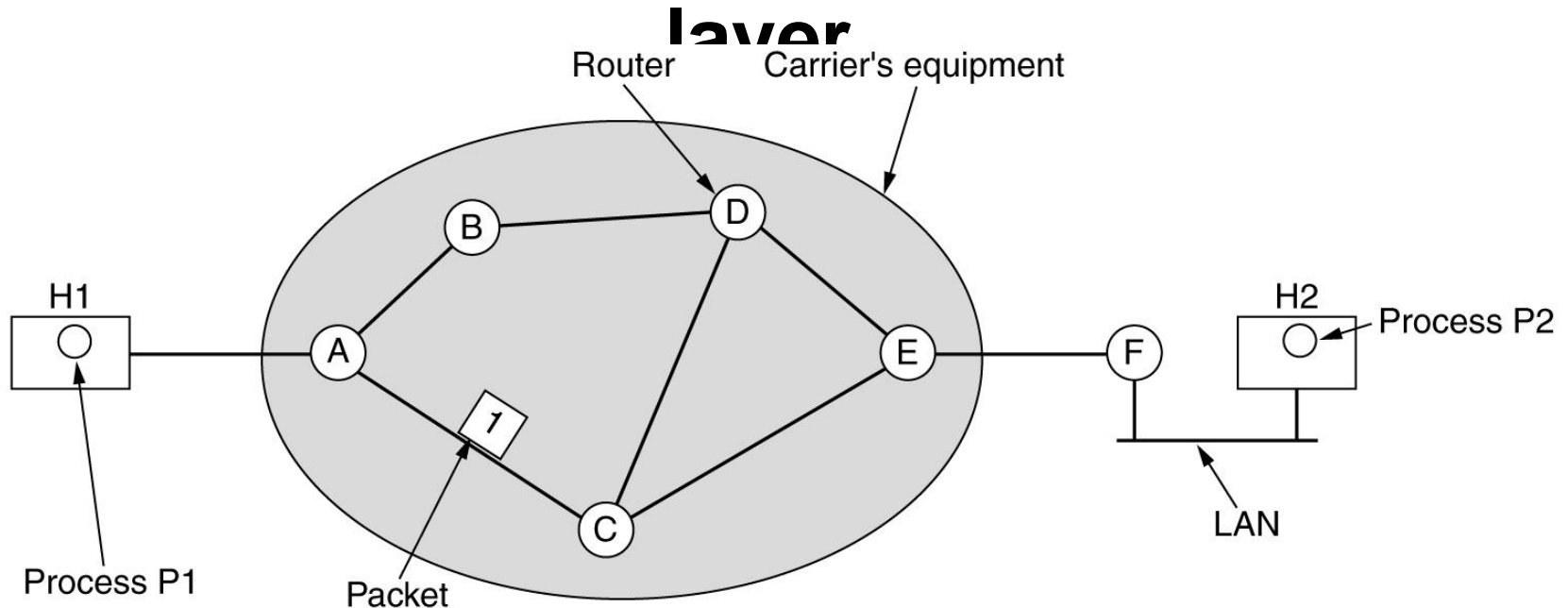
Design Issues & Routing Algorithms

Network Layer Design Issues

- Store-and-Forward Packet Switching
- Services Provided to the Transport Layer
- Implementation of Connectionless Service
- Implementation of Connection-Oriented Service
- Comparison of Virtual-Circuit and Datagram Subnets

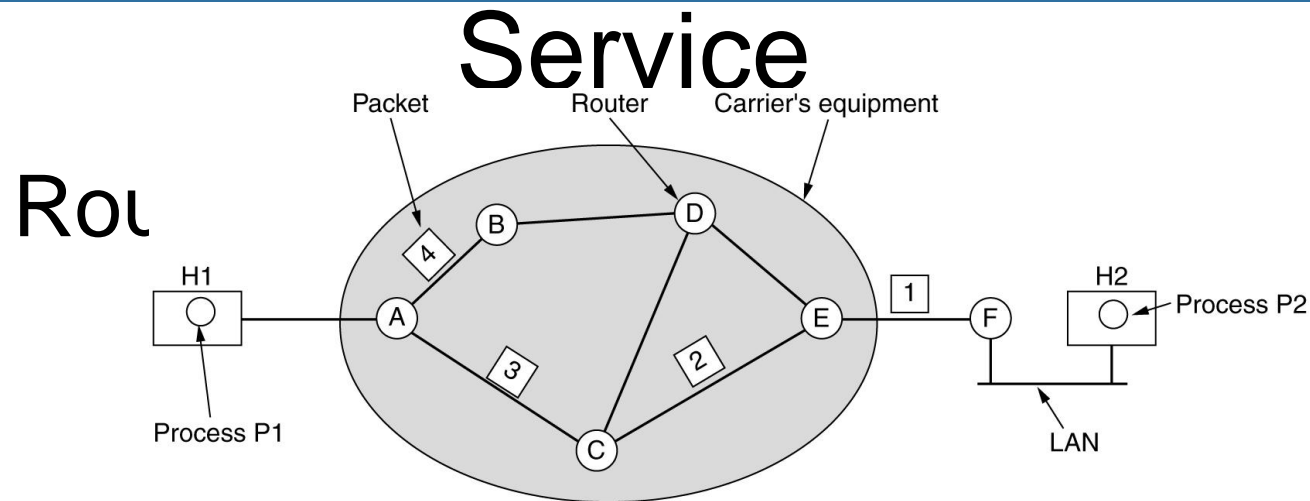
Store-and-Forward Packet Switching

The environment of the network



1. The services should be independent of the router technology
2. The transport layer should be shielded from the number, type and topology of the routers present
3. The network addresses made available to the transport layer should use a uniform numbering plan, even across LANs and WANs

Implementation of Connectionless Service



A's table

	initially	later
A	-	-
B	B	B
C	C	C
D	B	B
E	C	B
F	C	B

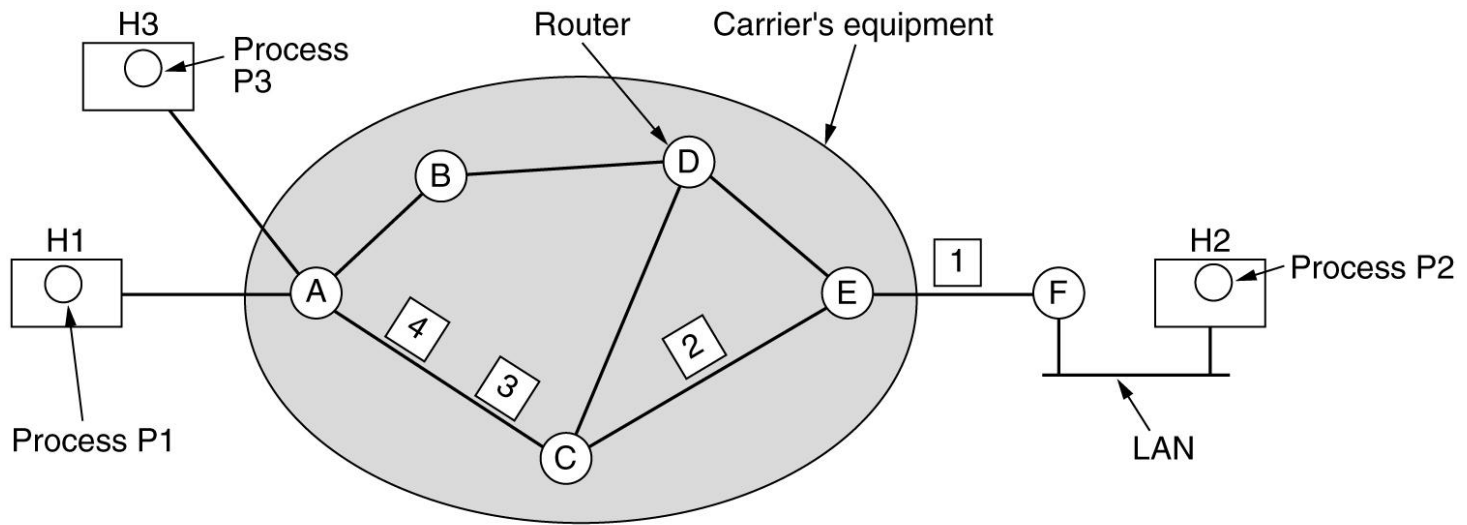
Dest. Line

C's table

A	A
B	A
C	-
D	D
E	E
F	E

E's table

A	C
B	D
C	C
D	D
E	-
F	F



A's table				C's table				E's table			
H1	1	C	1	A	1	E	1	C	1	F	1
H3	1	C	2	A	2	E	2	C	2	F	2
In											

Routing within a virtual-circuit subnet.

Issue	Datagram subnet	Virtual-circuit subnet
Circuit setup	Not needed	Required
Addressing	Each packet contains the full source and destination address	Each packet contains a short VC number
State information	Routers do not hold state information about connections	Each VC requires router table space per connection
Routing	Each packet is routed independently	Route chosen when VC is set up; all packets follow it
Effect of router failures	None, except for packets lost during the crash	All VCs that passed through the failed router are terminated
Quality of service	Difficult	Easy if enough resources can be allocated in advance for each VC
Congestion control	Difficult	Easy if enough resources can be allocated in advance for each VC

Routing Algorithms

- The Optimality Principle
- Shortest Path Routing
- Flooding
- Distance Vector Routing
- Link State Routing
- Hierarchical Routing
- Broadcast Routing
- Multicast Routing
- Routing for Mobile Hosts
- Routing in Ad Hoc Networks

Desirable Properties (Elaborate)

Correctness

Simplicity

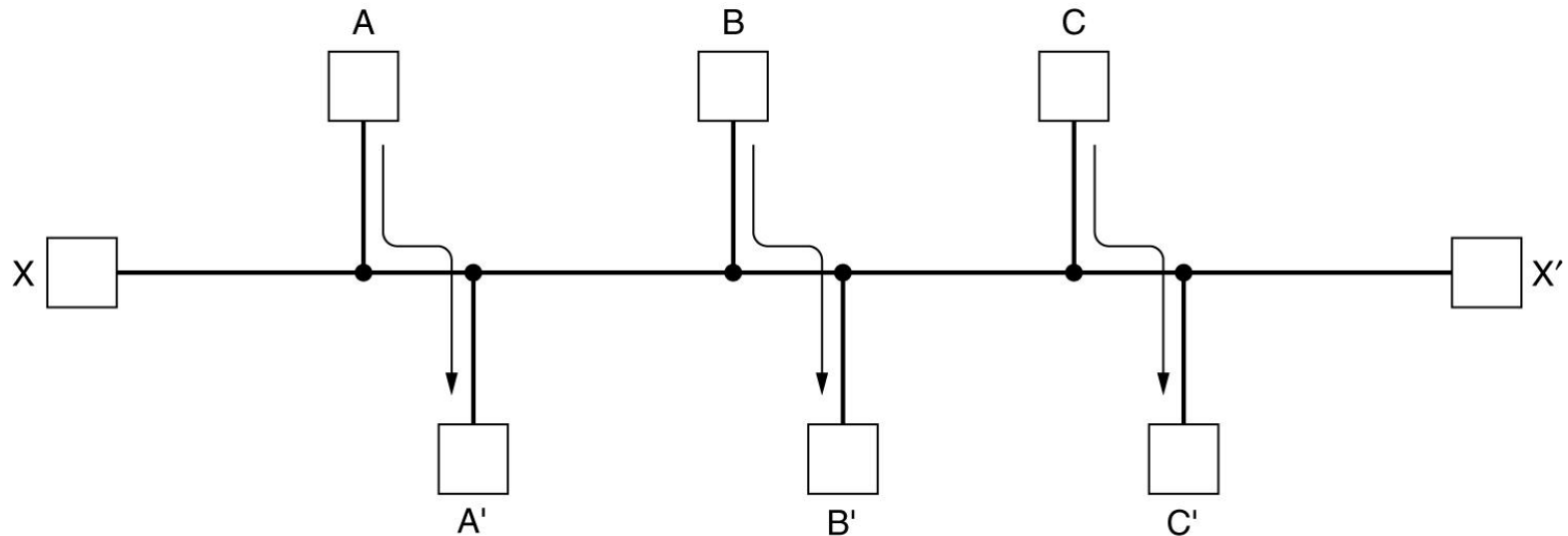
Robustness – System will be in place for years with small failures

Stability – Fast convergence

Fairness,

Efficiency.

Routing Algorithms (2)

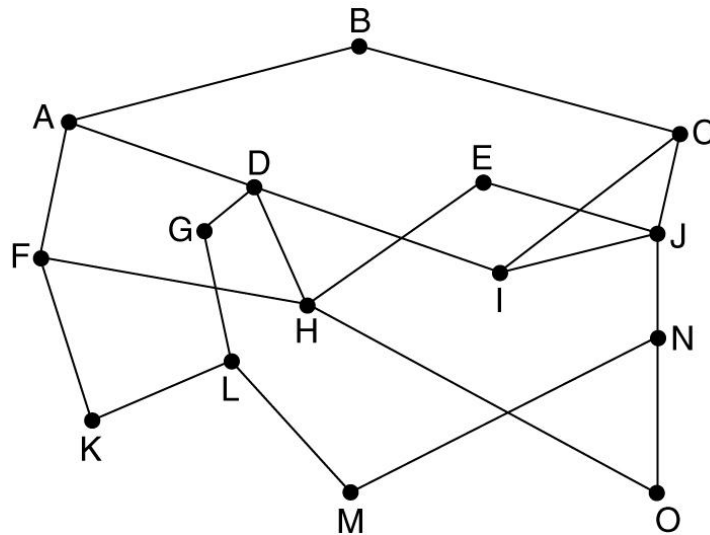


A – A', B – B', C – C', can fill the channel, then X-X' doesn't get a chance

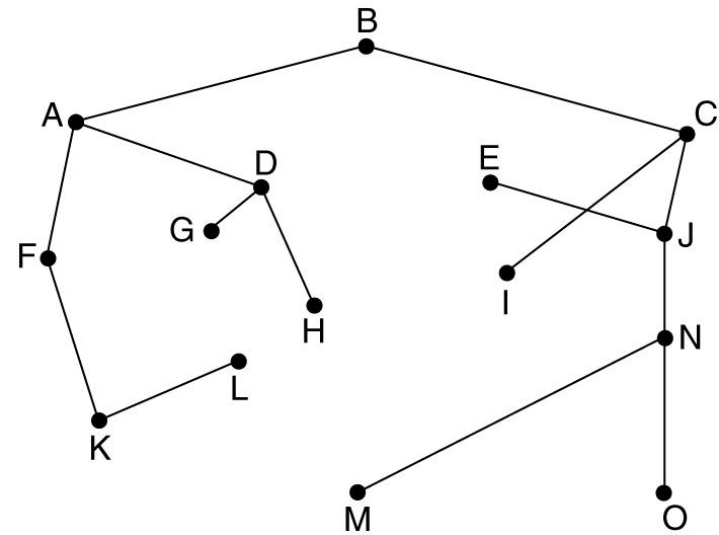
Conflict between fairness and optimality.

Minimizing the mean packet delay is an obvious candidate to send traffic through the network

The Optimality Principle



(a)



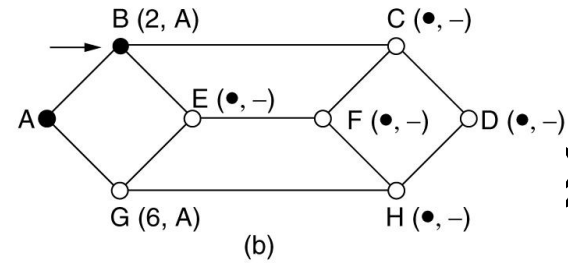
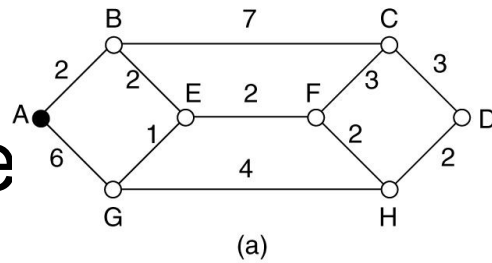
(b)

Optimality Principle – If router J is on the optimal path from router I to router K then the optimal path from J to K also falls along the same route.

(a) A subnet. (b) A sink tree for router B.

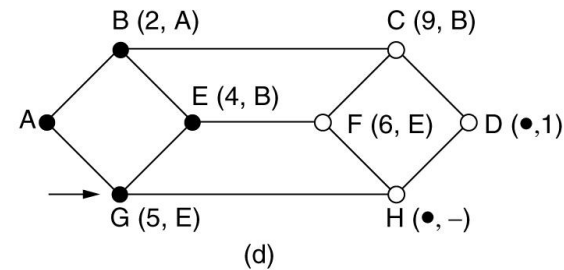
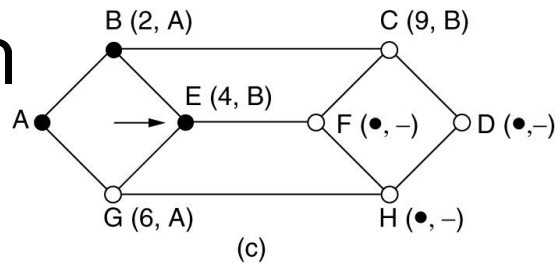
Shortest Path Routing

The

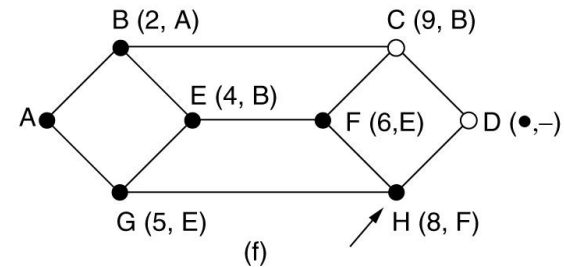
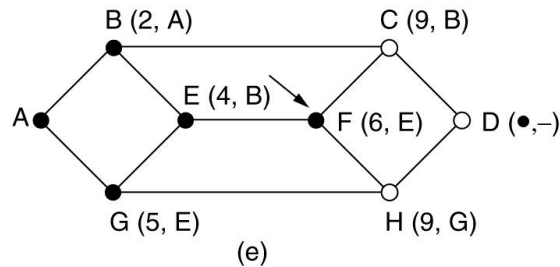


, the

Th



node.



Dijkstra

```
#define MAX_NODES 1024          /* maximum number of nodes */
#define INFINITY 1000000000     /* a number larger than every maximum path */
int n, dist[MAX_NODES][MAX_NODES]; /* dist[i][j] is the distance from i to j */

void shortest_path(int s, int t, int path[])
{ struct state {                /* the path being worked on */
    int predecessor;            /* previous node */
    int length;                 /* length from source to this node */
    enum {permanent, tentative} label; /* label state */
} state[MAX_NODES];

int i, k, min;
struct state *p;

for (p = &state[0]; p < &state[n]; p++) { /* initialize state */
    p->predecessor = -1;
    p->length = INFINITY;
    p->label = tentative;
}
state[t].length = 0; state[t].label = permanent;
k = t;                /* k is the initial working node */
```

Dijkstra's algorithm to compute the shortest path

Dijkstra

```

do {
    for (i = 0; i < n; i++)
        if (dist[k][i] != 0 && state[i].label == tentative) {
            if (state[k].length + dist[k][i] < state[i].length) {
                state[i].predecessor = k;
                state[i].length = state[k].length + dist[k][i];
            }
        }

    /* Find the tentatively labeled node with the smallest label. */
    k = 0; min = INFINITY;
    for (i = 0; i < n; i++)
        if (state[i].label == tentative && state[i].length < min) {
            min = state[i].length;
            k = i;
        }
    state[k].label = permanent;
} while (k != s);

/* Copy the path into the output array. */
i = 0; k = s;
do {path[i++] = k; k = state[k].predecessor; } while (k >= 0);
}

```

Dijkstra's algorithm to compute the shortest path

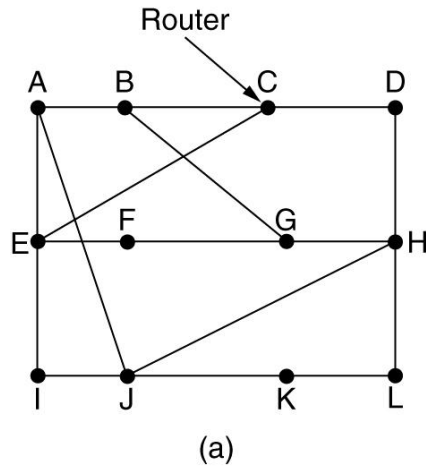
Flooding

Robust but costly.

**TTL and keep
track...**

- Used in military application
- Wireless Networks
- Distributed Database
- Metrics against which other routing algorithms are compared.

Distance Vector Routing



To	A	I	H	K	New estimated delay from J	Line
A	0	24	20	21	8	A
B	12	36	31	28	20	A
C	25	18	19	36	28	I
D	40	27	8	24	20	H
E	14	7	30	22	17	I
F	23	20	19	40	30	I
G	18	31	6	31	18	H
H	17	20	0	19	12	H
I	21	0	14	22	10	I
J	9	11	7	10	0	–
K	24	22	22	0	6	K
L	29	33	9	9	15	K

JA delay is 8	JI delay is 10	JH delay is 12	JK delay is 6
---------------	----------------	----------------	---------------

Vectors received from J's four neighbors

New routing table for J

(b)

(a) A subnet. (b) Input from A, I, H, K, and the new

Distance Vector Routing (2)

A	B	C	D	E	
•	•	•	•	•	Initially
	•	•	•	•	After 1 exchange
	1	•	•	•	After 2 exchanges
	1	2	•	•	After 3 exchanges
	1	2	3	•	After 4 exchanges
	1	2	3	4	

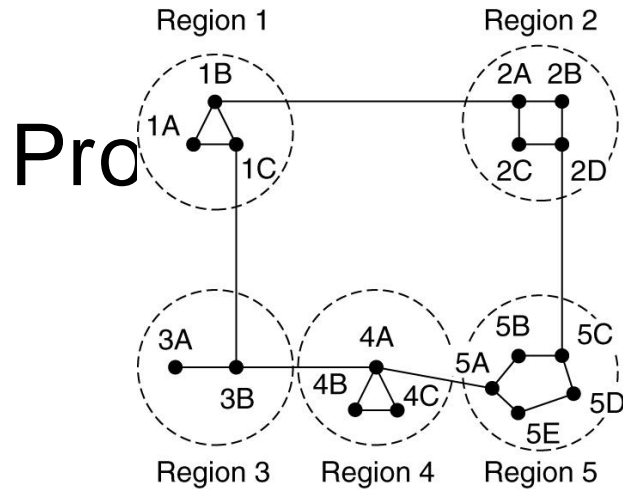
(a)

A	B	C	D	E	
•	•	•	•	•	Initially
	1	2	3	4	After 1 exchange
	3	2	3	4	After 2 exchanges
	3	4	3	4	After 3 exchanges
	5	4	5	4	After 4 exchanges
	5	6	5	6	After 5 exchanges
	7	6	7	6	After 6 exchanges
	7	8	7	8	
	•	•	•	•	

(b)

The count-to-infinity problem.

Hierarchical Routing



(a)

Full table for 1A

Dest.	Line	Hops
1A	—	—
1B	1B	1
1C	1C	1
2A	1B	2
2B	1B	3
2C	1B	3
2D	1B	4
3A	1C	3
3B	1C	2
4A	1C	3
4B	1C	4
4C	1C	4
5A	1C	4
5B	1C	5
5C	1B	5
5D	1C	6
5E	1C	5

(b)

Hierarchical table for 1A

Dest.	Line	Hops
1A	—	—
1B	1B	1
1C	1C	1
2	1B	2
3	1C	2
4	1C	3
5	1C	4

(c)

ed

Hierarchical Routing (2)

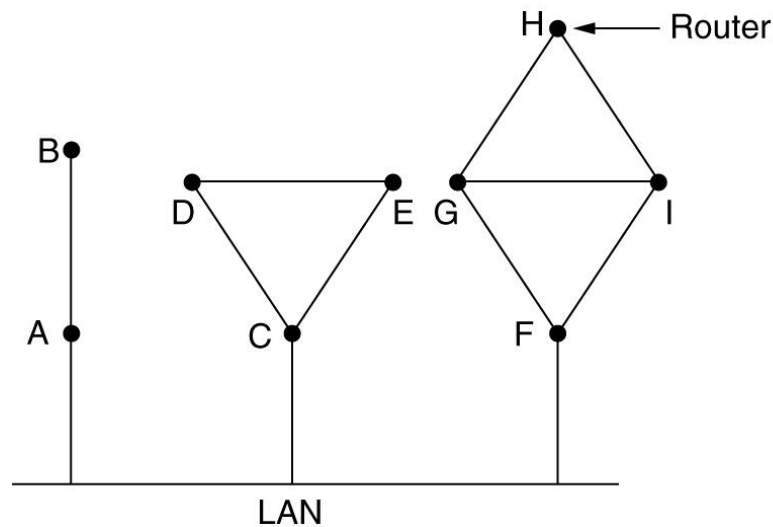
- How many levels of hierarchy?
- 720 routers.
- 720 routers in 24 regions.
- Three levels of hierarchy – 8 clusters each containing 9 regions of 10 clusters.

Link State Routing

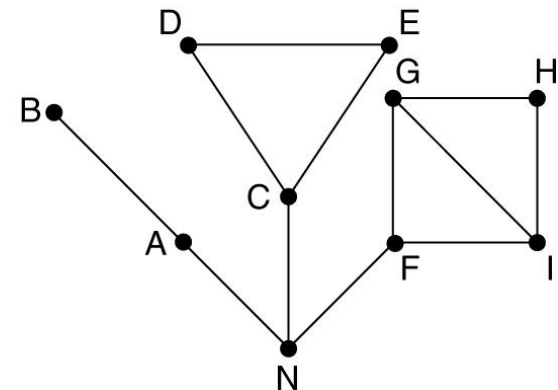
Each router must do the following:

1. Discover its neighbors, learn their network address.
2. Measure the delay or cost to each of its neighbors.
3. Construct a packet telling all it has just learned.
4. Send this packet to all other routers.
5. Compute the shortest path to every other router.

Learning about the Neighbors



(a)



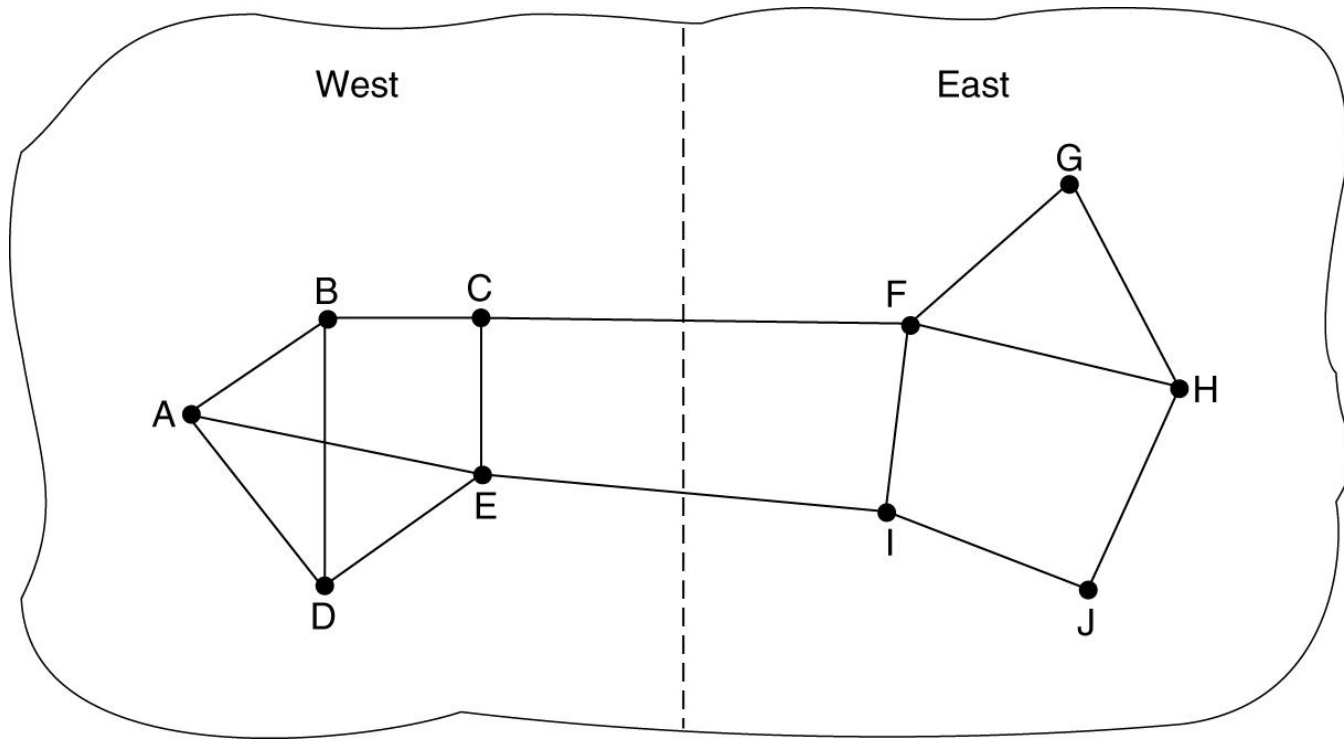
(b)

(a) Nine routers and a LAN. (b) A graph model of (a). All routers

Setting Link Cost

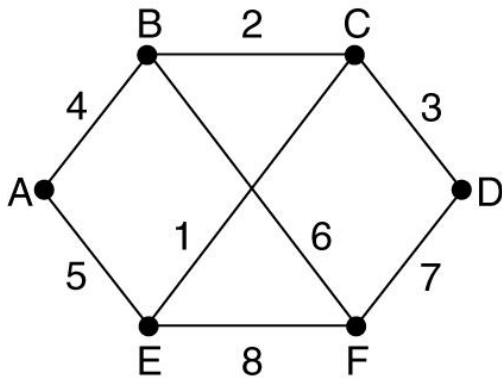
- Bandwidth
- Delay – measured by sending special ECHO
- Geographically spread out links

Measuring Line Cost



- A subnet in which the East and West parts are connected by 2 lines.
- Including queuing delay may lead to a lot of

Building Link State Packets



(a)

		Link		State		Packets	
A		B		C		D	
Seq.		Seq.		Seq.		Seq.	
Age		Age		Age		Age	
B	4	A	4	B	2	C	3
E	5	C	2	D	3	F	7
		F	6	E	1		

E		F	
Seq.		Seq.	
Age		Age	
A	5	B	6
C	1	D	7
F	8	E	8

(b)

(a) A subnet. (b) The link state packets for this subnet.

Few Problems

Algorithm – Sequence number less means obsolete

- If sequence numbers wrap around, confusion will reign
- Router crashes, sequence number is lost
- Sequence number gets corrupted
- Aging and then dropping the packet.

Source	Seq.	Age	Send flags			ACK flags			Data
			A	C	F	A	C	F	
A	21	60	0	1	1	1	0	0	
F	21	60	1	1	0	0	0	1	
E	21	59	0	1	0	1	0	1	
C	20	60	1	0	1	0	1	0	
D	21	59	1	0	0	0	1	1	

- The packet buffer for router B in the previous slide (Fig. 5-13).
- E has arrived twice.

OSPF (Open Shortest Path First) IS-IS (intermediate System- Intermediate System)

- Refreshed every 60 seconds.
- Hardware problem, router getting corrupt, etc.

Broadcast Routing

Multidimensional Routing

- Each packet contains a list of destinations.
- On arrival of a packet, router checks the set of destinations, and sends copies of packet along outgoing links to those destinations.

Flooding

- Flood with a sequence number per source.

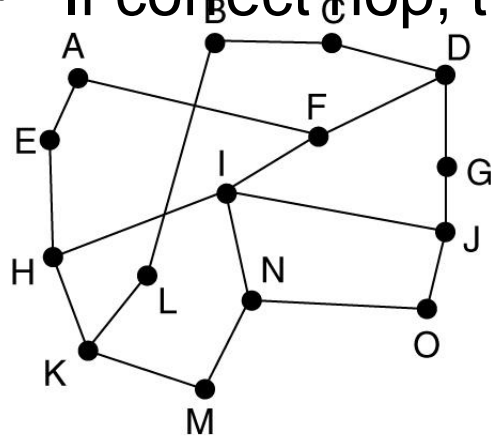
Spanning Tree

- Build spanning tree (such as, a sink tree).
- Forward packet along all links of spanning tree except the one from which packet is received.

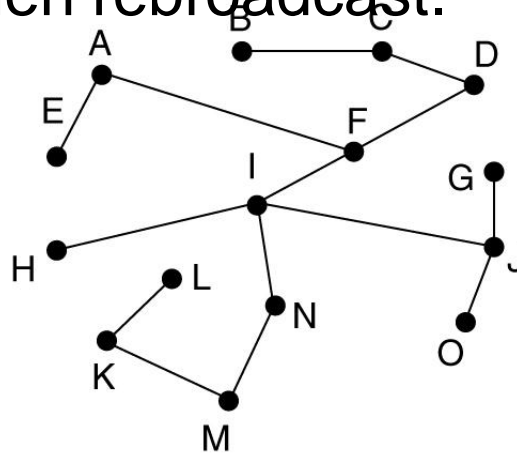
Reverse Path Forwarding

Broadcast – Reverse Path Forwarding

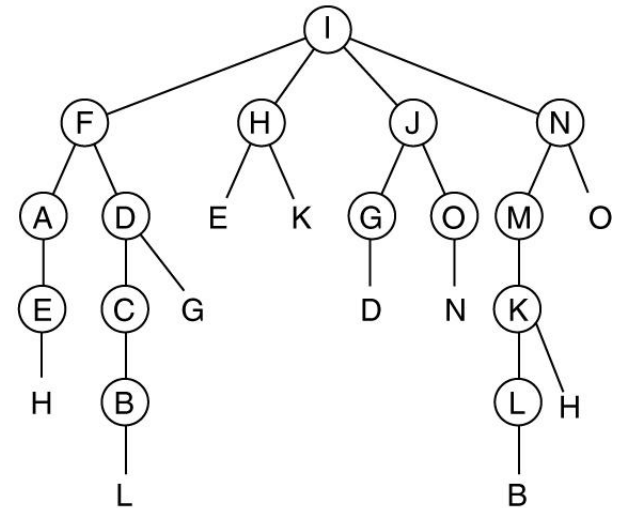
- Broadcast.
- Check if the packet has arrived following the correct hop or not.
- If correct hop, then rebroadcast.



(a)



(b)



(c)

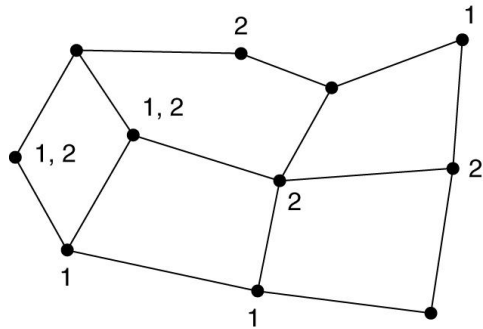
Reverse path forwarding.

(a) A subnet.

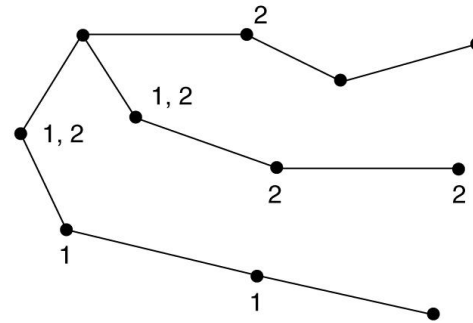
(b) a Sink tree.

(c) The tree built by reverse path forwarding

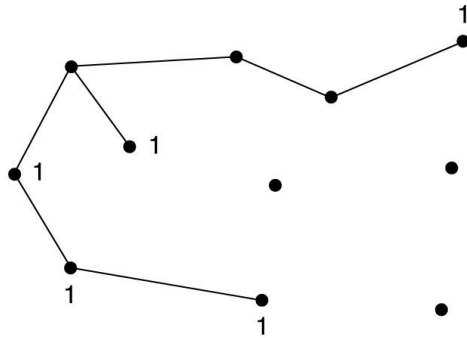
Multicast Routing



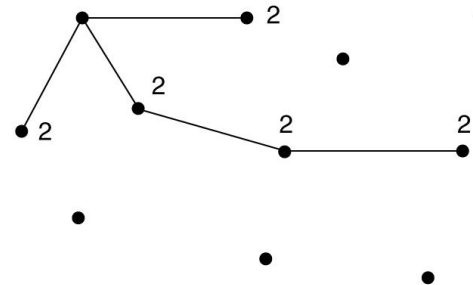
(a)



(b)



(c)



(d)

- (a) A network. (b) A spanning tree for the leftmost router.
(c) A multicast tree for group 1. (d) A multicast tree for group 2.

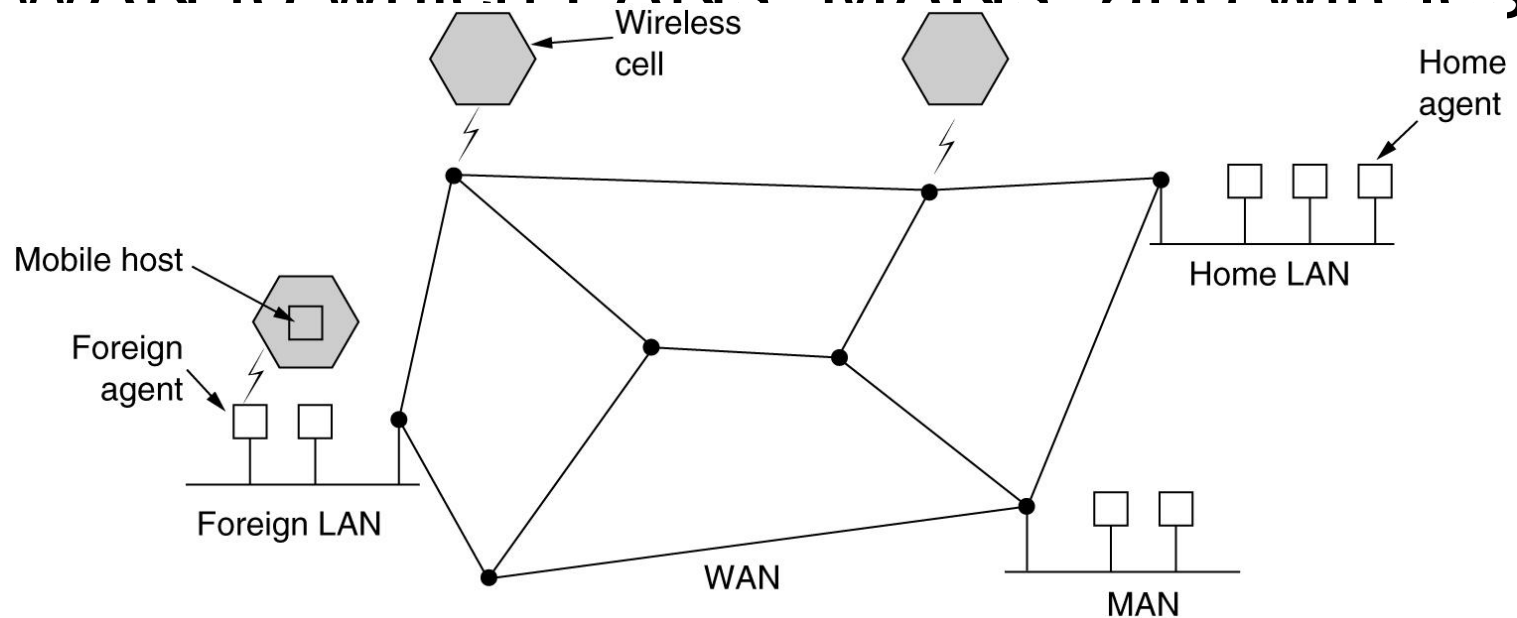
Typically done in Multi-state routing

Anycast Routing

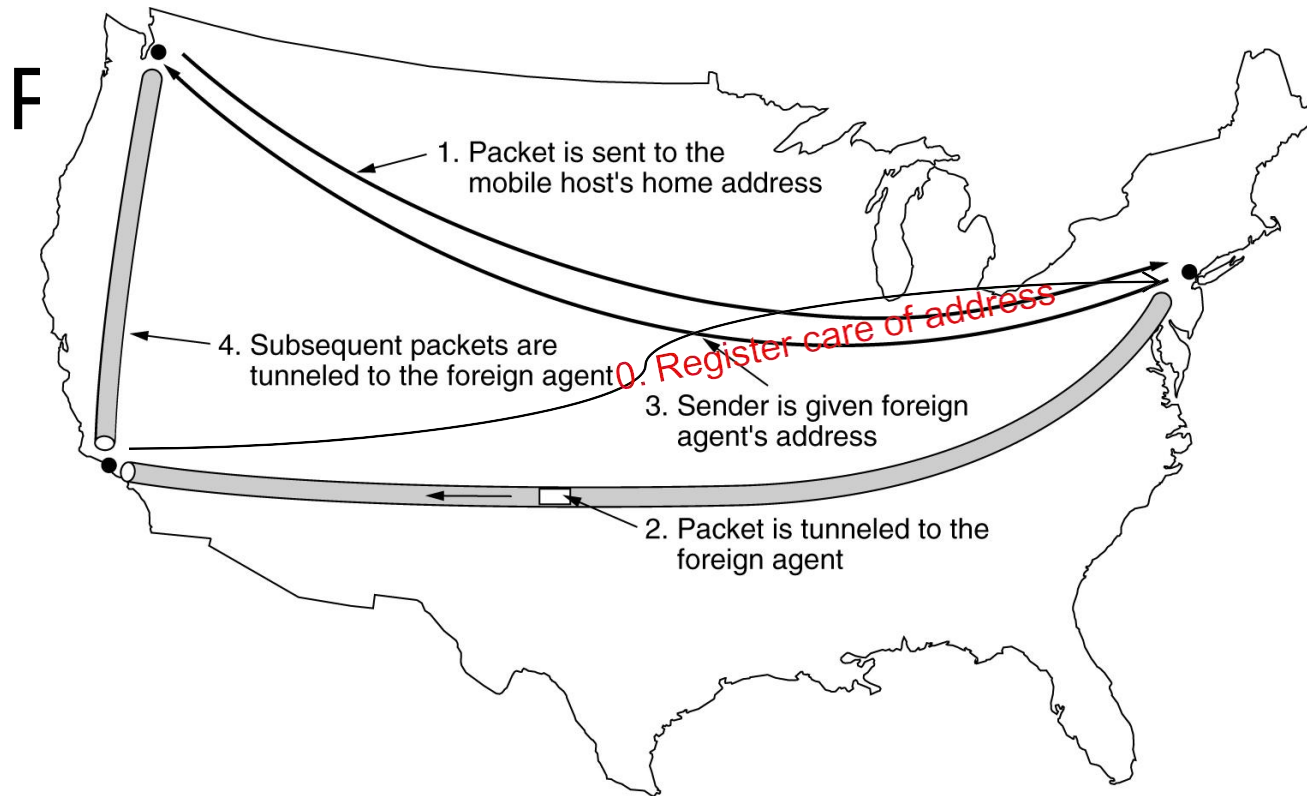
- Reaching any one of the servers in the group
- DNS server

Routing for Mobile Hosts

A WAN to which LANs, MANs and wireless



Routing for Mobile Hosts (2)



Routing in Ad Hoc Networks

Possibilities when the routers are mobile: 1. Military vehicles on battlefield.

- No infrastructure.

2. A fleet of ships at sea.

- All moving all the time

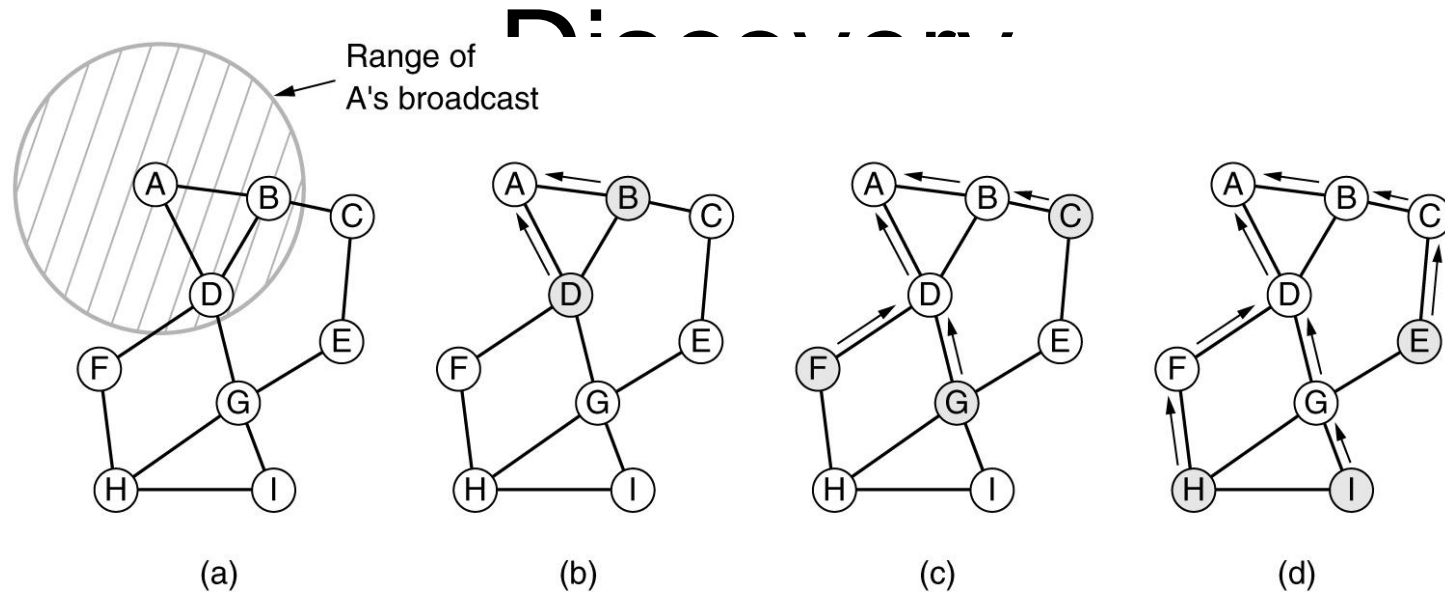
3. Emergency works at earthquake.

- The infrastructure destroyed.

4. A gathering of people with notebook computers.

- In an area lacking 802.11.

Ad Hoc Networks: Route



(a) Range of A's broadcast.

(b) After B and D have received A's broadcast.

(c) After C, F, and G have received A's broadcast.

Route Discovery (2)

Source address	Request ID	Destination address	Source sequence #	Dest. sequence #	Hop count
-------------------	---------------	------------------------	----------------------	---------------------	--------------

Format of a ROUTE REQUEST packet.

Route Discovery (3)

- The (Source Address, Request ID) pair is looked up in a local history table
- Receiver looks up the destination in its route table. If a fresh route is known, then a ROUTE REPLY is sent.
- Destination sequence number is higher than the Destination sequence in the Route Discovery Packet
- Increments Hop count and rebroadcasts ROUTE REPLY
- Stores the data in a new entry in its reverse route table.

Route Discovery (4)

Source address	Destination address	Destination sequence #	Hop count	Lifetime
-------------------	------------------------	---------------------------	--------------	----------

Format of a ROUTE REPLY packet.

Route Discovery (5)

IN response

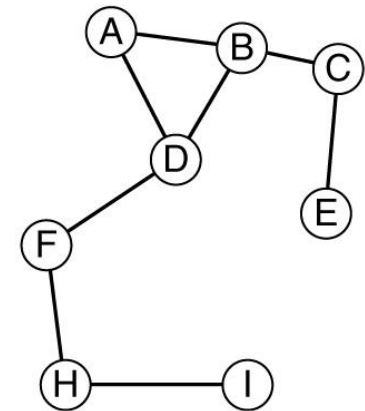
- Source addr., destination addr. and Hop Count copied but Dest. Seq. number taken from its counter.
- Hopcount is set to 0, Lifetime field controls how long the route is valid.

At each intermediate node:

- 1.No route to I is known,
- 2.Sequence number of I in the ROUTE REPLY packet is greater than the value in the routing table
- 3.The sequence numbers are equal but the new route is shorter
- 4.Hop Count incremented
- 5.In large network, discovery increases with Time to Live incrementally being increased from 1, 2, 3, ...

Dest.	Next hop	Distance	Active neighbors	Other fields
A	A	1	F, G	
B	B	1	F, G	
C	B	2	F	
E	G	2		
F	F	1	A, B	
G	G	1	A, B	
H	F	2	A, B	
I	G	2	A, B	

(a)



(b)

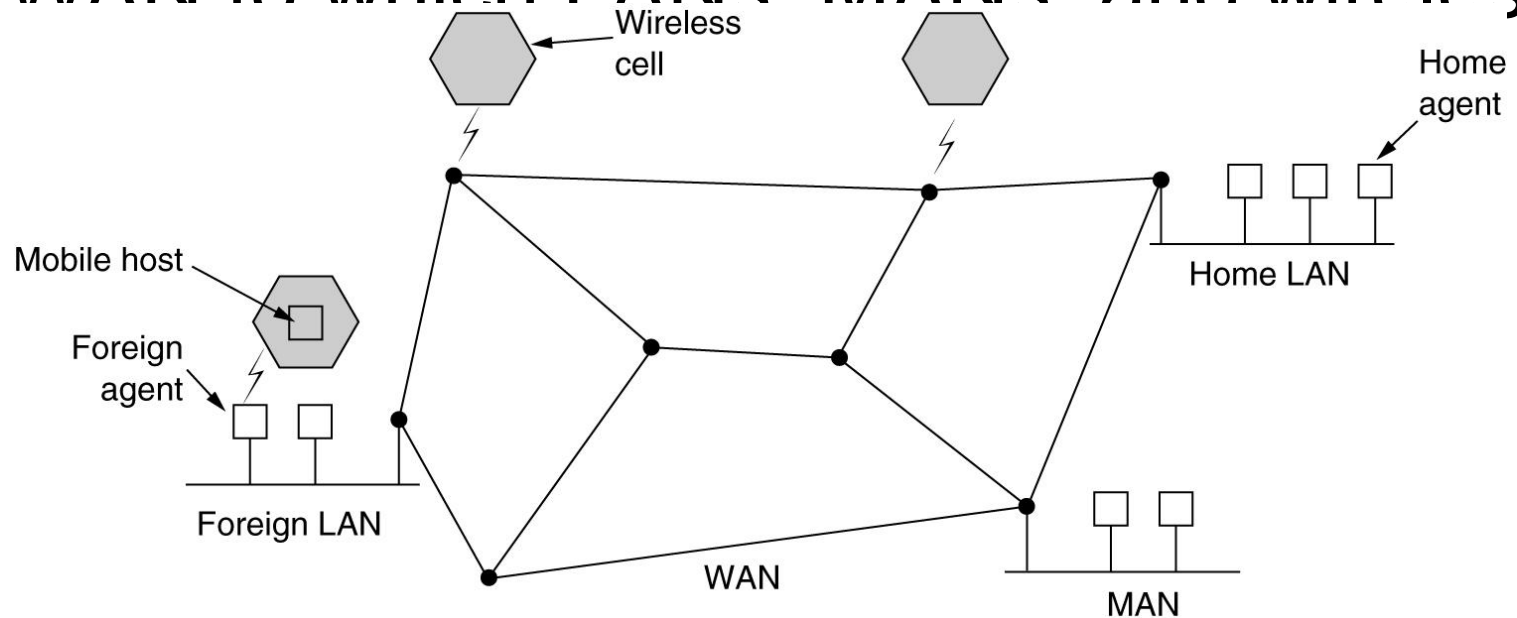
Active Neighbours that have fed in those destinations to A in last t seconds

(a) D's routing table before G goes down.

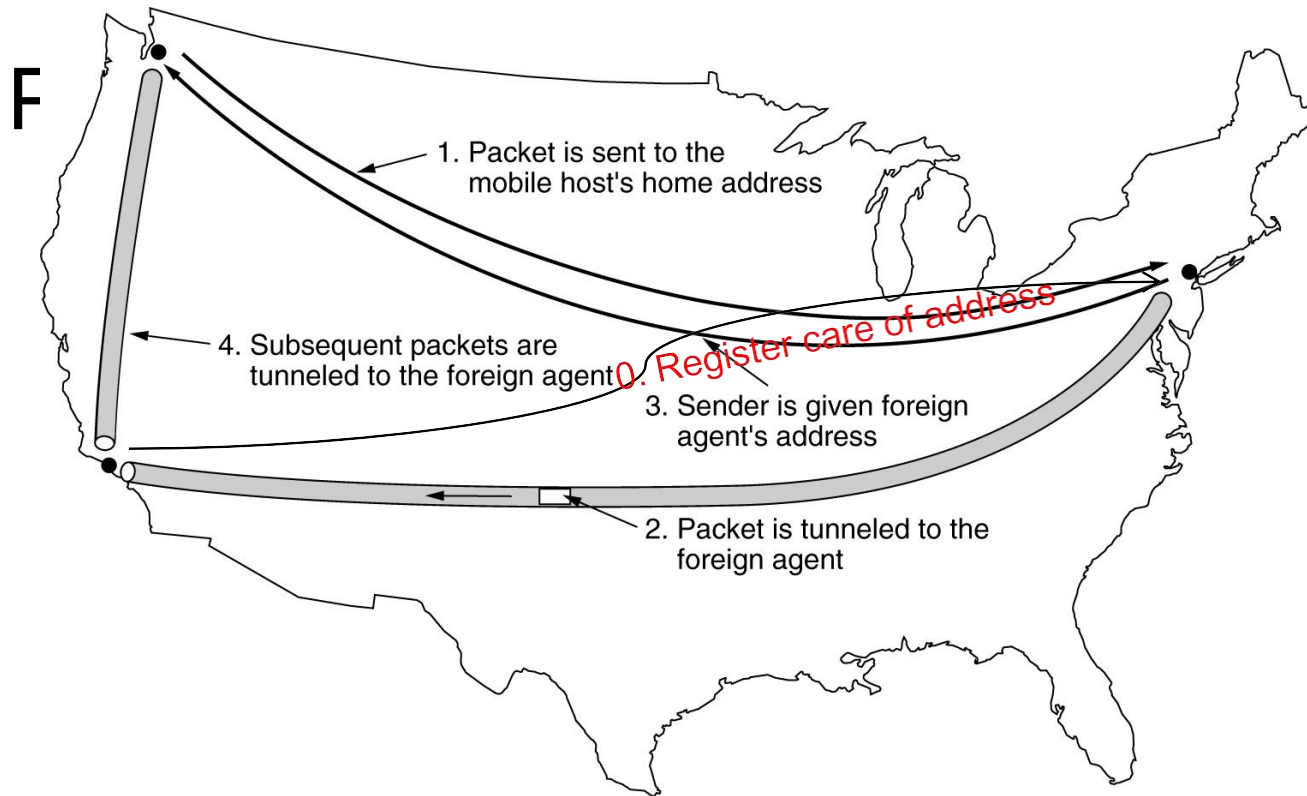
(b) The graph after G has gone down.

Routing for Mobile Hosts

A WAN to which LANs, MANs and wireless



Routing for Mobile Hosts (2)



Routing in Ad Hoc Networks

Possibilities when the routers are mobile: 1. Military vehicles on battlefield.

- No infrastructure.

2. A fleet of ships at sea.

- All moving all the time

3. Emergency works at earthquake .

- The infrastructure destroyed.

4. A gathering of people with notebook computers.

- In an area lacking 802.11.

Route Discovery (2)

Source address	Request ID	Destination address	Source sequence #	Dest. sequence #	Hop count
-------------------	---------------	------------------------	----------------------	---------------------	--------------

Format of a ROUTE REQUEST packet.

Route Discovery (3)

- The (Source Address, Request ID) pair is looked up in a local history table
- Receiver looks up the destination in its route table. If a fresh route is known, then a ROUTE REPLY is sent.
- Destination sequence number is higher than the Destination sequence in the Route Discovery Packet
- Increments Hop count and rebroadcasts ROUTE REPLY
- Stores the data in a new entry in its reverse route table.

Route Discovery (4)

Source address	Destination address	Destination sequence #	Hop count	Lifetime
-------------------	------------------------	---------------------------	--------------	----------

Format of a ROUTE REPLY packet.

Route Discovery (5)

IN response

- Source addr., destination addr. and Hop Count copied but Dest. Seq. number taken from its counter.
- Hopcount is set to 0, Lifetime field controls how long the route is valid.

At each intermediate node:

- 1.No route to I is known,
- 2.Sequence number of I in the ROUTE REPLY packet is greater than the value in the routing table
- 3.The sequence numbers are equal but the new route is shorter
- 4.Hop Count incremented
- 5.In large network, discovery increases with Time to Live incrementally being increased from 1, 2, 3, ...

Basics of Chord

- Uses a hash function such as SHA-1.
- SHA-1 converts a variable length input into a highly random 160 bit value
- Using SHA-1, Chord hashes:

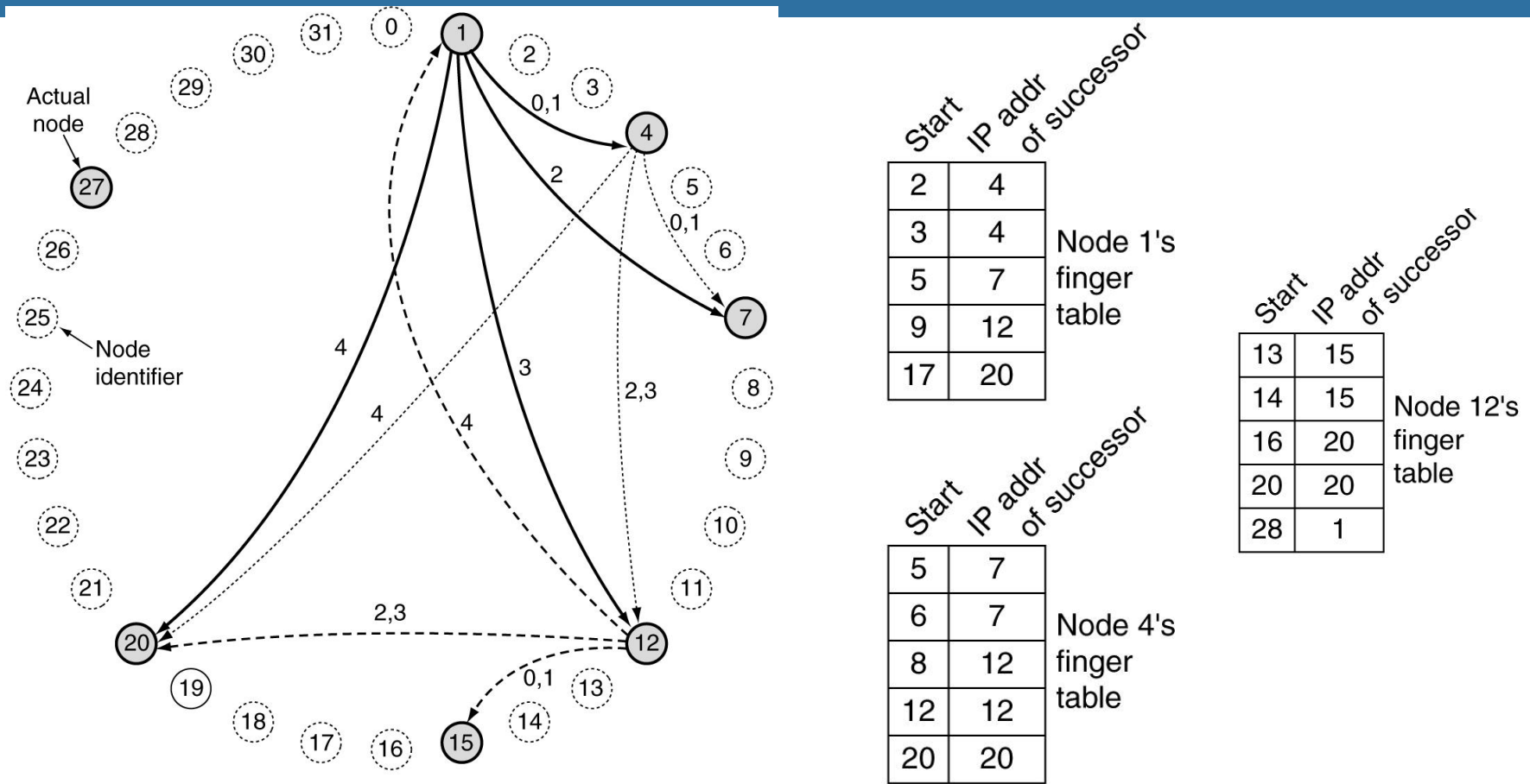
node IP addresses
(bits)

names of records
(bits)

node identifiers (160

keys (160

Storing Records



(a) A set of 32 node identifiers arranged in a circle. The shaded ones correspond to actual machines. The arcs show the fingers from nodes 1, 4, and 12. The labels on the

Storing records

- *successor(k)* is the first *real node* after *k*.
- To store data *name*, a node *N* creates a tuple (*name*, *N*'s IP address) and stores the tuple at *successor(hash(name))*. The original data remain at *N*, just the tuple is stored at *successor(hash(name))*.
- If *hash(name) = 22*, then the tuple is stored at node 27.
- To find information *name*, a node does *key = hash(name)*, then gets the record tuple from *successor(key)*.
- Simple? Mostly, except for implementing *successors(key)* efficiently.

Finding records

Each node needs to store the IP addresses of its successor.

Initially, the network start out with just a few nodes:

1. All nodes know each other.
2. They can easily arrange themselves into a the Chord ring.
3. $successor(k)$ can be computed.

When a node tries to join:

1. It calculates its node ID say p .
2. Then asks any node already in the ring to find $successor(p)$.
3. Asks $successor(p)$ for $successor(p)$'s predecessor and inserts itself between them.

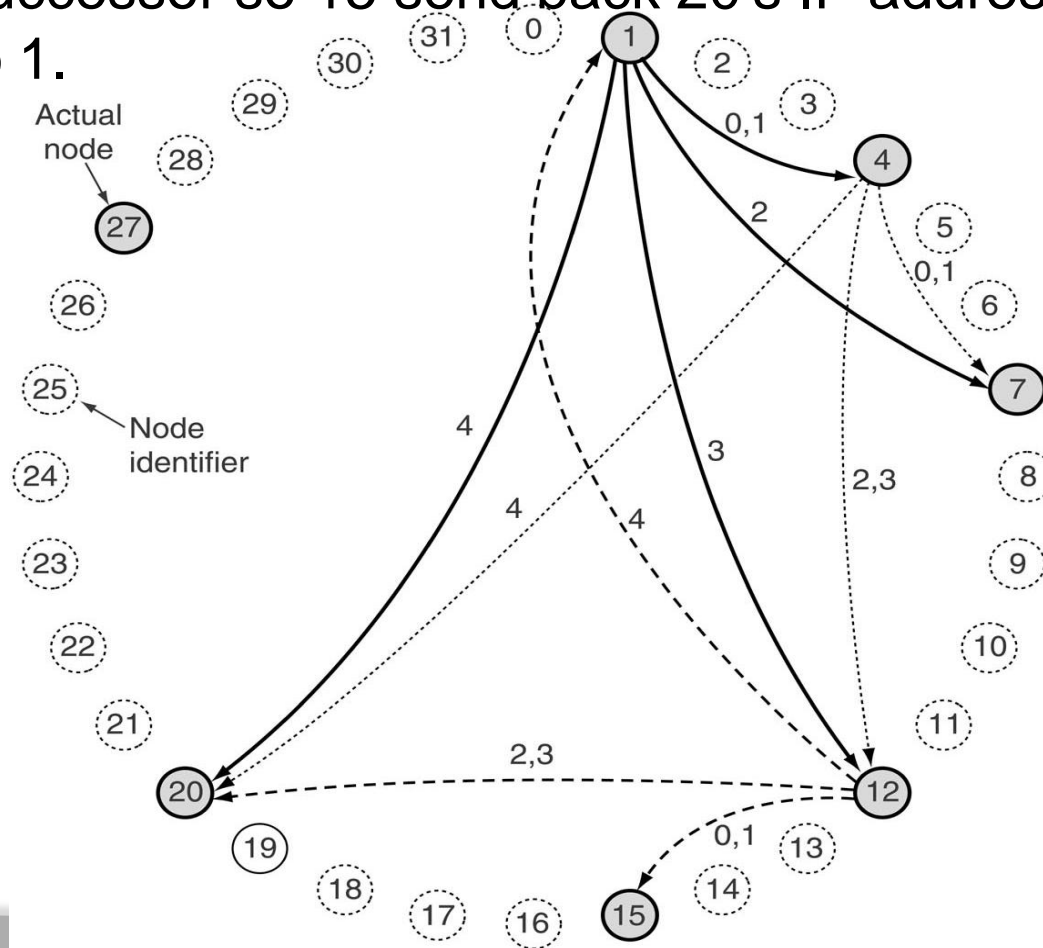
Any node in the ring can find $successor(k)$ by propagating the query around the ring starting with its successor.

Finger table

- Even if both successor and predecessor pointers are used, a sequential search will take time on average $O(n/2)$ [n is the number of nodes].
- Chord reduces this search time using a finger table at each node.
- The finger table contains up to m entries where each entry i consists of IP address of $\text{successor}(\text{start}[i])$
- $\text{Start}[i] = k + 2^i \pmod{2^m}$
- To find a record for key k , a node can directly jump to the closest predecessor of k .
- Average time can be reduced to $O(\log n)$.

Looking up key 16 at node 1

1. Nearest pred. 9 so query sent to 12
2. At 12 nearest pred. of 16 is 14 so query sent to 15
3. 15 knows that 16 is between itself and its successor so 15 send back 20's IP address to 1.



Start	IP addr of successor
2	4
3	4
5	7
9	12
17	20

Node 1's finger table

Start	IP addr of successor
5	7
6	7
8	12
12	12
20	20

Node 4's finger table

Start	IP addr of successor
13	15
14	15
16	20
20	20
28	1

Node 12's finger table

Maintaining finger table

- Maintaining the finger table does not come for free.
- Every time a new node is added a few successors and predecessor entries will change.

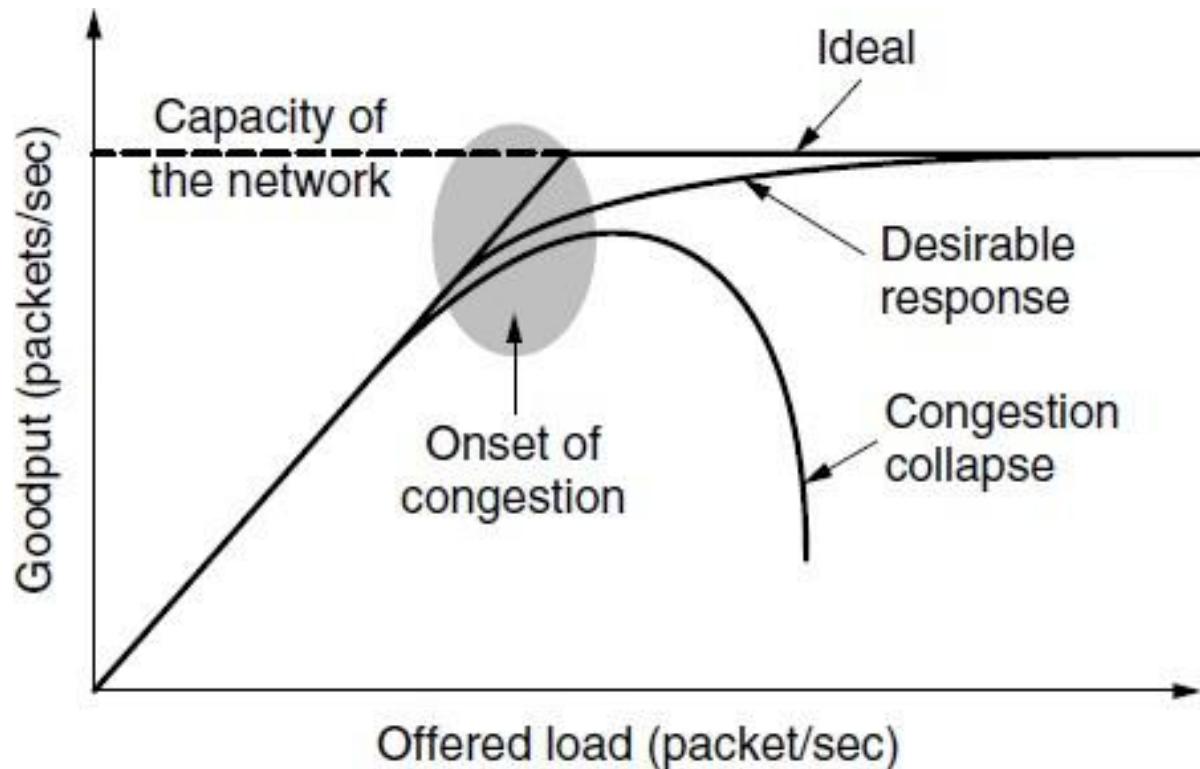
The Network Layer

Congestion Control Algorithms & Quality-of- Service

Congestion Control Algorithms

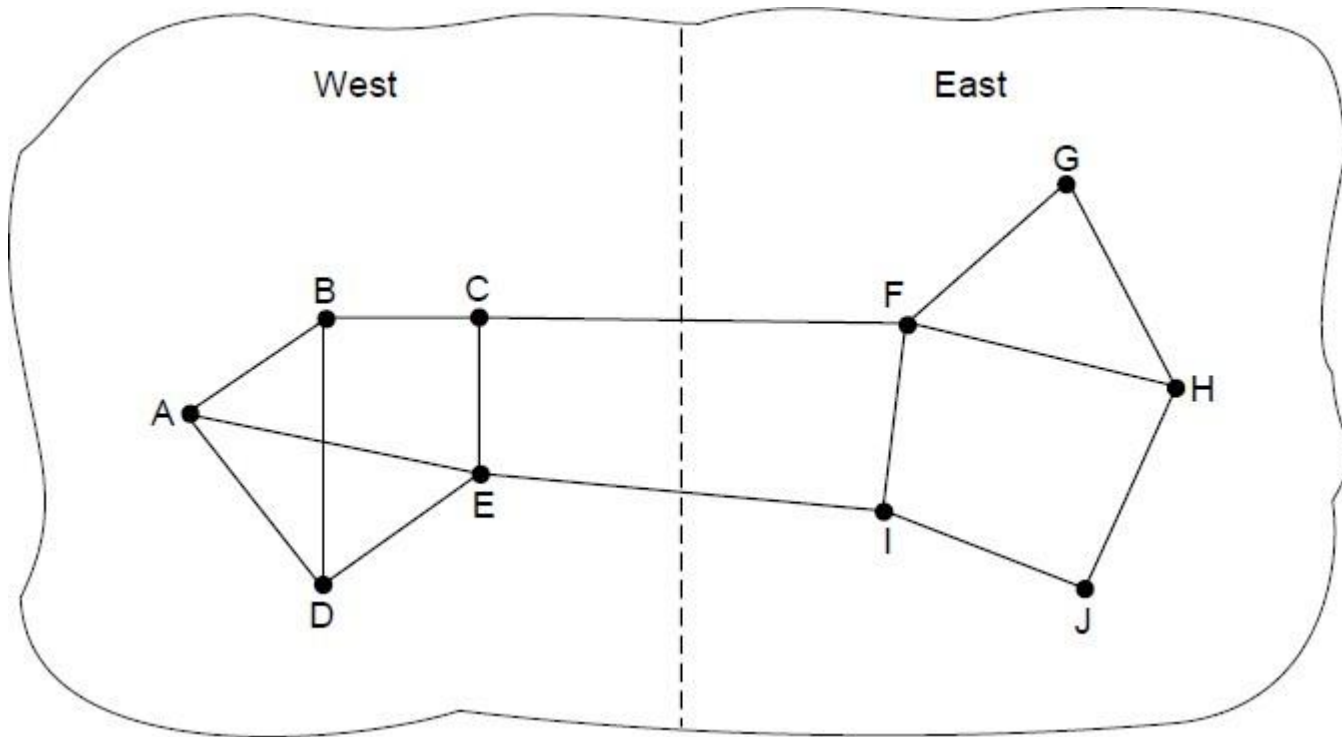
- Approaches to Congestion Control
- Traffic-Aware Routing
- Admission Control
- Traffic Throttling
- Load Shedding

Congestion



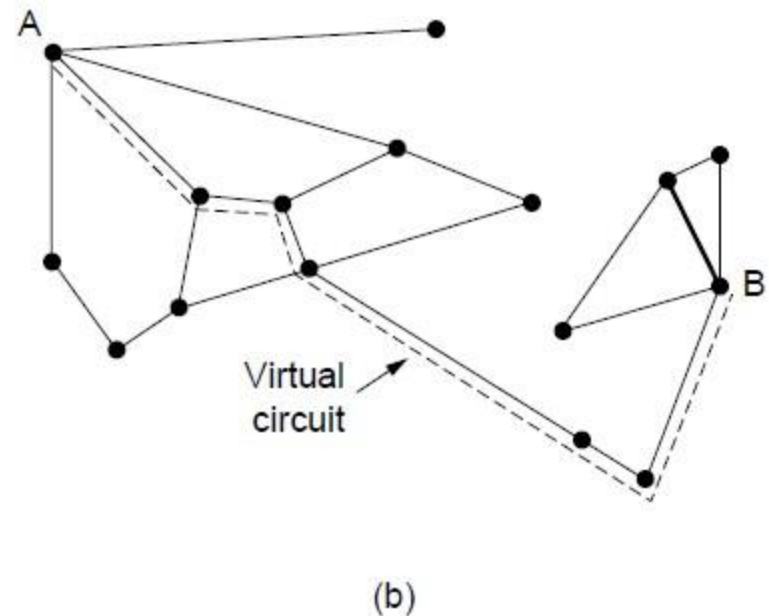
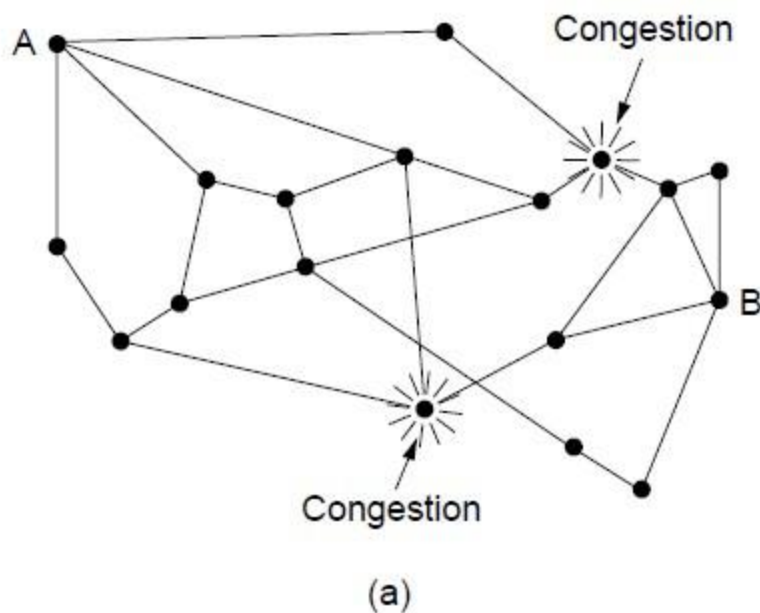
When too much traffic is offered, congestion sets in and performance degrades sharply.

Traffic-Aware Routing



A network in which the East and West parts are connected by two links.

Admission Control



(a) A congested network. (b) The portion of the network that is not congested. A virtual circuit from A to B is also shown.

Problem is in virtual circuits – there may be

Traffic

Throttling & Congestion

- Routers must detect new congestion as it is approaching, ideally before it has arrived.
- Each router can continuously monitor the resources it is using.
- 3 possibilities:
 1. utilization of the output links
 2. buffering of queued packets inside the router (most useful)
 3. no. of packets that are lost due to insufficient buffering

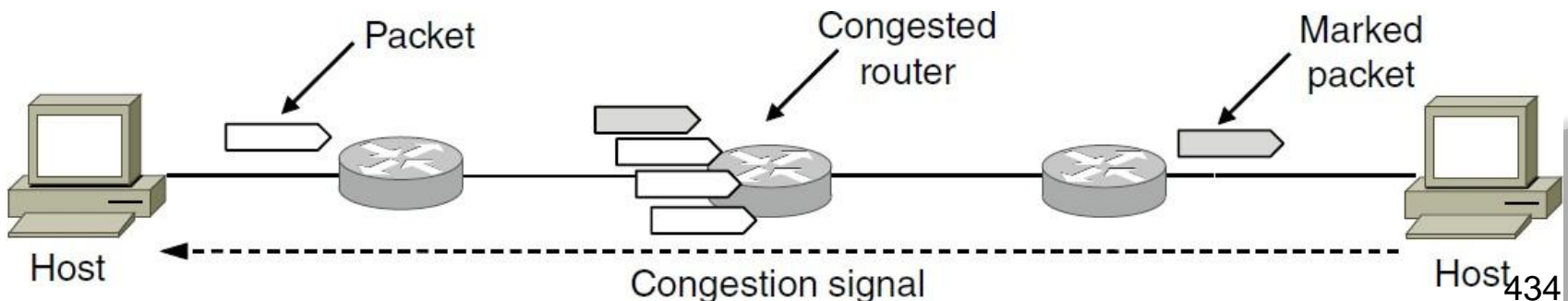
EWMA (Exponentially Weighted Moving Average)

- $d_{\text{new}} = \alpha d_{\text{old}} + (1 - \alpha)s$,
where, the constant α determines how fast the router forgets recent history.

Traffic Throttling: Feedback

- Routers must deliver timely feedback to the senders that are causing the congestion.
- The router must identify the appropriate senders.
- It must then warn them carefully, without sending many more packets into the already congested network.
- Many feedback mechanisms:

Mechanism 1: Explicit Congestion Notification (ECN)

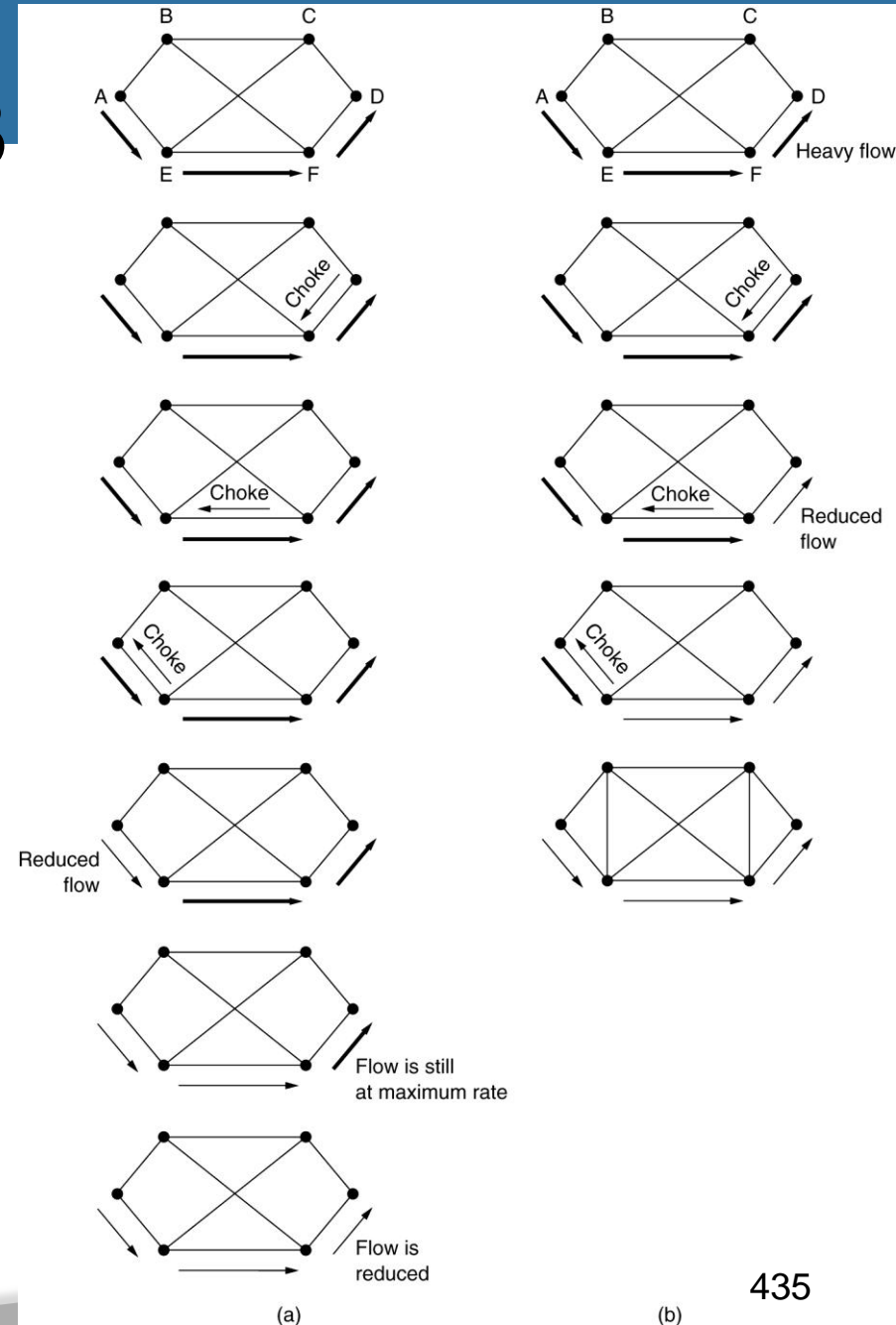


Mechanisms 2 & 3

Direct Choke Packets, Hop-by-Hop Backpressure

M-2: A choke packet that affects only the source.

M-3: A choke packet that affects each hop it passes through.



Mechanism 4: Load Shedding

- Performed when all other strategies fail.
- Cause blackout in some areas to save the entire network from failing.
- Intelligent packet drop policy desired.
- Which packets to discard may depend on application

Multimedia – old packets (full frame not to be discarded)

Text – Recent Packets

- Packet's importance can be marked in the beginning (application layer), then decision on which packets to discard can be taken

Quality of Service

- Requirements
 - Minimum throughput and maximum latency
- Techniques for Achieving Good Quality of Service
- Integrated Services
- Differentiated Services
- Label Switching and MPLS

Requirements

Application	Reliability	Delay	Jitter	Bandwidth
E-mail	High	Low	Low	Low
File transfer	High	Low	Low	Medium
Web access	High	Medium	Low	Medium
Remote login	High	Medium	Medium	Low
Audio on demand	Low	Low	High	Medium
Video on demand	Low	Low	High	High
Telephony	Low	High	High	Low
Videoconferencing	Low	High	High	High

Categories of QoS and Examples

1. Constant bit rate

- Telephony

2. Real-time variable bit rate

- Compressed videoconferencing

3. Non-real-time variable bit rate

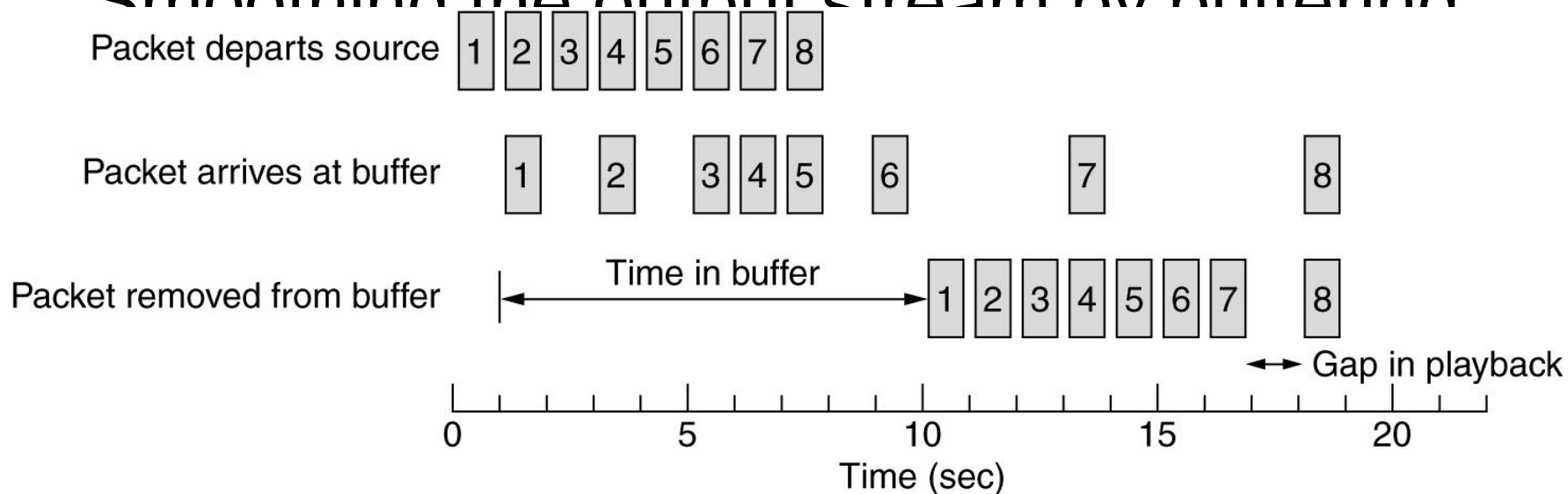
- Watching a movie on demand

4. Available bit rate

- File transfer

Buffering

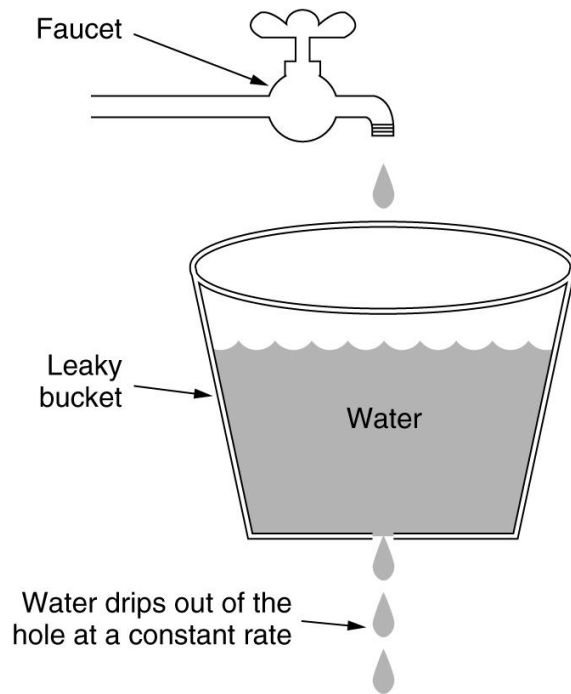
Smoothing the output stream by buffering



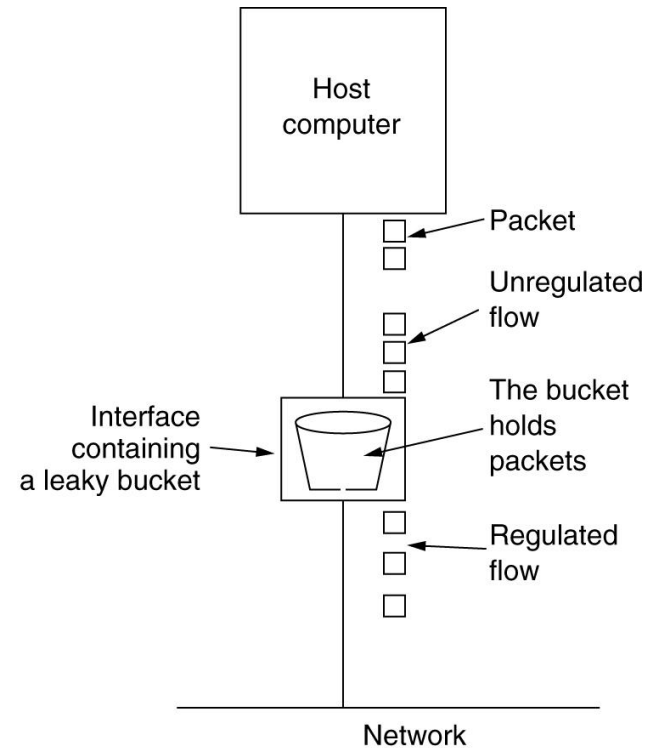
Traffic Shaping

- Traffic in data networks is **bursty** – typically arrives at non-uniform rates as the traffic rate varies.
- **Traffic shaping** is a technique for regulating the average rate and burstiness of a flow of data that enters the network.
- When a flow is set up, the user and the network agree on a certain traffic pattern (shape).
- Sometimes this agreement is called an **SLA (Service Level Agreement)**.

The Leaky Bucket Algorithm



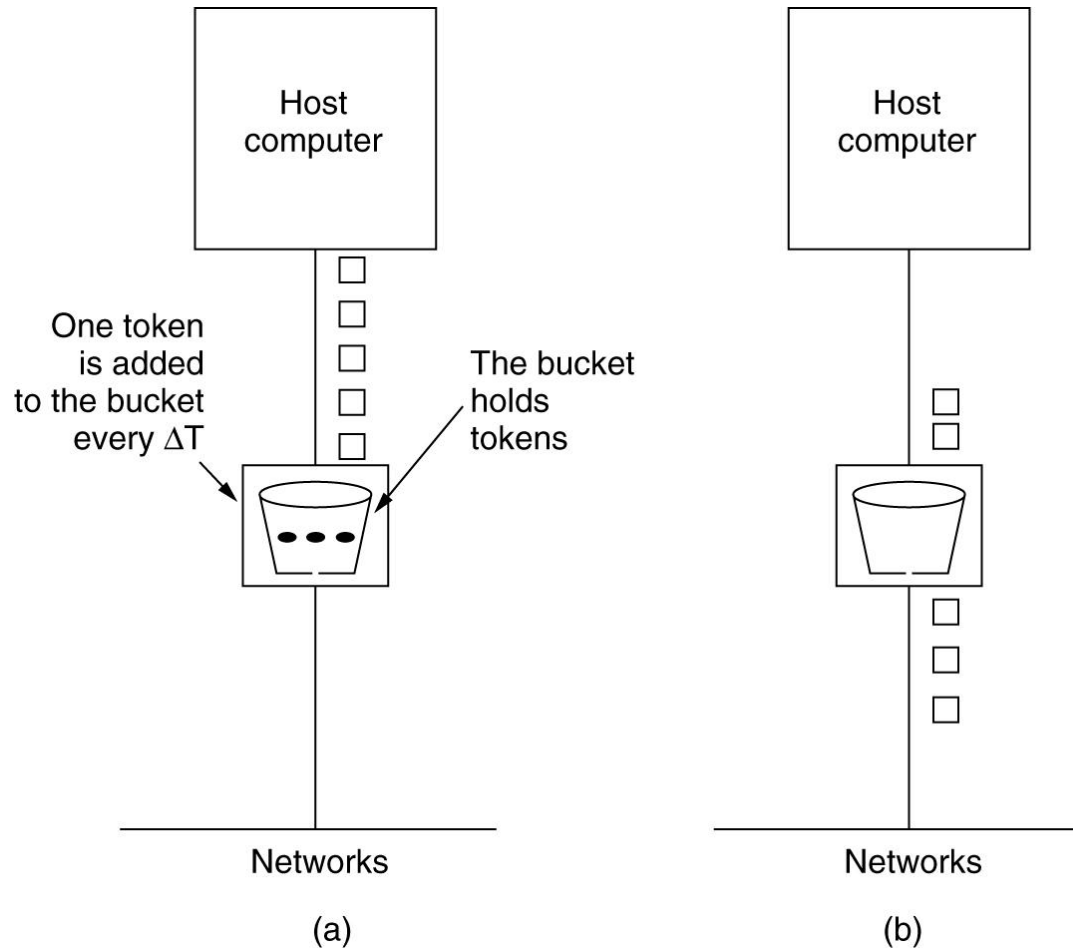
(a)



(b)

(a) A leaky bucket with water. (b) a leaky bucket with packets.

The Token Bucket Algorithm

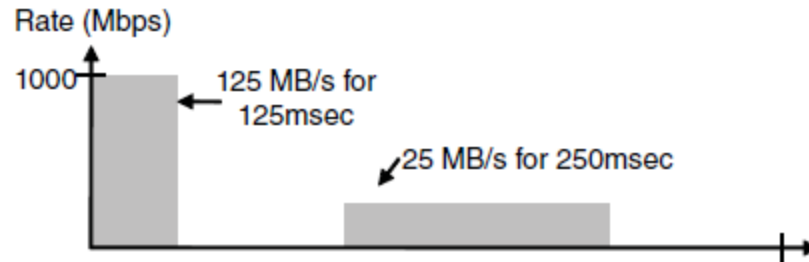


(a) Before. (b) After.

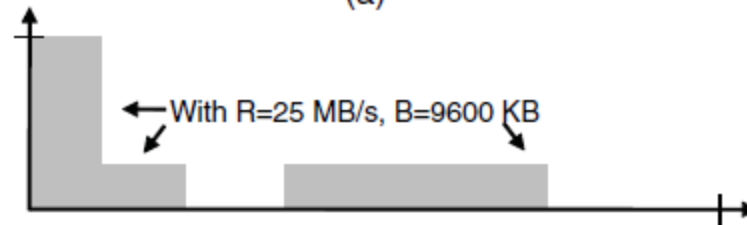
Token Bucket Algorithm (2)

- Burst length – S sec.
- Maximum output rate – M bytes/sec
- Token bucket capacity – B bytes
- Token arrival rate – R bytes/sec
- An output burst contains a maximum of $(B + RS)$ bytes.
- The number of bytes in a maximum speed burst of length S seconds is MS .
- Hence, we have: $B + RS = MS$
- This equation can be solved to get $S = B / (M - R)$

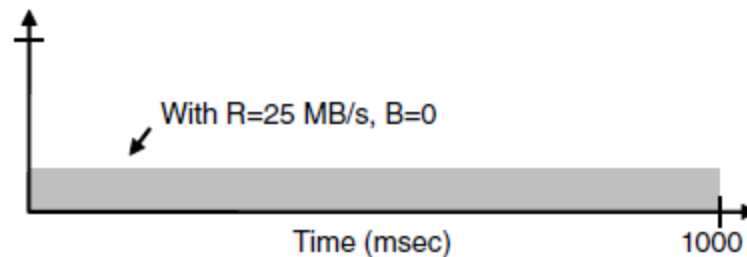
Traffic Sh



(a)



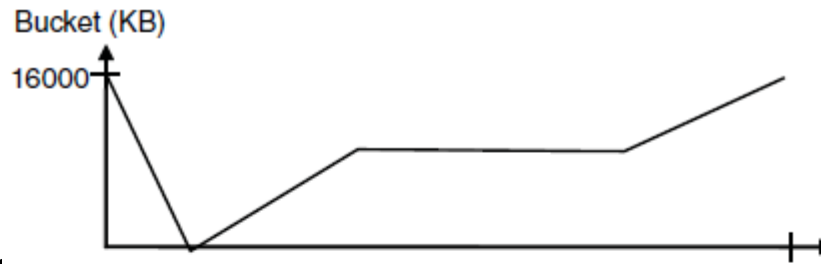
(b)



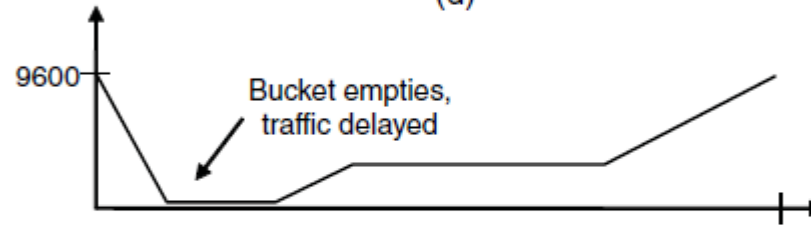
(c)

(a) Traffic from a host. Output shaped by a token bucket of rate 200 Mbps and capacity 447 (b)

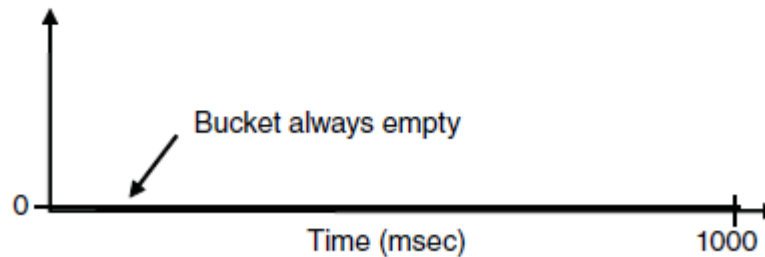
Traffic Shaping



(d)



(e)



(f)

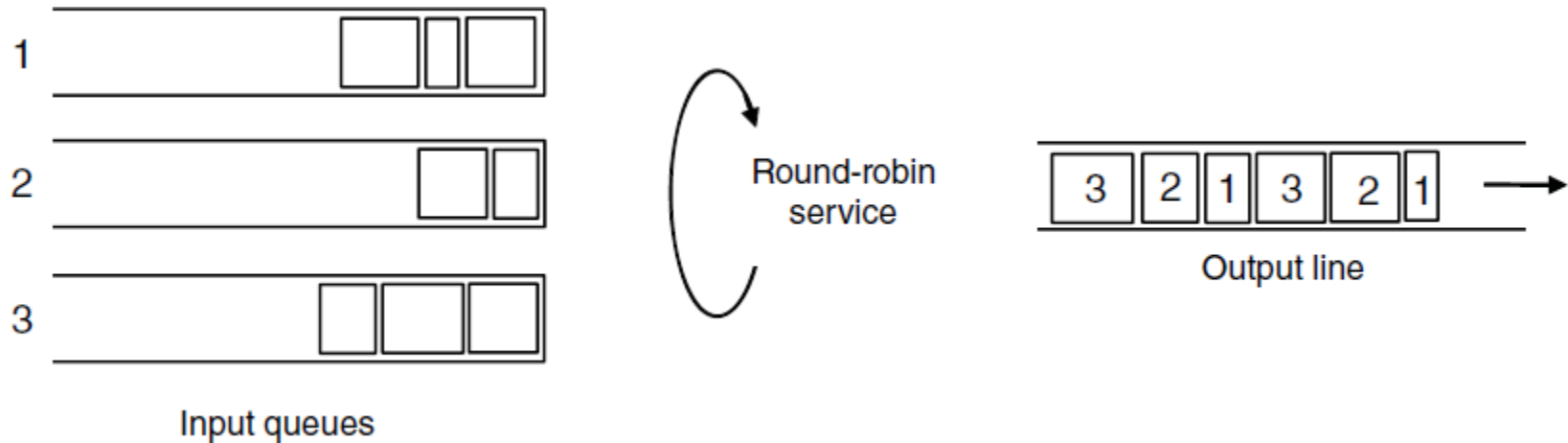
Token bucket level for shaping with rate 200 Mbps and capacity (d) 16000 KB, (e) 9600⁴⁴⁸

Packet Scheduling

Kinds of resources that can potentially be reserved for different flows:

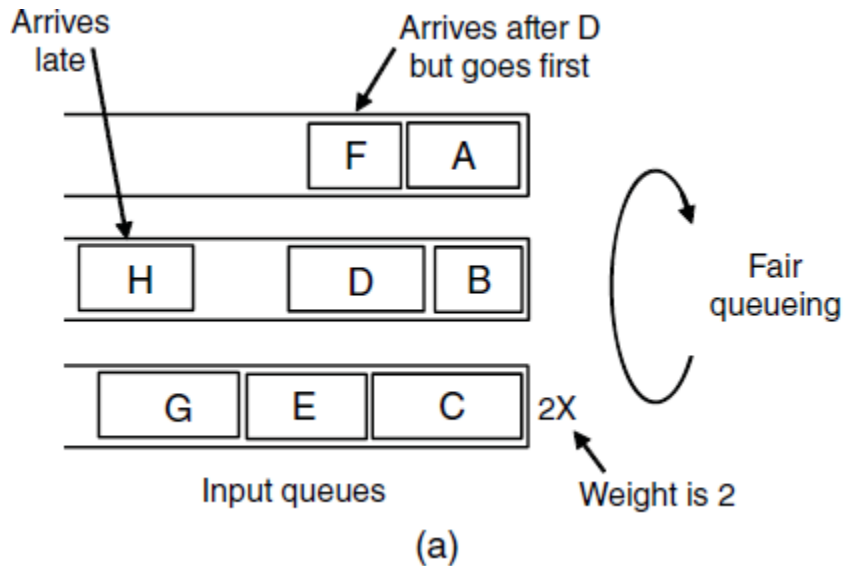
1. Bandwidth.
2. Buffer space.
3. CPU cycles.

Packet Scheduling (2)



Round-robin Fair Queuing

Packet Scheduling (3)



Packet	Arrival time	Length	Finish time	Output order
A	0	8	8	1
B	5	6	11	3
C	5	10	10	2
D	8	9	20	7
E	8	8	14	4
F	10	6	16	5
G	11	10	19	6
H	20	8	28	8

(b)

(a) Weighted Fair Queueing.

(b) Finishing times for the packets.

Admission Control (1)

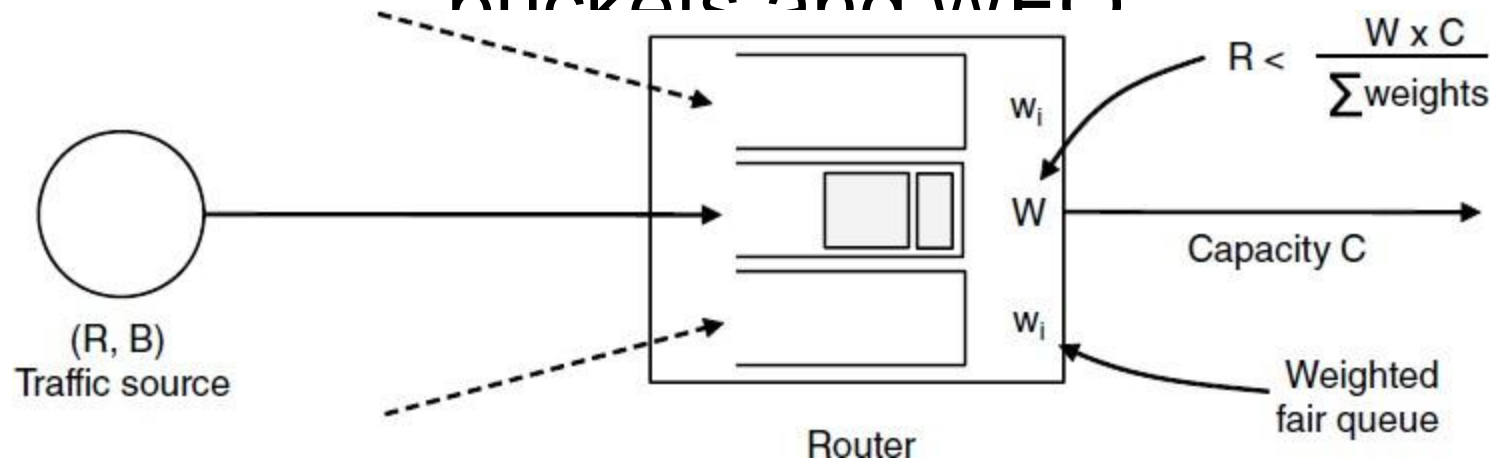
Parameter	Unit
Token bucket rate	Bytes/sec
Token bucket size	Bytes
Peak data rate	Bytes/sec
Minimum packet size	Bytes
Maximum packet size	Bytes

An example flow specification

- $T = 1/\mu \times 1/(1-\lambda/\mu) \text{ -- } \lambda$
 $= 0.95 \text{Mpackets/sec}$
- $\mu = 1 \text{Mb packets/sec}$

Admission Control (2)

Bandwidth and delay guarantees with token buckets and WRED



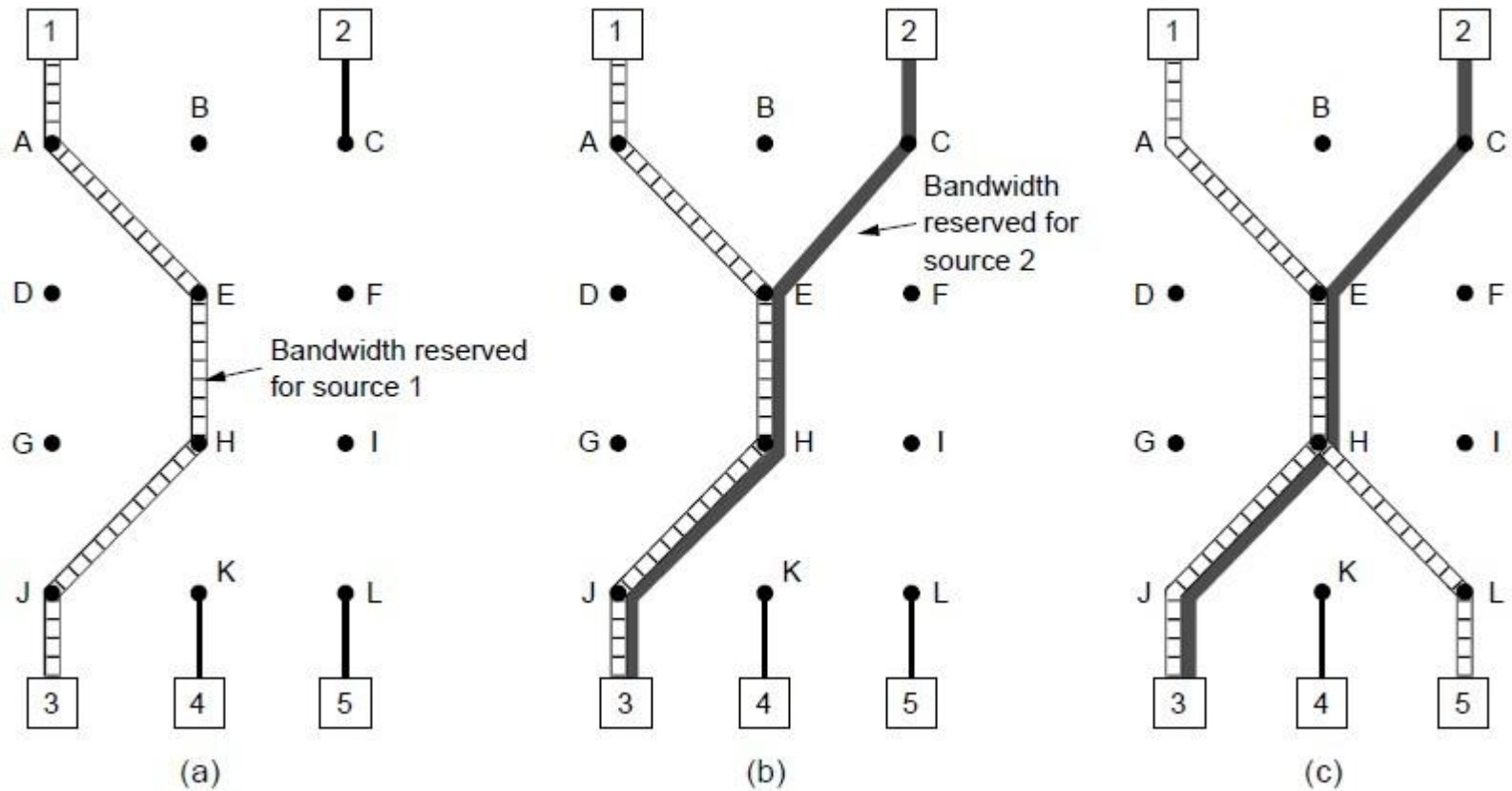
Hosts 1 and 2 are multicast sender

3,4, 5 are multicast receiver

Host 3 reserves for Host 1 and the Host 2

Host 5 reserves Host 1 (so the common path is utilized). However depending on need (Host 5 may be a bigger TV) – provision is made for the greediest part

RSVP (2)



(a) Host 3 requests a channel to host 1. (b) Host 3 then requests a second channel, to host 2.

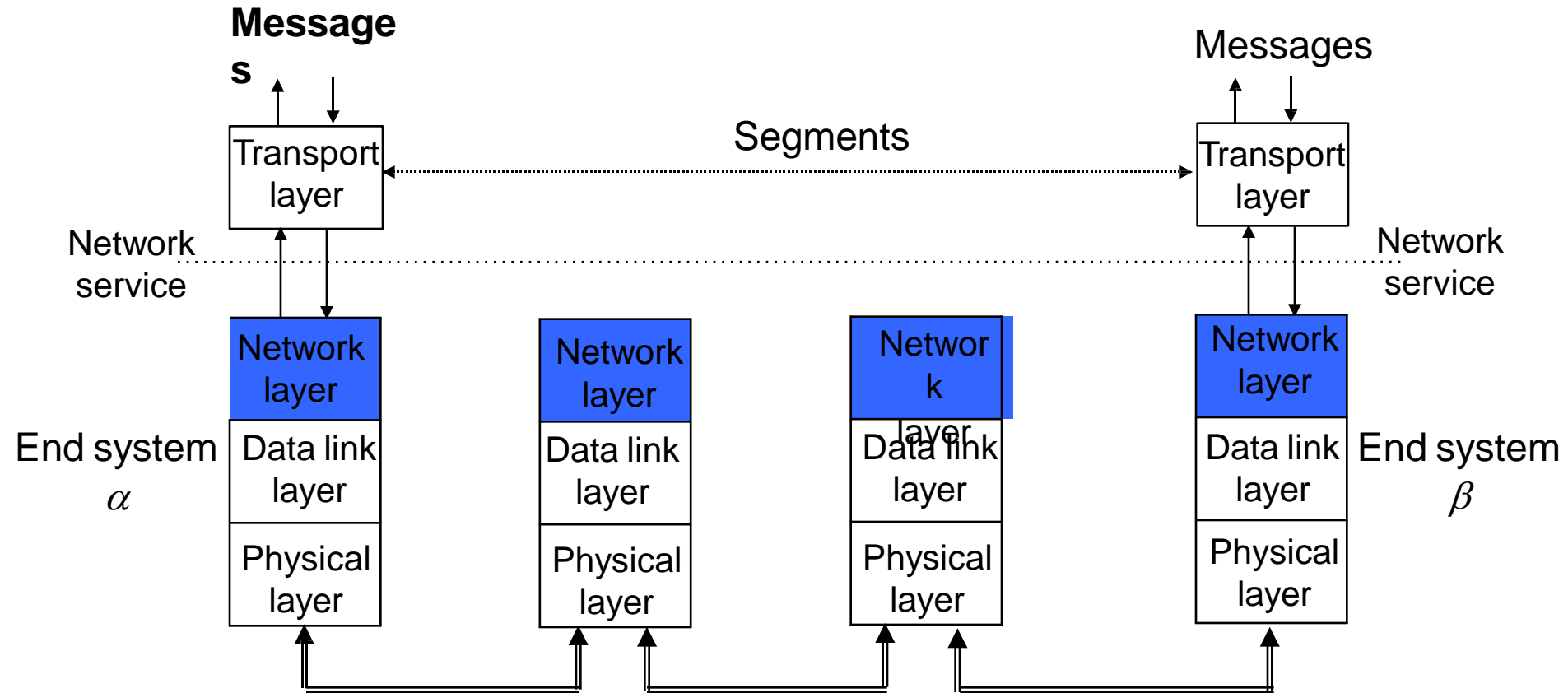
(c) Host 5 requests a channel to host 1

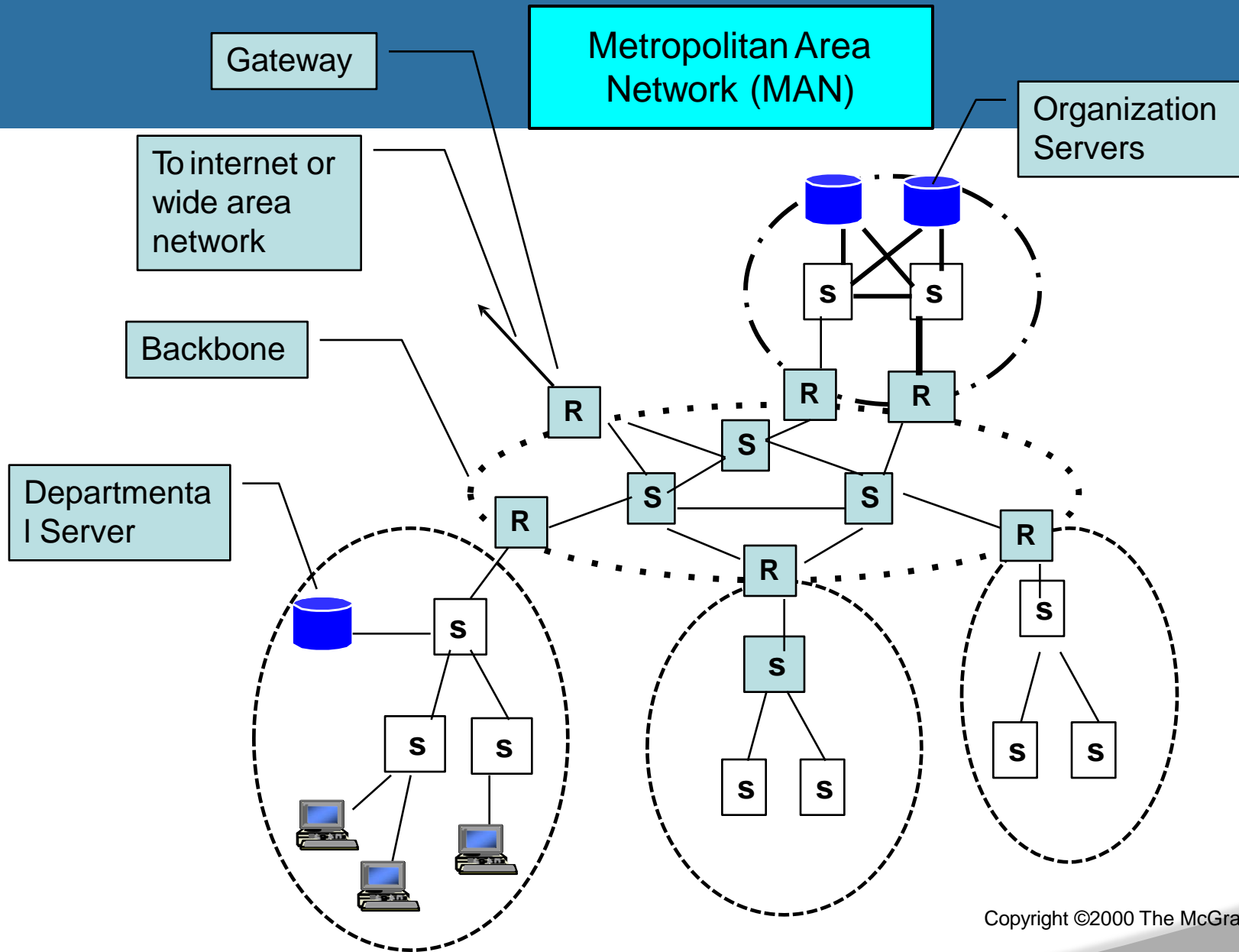
Network Layer

- Concerned with getting packets from source to destination.
- The network layer must know the topology of the subnet and choose appropriate paths through it.
- When source and destination are in *different networks*, the network layer **(IP)** must deal with these differences.
- * **Key issue:** *what service does the network layer provide to the transport layer (connection-oriented or connectionless).*

Network Layer Design Goals

1. The services provided by the network layer should be **independent** of the subnet topology.
2. The Transport Layer should be shielded from the number, type and topology of the subnets present.
3. The network addresses available to the Transport Layer should use a uniform numbering plan (even across LANs and WANs).





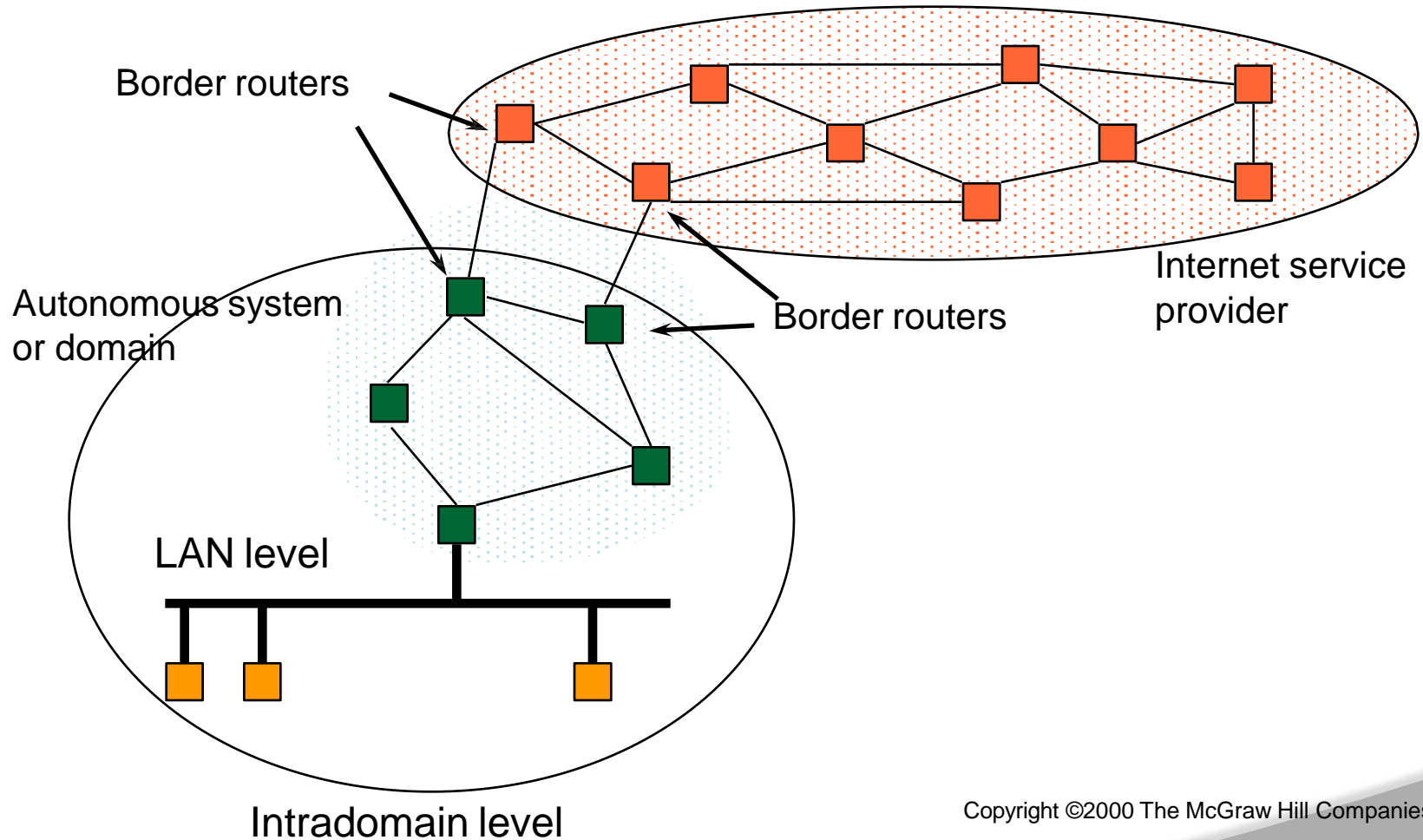
Copyright ©2000 The McGraw Hill Companies

Leon-Garcia & Widjaja: *Communication Networks*

Figure 7.6

Wide Area Network (WAN)

Interdomain level

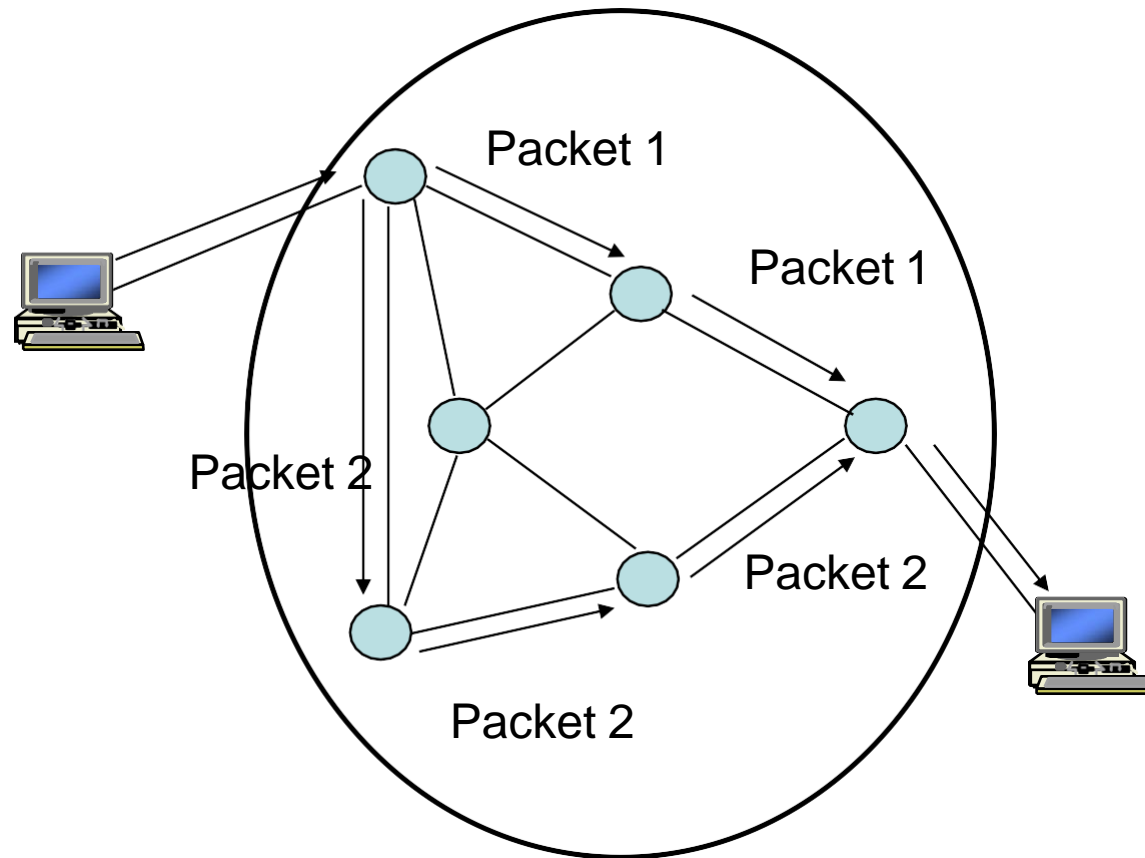


Copyright ©2000 The McGraw Hill Companies

Leon-Garcia & Widjaja: *Communication Networks*

Figure 7.7

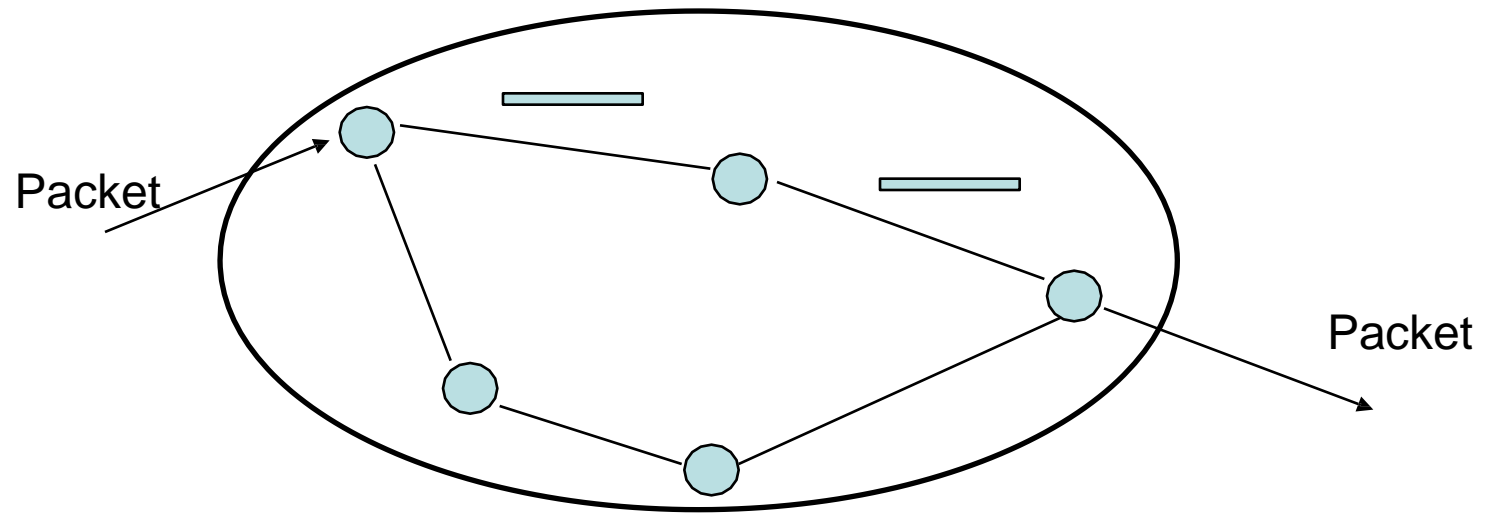
Datagram Packet Switching



Routing Table in Datagram Network

Destination address	Output port
0785	7
1345	12
1566	6
2458	12

Virtual Circuit Packet Switching



Routing Table in Virtual Circuit Network

Identifier Output
 port
 identifier

12	13	44
15	15	23
27	13	16
58	7	34

Entry for packets
with identifier 15



Routing

Routing algorithm:: that part of the Network Layer responsible for deciding on which output line to transmit an incoming packet.

- Remember: For virtual circuit subnets the routing decision is made ONLY at set up.

Algorithm properties:: correctness, simplicity, robustness, stability, fairness, optimality, and scalability.

Routing Classification

Adaptive Routing

- based on current measurements of traffic and/or topology.

1. centralized
2. isolated
3. distributed

Non-Adaptive Routing

- routing computed in advance and off-line

1. flooding
2. static routing using shortest path algorithms

Flooding

- *Pure flooding* :: every incoming packet to a node is sent out on **every outgoing line**.
 - Obvious adjustment – do not send out on arriving link (assuming full-duplex links).
 - The routing algorithm can use a hop counter (e.g., TTL) to **dampen the flooding**.
 - *Selective flooding* :: only send on those lines going “approximately” in the right direction.

Shortest Path Routing

1. Bellman-Ford Algorithm [Distance Vector]
2. Dijkstra's Algorithm [Link State]

What does it mean to be the shortest (or optimal) route?

Choices:

- a. Minimize the number of hops along the path.
- b. Minimize mean packet delay.
- c. Maximize the network throughput.

Metrics

- Set all link costs to 1.
 - Shortest hop routing.
 - Disregards delay.
- Measure the number of packets queued.
 - Did not work well.
- Timestamp **ArrivalTime** and **DepartTime*** and use link-level ACK to compute:

$$\text{Delay} = (\text{DepartTime} - \text{ArrivalTime}) + \text{TransmissionTime} + \text{Latency}$$

* Reset after retransmission

Metrics

- Unstable under heavy link load.
- Difficulty with granularity of the links.
- Revised ARPANET routing metric:
 - Compress dynamic range of the metric
 - Account for link type
 - Smooth variation of metric with time:
 - Delay transformed into link utilization
 - Utilization was averaged with last reported utilization.
 - Hard limit set on how much the metric could change per measurement cycle.

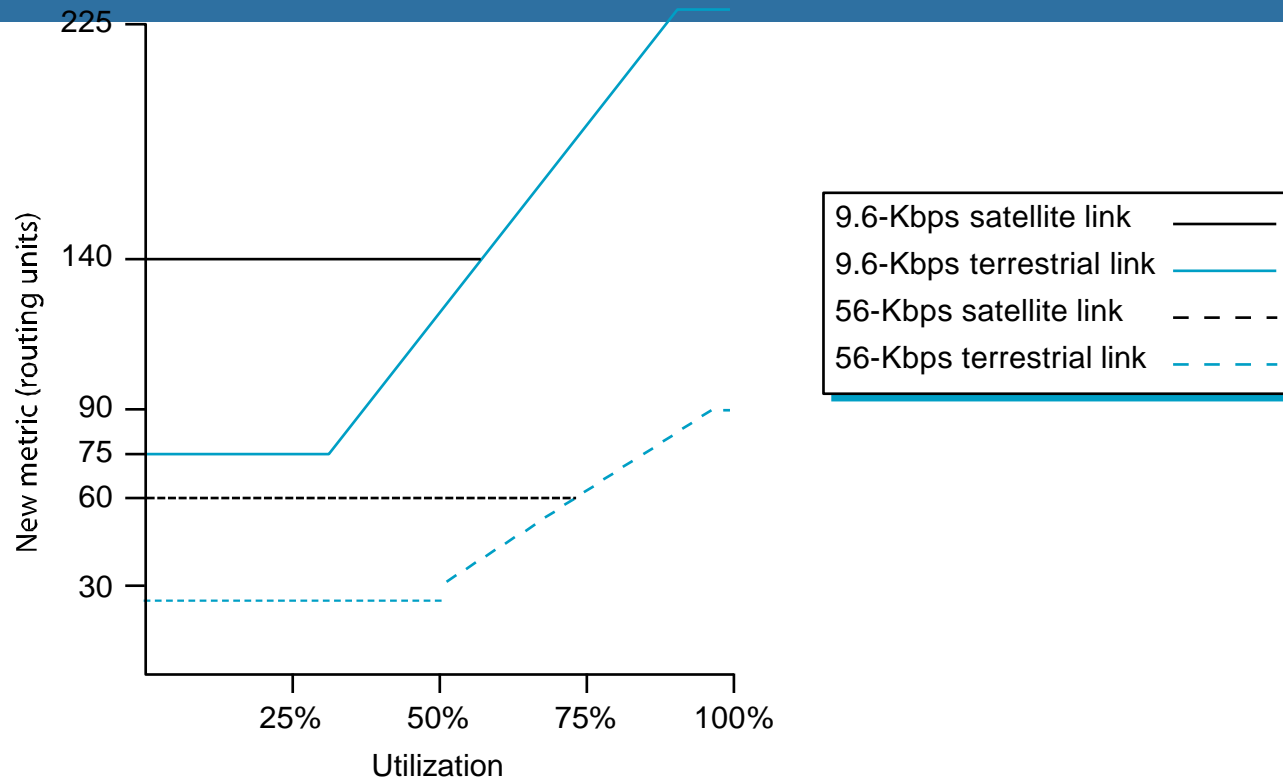
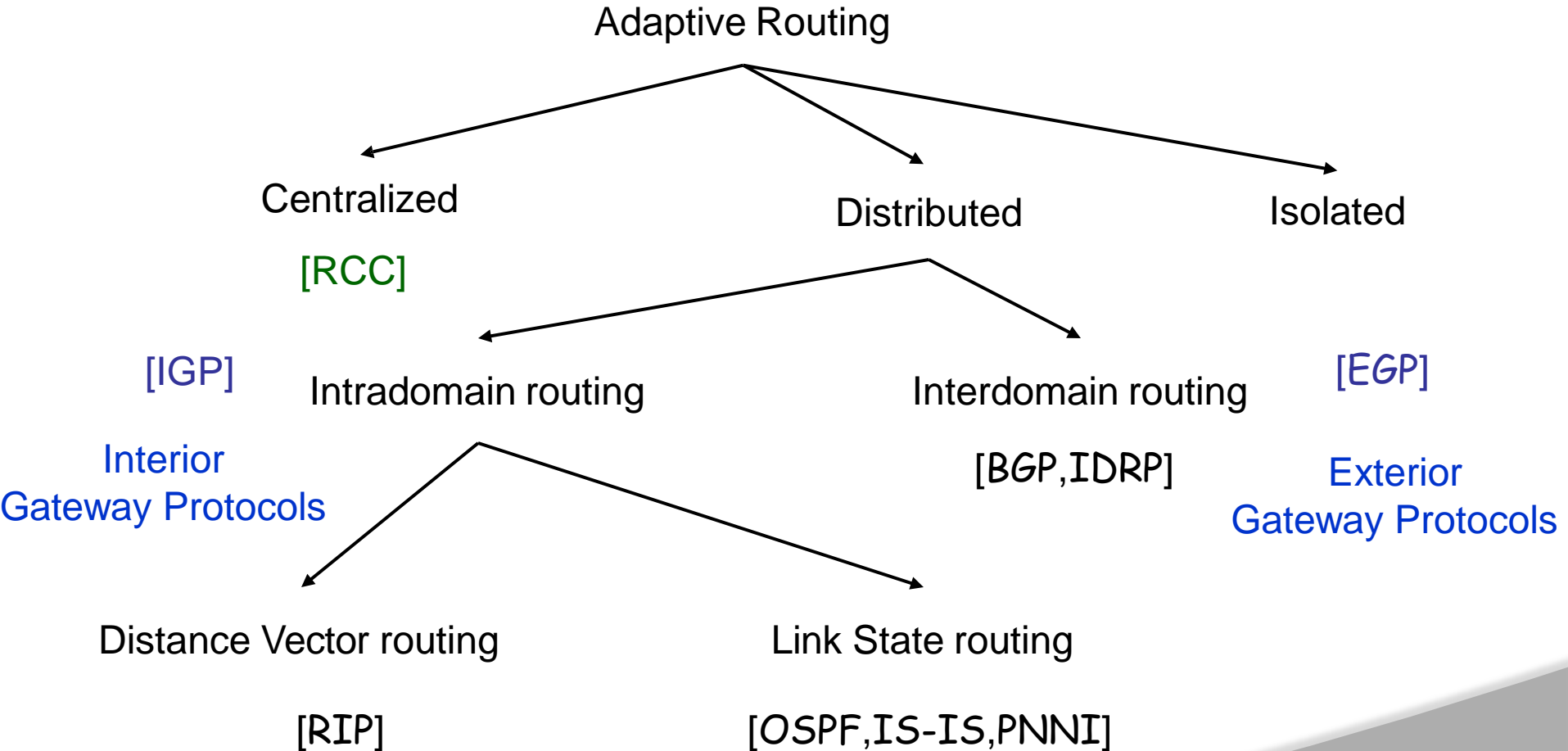


Figure 4.22 Revised ARPANET routing metric versus link utilization

P&D slide

Internetwork Routing [Halsall]



Adaptive Routing

Basic functions:

- 1.Measurement of pertinent network data.
- 2.Forwarding of information to where the routing computation will be done.
- 3.Compute the routing tables.
- 4.Convert the routing table information into a **routing decision** and then **dispatch** the data packet.

Adaptive Routing

Design Issues:

1. How much **overhead** is incurred due to gathering the routing information and sending **routing packets**?
2. What is the time frame (i.e, the frequency) for sending **routing packets** in support of adaptive routing?
3. What is the **complexity** of the routing strategy?

Distance Vector Routing

- Historically known as the **old** ARPANET routing algorithm {or known as *Bellman-Ford algorithm*}.

Basic idea: each network node maintains a Distance Vector table containing the ***distance*** between itself and **ALL possible destination nodes**.

- Distances are based on a chosen metric and are computed using information from the **neighbors'** distance vectors.

Metric: usually hops or delay

Distance Vector Routing

Information kept by DV router

1. each router has an ID
2. associated with each link connected to a router, there is a link cost (static or dynamic) **the metric issue!**

Distance Vector Table Initialization

Distance to itself = 0

Distance to ALL other routers = infinity number

Distance Vector Routing

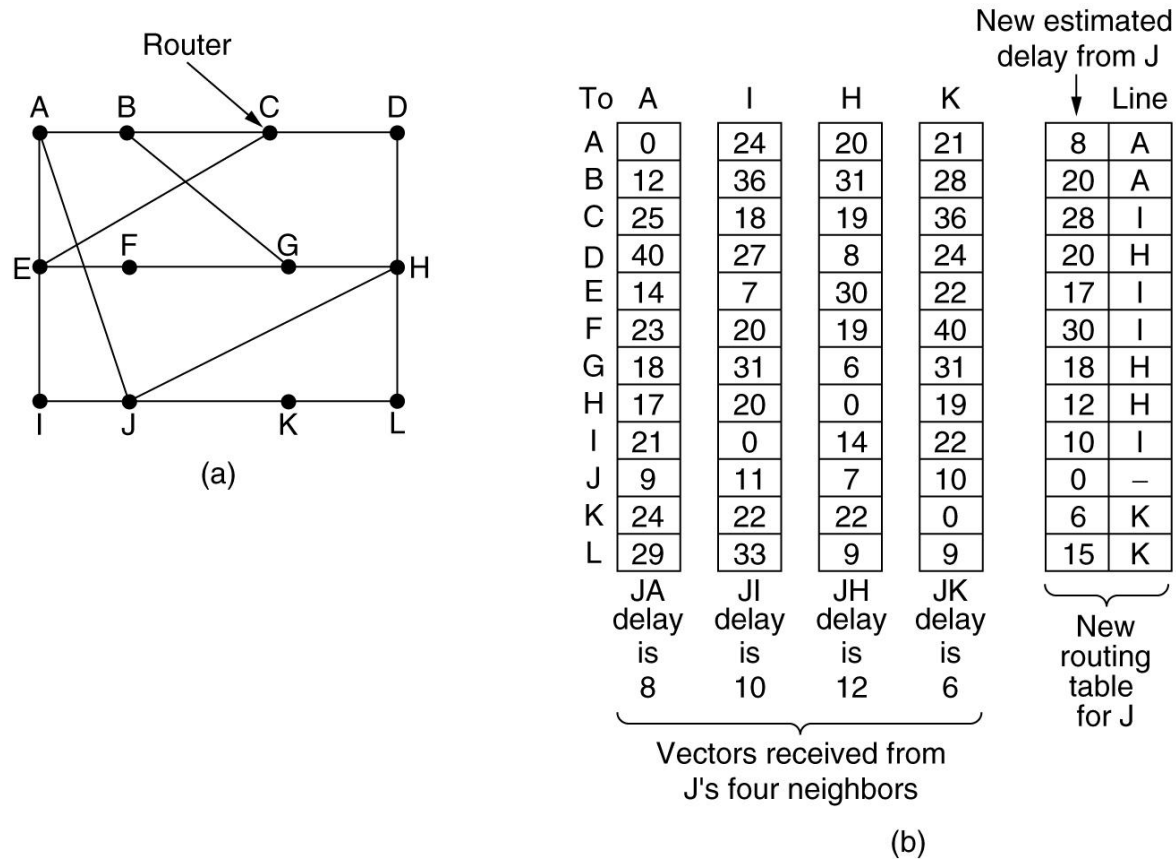
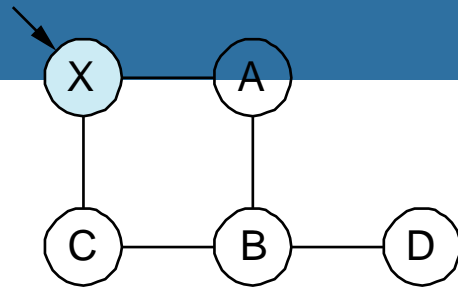


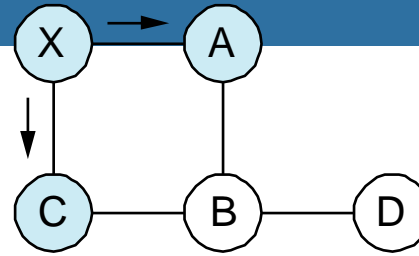
Figure 5-9.(a) A subnet. (b) Input from A, I, H, K, and the new routing table for J.

Link State Algorithm

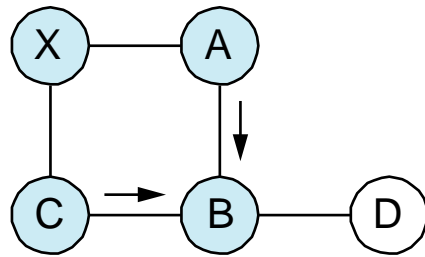
1. Each router is responsible for meeting its neighbors and learning their names.
2. Each router constructs a **link state packet (LSP)** which consists of a list of names and cost to reach each of its neighbors.
3. The **LSP** is transmitted to **ALL other routers**. Each router stores the most recently generated **LSP** from each other router.
4. Each router uses complete information on the network topology to compute the *shortest path route* to each destination node.



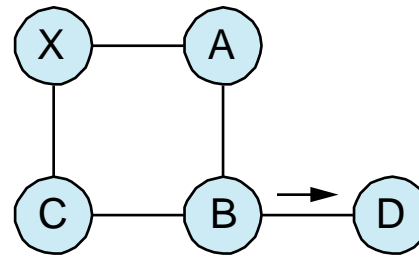
(a)



(b)



(c)



(d)

Figure 4.18 Reliable LSP Flooding

P&D slide

OSPF Terminology

Internal router :: a level 1 router.

Backbone router :: a level 2 router.

Area border router (ABR) :: a **backbone** router that attaches to more than one area.

AS border router :: (an interdomain router), namely, a router that attaches to routers from other ASs across AS boundaries.

OSPF LSA Types

- 1.Router link advertisement [Hello message]
- 2.Network link advertisement
- 3.Network summary link advertisement
- 4.AS border router's summary link advertisement
- 5.AS external link advertisement

LS Age		Options		Type=1
Link-state ID				
Advertising router				
LS sequence number				
LS checksum			Length	
0	Flags	0	Number of links	
Link ID				
Link data				
Link type	Num_TOS		Metric	
Optional TOS information				
More links				

Indicates
LSA type

Indicates
link cost

Figure 4.21 OSF Type 1 Link-State Advertisement

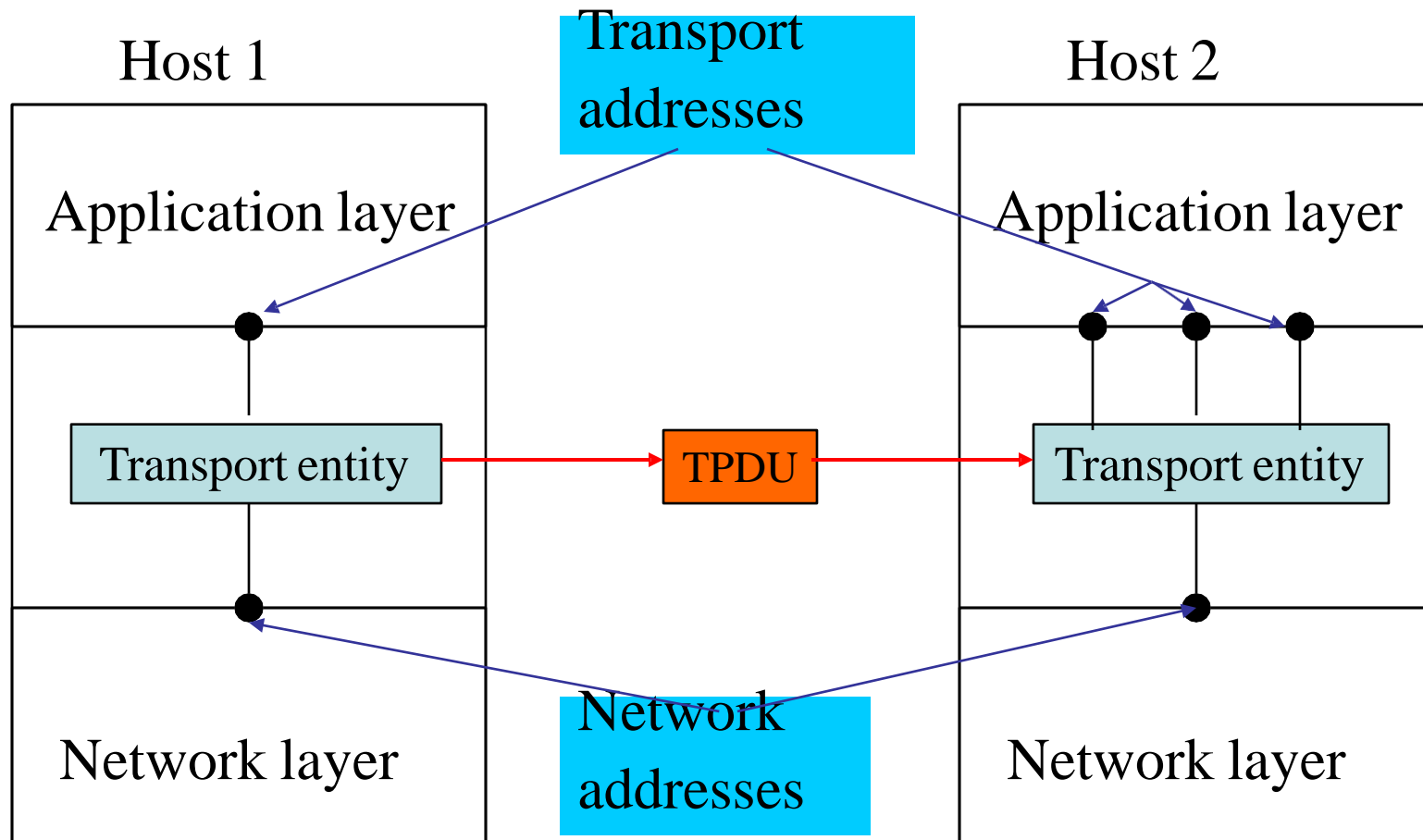
P&D slide

The Transport layer

Transport Layer

- **Services**
- Elements of transport protocol
- Simple transport protocol
- UDP
- Remote Procedure Call (see Distributed Systems)
- TCP

Layer overview



Services

- needed from network layer
 - packet transport between hosts
 - relationship network <> transport
 - Hosts <> processes
 - Transport service
 - independent network
 - more reliable
 - Network
 - run by carrier
 - part of communication subnet for WANs

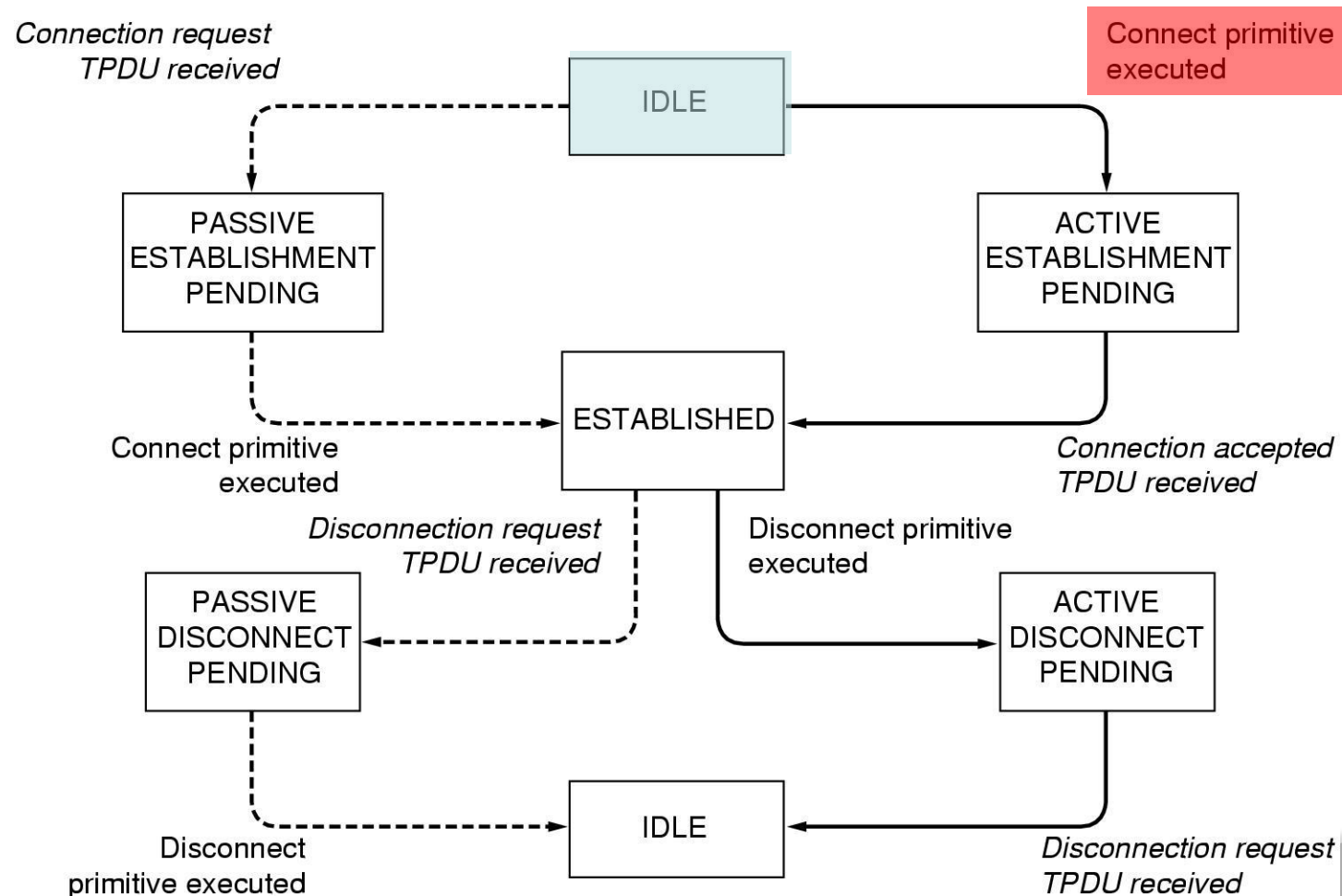
Simple service: primitives

- Simple primitives:
 - connect
 - send
 - receive
 - disconnect
- How to handle incoming connection request in server process?
 - ➔ Wait for connection request from client!
 - listen

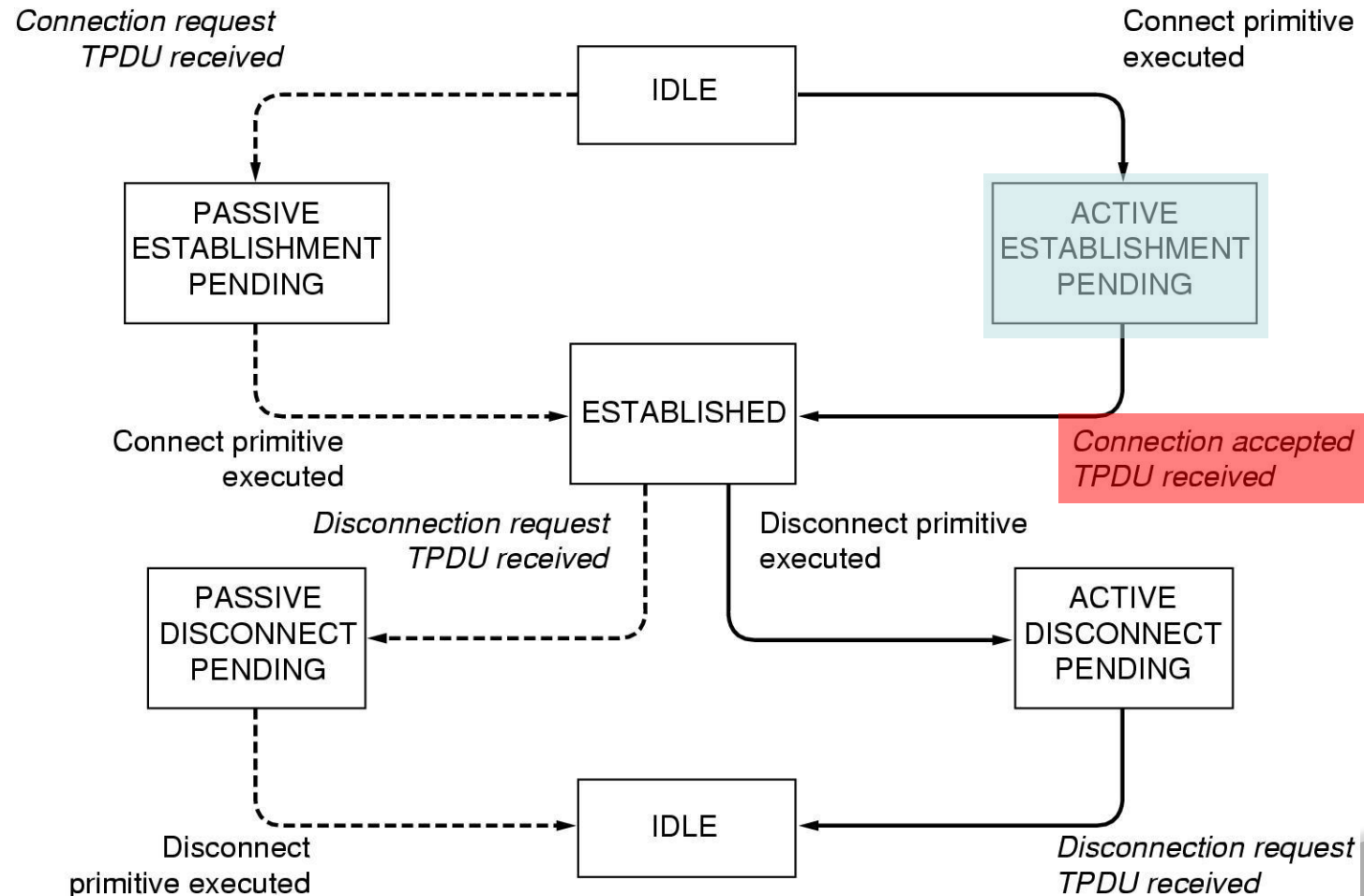
Simple service: primitives

listen	No TPDU
connect	Connection Request TPDU
send	Data TPDU
receive	No TPDU
disconnect	Disconnect TPDU

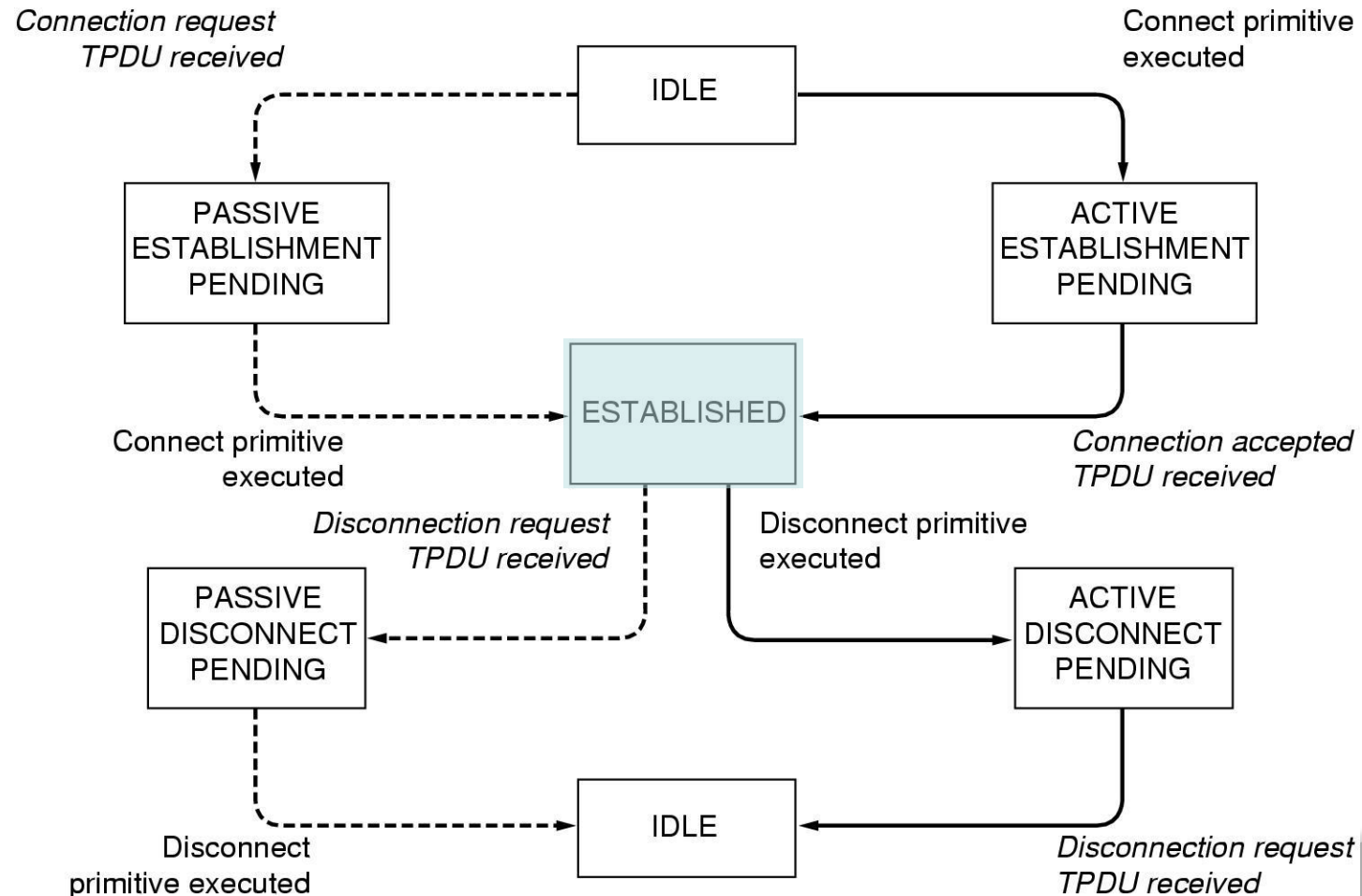
Simple service: state diagram



Simple service: state diagram



Simple service: state diagram



Berkeley service primitives

socket	create new communication end point
bind	attach a local address to a socket
listen	announce willingness to accept connections; give queue size
accept	block caller until a connection request arrives
connect	actively attempt to establish a connection
send	send some data over the connection
receive	receive some data from the connection
close	release the connection

Transport Layer

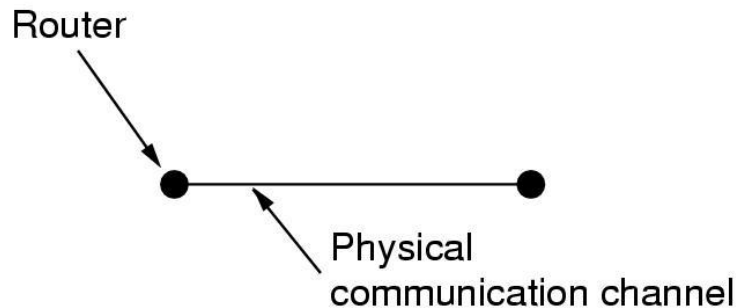
- Services
- Elements of transport protocol
- Simple transport protocol
- UDP
- Remote Procedure Call (see Distributed Systems)
- TCP

Elements of transport protocols (*etp*)

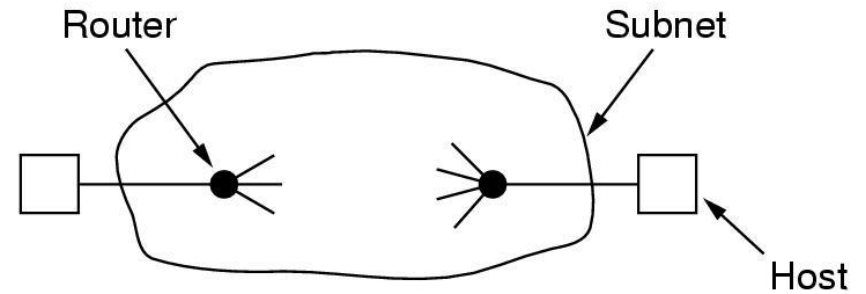
- Transport <> Data Link
- Addressing
- Establishing a connection
- Releasing a connection
- Flow control and buffering
- Multiplexing
- Crash recovery

etp: Transport \leftrightarrow data link

- Physical channel \leftrightarrow subnet



(a)



(b)

Explicit addressing

Connection establishment

Potential existence of storage capacity in subnet

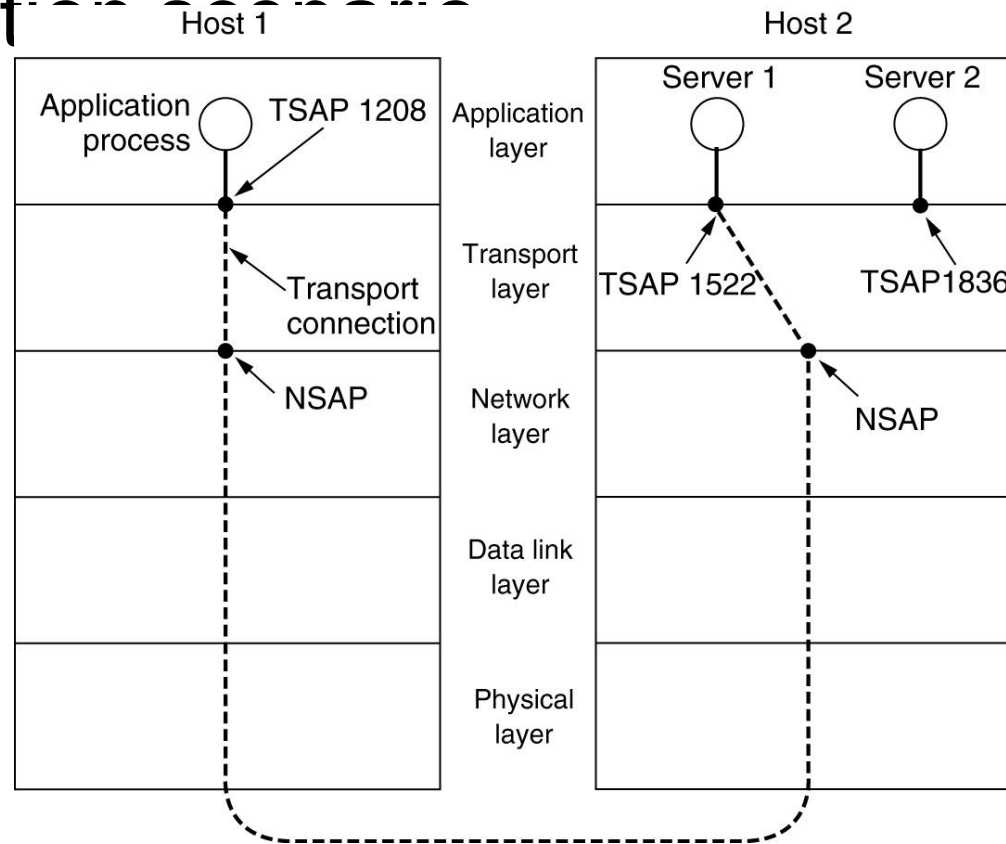
Dynamically varying number of connections

etp: Addressing

- TSAP = transport service access point
 - Internet: IP address + local port
 - ATM: AAL-SAPs
- *Connection scenario*
- *Getting TSAP addresses?*
- *From TSAP address to NSAP address?*

etp: Addressing

- Connect



etp: Addressing

- Connection scenario
 - Host 2 (server)
 - Time-of-day server attaches itself to TSAP 1522
 - Host 1 (client)
 - Connect from TSAP 1208 to TSAP 1522
 - Setup network connection to host 2
 - Send transport connection request
 - Host 2
 - Accept connection request

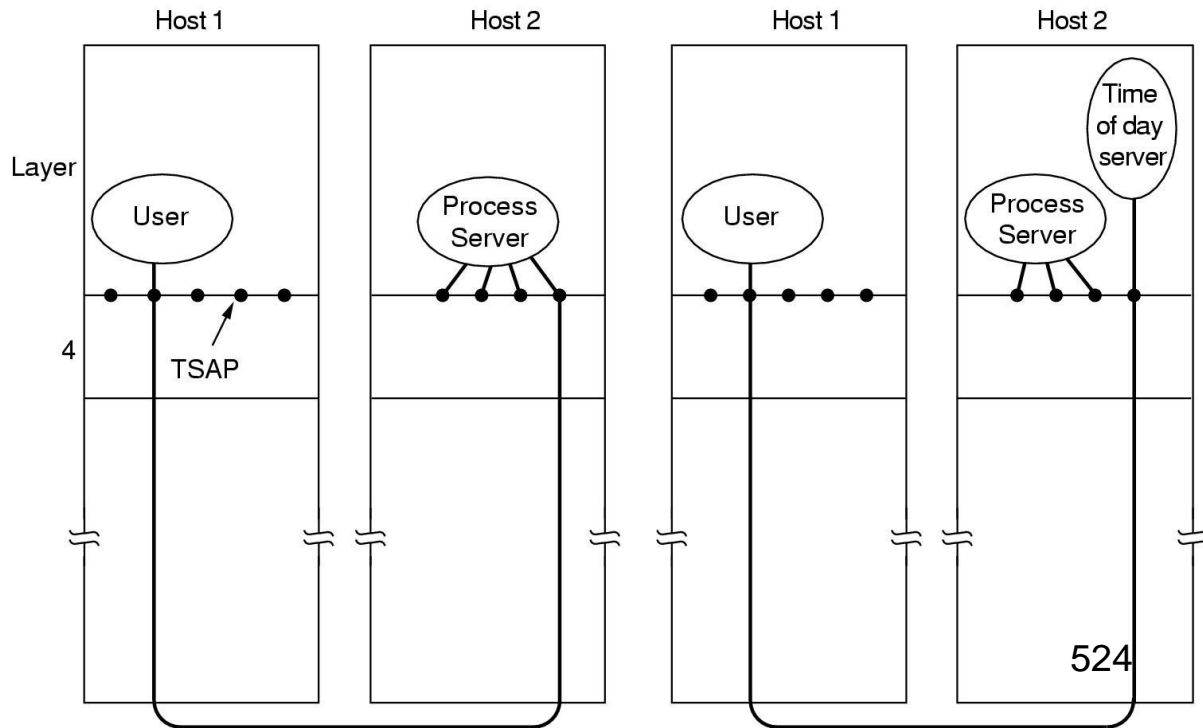
etp: Addressing

- Getting TSAP addresses?
 - Stable TSAP addresses
 - For key services
 - Not for user processes
 - active for a short time
 - number of addresses limited
 - Name servers
 - to find existing servers
 - map service name into TSAP address
 - *Initial connection protocol*

etp: Addressing

- Getting TSAP addresses?
 - Initial connection protocol
 - to avoid many waiting servers → one process server

- waits on
- creates r



etp: Addressing

- From TSAP address to NSAP address?
 - hierarchical addresses
 - address = <country> <network> <host> <port>
 - Examples: IP address + port
 - Telephone numbers (<> number portability?)
 - Disadvantages:
 - TSAP bound to host!
 - flat address space
 - Advantages:
 - Independent of underlying network addresses
 - TSAP address not bound to host
 - Mapping to network addresses:
 - Name server
 - broadcast

etp: Establishing a connection

- Problem: delayed duplicates!
- Scenario:
 - Correct bank transaction
 - connect
 - data transfer
 - disconnect
 - Problem: same packets are received in same order a second time Recognized?

etp: Establishing a connection

- Unsatisfactory solutions:
 - throwaway TSAP addresses
 - need unlimited number of addresses?
 - process server solution impossible
 - connection identifier
 - Never reused!
 - ➔ Maintain state in hosts
- Satisfactory solutions

etp: Establishing a connection

- Satisfactory solutions
 - Ensure limited packet lifetime (incl. Acks)
 - Mechanisms
 - prevent packets from looping + bound congestion delay
 - hopcounter in each packet
 - timestamp in each packet
 - Basic assumption
 - Maximum packet lifetime T
 - If we wait a time T after sending a packet all traces of it (including Acks) are gone

etp: Establishing a connection

- Tomlinson's method
 - requires: clock in each host
 - Number of bits $>$ number of bits in sequence number
 - Clock keeps running, even when a hosts crashes
 - Basic idea:

2 identically numbered TPDU's are never outstanding at the same time!

etp: Establishing a connection

- Tomlinson's rule: Never reuse a sequence number x within the lifetime T for the packet with x

– Problems to solve

- Selection of the initial sequence number for a new connection
- Wrap around of sequence numbers for an active connection
- Handle host crashes

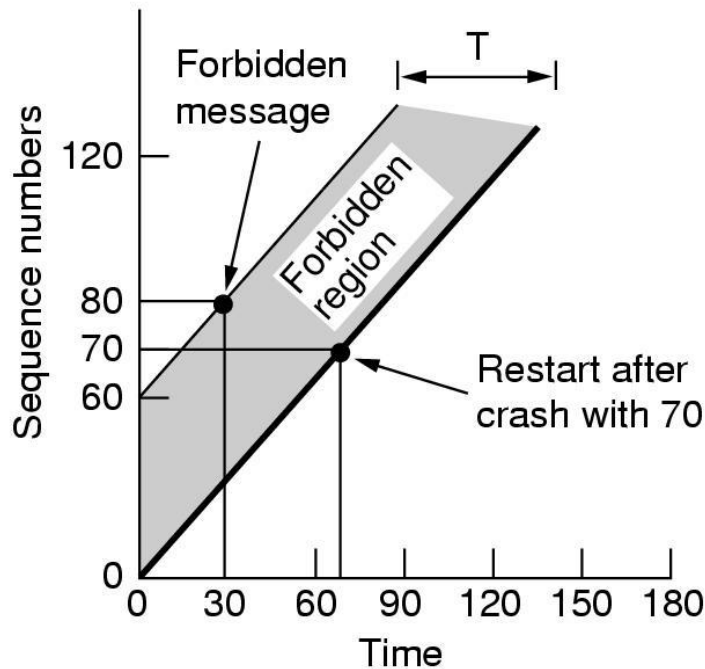
➔ Forbidden region

etp: Establishing a connection

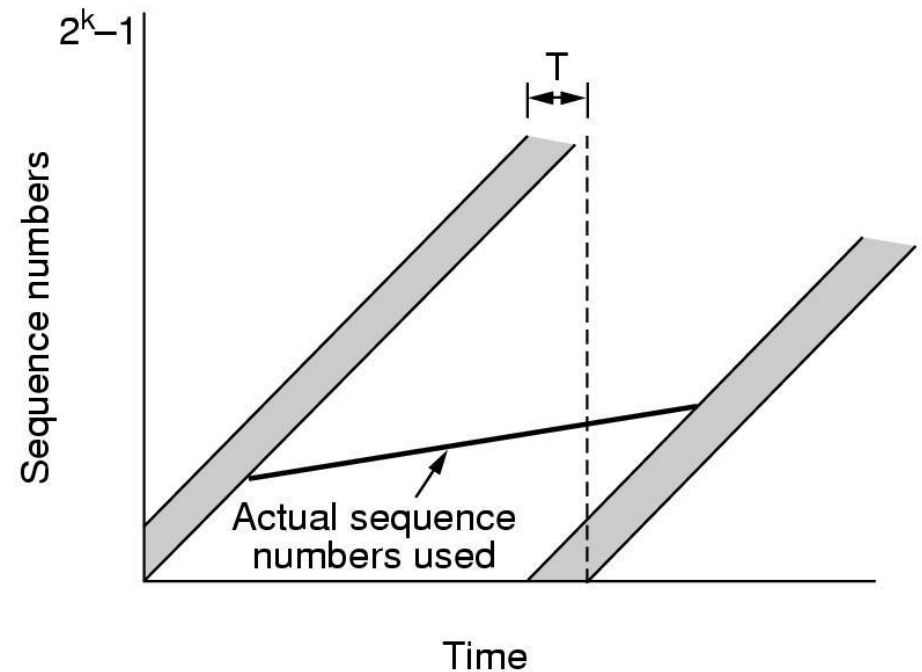
- Tomlinson's method
 - Initial sequence number
 - = lower order bits of clock
 - Ensure initial sequence numbers are always OK
 - forbidden region
 - Wrap around
 - Idle
 - Resynchronize sequence numbers

etp: Establishing a connection

- Tomlinson - forbidden region



(a)

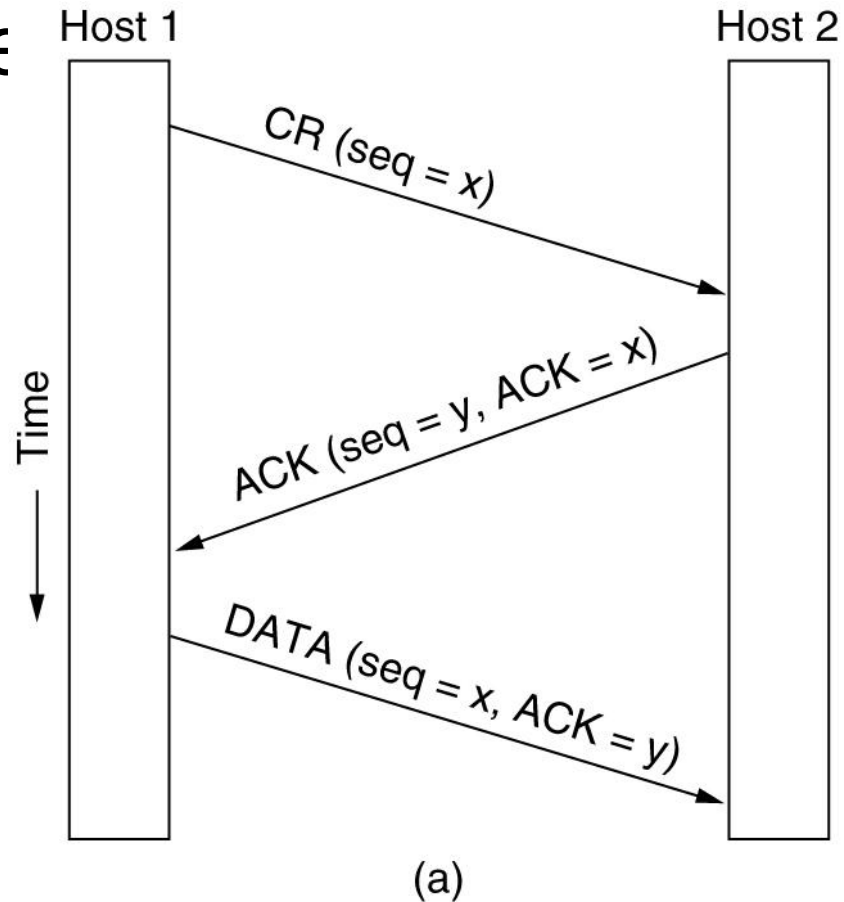


(b)

etp: Establishing a connection

- Tomlinson – $\text{thr} \in$

No combination of delayed packets can cause the protocol to fail

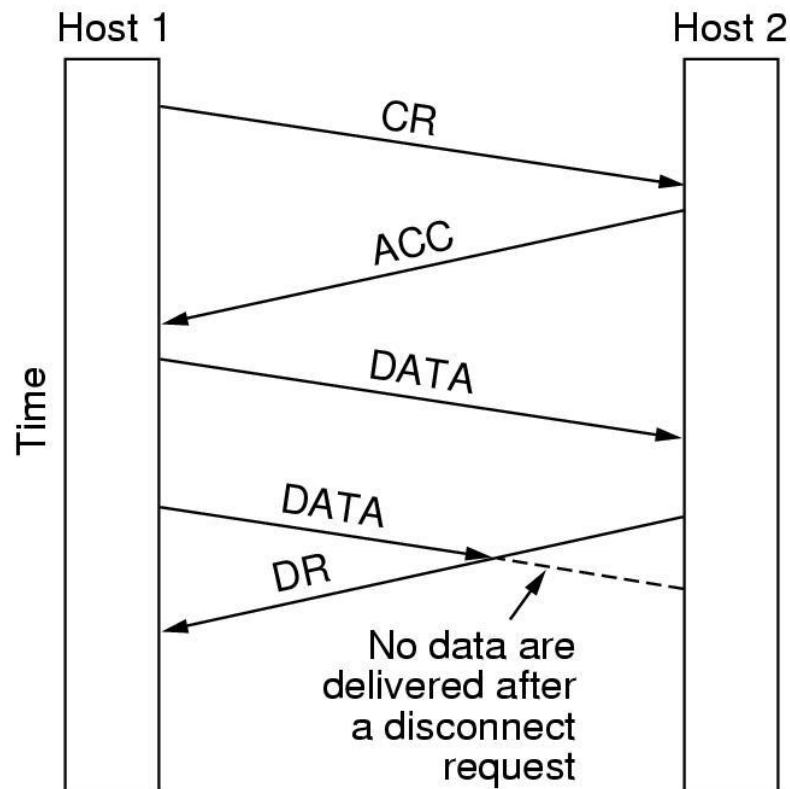


etp: Releasing a connection

- 2 styles:
 - Asymmetric
 - Connection broken when one party hangs up
 - Abrupt! → may result in **data loss**
 - Symmetric
 - Both parties should agree to release connection
 - How to reach agreement? **Two-army problem**
 - Solution: **three-way-handshake**
 - Pragmatic approach
 - Connection = 2 unidirectional connections
 - Sender can close unidirectional connection

etp: Releasing a connection

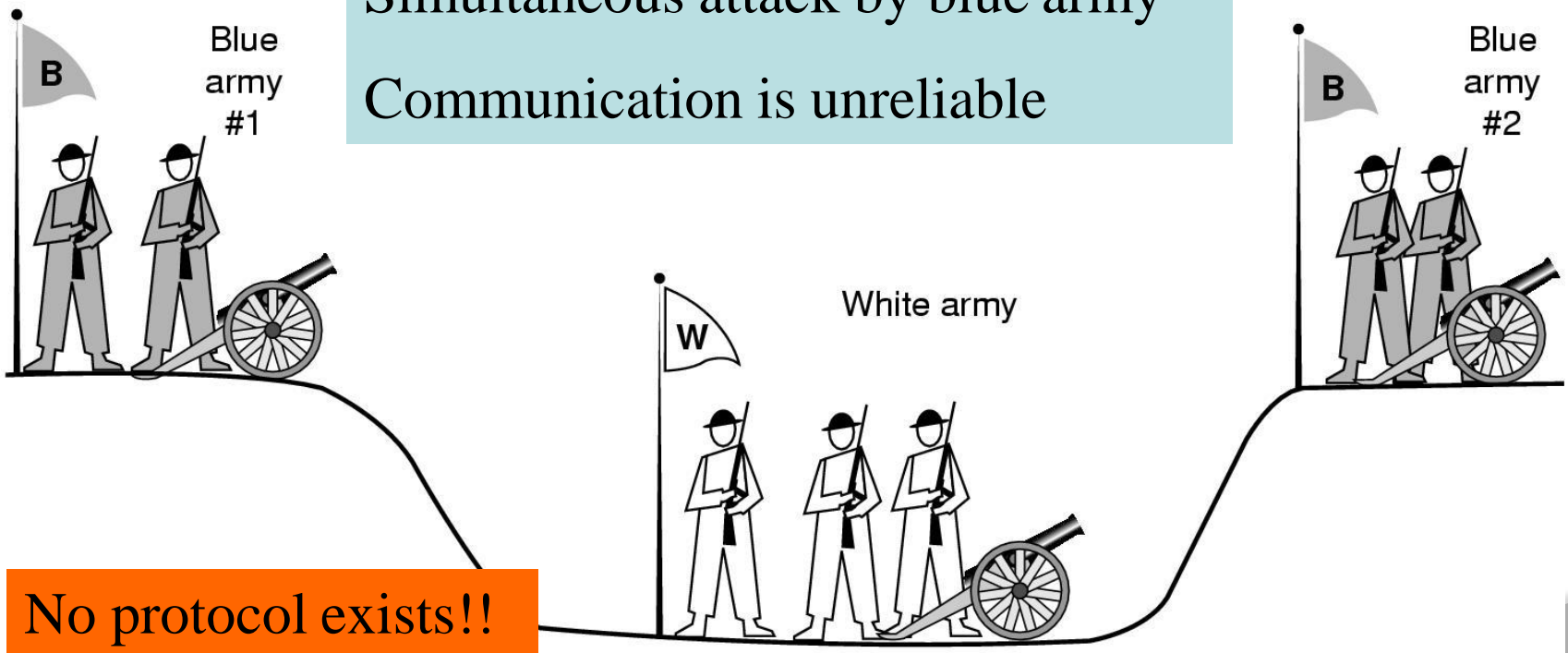
- Asymmetric: data loss



etp: Releasing a connection

- Symmetric: two-army-problem

Simultaneous attack by blue army
Communication is unreliable

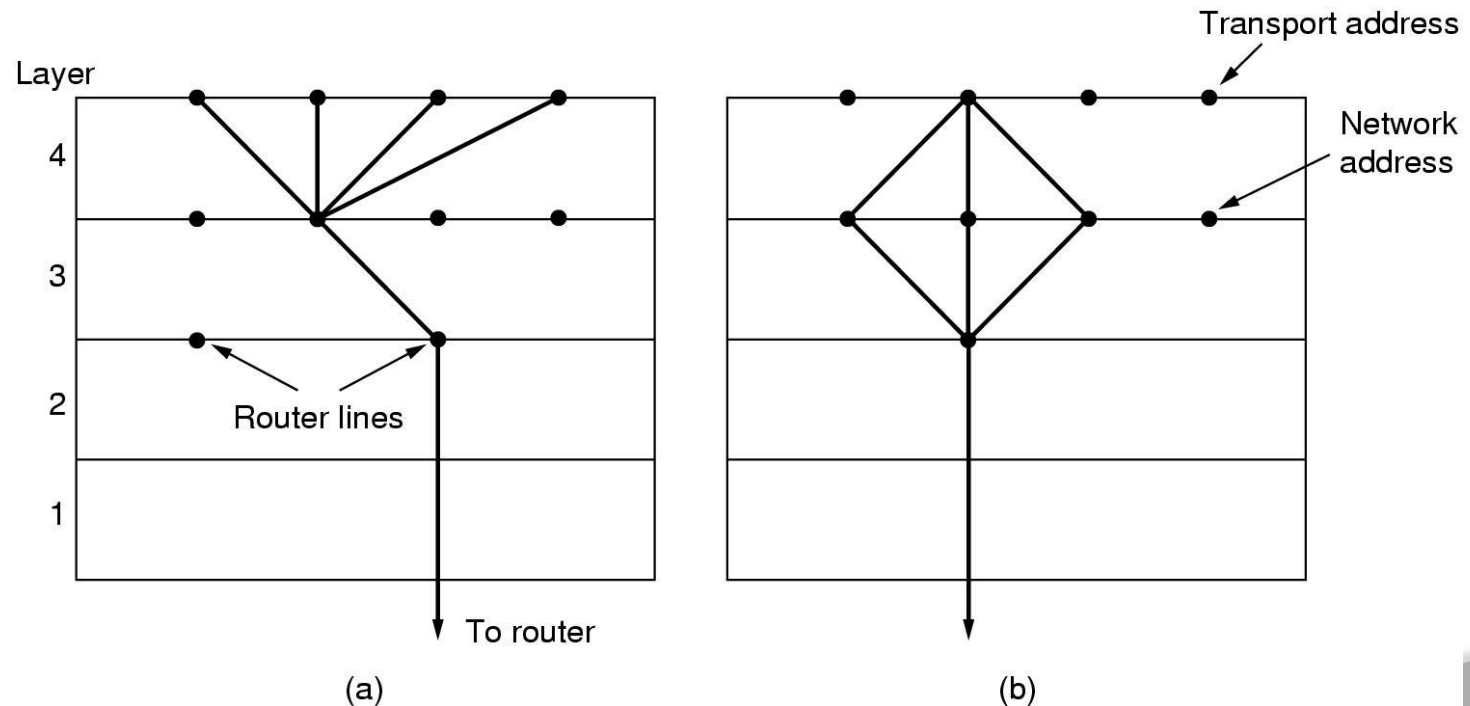


etp: Releasing a connection

- Three-way-handshake + timers
 - Send disconnection request
 - + start timer **RS** to resend (at most N times) the disconnection request
 - Ack disconnection request
 - + start timer **RC** to release connection

etp: Multiplexing

- Upward: reduce number of network connections to reduce cost
- Downward: increase bandwidth to avoid per connection limits



etp: Crash recovery

- recovery from network, router crashes?
 - No problem
 - Datagram network: loss of packet is always handled
 - Connection-oriented network: establish new connection + use state to continue service
- recovery from host crash?
 - server crashes, restarts: implications for client?
 - assumptions:

Recovery from a layer N crash can only

- be done by layer N+1 and only if the higher layer retains enough status information.

etp: Crash recovery

- Illustration of problem: File transfer:
 - Sender: 1 bit window protocol: states S0, S1
 - packet with seq number 0 transmitted; wait for ack
 - Receiver: actions
 - Ack packet
 - Write data to disk
 - Order?

Transport Layer

- Services
- Elements of transport protocol
- Simple transport protocol
- UDP
- Remote Procedure Call (see Distributed Systems)
- TCP

Simple transport protocol

- Service primitives:
 - connum = **LISTEN** (local)
 - Caller is willing to accept connection
 - Blocked till request received
 - connum = **CONNECT** (local, remote)
 - Tries to establish connection
 - Returns identifier (nonnegative number)
 - status = **SEND** (connum, buffer, bytes)
 - Transmits a buffer
 - Errors returned in status
 - status = **RECEIVE** (connum, buffer, bytes)
 - Indicates caller's desire to get data
 - status = **DISCONNECT** (connum)
 - Terminates connection

Example Transport Entity

```
(
#define MAX_CONN 32                /* max number of simultaneous connections */
#define MAX_MSG_SIZE 8192          /* largest message in bytes */
#define MAX_PKT_SIZE 512          /* largest packet in bytes */
#define TIMEOUT 20
#define CRED 1
#define OK 0

#define ERR_FULL -1
#define ERR_REJECT -2
#define ERR_CLOSED -3
#define LOW_ERR -3

typedef int transport_address;
typedef enum {CALL_REQ,CALL_ACC,CLEAR_REQ,CLEAR_CONF,DATA_PKT,CREDIT} pkt_type;
typedef enum {IDLE,WAITING,QUEUED,ESTABLISHED,SENDING,RECEIVING,DISCONN} cstate;

/* Global variables. */
transport_address listen_address; /* local address being listened to */
int listen_conn;                 /* connection identifier for listen */
unsigned char data[MAX_PKT_SIZE]; /* scratch area for packet data */

struct conn {
    transport_address local_address, remote_address;
    cstate state; /* state of this connection */
    unsigned char *user_buf_addr; /* pointer to receive buffer */
    int byte_count; /* send/receive count */
    int clr_req_received; /* set when CLEAR_REQ packet received */
    int timer; /* used to time out CALL_REQ packets */
    int credits; /* number of messages that may be sent */
} conn[MAX_CONN + 1]; /* slot 0 is not used */
```


Example Transport Entity (2)

```
void sleep(void);                /* prototypes */
void wakeup(void);
void to_net(int cid, int q, int m, pkt_type pt, unsigned char *p, int bytes);
void from_net(int *cid, int *q, int *m, pkt_type *pt, unsigned char *p, int *bytes);

int listen(transport_address t)
{ /* User wants to listen for a connection. See if CALL_REQ has already arrived. */
  int i, found = 0;

  for (i = 1; i <= MAX_CONN; i++)          /* search the table for CALL_REQ */
    if (conn[i].state == QUEUED && conn[i].local_address == t) {
      found = i;
      break;
    }

  if (found == 0) {
    /* No CALL_REQ is waiting. Go to sleep until arrival or timeout. */
    listen_address = t; sleep(); i = listen_conn ;
  }
  conn[i].state = ESTABLISHED;              /* connection is ESTABLISHED */
  conn[i].timer = 0;                        /* timer is not used */
}
```

Example Transport Entity (3)

```
listen_conn = 0;                                /* 0 is assumed to be an invalid address */
to_net(i, 0, 0, CALL_ACC, data, 0);              /* tell net to accept connection */
return(i);                                       /* return connection identifier */
}

int connect(transport_address l, transport_address r)
{ /* User wants to connect to a remote process; send CALL_REQ packet. */
  int i;
  struct conn *cptr;

  data[0] = r; data[1] = l;                      /* CALL_REQ packet needs these */
  i = MAX_CONN;                                  /* search table backward */
  while (conn[i].state != IDLE && i > 1) i = i - 1;
  if (conn[i].state == IDLE) {
    /* Make a table entry that CALL_REQ has been sent. */
    cptr = &conn[i];
    cptr->local_address = l; cptr->remote_address = r;
    cptr->state = WAITING; cptr->clr_req_received = 0;
    cptr->credits = 0; cptr->timer = 0;
    to_net(i, 0, 0, CALL_REQ, data, 2);
    sleep();                                     /* wait for CALL_ACC or CLEAR_REQ */
    if (cptr->state == ESTABLISHED) return(i);
    if (cptr->clr_req_received) {
      /* Other side refused call. */
      cptr->state = IDLE;                         /* back to IDLE state */
      to_net(i, 0, 0, CLEAR_CONF, data, 0);
      return(ERR_REJECT);
    }
  }
  } else return(ERR_FULL);                       /* reject CONNECT: no table space */
}
```

Example Transport Entity (4)

```
int send(int cid, unsigned char bufptr[], int bytes)
{ /* User wants to send a message. */
  int i, count, m;
  struct conn *cptr = &conn[cid];

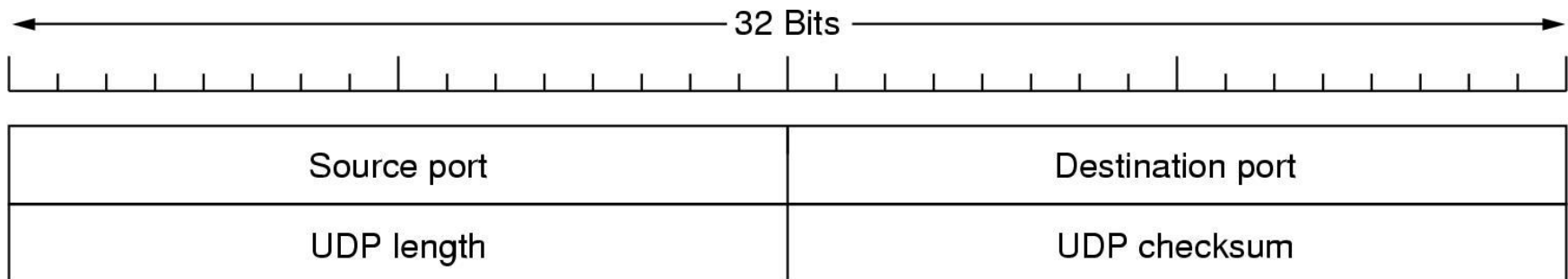
  /* Enter SENDING state. */
  cptr->state = SENDING;
  cptr->byte_count = 0; /* # bytes sent so far this message */
  if (cptr->clr_req_received == 0 && cptr->credits == 0) sleep();
  if (cptr->clr_req_received == 0) {
    /* Credit available; split message into packets if need be. */
    do {
      if (bytes - cptr->byte_count > MAX_PKT_SIZE) { /* multipacket message */
        count = MAX_PKT_SIZE; m = 1; /* more packets later */
      } else { /* single packet message */
        count = bytes - cptr->byte_count; m = 0; /* last pkt of this message */
      }
      for (i = 0; i < count; i++) data[i] = bufptr[cptr->byte_count + i];
      to_net(cid, 0, m, DATA_PKT, data, count); /* send 1 packet */
      cptr->byte_count = cptr->byte_count + count; /* increment bytes sent so far */
    } while (cptr->byte_count < bytes); /* loop until whole message sent */
  }
}
```

Transport Layer

- Services
- Elements of transport protocol
- Simple transport protocol
- **UDP**
- Remote Procedure Call (see Distributed Systems)
- TCP

UDP

- User Data Protocol
 - Datagram service between processes
 - No connection overhead
 - UDP header:
 - Ports – identification of end points



UDP

- Some characteristics
 - Supports broadcasting, multicasting (not in TCP)
 - Packet oriented (TCP gives byte stream)
 - Simple protocol

- Why needed above IP?

Transport Layer

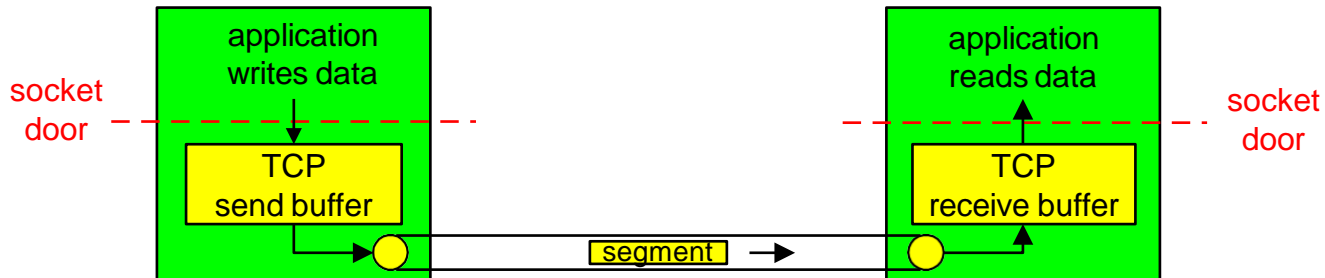
- Services
- Elements of transport protocol
- Simple transport protocol
- UDP
- Remote Procedure Call (see Distributed Systems)
- TCP

TCP service model

- point-to-point
 - one sender, one receiver
- reliable, in-order byte stream
 - no *message/packet boundaries*
- pipelined & flow controlled
 - window size set by TCP congestion and flow control algorithms
- connection-oriented
 - handshaking to get at initial state
- full duplex data
 - bi-directional data flow in same connection

TCP service model

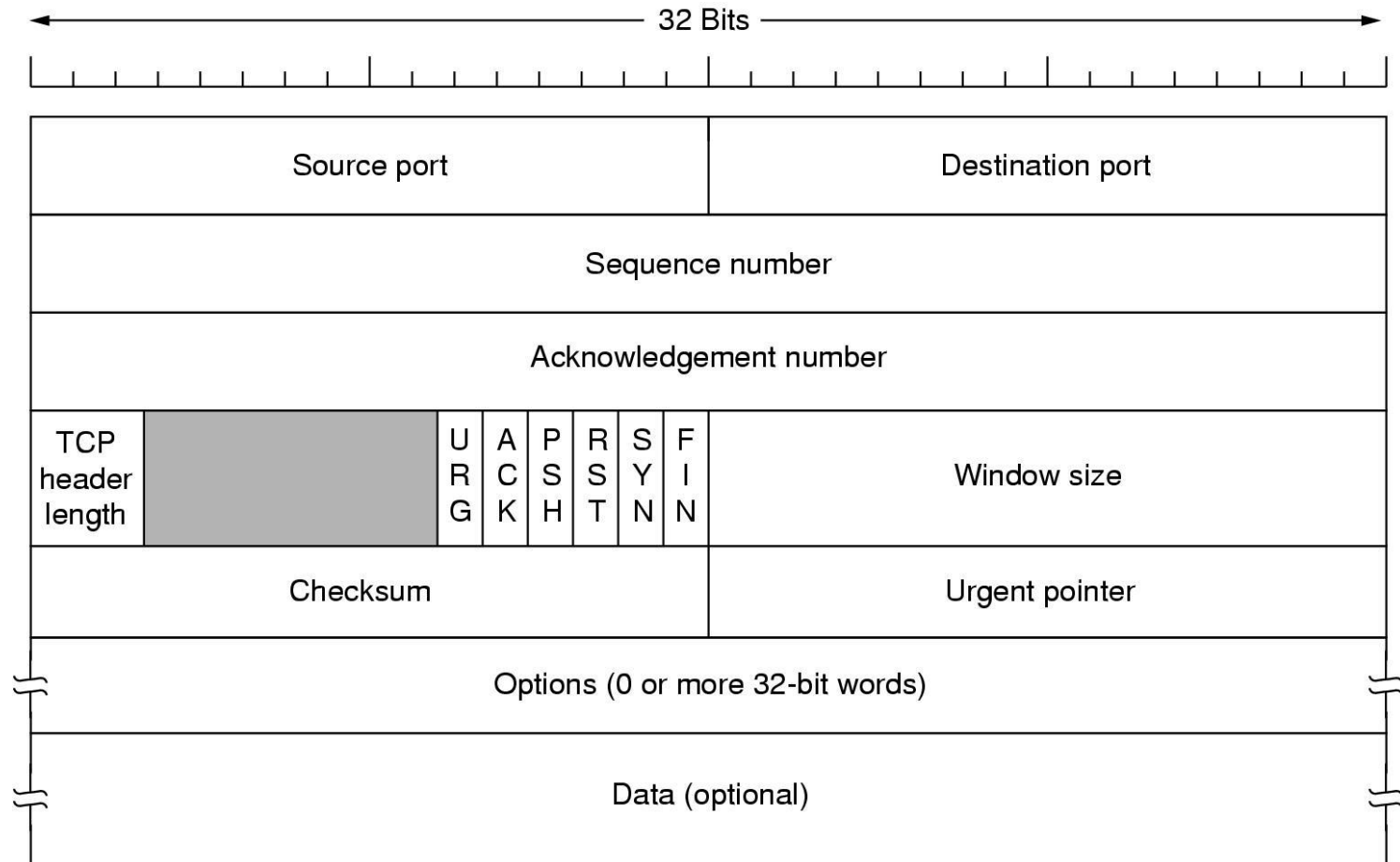
- ...
- send & receive buffers



TCP protocol

- Three-way handshake to set up connections
- Every byte has its own 32-bit sequence number
 - Wrap around
 - 32-bit Acks; window size in bytes
- Segment = unit of data exchange
 - 20-byte header + options + data
 - Limits for size
 - 64Kbyte
 - MTU, agreed upon for each direction
 - Data from consecutive writes may be accumulated in a single segment
 - Fragmentation possible
- Sliding window protocol

TCP header



TCP header

- **source & destination ports** (16 bit)
- **sequence number** (32 bit)
- **Acknowledgement number** (32 bit)
- **Header length** (4 bits) in 32-bit words
- **6 flags** (1 bit)
- **window size** (16 bit): number of bytes the sender is allowed to send starting at byte acknowledged
- **checksum** (16 bit)
- **urgent pointer** (16 bit) : byte position of urgent data

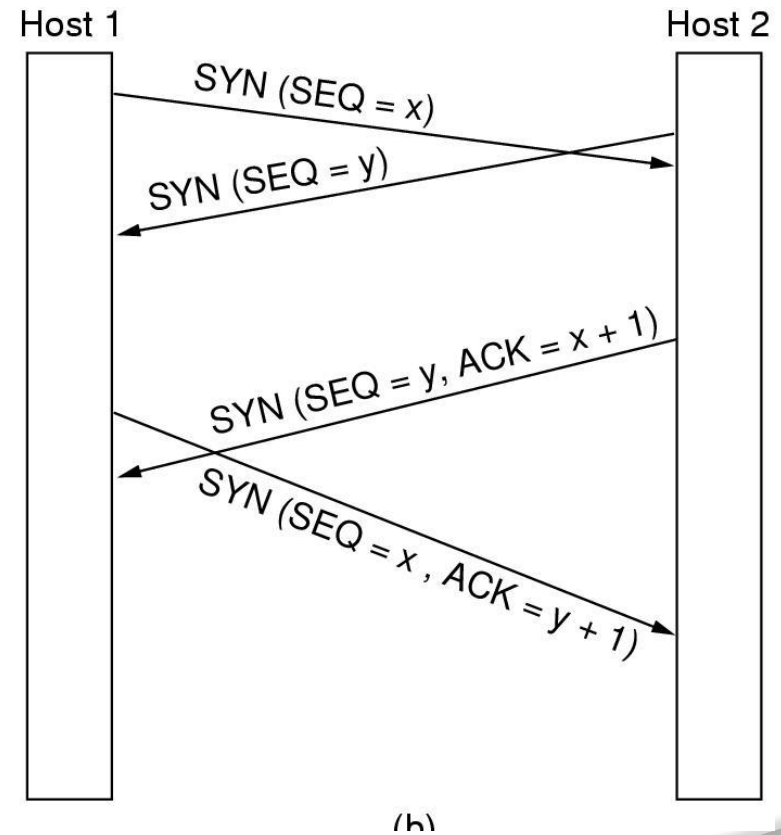
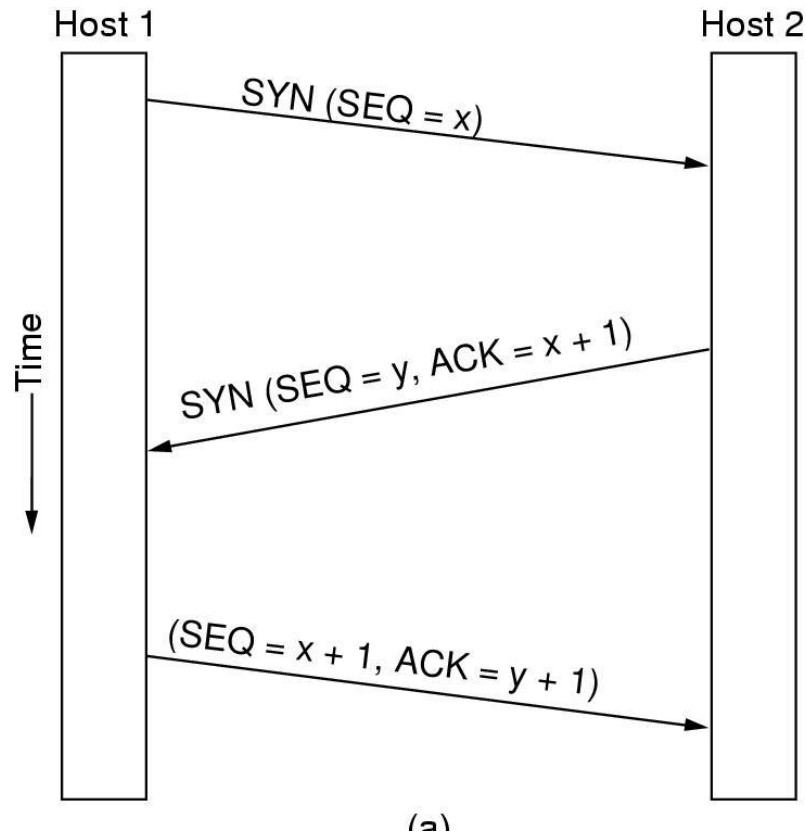
TCP header

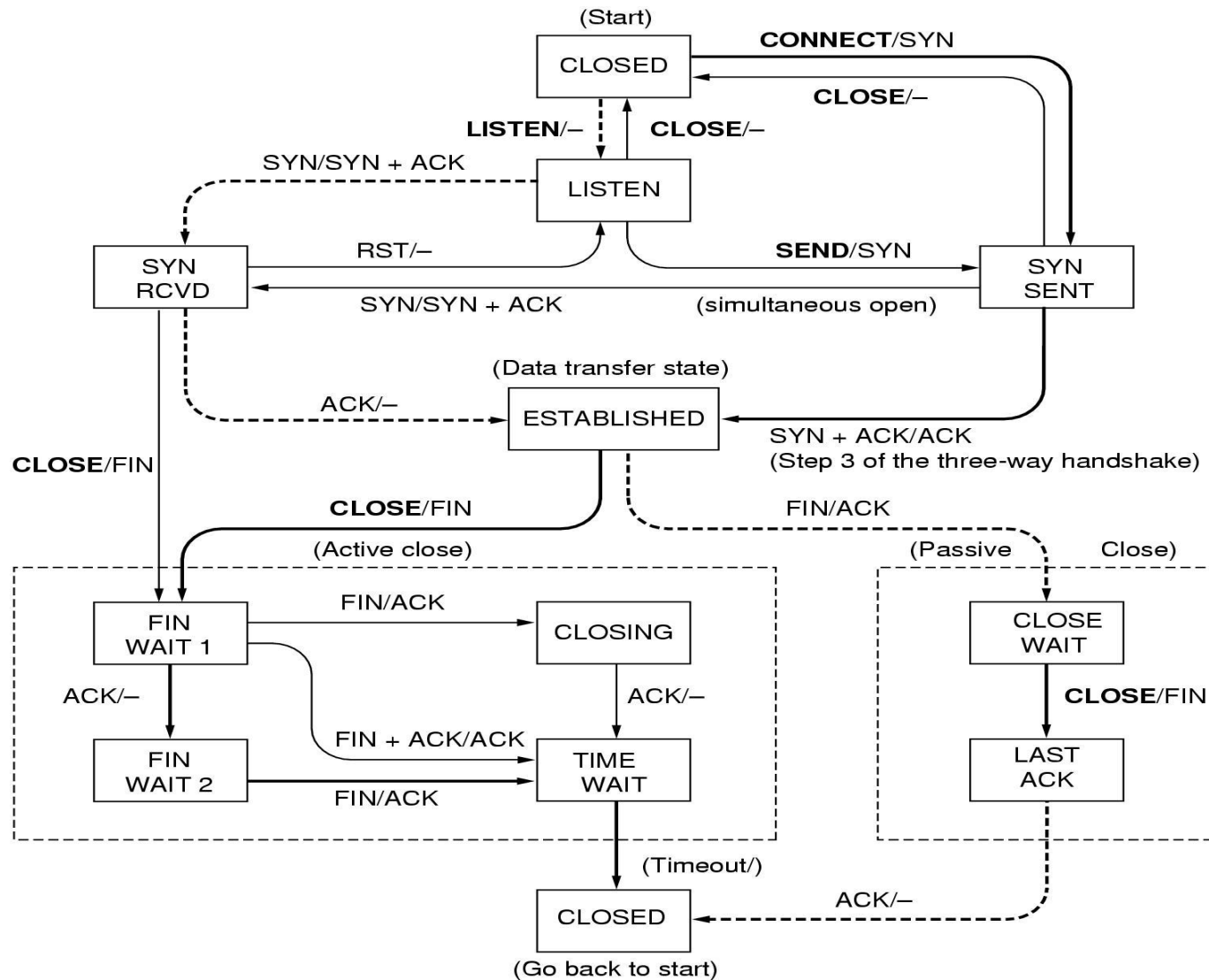
- Flags:
 - URG: urgent pointer in use
 - ACK: valid Acknowledgement number
 - PSH: receiver should deliver data without delay to user
 - RST: reset connection
 - SYN: used when establishing connections
 - FIN: used to release connection
- Options:
 - Maximum payload a host is willing to receive
 - Scale factor window size
 - Use selective repeat instead of go back n

TCP connection management

- Three-way handshake
 - Initial sequence number: clock based
 - No reboot after crash for T (maximum packet lifetime=120 sec)
 - Wrap around?
- Connection identification
 - Pair of ports of end points
- Connection release
 - Both sides are closed separately
 - No response to FIN: release after $2 * T$
 - Both sides closed: wait for time $2 * T$

TCP connection management





TCP connection

State	Description
management	
Closed	No connection is active or pending
Listen	The server is waiting for an incoming call
SYN rcvd	A connection request has arrived; wait for ACK
SYN sent	The application has started to open a connection
Established	The normal data transfer state
FIN wait 1	The application has said it is finished
FIN wait 2	The other side has agreed to release
Timed wait	Wait for all packets to die off
Closing	Both sides have tried to close simultaneously
Close wait	The other side has initiated a release
Last Ack	Wait for all packets to die off

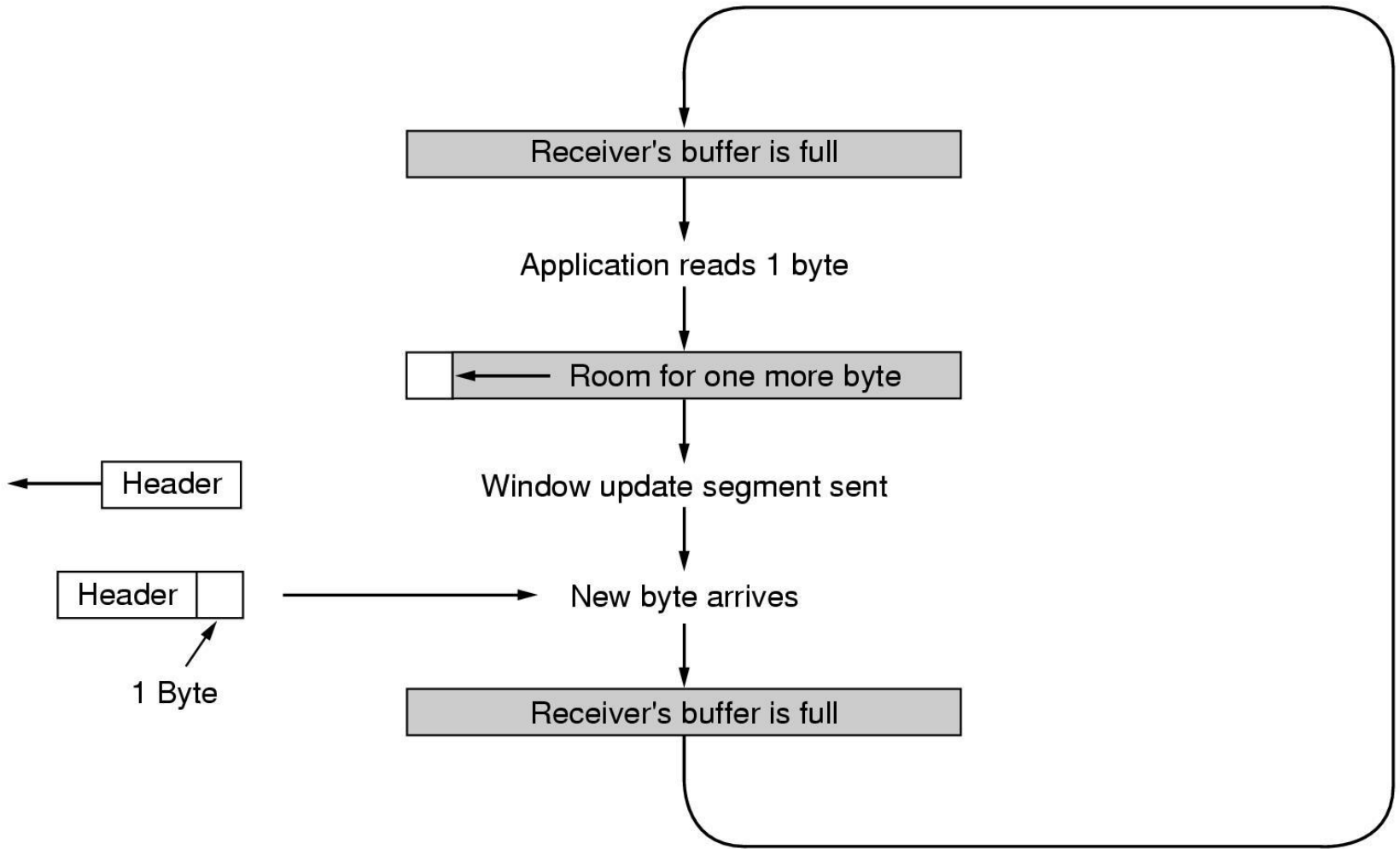
TCP transmission policy

- **Telnet scenario:** interactive editor reacting on each keystroke:
One character typed
 - ➔ 21 byte segment or 41 byte IP packet
 - ⬅ (packet received) 20 byte segment with Ack
 - ⬅ (editor has read byte) 20 byte segment with window update
 - ⬅ (editor has processed byte; sends echo) 21 byte segment
 - ➔ (client gets echo) 20 byte segment with Ack
- **Solutions:**
 - delay acks + window updates for 500 msec
 - Nagle's algorithm:
 - Receive one byte from user; send it in segment
 - Buffer all other chars till Ack for first char arrives
 - Send other chars in a single segment
 - **Disable** algorithm for X-windows applications (do not delay sending of mouse movements)

TCP transmission policy

- Silly window syndrome
 - Problem:
 - Sender transmits data in large blocks
 - Receiver reads data 1 byte at a time
 - Scenario: next slide
 - Solution:
 - Do not send window update for 1 byte
 - Wait for window update till
 - Receiver can accept MTU
 - Buffer is half empty

TCP transmission policy



TCP transmission policy

- General approach:
 - Sender should not send small segments
 - Nagle: buffer data in TCP send buffer
 - Receiver should not ask for small segments
 - Silly window: do window updates in large units

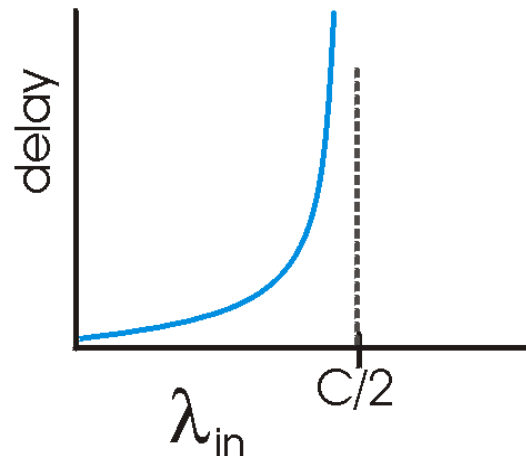
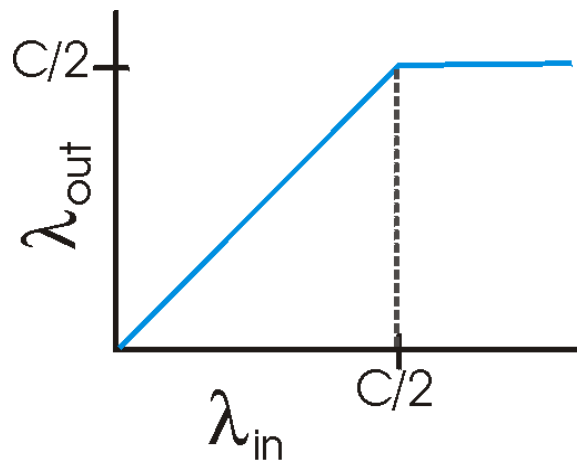
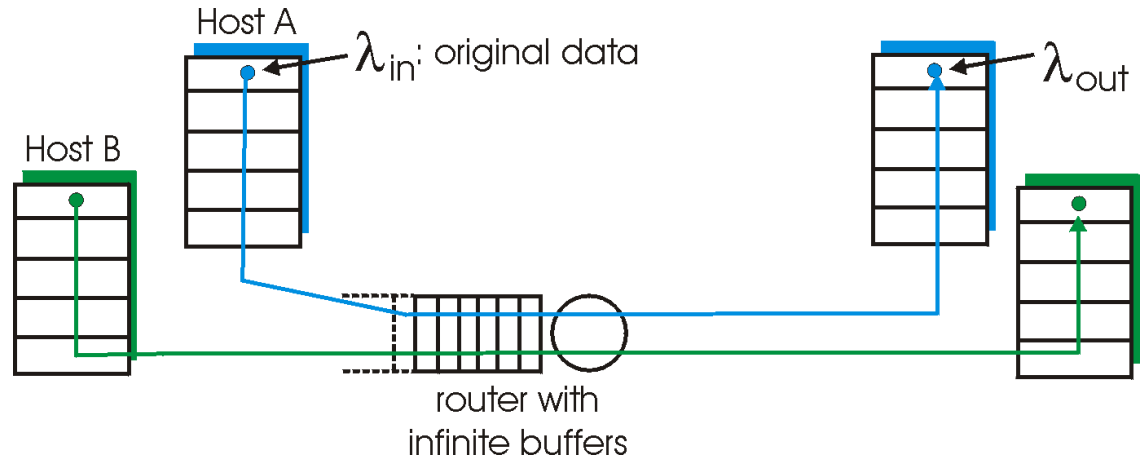
Principles of Congestion Control

Congestion:

- informally: “too many sources sending too much data too fast for *network* to handle”
- different from flow control!
 - = end-to-end issue!
- manifestations:
 - lost packets (buffer overflow at routers)
 - long delays (queue-ing in router buffers)
- a top-10 problem!

Causes/costs of congestion: scenario

- two senders, two receivers
- one router, infinite buffers
- no retransmission



- large delays when congested
- maximum achievable throughput

Approaches towards congestion control

Two broad approaches towards congestion control:

end-to-end congestion control:

- no explicit feedback from network
- congestion inferred from end-system observed loss, delay
- approach taken by TCP

Network-assisted congestion control:

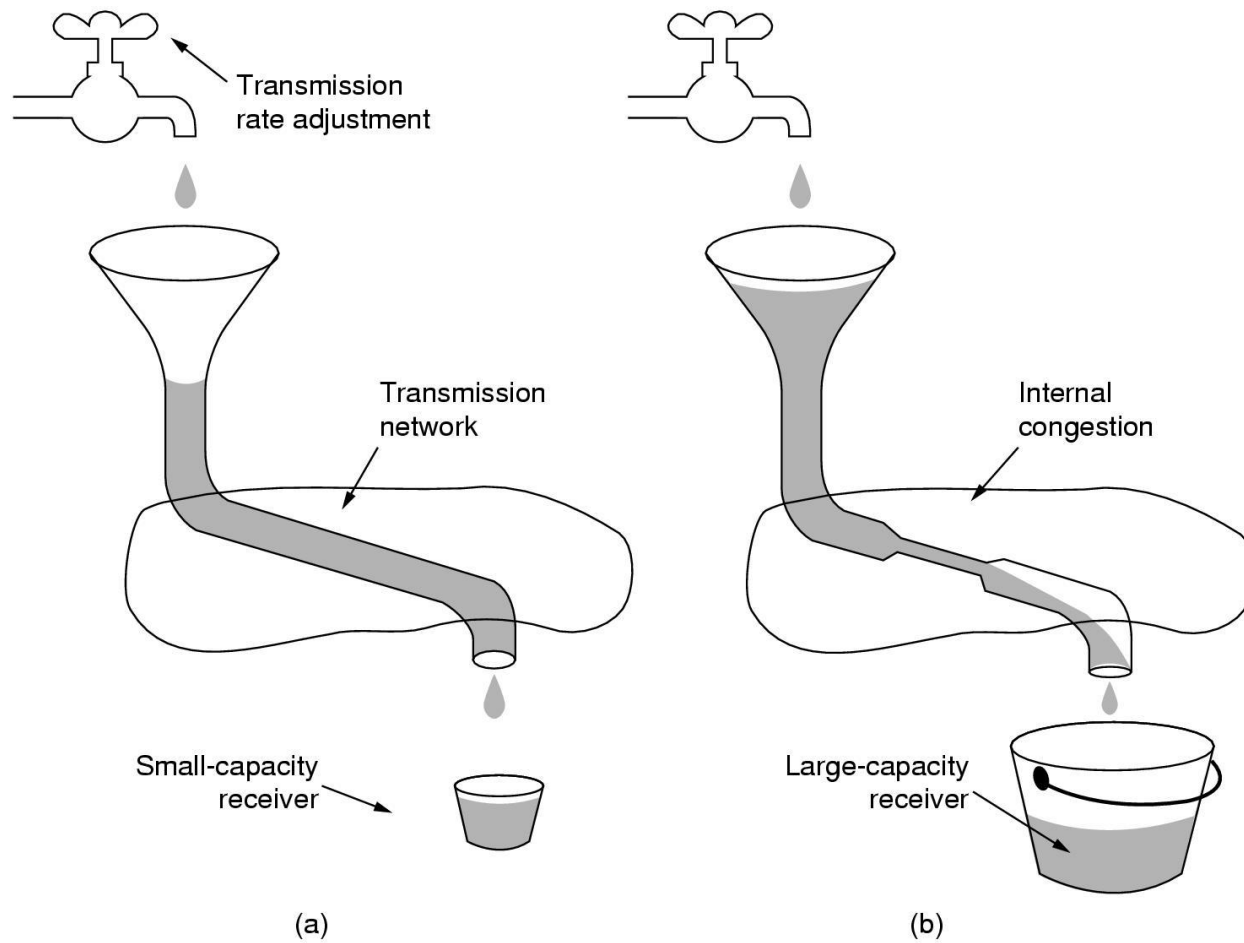
- routers provide feedback to end systems
 - single bit indicating congestion (SNA, ATM)
 - explicit rate sender should send at

TCP Congestion Control

- How to detect congestion?
- Timeout caused by packet loss reasons
 - Transmission errors
 - Packet discarded at congested router

Hydraulic example illustrating two limitations for sender!

TCP congestion control



TCP Congestion Control

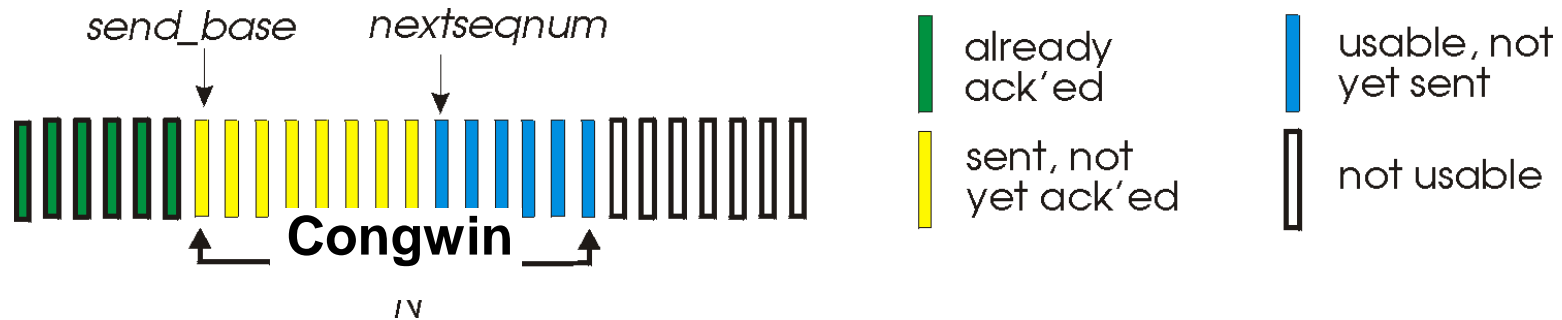
- How to detect congestion?
- Timeout caused by ^{Rare} packet loss: reasons
 - Transmission errors
 - Packets discarded at congested router ^{Packet loss → congestion}

Approach: 2 windows for sender
Receiver window

Minimum of { Congestion window

TCP Congestion Control

- end-end control (no network assistance)
- transmission rate limited by congestion window size, **Congwin**, over segments:



□ w segments, each with MSS bytes sent in one RTT:

$$\text{throughput} = \frac{w * \text{MSS}}{\text{RTT}} \text{ Bytes/sec}$$

TCP Congestion Control:

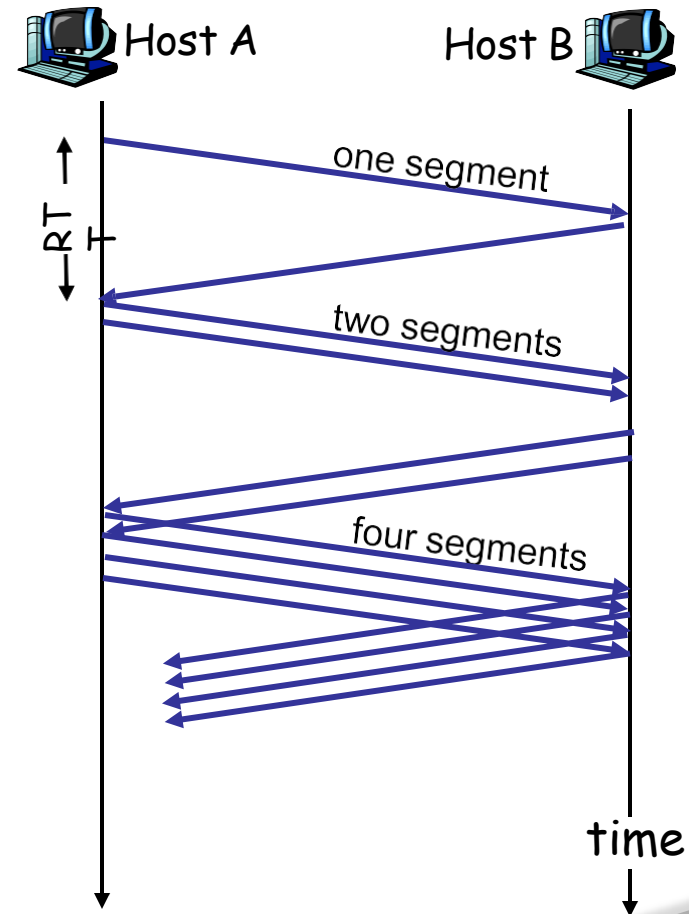
- “**probing**” for usable bandwidth:
 - **ideally**: transmit as fast as possible (**Congwin** as large as possible) without loss
 - *increase* **Congwin** until loss (congestion)
 - loss: *decrease* **Congwin**, then begin probing (increasing) again
- two “phases”
 - **slow start**
 - **congestion avoidance**
- important variables:
 - **Congwin**
 - **threshold**: defines threshold between two phases:
 - slow start phase
 - congestion control phase

TCP Slow start

Slow start algorithm

initialize: Congwin = 1
for (each segment ACKed)
 Congwin++
until (loss event OR
 CongWin > threshold)

- exponential increase (per RTT) in window size (not so slow!)
- loss event: timeout (Tahoe TCP) and/or three duplicate ACKs (Reno TCP)

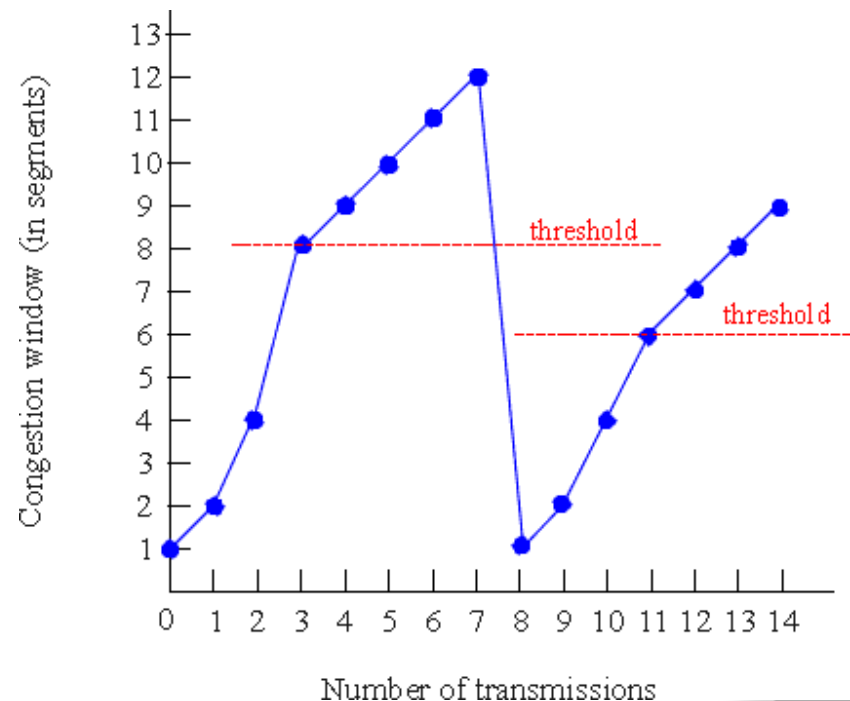


TCP Congestion Avoidance

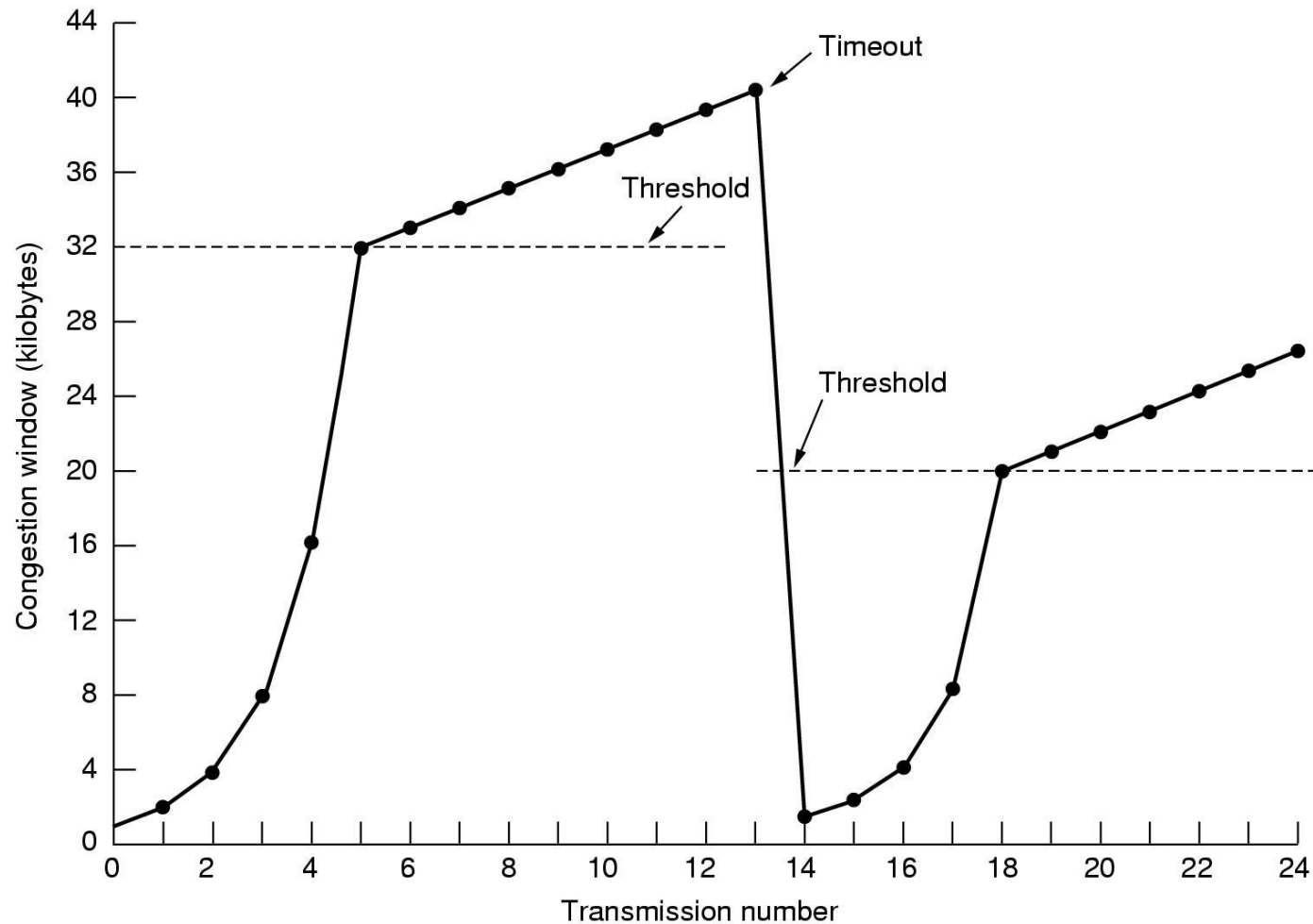
Congestion avoidance

```
/* slowstart is over */
/* Congwin > threshold */
Until (loss event) {
    every w segments ACKed:
        Congwin++
}
threshold = Congwin/2
Congwin = 1
perform slowstart1
```

1: TCP Reno skips slowstart (fast recovery) after three duplicate ACKs



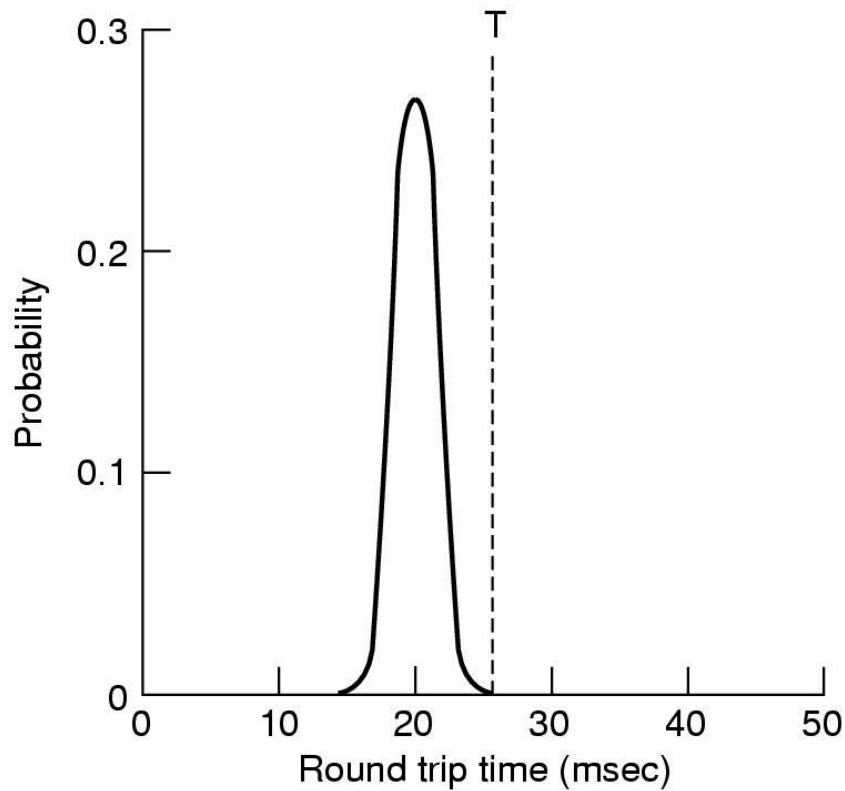
TCP congestion control



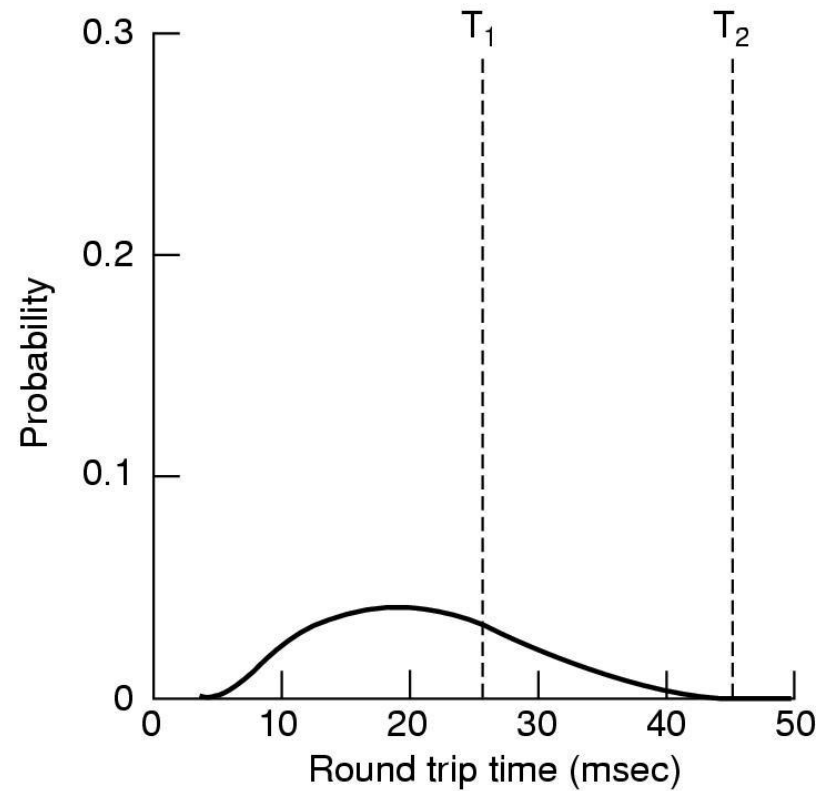
TCP timer management

- How long should the timeout interval be?
 - Data link: expected delay predictable
 - Transport: different environment; impact of
 - Host
 - Network (routers, lines)unpredictable
- Consequences
 - Too small: unnecessary retransmissions
 - Too large: poor performance
- Solution: **adjust timeout interval** based on continuous measurements of network performance

TCP timer management



Data link layer



Transport layer

TCP timer management

- Algorithm of Jacobson $\text{Timeout} = \text{RTT} + 4 * D$
 - RTT = best current estimate of the round-trip time
 - D = mean deviation (cheap estimator of the standard variance)
 - 4?
 - Less than 1% of all packets come in more than 4 standard deviations late
 - Easy to compute

TCP timer management

- Algorithm of Jacobson: $\text{Timeout} = \text{RTT} + 4 * D$
 - $\text{RTT} = \alpha \text{RTT} + (1 - \alpha) M$ $\alpha = 7/8$
M = last measurement of round-trip time
 - $D = \alpha D + (1 - \alpha) / \text{RTT} - M$
- Karn's algorithm: how handle retransmitted segments?
 - Do not update RTT for retransmitted segments
 - Double timeout

TCP timer management

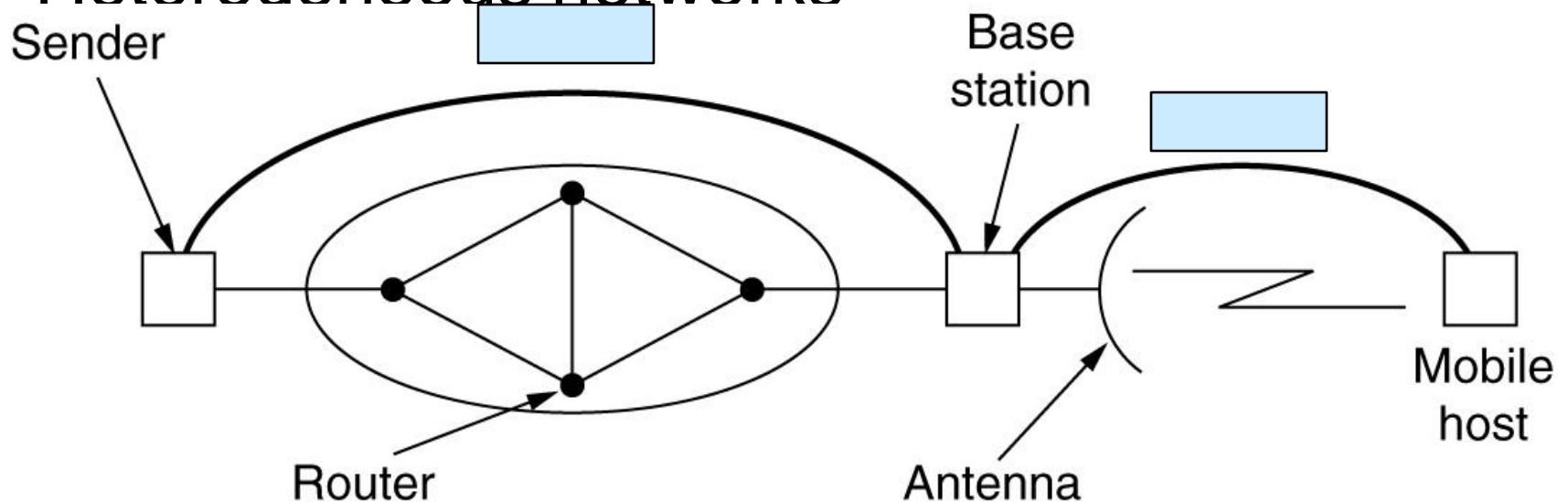
- Other timers:
 - Persistence timer
 - Problem: lost window update packet when window is 0
 - Sender transmits probe; receiver replies with window size
 - Keep alive timer
 - Check whether other side is still alive if connection is idle for a long time
 - No response: close connection
 - Timed wait
 - Make sure all packets are died off when connection is closed
 - $= 2 T$

Wireless TCP & UDP

- Transport protocols
 - **Independent** of underlying network layer
 - **BUT**: carefully optimized for wired networks
 - Assumption:
 - Packet loss caused by congestion
 - Invalid for wireless networks: always loss of packets
 - Congestion algorithm:
 - Timeout (= congestion) → slowdown
 - Solution for congestion control:
 - Retransmit asap
- Wireless: Lower throughput – same loss → **NO** solution

Wireless TCP

- Heterogeneous networks



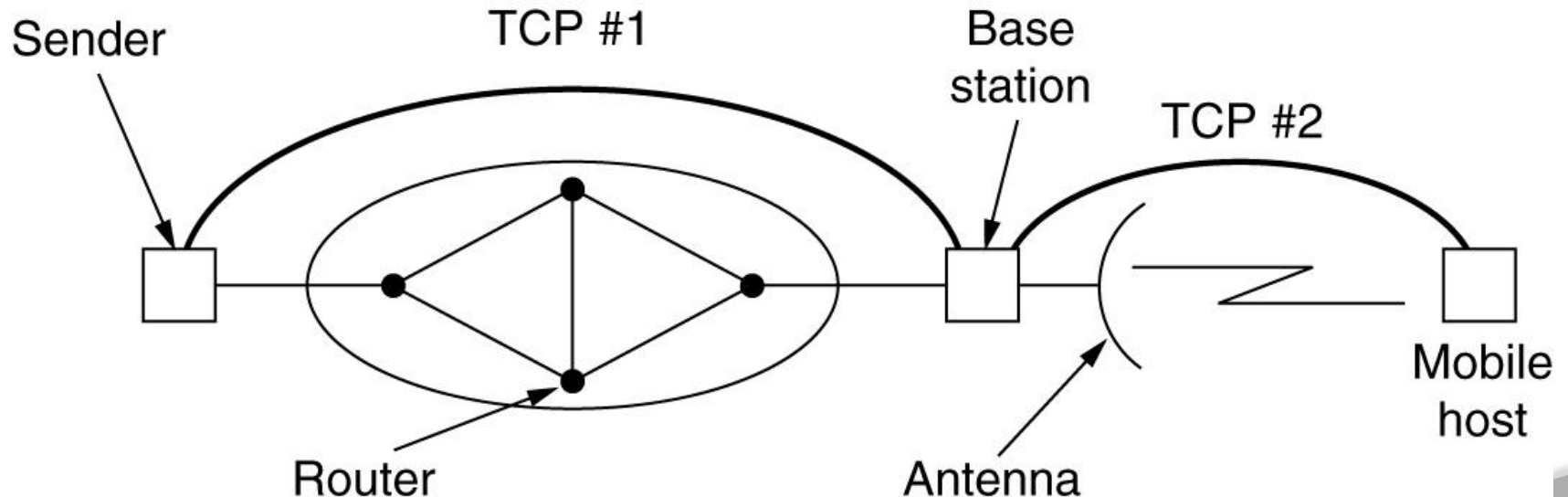
- Solutions?

- Retransmissions can cause congestion in wired network

Wireless TCP

- Solutions for heterogeneous networks
 - Indirect TCP

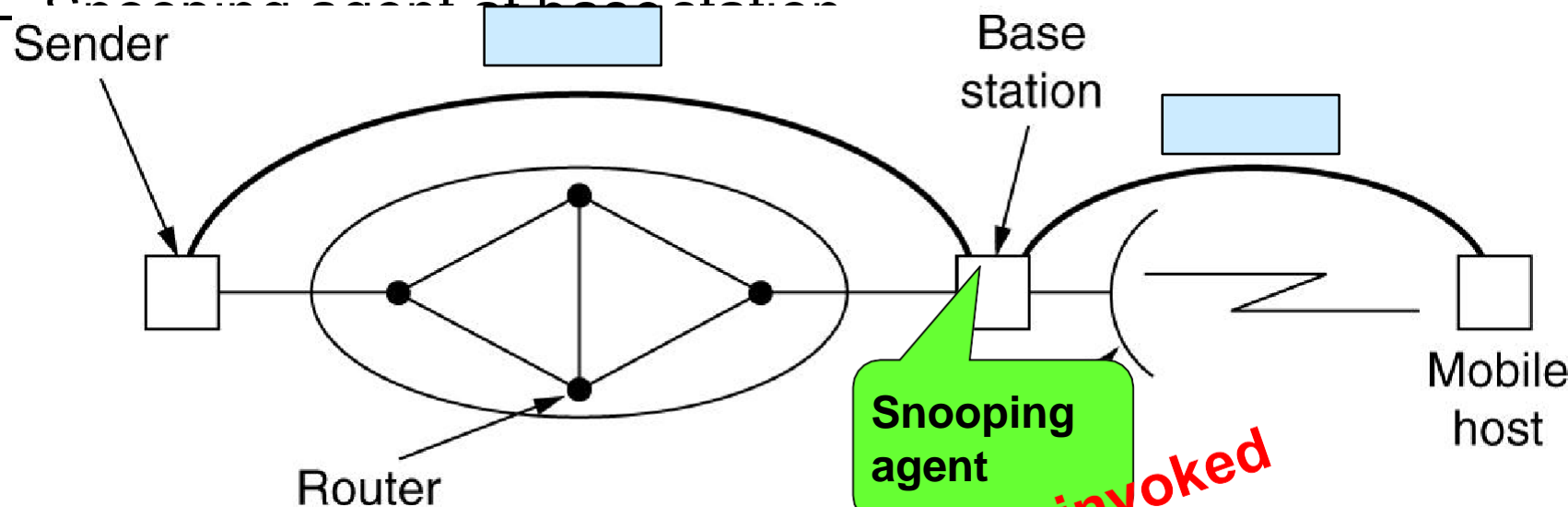
+ 2 homogeneous connections



Wireless TCP

- Solutions for heterogeneous networks

- Snooping agent at base station



- Retransmits segment if ack is missing
- Removes duplicate acks
- Generates selective repeat requests for segments originating at mobile host

Wireless UDP

- UDP = datagram service → loss permitted
no problems?
 - Programs using UDP expect it to be highly reliable
 - Wireless UDP: far from perfect!!!
- Implications for programs recovering from lost UDP messages

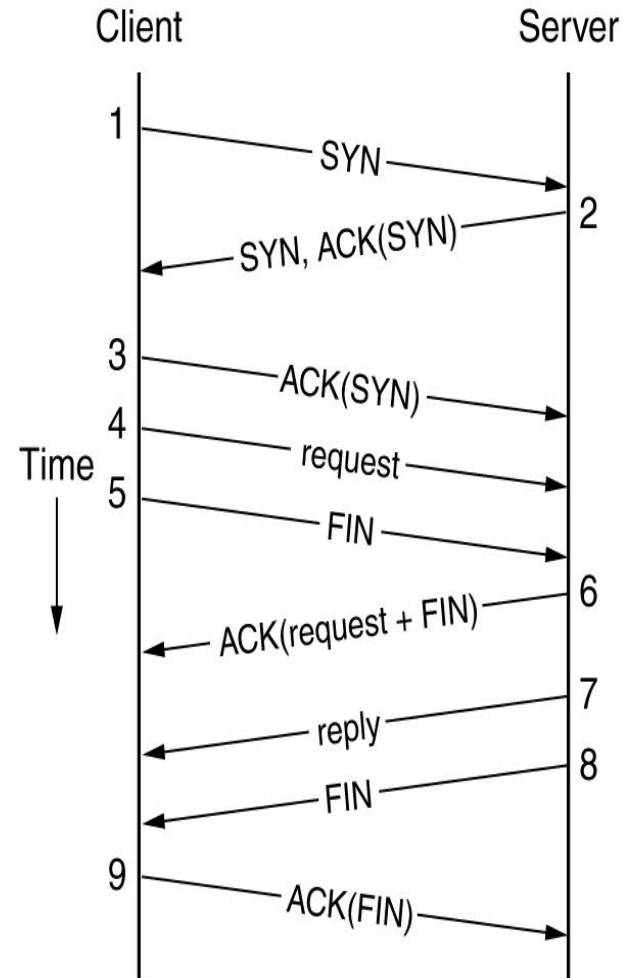
Transactional TCP

- How to implement RPC?
 - On top of UDP?
 - Yes if
 - Request and reply fit in a single packet
 - Operations are idempotent
 - Otherwise
 - Reimplementation of reliability
 - On top of TCP?

Transactional TCP

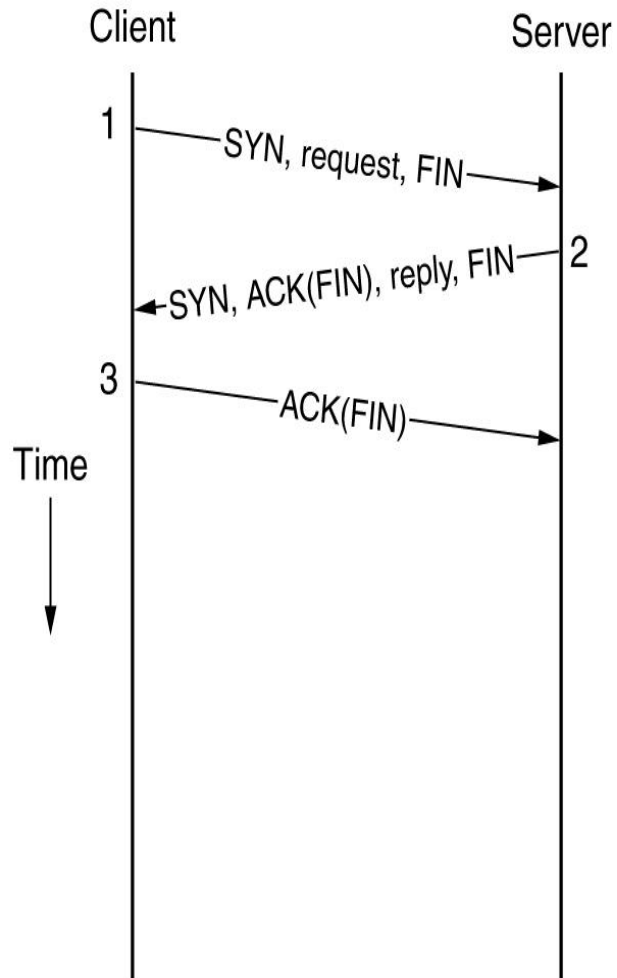
How to implement RPC?

- On top of UDP?
 - Yes if
 - Request and reply fit in a single packet
 - Operations are idempotent
 - Otherwise
 - Reimplementation of reliability
- On top of TCP?
 - Unattractive because of connection set up
- Solution: transactional TCP



(a)

Transactional TCP



(b)

How to implement RPC?

- On top of UDP?
 - Problems with reliability
- On top of TCP?
 - Overhead of connection set up
- Solution: transactional TCP
 - Allow data transfer during setup
 - Immediate close of stream

INTRODUCTION TO APPLICATION LAYER

Application layer

- 2.1 Principles of network applications
- 2.2 Web and HTTP
- 2.3 FTP
- 2.4 Electronic Mail
 - SMTP, POP3, IMAP
- 2.5 DNS
- 2.6 P2P applications
- 2.7 Socket programming with TCP
- 2.8 Socket programming with UDP

Processes communicating

Process:

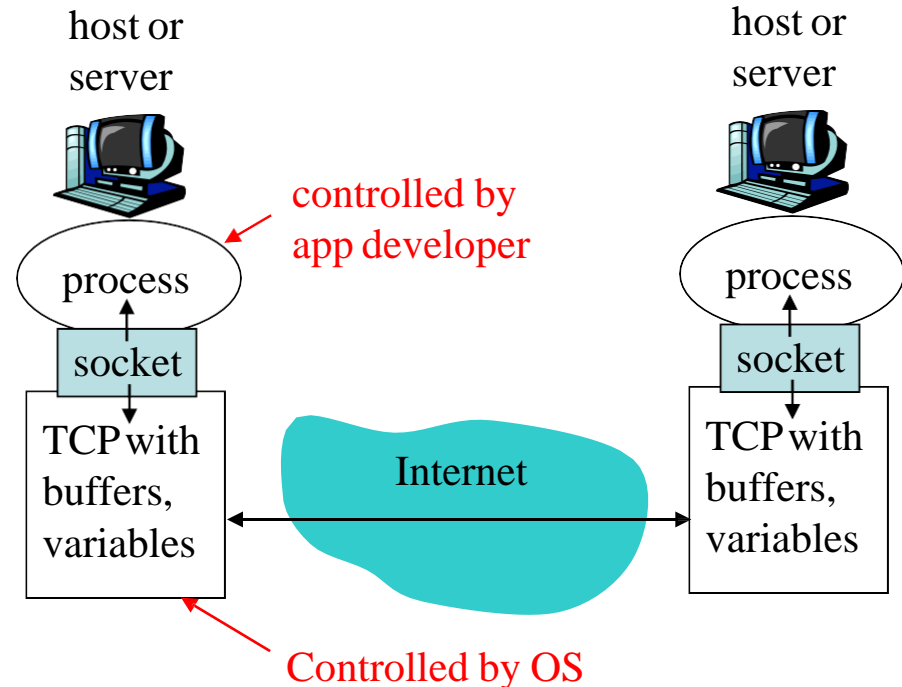
program running within a host

Client process:

initiates communication

Server process:

waits to be contacted



process sends/receives messages to/from its **socket**

identifier includes both **IP address** and **port numbers** associated with process on host.

App-layer protocol defines

- Types of messages exchanged,
 - e.g., request, response
- Message syntax:
 - what fields in messages & how fields are delineated
- Message semantics
 - meaning of information in fields
- Rules for when and how processes send & respond to messages

Public-domain protocols:

- ◆ defined in RFCs
- ◆ allows for interoperability
- ◆ e.g., HTTP, SMTP

Proprietary protocols:

- ◆ e.g., Skype

Transport service requirements of common apps

Application	Data loss	Throughput	Time Sensitive
file transfer	no loss	elastic	no
e-mail	no loss	elastic	no
Web documents	no loss	elastic	no
real-time audio/video	loss-tolerant	audio: 5kbps-1Mbps video: 10kbps-5Mbps	yes, 100's msec
stored audio/video	loss-tolerant	same as above	yes, few secs
interactive games	loss-tolerant	few kbps up	yes, 100's msec
instant messaging	no loss	elastic	yes and no

Internet transport protocols services

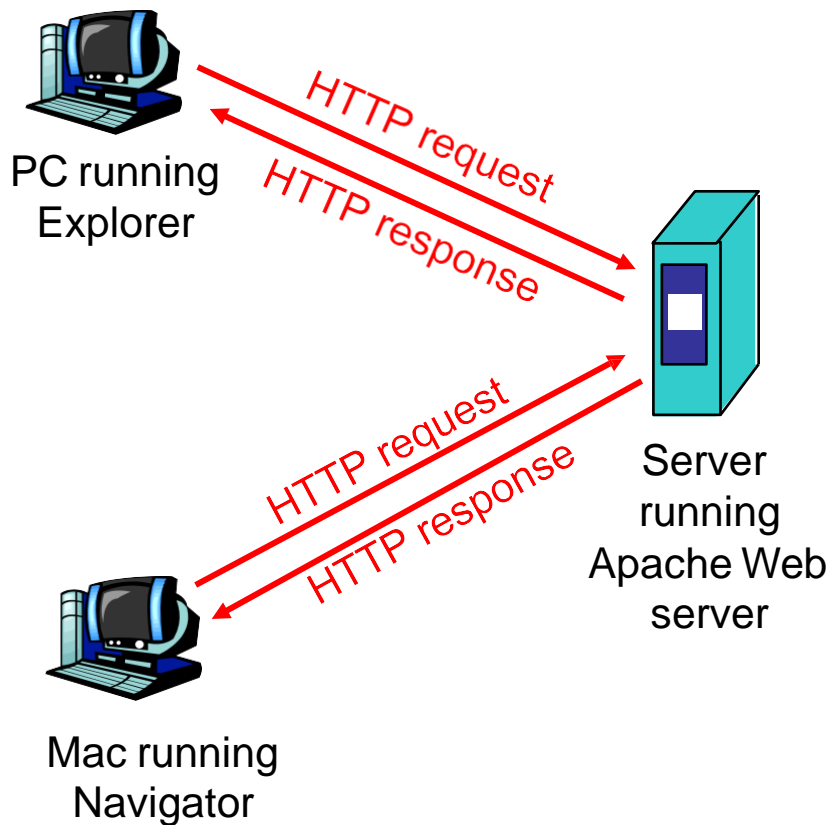
TCP service:

- *connection-oriented*: setup required between client and server processes
- *reliable transport* between sending and receiving process
- *flow control*: sender won't overwhelm receiver
- *congestion control*: throttle sender when network overloaded
- *does not provide*: timing, minimum throughput guarantees, security

UDP service:

- unreliable data transfer between sending and receiving process
- does not provide: connection setup, reliability, flow control, congestion control, timing, throughput guarantee, or security

HTTP overview



- ❑ Web page consists of **base HTML-file** which includes several referenced **objects**

- ❑ Each object is addressable by a **URL**

HTTP: hypertext transfer protocol

- ❑ Web's application layer protocol

- ❑ client/server model

- ❖ **client**: browser that requests, receives, "displays" Web objects

- ❖ **server**: Web server sends objects in response to requests

- ❑ uses TCP

- ❑ is "**stateless**"

HTTP connections

Nonpersistent HTTP

- At most one object is sent over a TCP connection.

Persistent HTTP

- Multiple objects can be sent over single TCP connection between client and server.

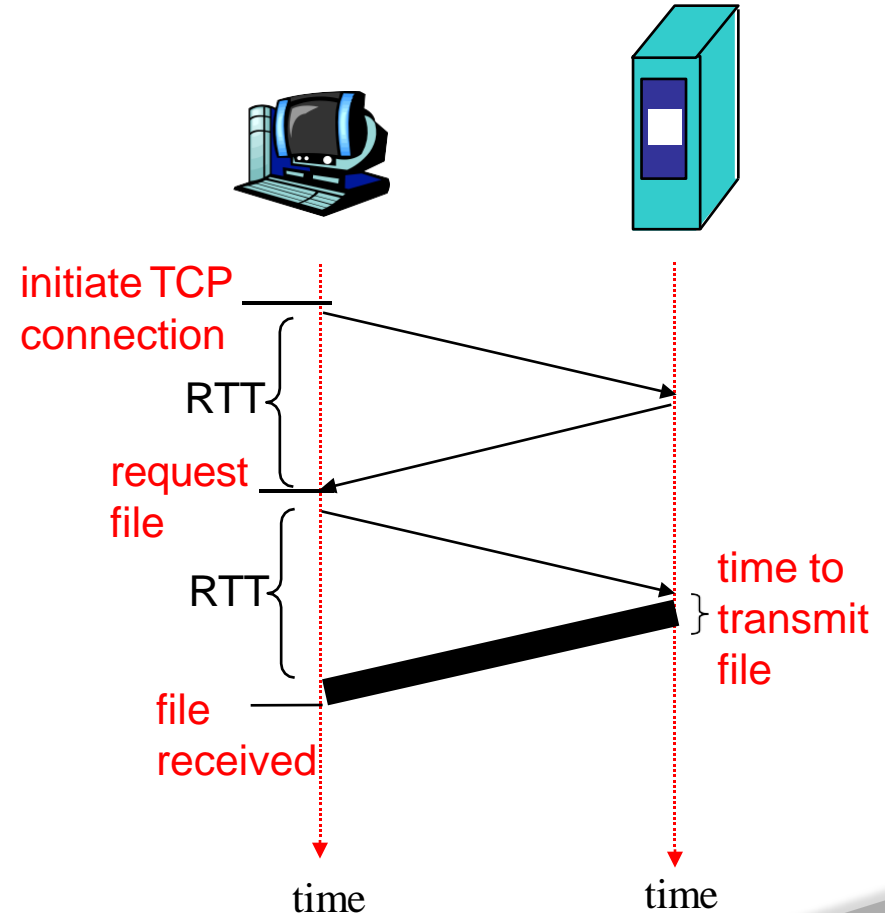
Non-Persistent HTTP: Response time

Definition of RTT: time for a small packet to travel from client to server and back.

Response time:

- one RTT to initiate TCP connection
- one RTT for HTTP request and first few bytes of HTTP response to return
- file transmission time

total = $2RTT + \text{transmit time}$



Persistent HTTP

Nonpersistent HTTP issues:

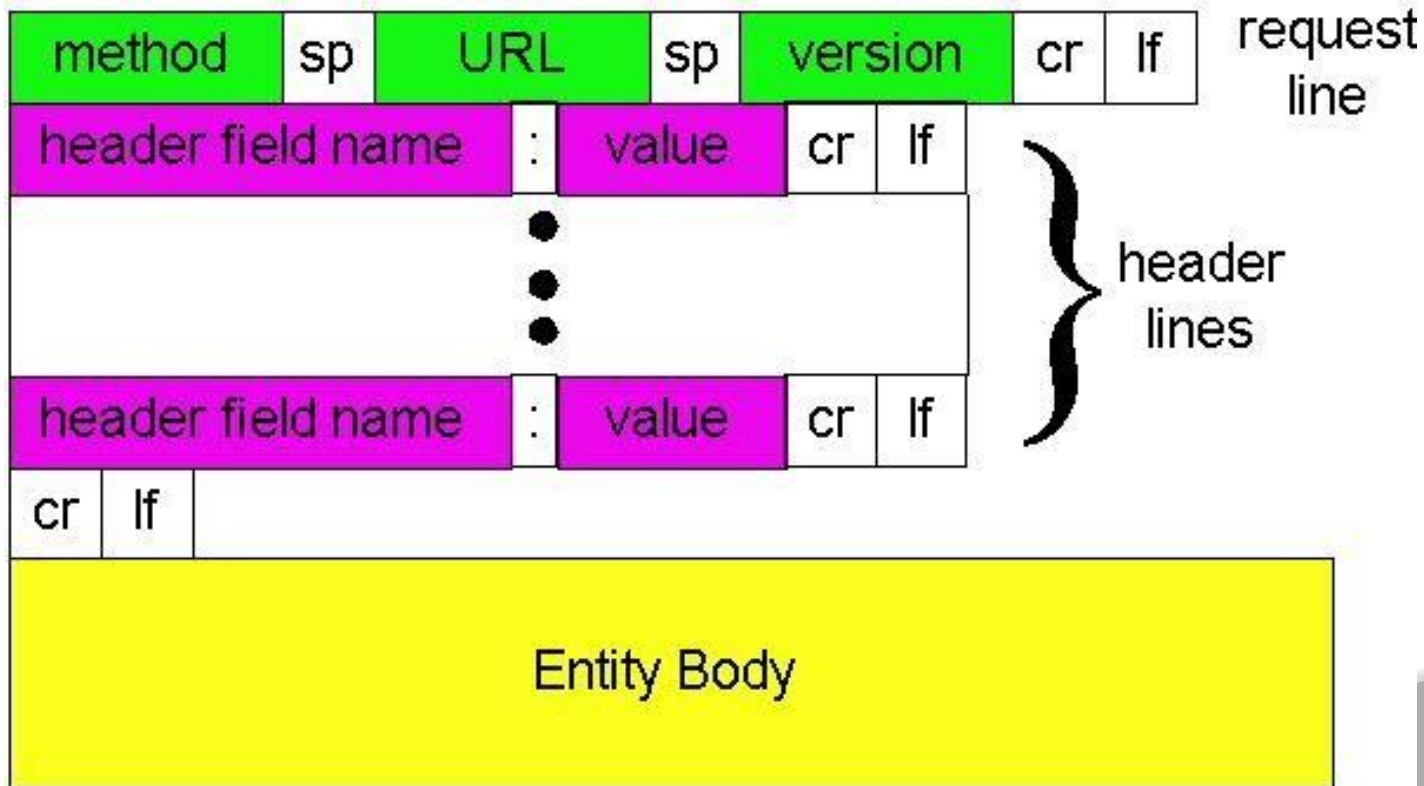
- requires 2 RTTs per object
- OS overhead for *each* TCP connection
- browsers often open parallel TCP connections to fetch referenced objects

Persistent HTTP

- server leaves connection open after sending response
- subsequent HTTP messages between same client/server sent over open connection
- client sends requests as soon as it encounters a referenced object
- as little as one RTT for all the referenced objects

HTTP messages

- ❑ two types of HTTP messages: *request*, *response*
- ❑ HTTP request message:
 - ❖ ASCII (human-readable format)



Method types

HTTP/1.0

- GET
 - request an object from server
- POST
 - upload information using forms
- HEAD
 - asks server to leave requested object out of response

HTTP/1.1

- GET, POST, HEAD
- PUT
 - uploads file in entity body to path specified in URL field
- DELETE
 - deletes file specified in the URL field

Cookies: Keeping state

What cookies can bring:

- authorization
- shopping carts
- recommendations
- user session state (Web e-mail)

- aside
- ### Cookies and privacy:
- ☐ cookies permit sites to learn a lot about you
 - ☐ you may supply name and e-mail to sites

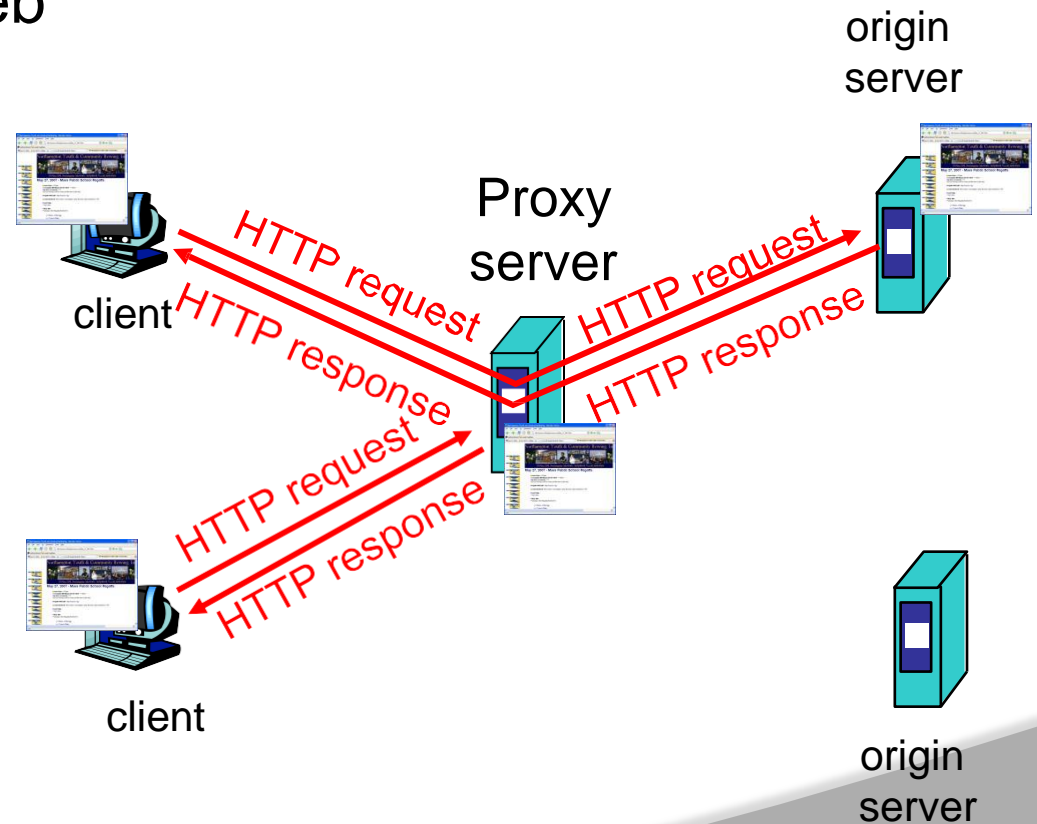
How to keep “state”:

- ☐ protocol endpoints: maintain state at sender/receiver over multiple transactions
- ☐ cookies: http messages carry state

Web caches (proxy server)

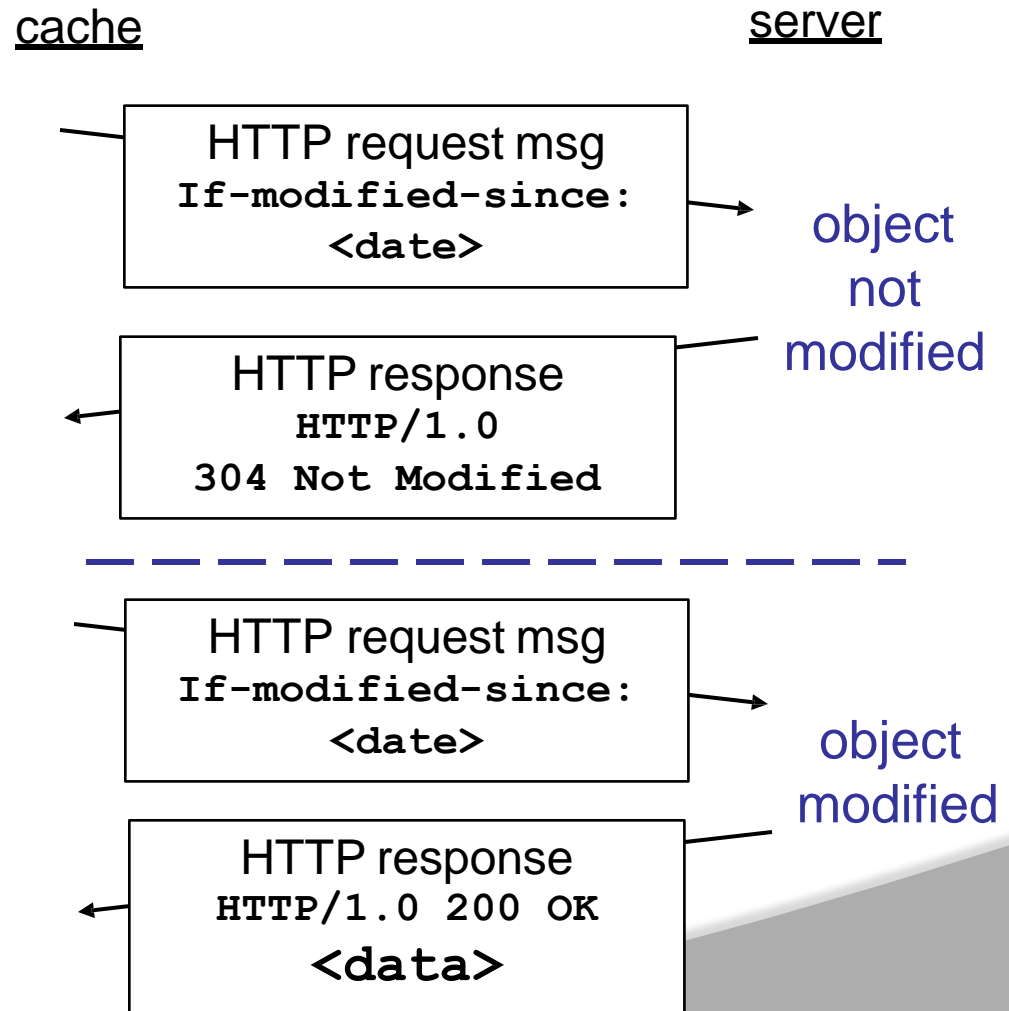
Goal: satisfy client request without involving origin server

- user sets browser: Web accesses via cache
- browser sends all HTTP requests to cache
- Why Web caching?
 - reduce response time for client request
 - reduce traffic on an institution's access link.
 - enables “poor” content providers to effectively deliver content



Conditional GET

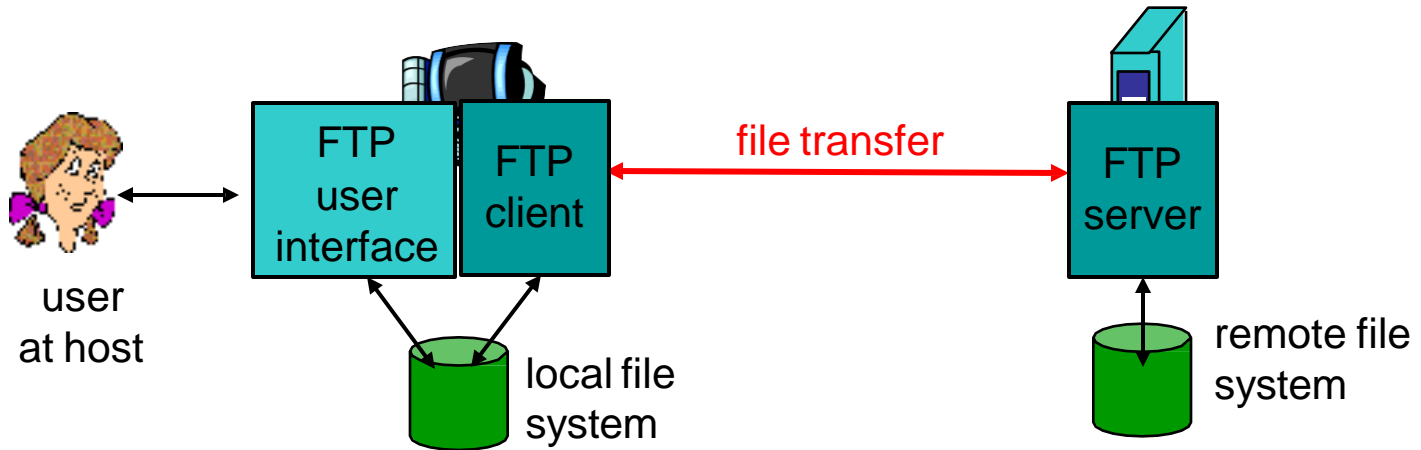
- **Goal:** don't send object if cache has up-to-date cached version
- **cache:** specify date of cached copy in HTTP request
`If-modified-since: <date>`
- **server:** response contains no object if cached copy is up-to-date:
`HTTP/1.0 304 Not Modified`



Lecture 5: Outline

- 2.1 Principles of network applications
- 2.2 Web and HTTP
- **2.3 FTP**
- 2.4 Electronic Mail
 - SMTP, POP3, IMAP
- 2.5 DNS

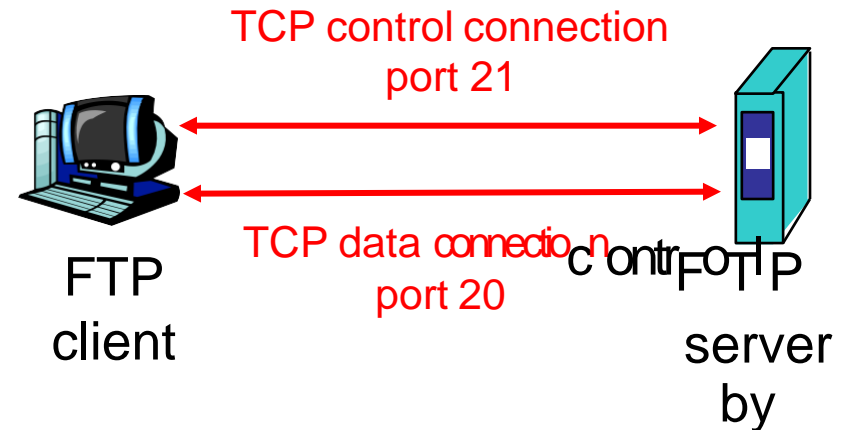
FTP: the file transfer protocol



- transfer file to/from remote host
- client/server model
 - *client*: side that initiates transfer (either to/from remote)
 - *server*: remote host
- ftp: RFC 959
- ftp server: port 21

FTP: separate control, data connections

- FTP client contacts FTP server port 21
- client authorized over connection
- client browses remote directory sending commands over control connection.
- when server receives file transfer command, server opens 2nd TCP connection (for file) to client
- after transferring one file, server closes data connection.
- server opens another TCP data connection to transfer another file.
- control connection: “out of band”
- FTP server maintains “state”: current directory, earlier authentication



FTP commands, responses

Sample commands:

- sent as ASCII text over control channel
- **USER *username***
- **PASS *password***
- **LIST** return list of file in current directory
- **RETR *filename*** retrieves (gets) file
- **STOR *filename*** stores (puts) file onto remote host

Sample return codes

- status code and phrase (as in HTTP)
- **331 Username OK, password required**
- **125 data connection already open; transfer starting**
- **425 Can't open data connection**
- **452 Error writing file**

FTP issues

- Multiple connections are used
 - for each directory listing and file transmission
- No integrity check at receiver
- Messages are sent in clear text
 - including Passwords and file contents
 - can be sniffed by eavesdroppers
- Solution
 - Secure FTP (SSH FTP)
 - allows a range of operations on remote files
 - FTPS (FTP over Secure Sockets Layer (SSL))
 - Transport Layer Security (TLS) encryption

Lecture 5: Outline

- 2.1 Principles of network applications
- 2.2 Web and HTTP
- 2.3 FTP
- 2.4 Electronic Mail
 - SMTP
 - POP3
 - IMAP
- 2.5 DNS

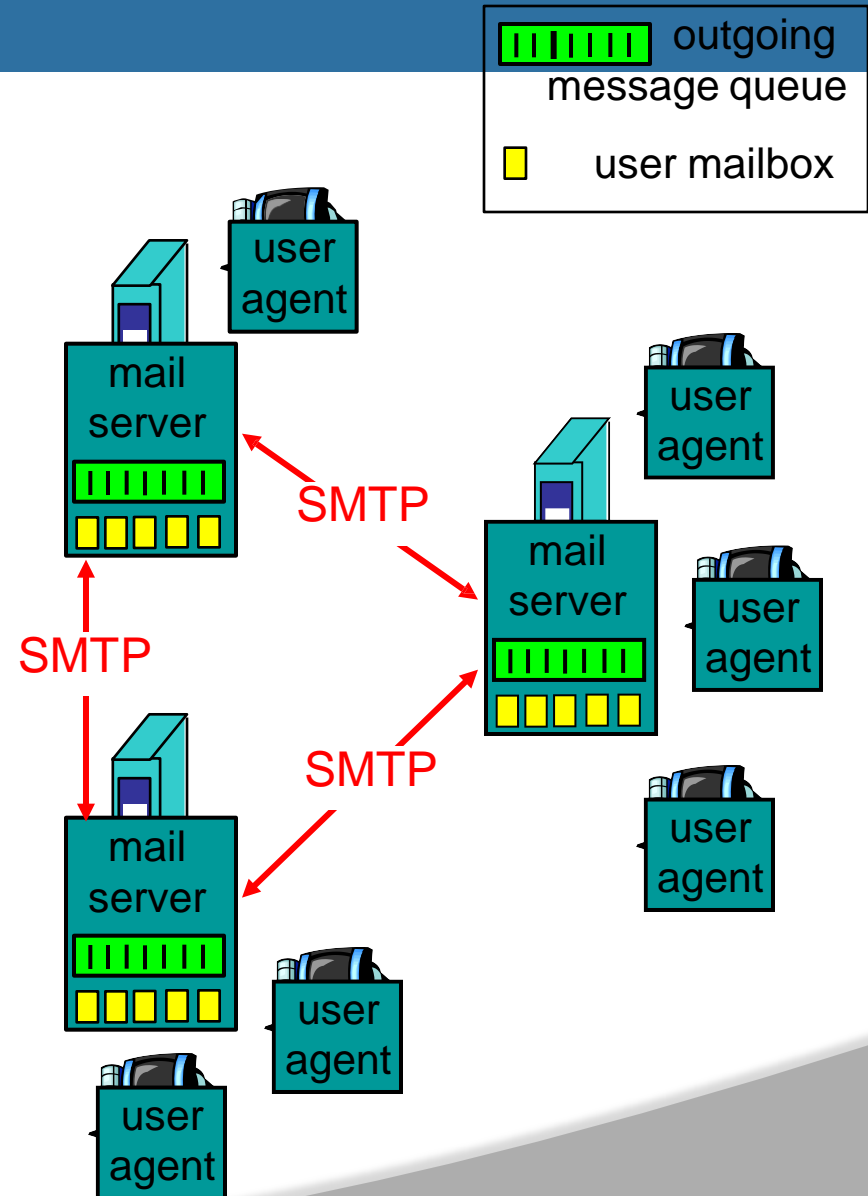
Electronic Mail

Three major components:

- user agents
- mail servers
- simple mail transfer protocol: SMTP

User Agent

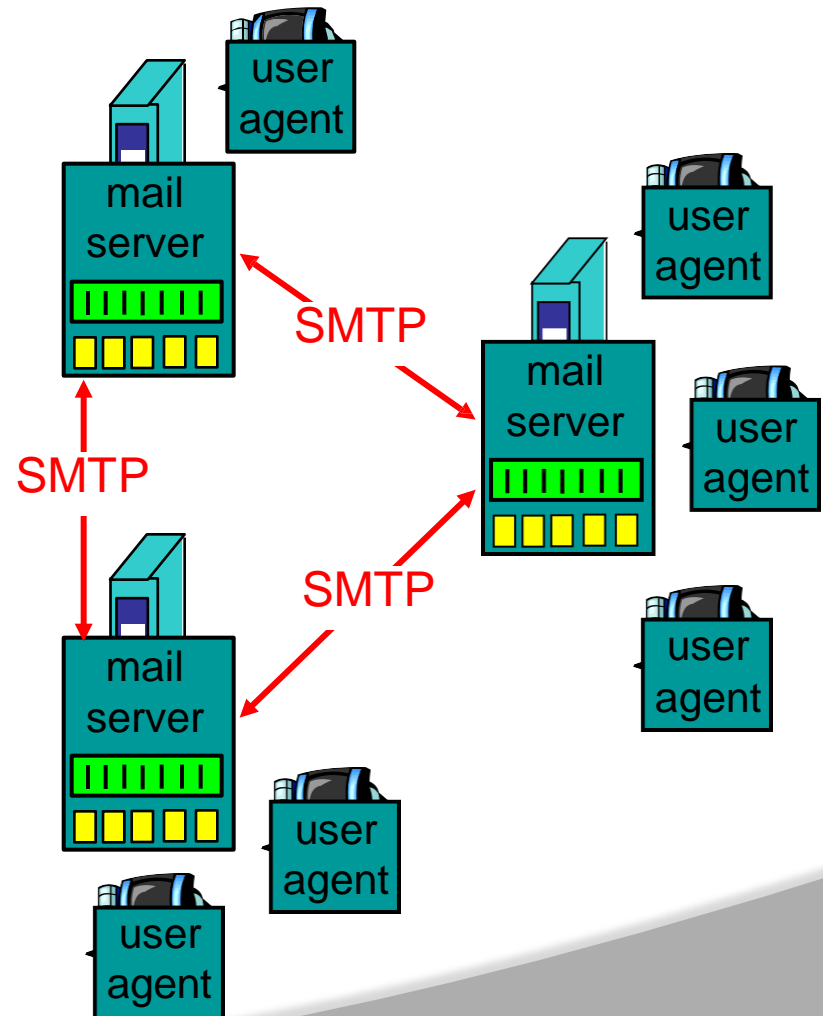
- a.k.a. “mail reader”
- composing, editing, reading mail messages
- e.g., Eudora, Outlook, elm, Mozilla Thunderbird
- outgoing, incoming messages stored on server



Electronic Mail: mail servers

Mail Servers

- **mailbox** contains incoming messages for user
- **message queue** of outgoing (to be sent) mail messages
- **SMTP protocol** between mail servers to send email messages
 - client: sending mail server
 - “server”: receiving mail server

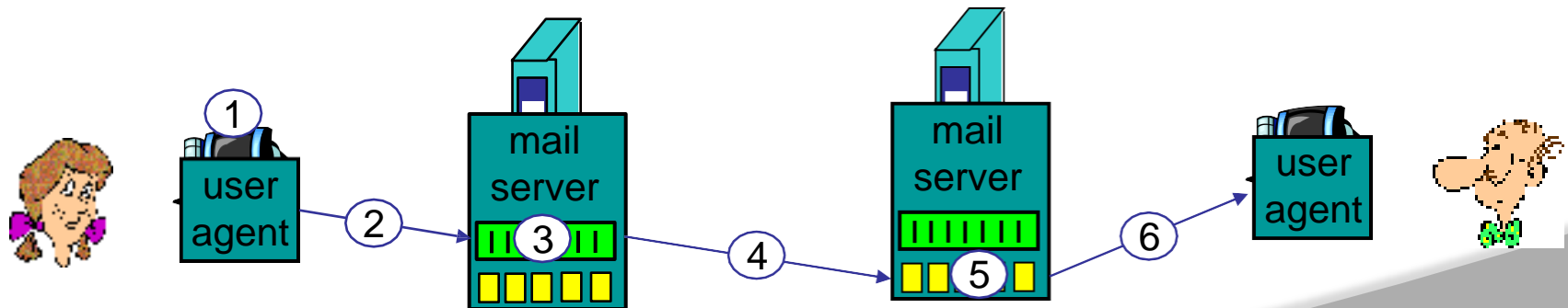


Electronic Mail: SMTP [RFC 2821]

- uses TCP to reliably transfer email message from client to server (port 25)
- direct transfer: sending server to receiving server
- three phases of transfer
 - handshaking (greeting)
 - transfer of messages
 - closure
- command/response interaction
 - **commands**: ASCII text
 - **response**: status code and phrase
- messages must be in 7-bit ASCII

Scenario: Alice sends message to Bob

- 1) Alice uses UA to compose message and “to” bob@someschool.edu
- 2) Alice’s UA sends message to her mail server; message placed in message queue
- 3) Client side of SMTP opens TCP connection with Bob’s mail server
- 4) SMTP client sends Alice’s message over the TCP connection
- 5) Bob’s mail server places the message in Bob’s mailbox
- 6) Bob invokes his user agent to read message



Sample SMTP interaction

```
S: 220 hamburger.edu
C: HELO crepes.fr
S: 250 Hello crepes.fr, pleased to meet you
C: MAIL FROM: <alice@crepes.fr>
S: 250 alice@crepes.fr... Sender ok
C: RCPT TO: <bob@hamburger.edu>
S: 250 bob@hamburger.edu ... Recipient ok
C: DATA
S: 354 Enter mail, end with "." on a line by itself
C: Do you like ketchup?
C: How about pickles?
C: .
S: 250 Message accepted for delivery
C: QUIT
S: 221 hamburger.edu closing connection
```

SMTP: final words

- SMTP uses persistent connections
- SMTP requires message (header & body) to be in 7-bit ASCII
- SMTP server uses `CRLF.CRLF` to determine end of message

Comparison with HTTP:

- HTTP: pull
- SMTP: push
- both have ASCII command/response interaction, status codes
- HTTP: each object encapsulated in its own response msg
- SMTP: multiple objects sent in multipart msg

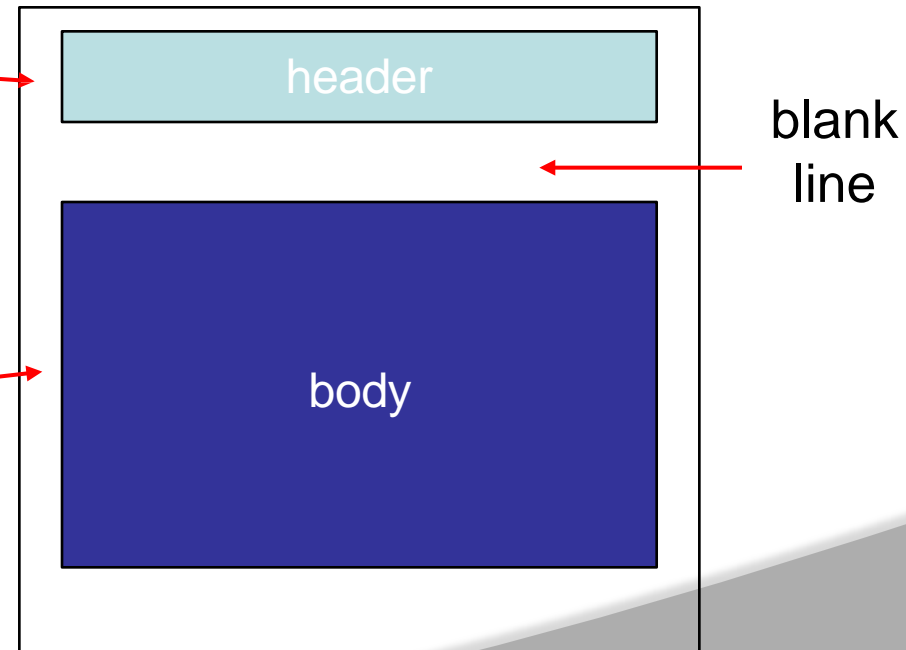
Mail message format

SMTP: protocol for exchanging email msgs

RFC 822: standard for text message format:

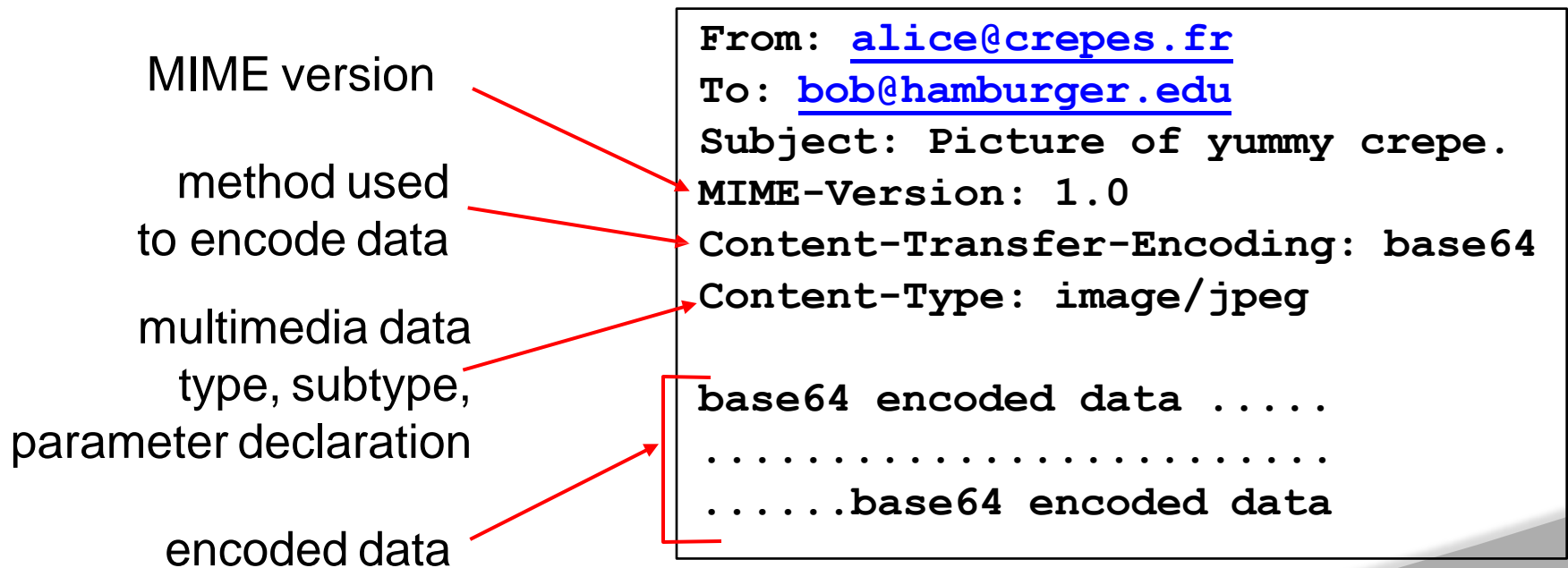
- header lines, e.g.,
 - To:
 - From:
 - Subject:*different from SMTP commands!*

- body
 - the “message”,
ASCII characters only

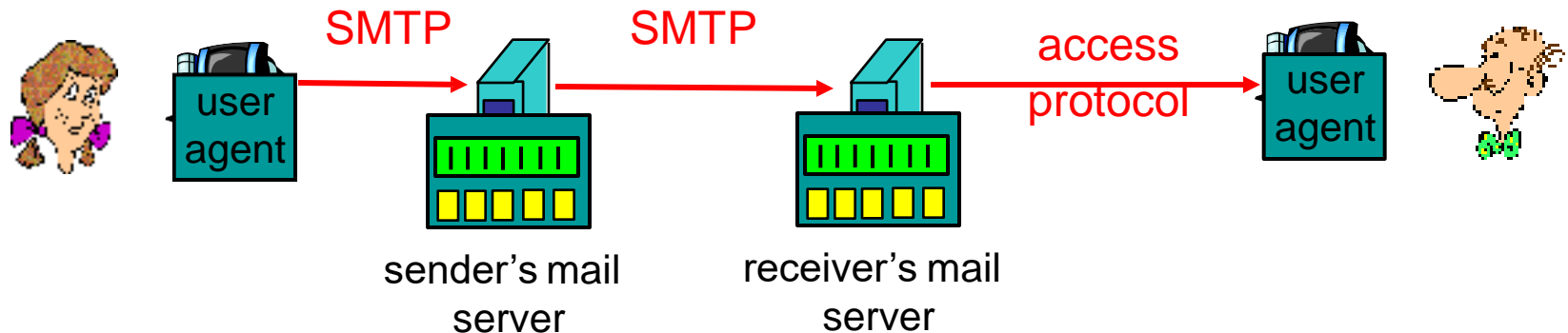


Message format: multimedia extensions

- MIME: multimedia mail extension, RFC 2045, 2056
- additional lines in msg header declare MIME content type



Mail access protocols



- SMTP: delivery/storage to receiver's server
- Mail access protocol: retrieval from server
 - POP: Post Office Protocol [RFC 1939]
 - authorization (agent <-->server) and download
 - IMAP: Internet Mail Access Protocol [RFC 1730]
 - more features (more complex)
 - manipulation of stored msgs on server
 - HTTP: gmail, Hotmail, Yahoo! Mail, etc.

POP3 (more) and IMAP

More about POP3

- Previous example uses “download and delete” mode.
- Bob cannot re-read e-mail if he changes client
- “Download-and-keep”: copies of messages on different clients
- POP3 is stateless across sessions

IMAP

- Keep all messages in one place: the server
- Allows user to organize messages in folders
- IMAP keeps user state across sessions:
 - names of folders and mappings between message IDs and folder name

SMTP interaction for yourself

- `telnet servername 25`
- see 220 reply from server
- enter HELO, MAIL FROM, RCPT TO, DATA, QUIT commands

above lets you send email without using email client (reader)