# Chapter 5
# Link Layer and LANs

# Chapter 5: The Data Link Layer

## Our goals:

□ understand principles behind data link layer services:

- ○ error detection, correction
- ○ sharing a broadcast channel: multiple access
- ○ link layer addressing
- ○ reliable data transfer, flow control: *done!*

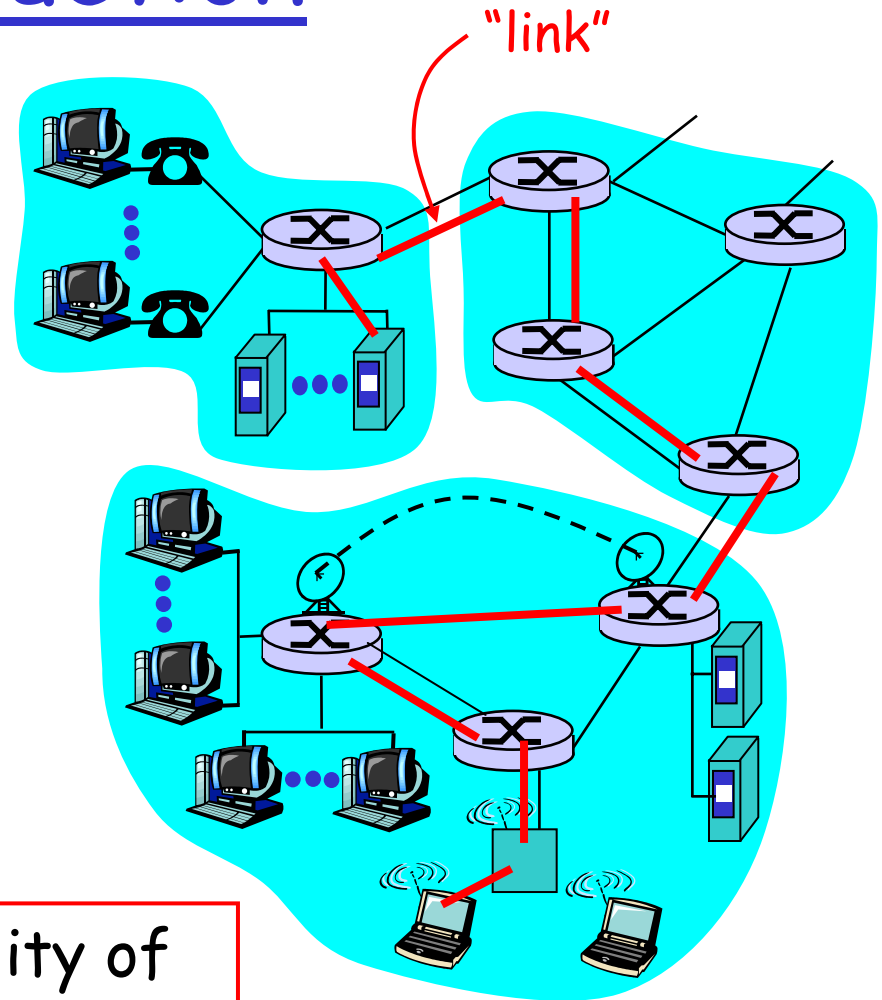□ instantiation and implementation of various link layer technologies

# Link Layer

# Link Layer: Introduction

## Some terminology:

□ hosts and routers are **nodes**

□ communication channels that connect adjacent nodes along communication path are **links**
- ○ wired links
- ○ wireless links
- ○ LANs

□ layer-2 packet is a **frame**, encapsulates datagram

"link"

**data-link layer** has responsibility of transferring datagram from one node to 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
- intermediate trips= 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:
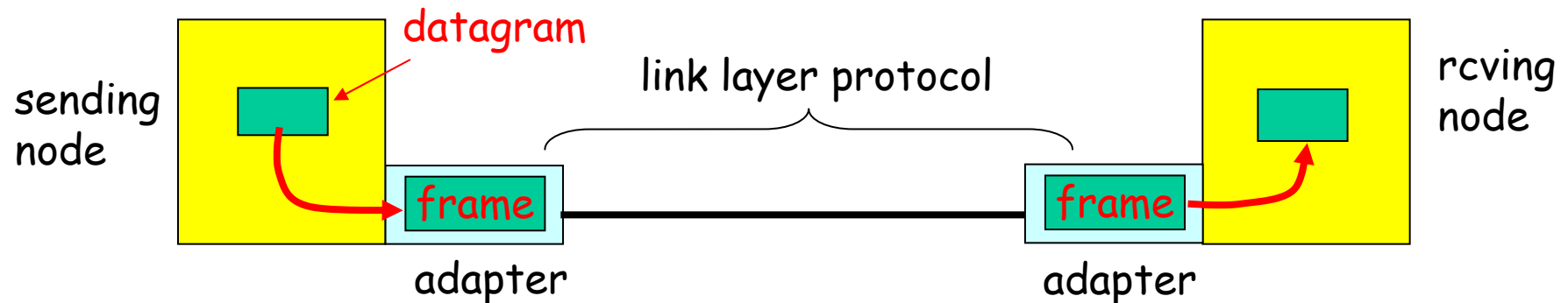  - • signals sender for retransmission or drops frame

□ 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

# Adaptors Communicating



- link layer implemented in "adaptor" (aka NIC)
  - Ethernet card, PCMCI card, 802.11 card
- sending side:
  - encapsulates datagram in a frame
  - adds error checking bits, rdt, flow control, etc.
- receiving side
  - looks for errors, rdt, flow control, etc
  - extracts datagram, passes to rcving node
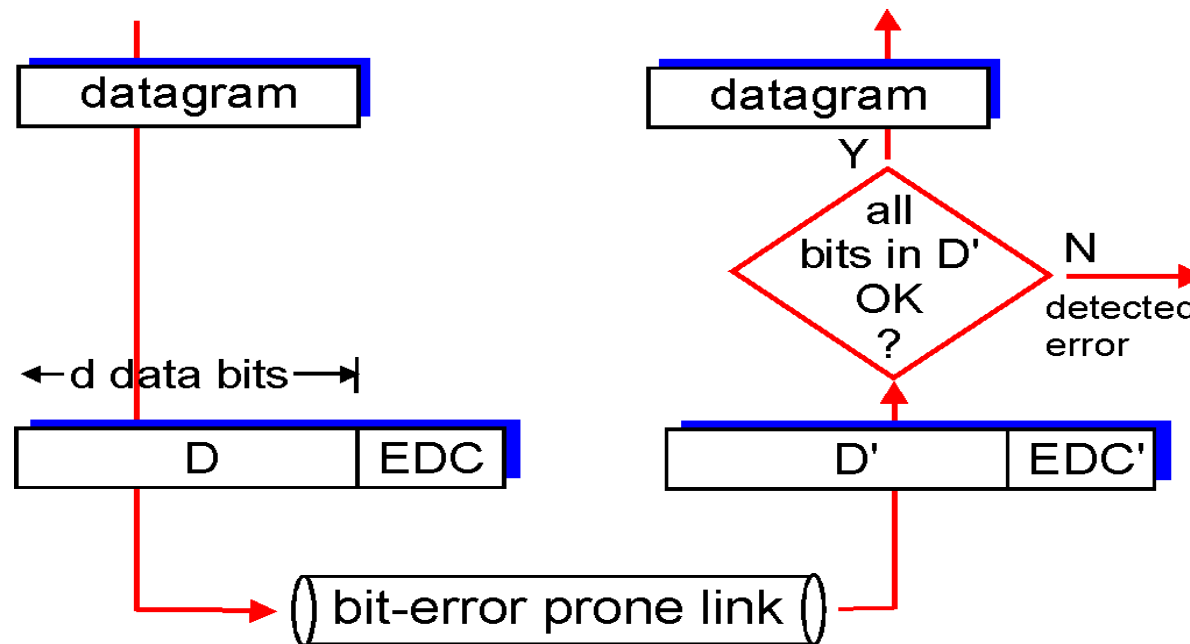- adapter is semi-autonomous
- link & physical layers

# Link Layer

# Error Detection

EDC= Error Detection and Correction bits (redundancy)
D   = Data protected by error checking, may include header fields

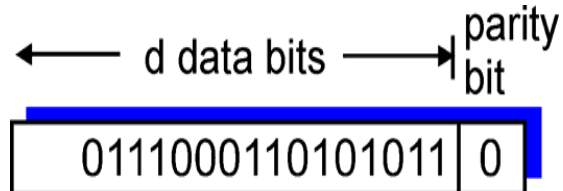- Error detection not 100% reliable!
    - protocol may miss some errors, but rarely
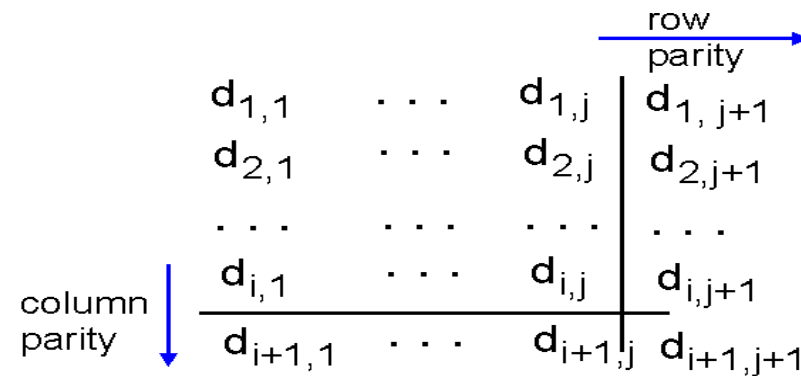    - larger EDC field yields better detection and correction

```
    datagram                          datagram
                                          Y
                                        ┌─────┐
                                        │ all │
                                    ┌───┤bits in D'├───┐ N
                                        │ OK  │          → detected
                                        │  ?  │            error
                                        └─────┘
   ←d data bits→
  ┌─────────┬──────┐              ┌──────────┬───────┐
  │    D    │ EDC  │              │    D'    │ EDC'  │
  └─────────┴──────┘              └──────────┴───────┘

              ┌─ bit-error prone link ─┐
```

# Parity Checking

## Single Bit Parity:
**Detect single bit errors**



d data bits | parity bit

011100110101011 | 0

## Two Dimensional Bit Parity:
**Detect *and correct* single bit errors**

row parity →

$$d_{1,1} \quad \cdots \quad d_{1,j} \quad | \quad d_{1,j+1}$$
$$d_{2,1} \quad \cdots \quad d_{2,j} \quad | \quad d_{2,j+1}$$
$$\cdots \qquad \cdots \qquad \cdots \qquad \cdots$$
$$d_{i,1} \quad \cdots \quad d_{i,j} \quad | \quad d_{i,j+1}$$

column parity ↓

$$d_{i+1,1} \quad \cdots \quad d_{i+1,j} \quad d_{i+1,j+1}$$

```
1 0 1 0 1 | 1          1 0 1 0 1 | 1
1 1 1 1 0 | 0          1 0 1 1 0 | 0  → parity error
0 1 1 1 0 | 1          0 1 1 1 0 | 1
0 0 1 0 1 | 0          0 0 1 0 1 | 0
```

*no errors*          ↓ parity error

*correctable single bit error*

# Internet checksum

**Goal:** detect "errors" (e.g., flipped bits) in transmitted segment (note: used at transport layer *only*)

## Sender:
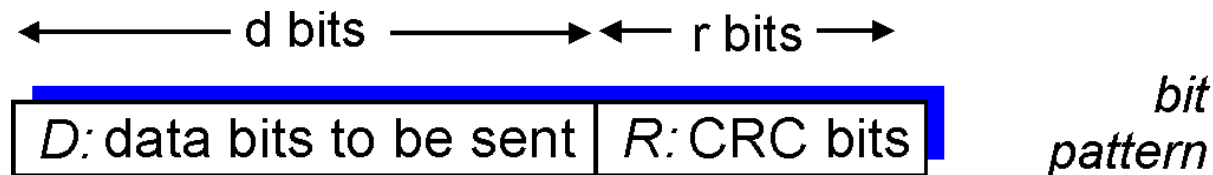
□ treat segment contents as sequence of 16-bit integers

□ 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?* More later ....

# Checksumming: Cyclic Redundancy Check

- view data bits, D, as a binary number
- choose r+1 bit pattern (generator), G
- goal: choose r CRC bits, R, such that
  - <D,R> exactly divisible by G (modulo 2)
  - receiver knows G, divides <D,R> by G.  If non-zero remainder: error detected!
  - can detect all burst errors less than r+1 bits
- widely used in practice (ATM, HDCL)



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

# CRC Example

Want:

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

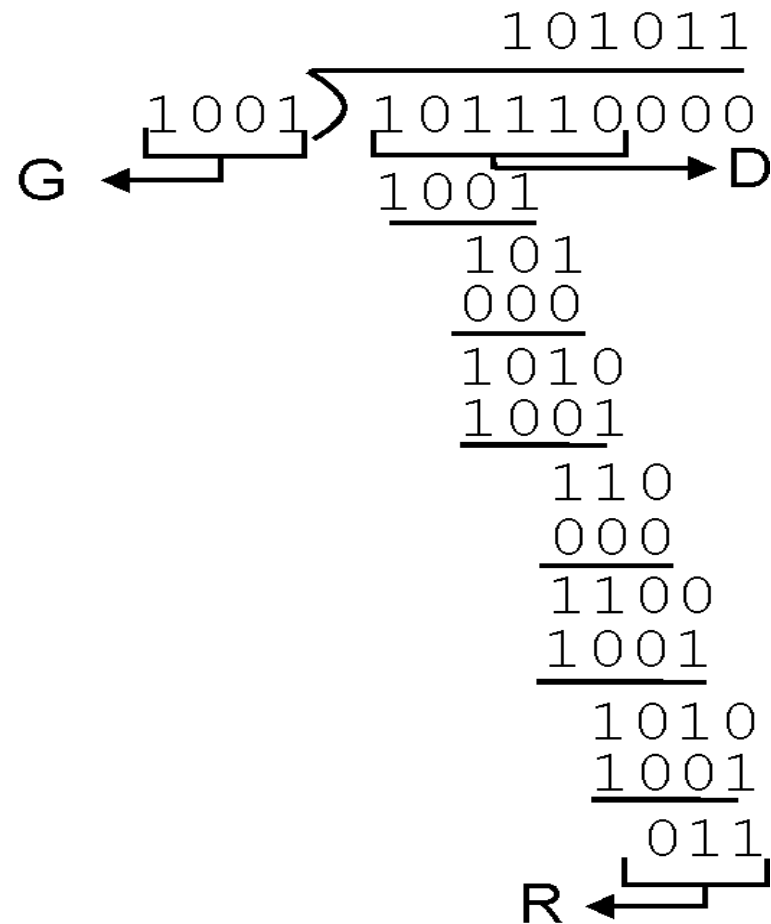*equivalently:*

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

*equivalently:*

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

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

```
                    101011
           1001 ) 101110000
  G <-----        1001              -----> D
                   101
                   000
                   1010
                   1001
                    110
                    000
                    1100
                    1001
                     1010
                     1001
                      011
  R <-----
```

# Link Layer

- 5.1 Introduction and services
- 5.2 Error detection and correction
- <span style="color:red">5.3 Multiple access protocols</span>
- 5.4 Link-Layer Addressing
- 5.5 Ethernet

- 5.6 Hubs and switches
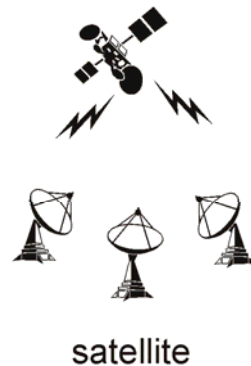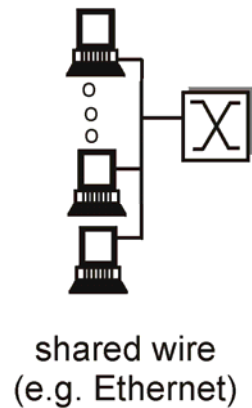- 5.7 PPP

# Multiple Access Links and Protocols

Two types of "links":

- □ point-to-point
  - ○ PPP for dial-up access
  - ○ point-to-point link between Ethernet switch and host
- □ broadcast (shared wire or medium)
  - ○ traditional Ethernet
  - ○ upstream HFC
  - ○ 802.11 wireless LAN

shared wire (e.g. Ethernet)

shared wireless (e.g. Wavelan)

satellite

Blah, blah, blah

ZZZzzzzzzzzzz

cocktail party

# 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

# Ideal Multiple Access Protocol

**Broadcast channel of rate R bps**

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: a 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

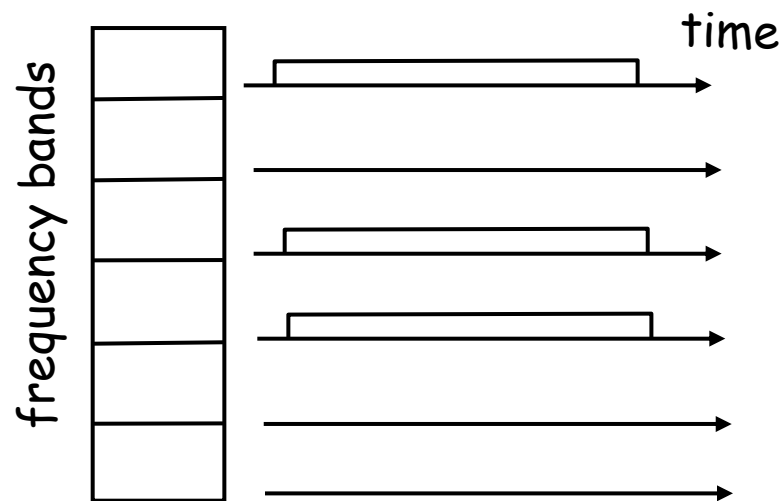# Channel Partitioning MAC protocols: TDMA

**TDMA: time division multiple access**

- access to channel in "rounds"
- each station gets fixed length slot (length = pkt trans time) in each round
- unused slots go idle
- example: 6-station LAN, 1,3,4 have pkt, slots 2,5,6 idle

# Channel Partitioning MAC protocols: FDMA

## FDMA: frequency division multiple access

- channel spectrum divided into frequency bands
- each station assigned fixed frequency band
- unused transmission time in frequency bands go idle
- example: 6-station LAN, 1,3,4 have pkt, frequency bands 2,5,6 idle

# 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
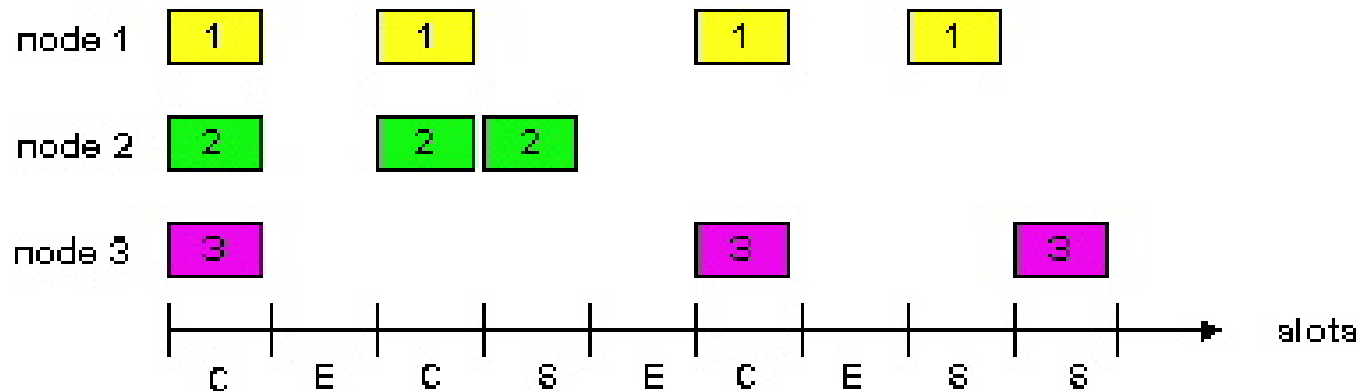  ○ ALOHA
  ○ CSMA, CSMA/CD, CSMA/CA

# Slotted ALOHA

## Assumptions

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

## Operation

- when node obtains fresh frame, it transmits in next slot
- 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

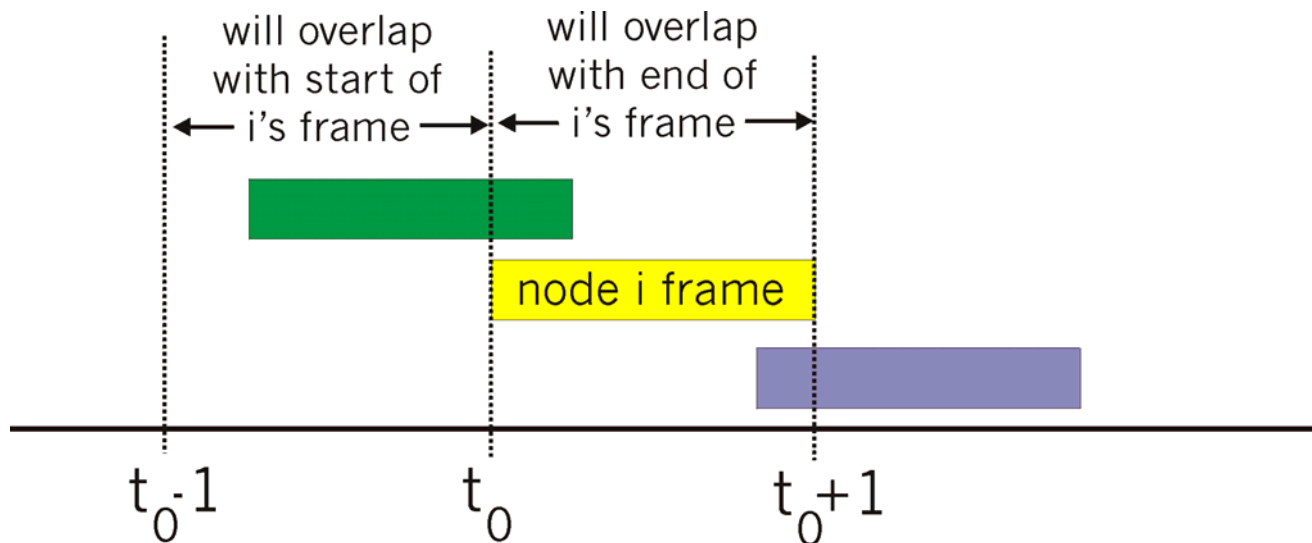**Efficiency** is the long-run fraction of successful slots when there are many nodes, each with many frames to send

- Suppose N nodes with many frames to send, each transmits in slot with probability $p$

- prob that node 1 has success in a slot
  = $p(1-p)^{N-1}$

- prob that any node has a success = $Np(1-p)^{N-1}$

- For max efficiency with N nodes, 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 $1/e = .37$

*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 in $[t_0-1, t_0+1]$

will overlap with start of ← i's frame →

will overlap with end of ← i's frame →

node i frame

$t_0-1$          $t_0$          $t_0+1$

# Pure Aloha efficiency

P(success by given node) = P(node transmits) ·

$\qquad$ P(no other node transmits in $[t_0-1,t_0]$) ·

$\qquad$ P(no other node transmits in $[t_0,t_0+1]$

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

$\qquad$ = p · $(1-p)^{2(N-1)}$

$\qquad$ … choosing optimum p and then letting n -> infty …

Even worse !$\qquad$ = 1/(2e) = .18

# CSMA (Carrier Sense Multiple Access)

**CSMA**: listen before transmit:

If channel sensed idle: transmit entire frame

□ If channel sensed busy, defer transmission


□ Human analogy: don't interrupt others!

# CSMA collisions

**collisions *can* still occur:**

propagation delay means
two nodes may not hear
each other's transmission

**collision:**

entire packet transmission
time wasted

**note:**

role of distance & propagation
delay in determining collision
probability

spatial layout of nodes

# 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: receiver shut off while transmitting

☐ 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

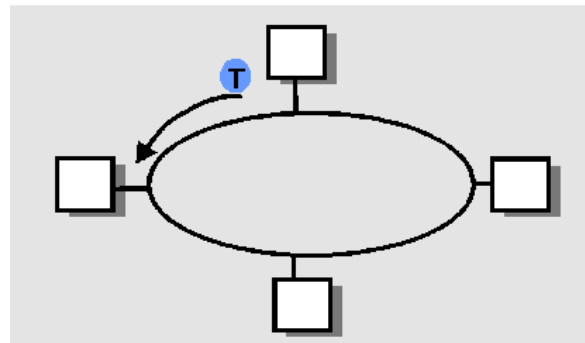"taking turns" protocols

   look for best of both worlds!

# "Taking Turns" MAC protocols

## Polling:

□ master node "invites" slave nodes to transmit in turn

□ concerns:
  ○ polling overhead
  ○ latency
  ○ single point of failure (master)

## 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

□ **What do you do with a shared media?**

  ○ Channel Partitioning, by time, frequency or code
    • Time Division, Frequency Division

  ○ Random partitioning (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 a central site, token passing

# LAN technologies

Data link layer so far:
- services, error detection/correction, multiple access

Next: LAN technologies
- addressing
- Ethernet
- hubs, switches
- PPP

# Link Layer

# MAC Addresses and ARP
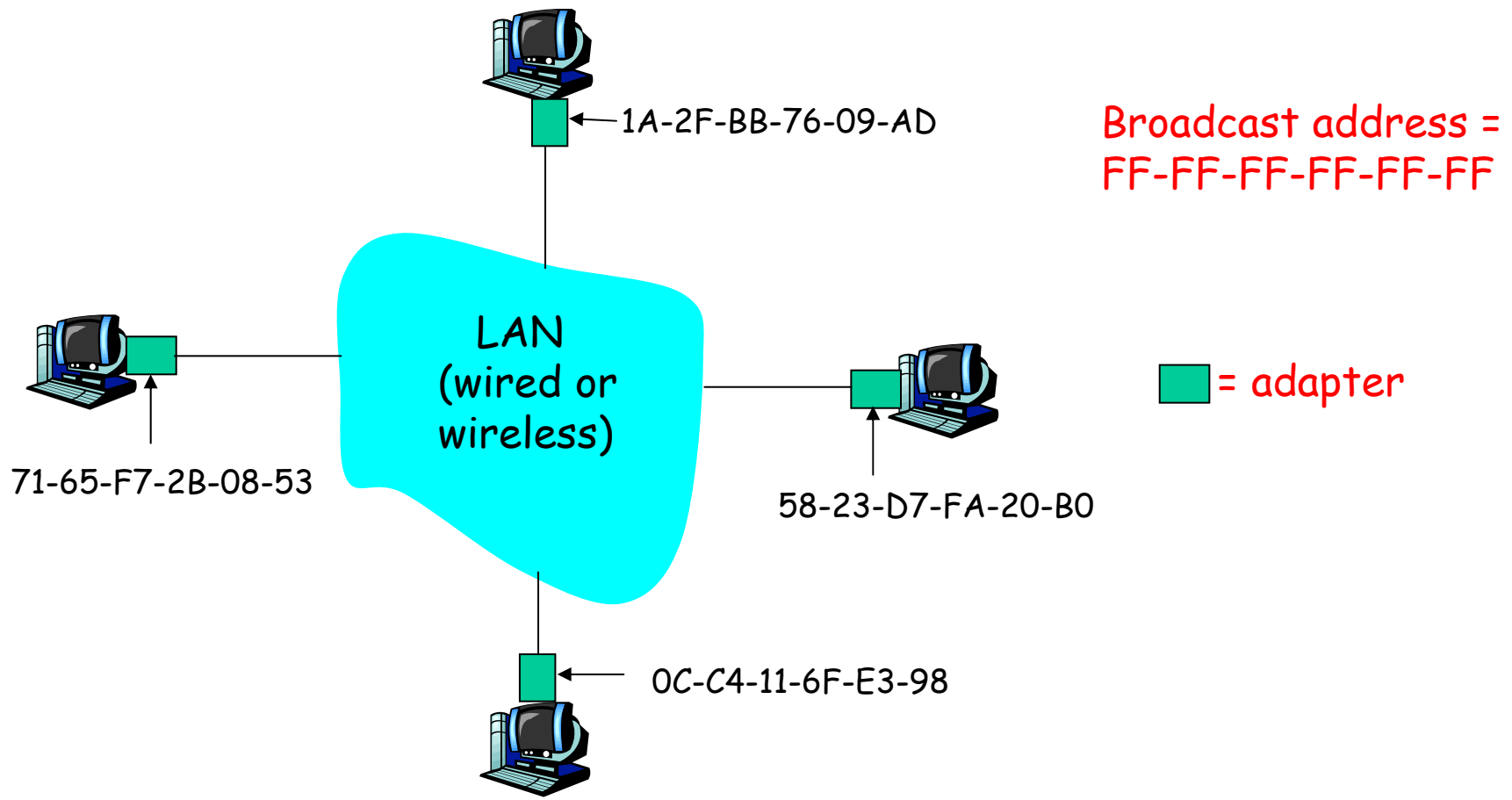
□ **32-bit IP address:**

  ○ *network-layer* address

  ○ used to get datagram to destination IP subnet

□ **MAC (or LAN or physical or Ethernet) address:**

  ○ used to get datagram from one interface to another physically-connected interface (same network)

  ○ 48 bit MAC address (for most LANs) burned in the adapter ROM

# LAN Addresses and ARP

Each adapter on LAN has unique LAN address

1A-2F-BB-76-09-AD

Broadcast address =
FF-FF-FF-FF-FF-FF

LAN
(wired or
wireless)

71-65-F7-2B-08-53

58-23-D7-FA-20-B0

= adapter

0C-C4-11-6F-E3-98

# LAN Address (more)

☐ MAC address allocation administered by IEEE
☐ manufacturer buys portion of MAC address space (to assure uniqueness)
☐ Analogy:

 (a) MAC address: like Social Security Number

 (b) IP address: like postal address

☐ MAC flat address ➜ portability

 ○ can move LAN card from one LAN to another

☐ IP hierarchical address NOT portable

 ○ depends on IP subnet to which node is attached
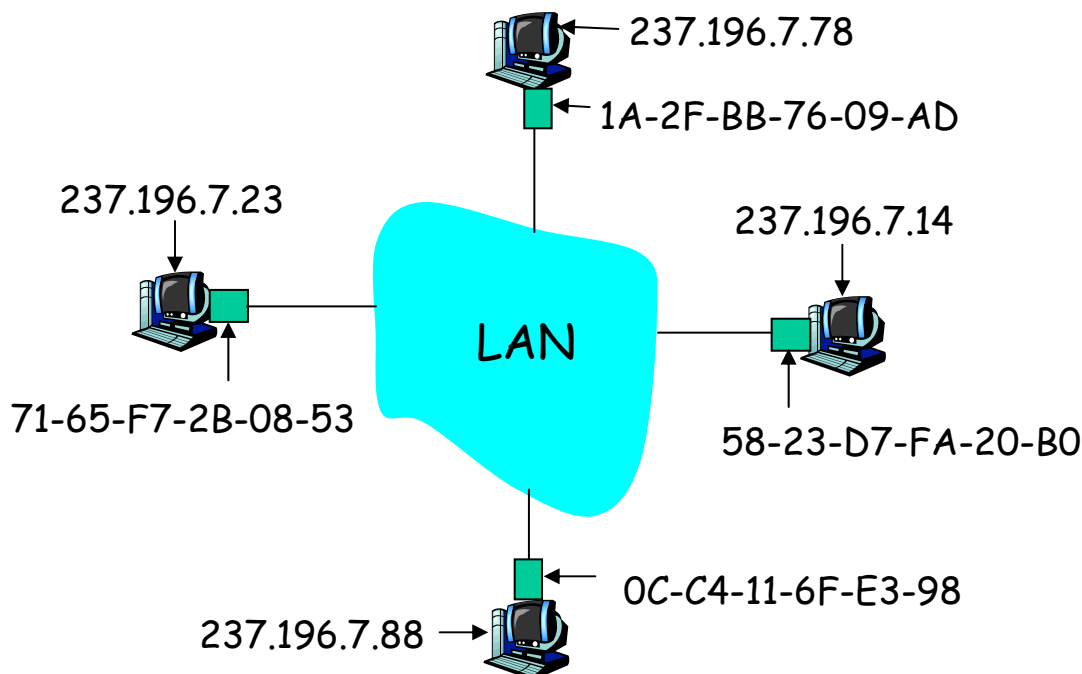
# ARP: Address Resolution Protocol

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

- Each IP node (Host, Router) on LAN has ARP table

- ARP Table: IP/MAC address mappings for some LAN nodes

  < IP address; MAC address; TTL>

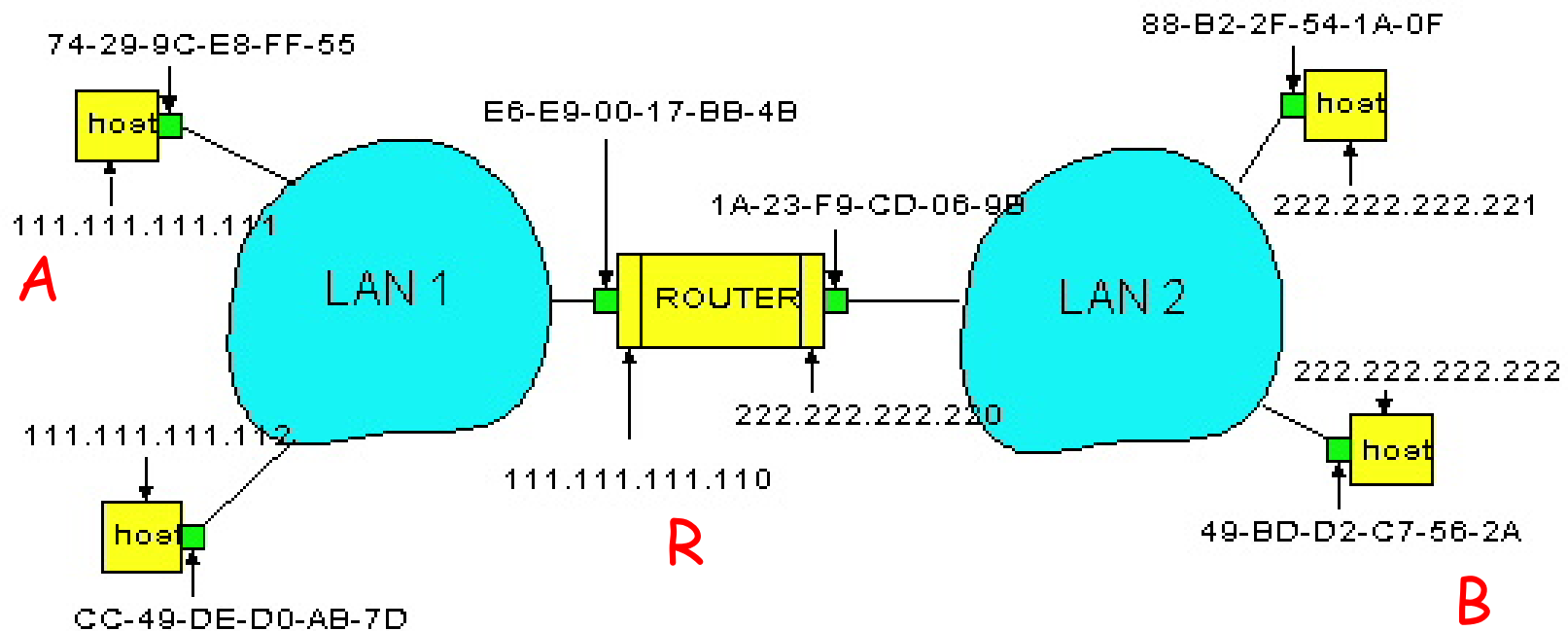  - TTL (Time To Live): time after which address mapping will be forgotten (typically 20 min)

237.196.7.78
1A-2F-BB-76-09-AD

237.196.7.23

237.196.7.14

71-65-F7-2B-08-53

LAN

58-23-D7-FA-20-B0

237.196.7.88

0C-C4-11-6F-E3-98

# ARP protocol: Same LAN (network)

□ A wants to send datagram to B, and 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 administrator
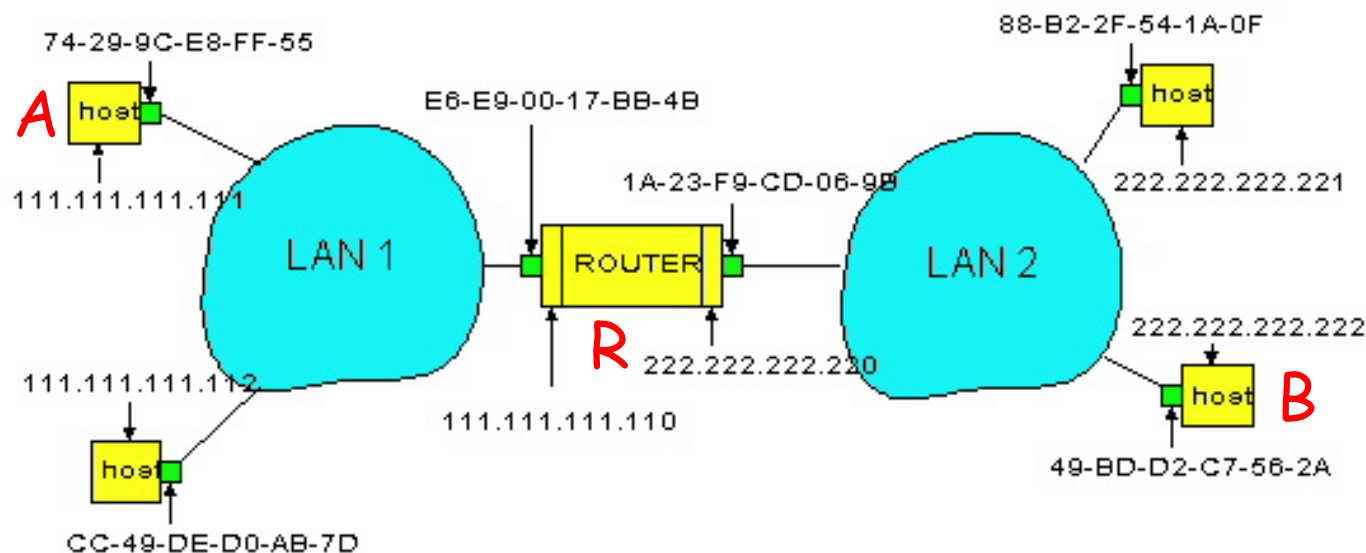
# Routing to another LAN

walkthrough: send datagram from A to B via R

assume A know's B IP address



- Two ARP tables in router R, one for each IP network (LAN)

- A creates datagram with source A, destination B
- A uses ARP to get R's MAC address for 111.111.111.110
- A creates link-layer frame with R's MAC address as dest, frame contains A-to-B IP datagram
- A's adapter sends frame
- R's adapter receives frame
- R removes IP datagram from Ethernet frame, sees its destined to B
- R uses ARP to get B's MAC address
- R creates frame containing A-to-B IP datagram sends to B

74-29-9C-E8-FF-55

88-B2-2F-54-1A-0F

A host

host

E6-E9-00-17-BB-4B

111.111.111.111

1A-23-F9-CD-06-9B

222.222.222.221

LAN 1

ROUTER

LAN 2

R

222.222.222.222

111.111.111.112

222.222.222.220

host

B

111.111.111.110

host

CC-49-DE-D0-AB-7D

49-BD-D2-C7-56-2A

# Dynamic Host Configuration Protocol (DHCP)

Src IP:     0.0.0.0

Dest IP:    255.255.255.255

MAC:        FF-FF-FF-FF-FF-FF

DHCP server

223.1.2.5

223.1.1.1

223.1.1.4    223.1.2.9    223.1.2.1

223.1.3.27

223.1.1.2

Arriving DHCP client

223.1.1.3

223.1.2.2

223.1.3.1    223.1.3.2

**Figure 5.20** ♦ DHCP client-server scenario

# DHCP



**Figure 5.21** ◆ DHCP client-server interaction

Use transaction ID to match the query request with the response!

# Link Layer

□ 5.1 Introduction and services

□ 5.2 Error detection and correction

□ 5.3Multiple access protocols

□ 5.4 Link-Layer Addressing

□ 5.5 Ethernet

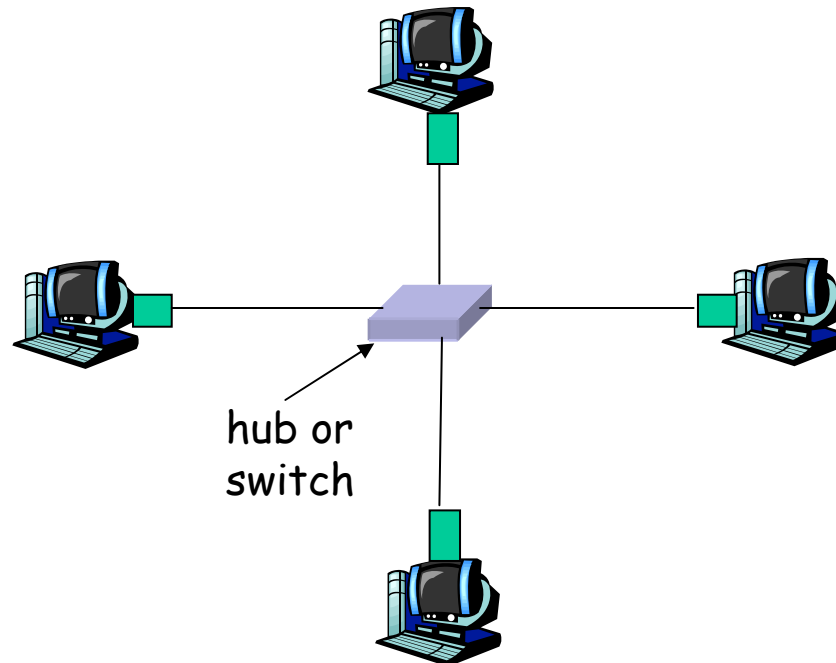□ 5.6 Hubs and switches

# Ethernet

"dominant" wired LAN technology:

☐ cheap $20 for 100Mbs!

☐ first widely used LAN technology

☐ Simpler, cheaper than token LANs and ATM

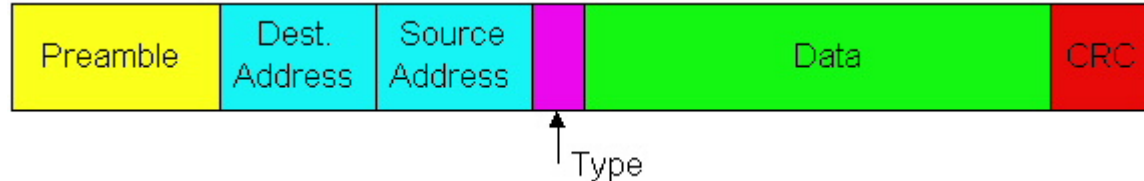☐ Kept up with speed race: 10 Mbps – 10 Gbps

Metcalfe's Ethernet sketch

# Star topology

- Bus topology popular through mid 90s
- Now star topology prevails
- Connection choices: hub or switch (more later)

hub or
switch

# Ethernet Frame Structure

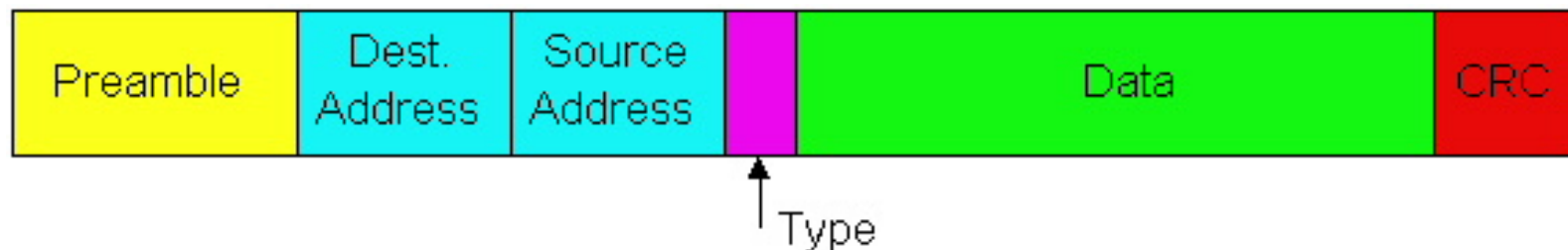Sending adapter encapsulates IP datagram (or other network layer protocol packet) in Ethernet frame



## 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 (eg ARP packet), it passes data in frame to net-layer protocol
  - otherwise, adapter discards frame
- **Type:** indicates the higher layer protocol (mostly IP but others may be supported such as Novell IPX and AppleTalk)
- **CRC:** checked at receiver, if error is detected, the frame is simply dropped

| Preamble | Dest. Address | Source Address | | Data | CRC |
|---|---|---|---|---|---|

Type

# Unreliable, connectionless service

□ **Connectionless:** No handshaking between sending and receiving adapter.

□ **Unreliable:** receiving adapter doesn't send acks or nacks to sending adapter
  ○ stream of datagrams passed to network layer can have gaps
  ○ gaps will be filled if app is using TCP
  ○ otherwise, app will see the gaps

# Ethernet uses CSMA/CD

□ No slots

□ adapter doesn't transmit
if it senses that some
other adapter is
transmitting, that is,
carrier sense

□ transmitting adapter
aborts when it senses
that another adapter is
transmitting, that is,
collision detection

□ Before attempting a
retransmission,
adapter waits a
random time, that is,
random access

# Ethernet CSMA/CD algorithm

1. Adaptor receives datagram from net layer & creates frame

2. If adapter senses channel idle, it starts to transmit frame. If it senses channel busy, waits until channel idle and then transmits

3. If adapter transmits entire frame without detecting another transmission, the adapter is done with frame !

4. If adapter detects another transmission while transmitting, aborts and sends jam signal

5. After aborting, adapter enters **exponential backoff**: after the mth collision, adapter chooses a K at random from {0,1,2,…,$2^m$-1}. Adapter waits K·512 bit times and returns to Step 2

# Ethernet's CSMA/CD (more)

Jam Signal: make sure all other transmitters are aware of collision; 48 bits

Bit time: .1 microsec for 10 Mbps Ethernet ;
for K=1023, wait time is about 50 msec

Exponential Backoff:

□ *Goal*: adapt retransmission attempts to estimated current load

    ○ heavy load: random wait will be longer

□ first collision: choose K from {0,1}; delay is K· 512 bit transmission times

□ after second collision: choose K from {0,1,2,3}...

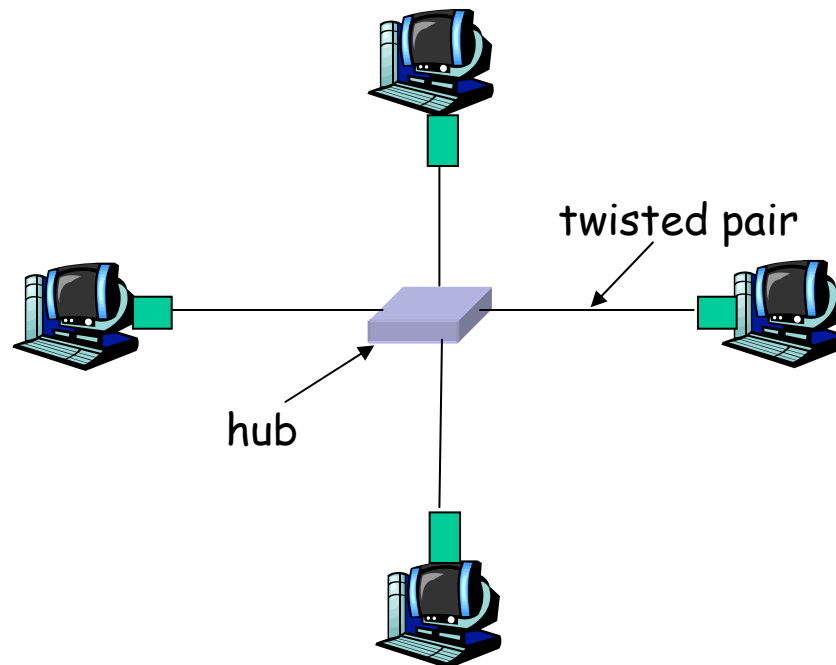□ after ten collisions, choose K from {0,1,2,3,4,...,1023}

# CSMA/CD efficiency

- $T_{prop}$ = max prop between 2 nodes in LAN
- $t_{trans}$ = time to transmit max-size frame

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

- Efficiency goes to 1 as $t_{prop}$ goes to 0
- Goes to 1 as $t_{trans}$ goes to infinity
- Much better than ALOHA, but still decentralized, simple, and cheap
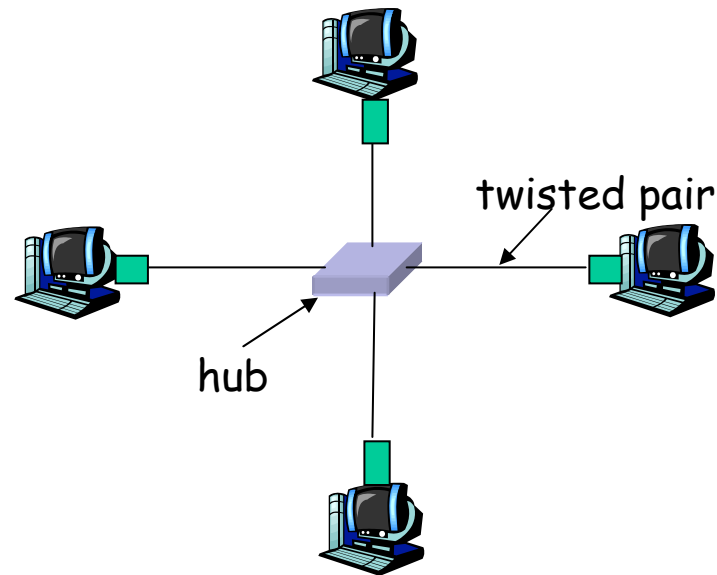
# 10BaseT and 100BaseT

□ 10/100 Mbps rate; latter called "fast ethernet"

□ T stands for Twisted Pair

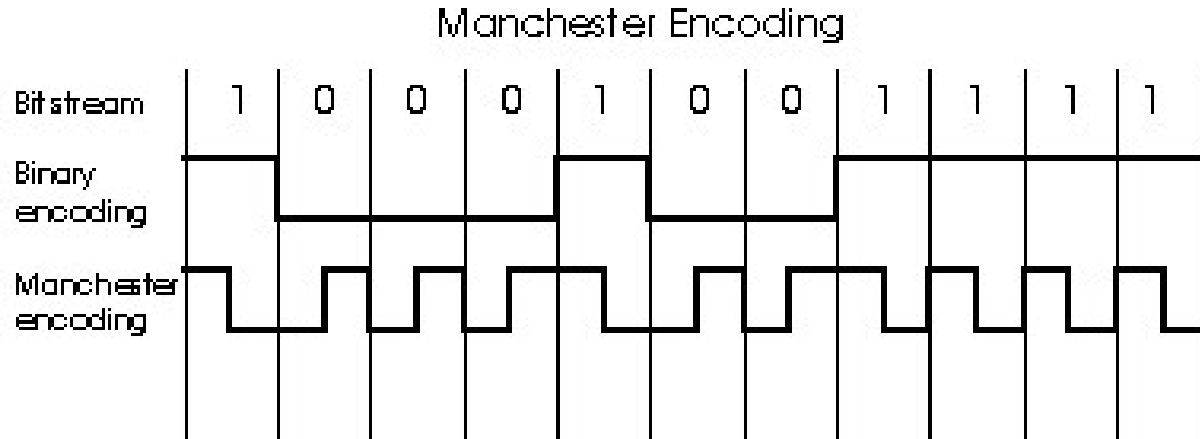□ Nodes connect to a hub: "star topology"; 100 m max distance between nodes and hub

twisted pair

hub

# Hubs

Hubs are essentially physical-layer repeaters:
- ❍ bits coming from one link go out all other links
- ❍ at the same rate
- ❍ no frame buffering
- ❍ no CSMA/CD at hub: adapters detect collisions
- ❍ provides net management functionality

twisted pair

hub

# Manchester encoding

Manchester Encoding

| Bit stream | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 |

Binary encoding

Manchester encoding

□ Used in 10BaseT

□ Each bit has a transition

□ Allows clocks in sending and receiving nodes to synchronize to each other
  ○ no need for a centralized, global clock among nodes!

□ Hey, this is physical-layer stuff!

# Gbit Ethernet
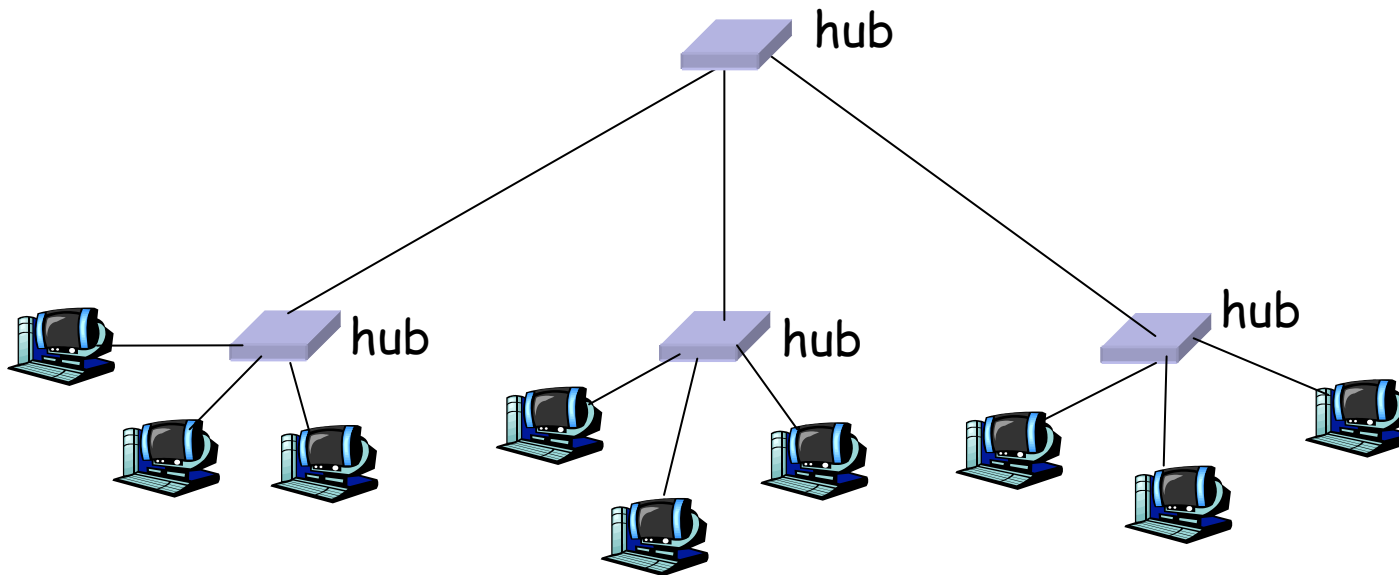
❑ uses standard Ethernet frame format

❑ allows for point-to-point links and shared broadcast channels

❑ in shared mode, CSMA/CD is used; short distances between nodes required for efficiency

❑ uses hubs, called here "Buffered Distributors"

❑ Full-Duplex at 1 Gbps for point-to-point links

❑ 10 Gbps now !

# Link Layer

- 5.1 Introduction and services
- 5.2 Error detection and correction
- 5.3 Multiple access protocols
- 5.4 Link-Layer Addressing
- 5.5 Ethernet

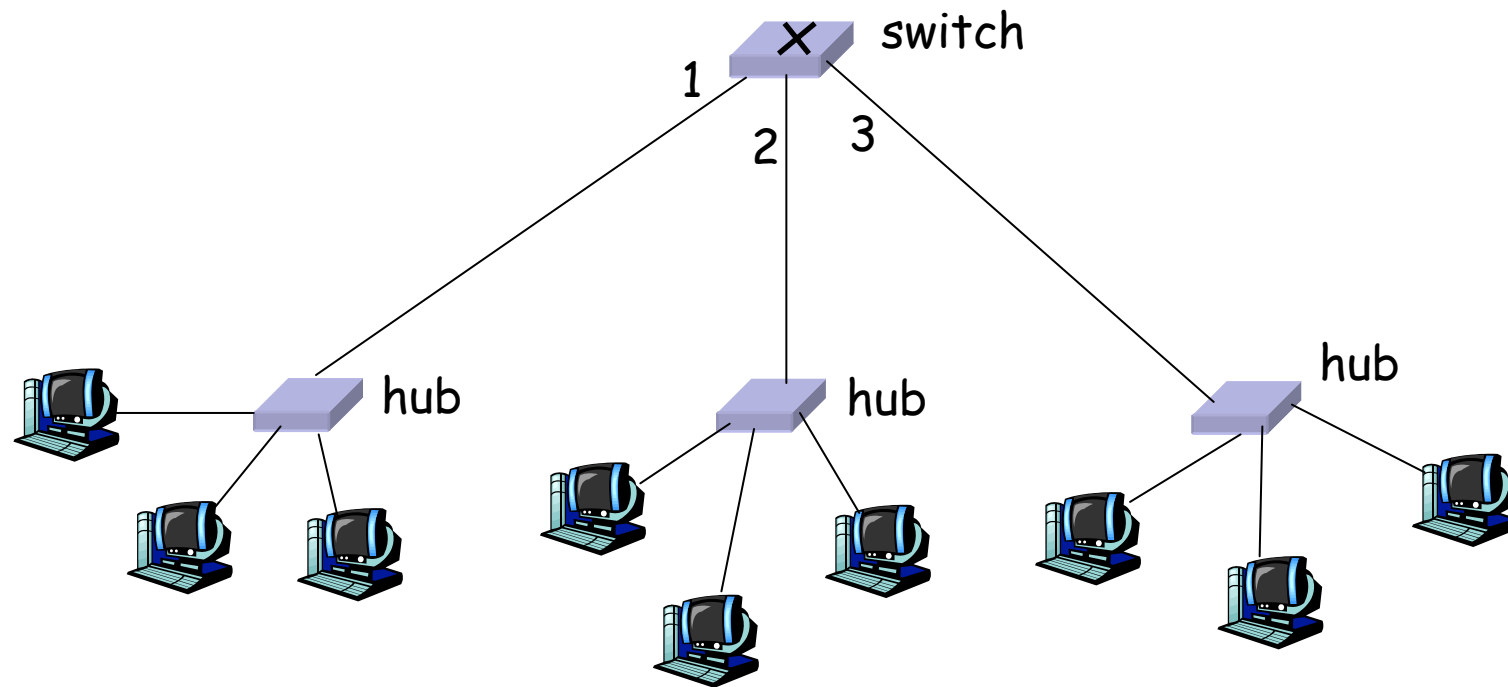- 5.6 Interconnections: Hubs and switches

# Interconnecting with hubs

□ Backbone hub interconnects LAN segments

□ Extends max distance between nodes

□ But individual segment collision domains become one large collision domain

□ Can't interconnect 10BaseT & 100BaseT

# Switch

□ Link layer device
  ○ stores and forwards Ethernet frames
  ○ examines frame header and selectively forwards frame based on MAC dest address
  ○ 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

# Forwarding



- How do determine onto which LAN segment to forward frame?
- Looks like a routing problem...

# Self learning

□ A switch has a <span style="color:red">switch table</span>

□ entry in switch table:
  ○ (MAC Address, Interface, Time Stamp)
  ○ stale entries in table dropped (TTL can be 60 min)

□ 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

# Filtering/Forwarding

**When switch receives a frame:**

index switch table using MAC dest address
**if** entry found for destination
   **then**{
      **if** dest on segment from which frame arrived
        **then** drop the frame
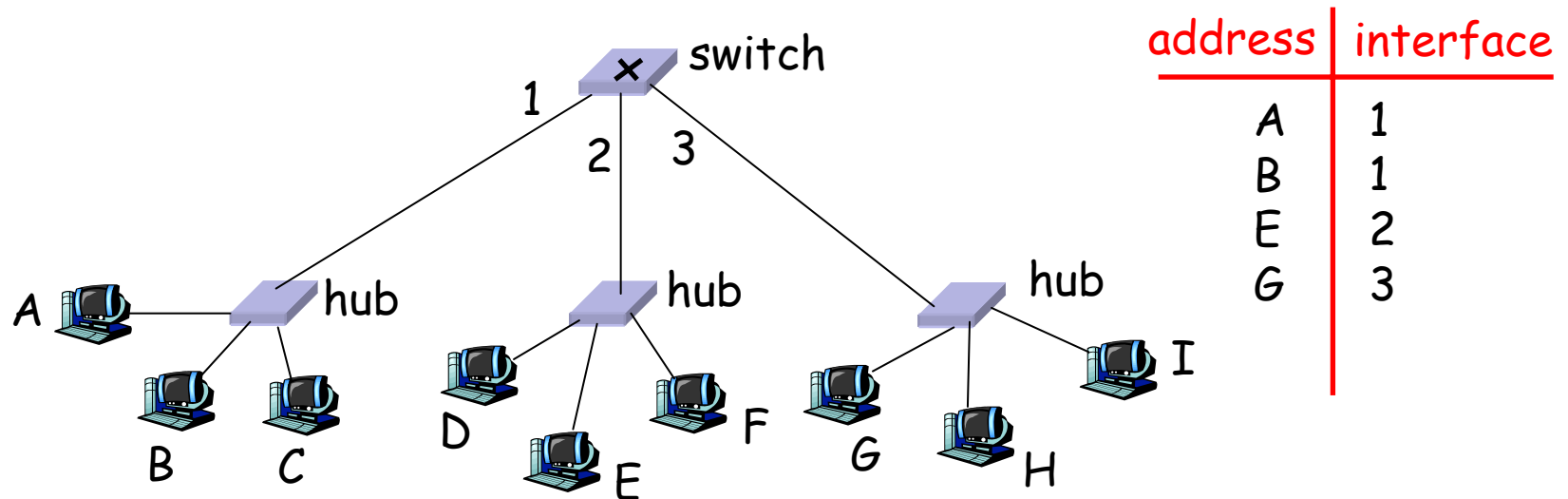        **else** forward the frame on interface indicated
    }
  **else** flood

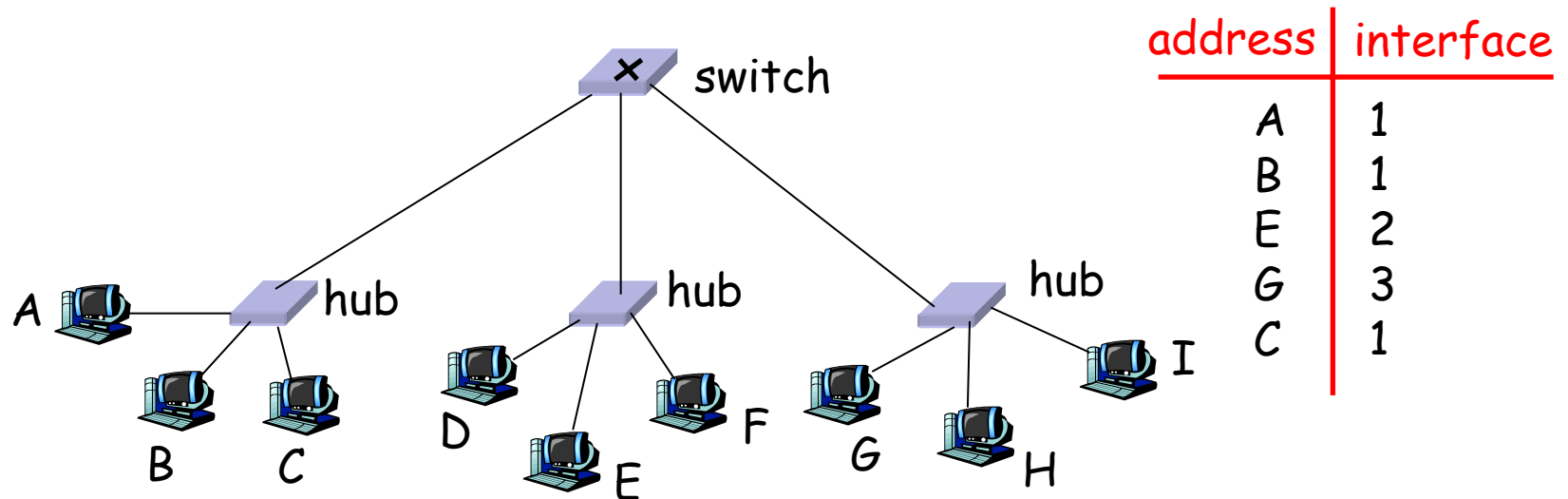*forward on all but the interface on which the frame arrived*

# Switch example

Suppose C sends frame to D



| address | interface |
|---------|-----------|
| A | 1 |
| B | 1 |
| E | 2 |
| G | 3 |

- ❒ Switch receives frame from C
  - ○ notes in bridge table that C is on interface 1
  - ○ because D is not in table, switch forwards frame into interfaces 2 and 3
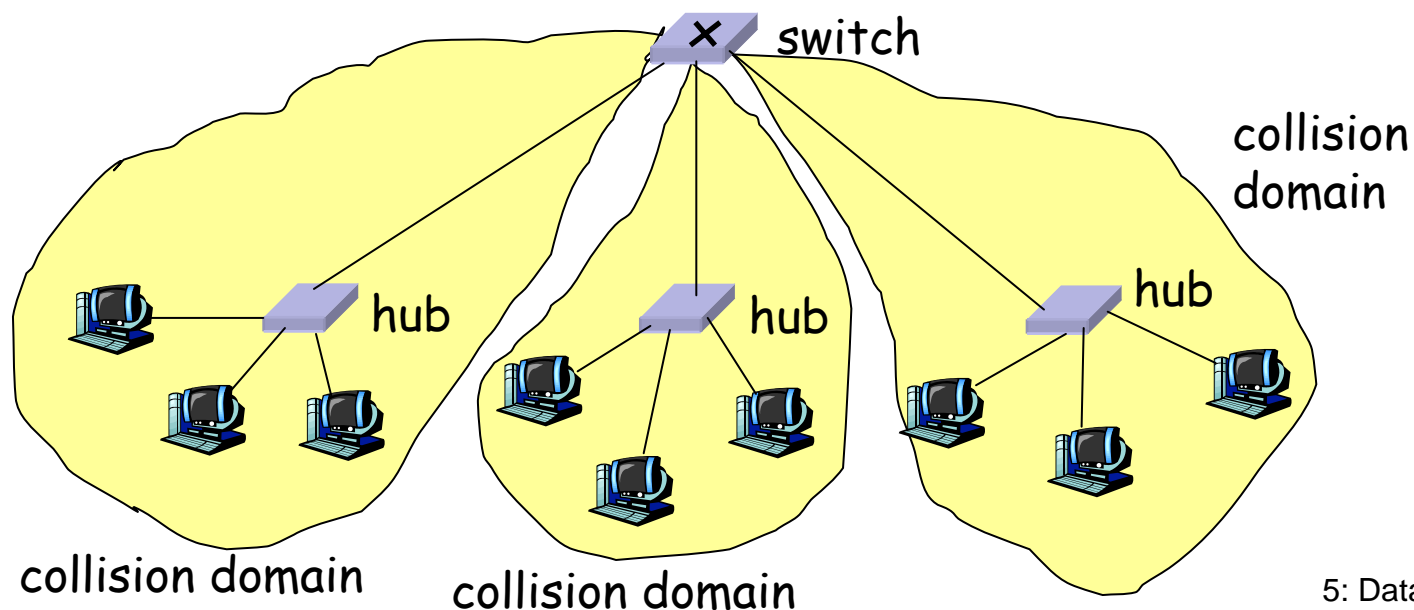- ❒ frame received by D

# Switch example

Suppose D replies back with frame to C.



| address | interface |
|---------|-----------|
| A | 1 |
| B | 1 |
| E | 2 |
| G | 3 |
| C | 1 |

□ Switch receives frame from from D
  ○ notes in bridge table that D is on interface 2
  ○ because C is in table, switch forwards frame only to interface 1

□ frame received by C
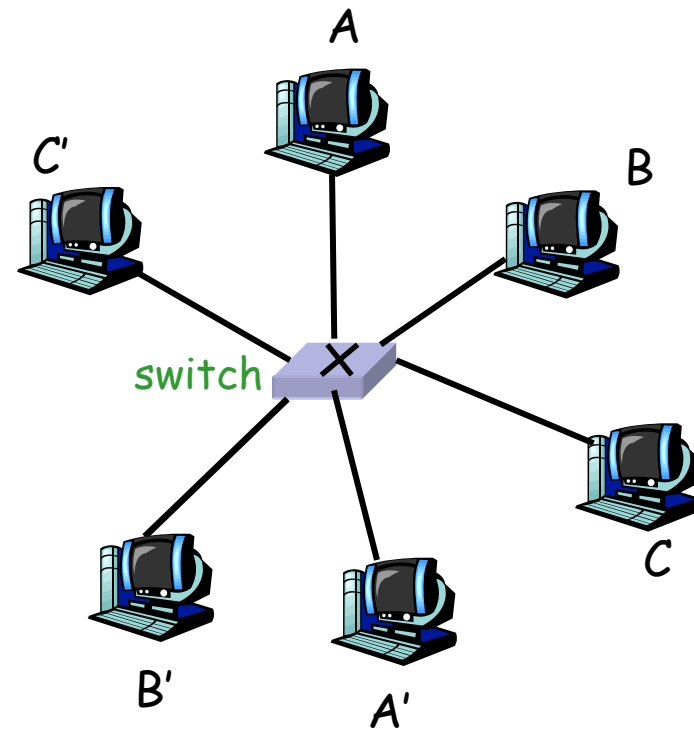
# Switch: traffic isolation

□ switch installation breaks subnet into LAN segments

□ switch filters packets:

  ○ same-LAN-segment frames not usually forwarded onto other LAN segments

  ○ segments become separate collision domains



switch

collision domain

hub

hub

hub

collision domain

collision domain

# Switches: dedicated access

- Switch with many interfaces
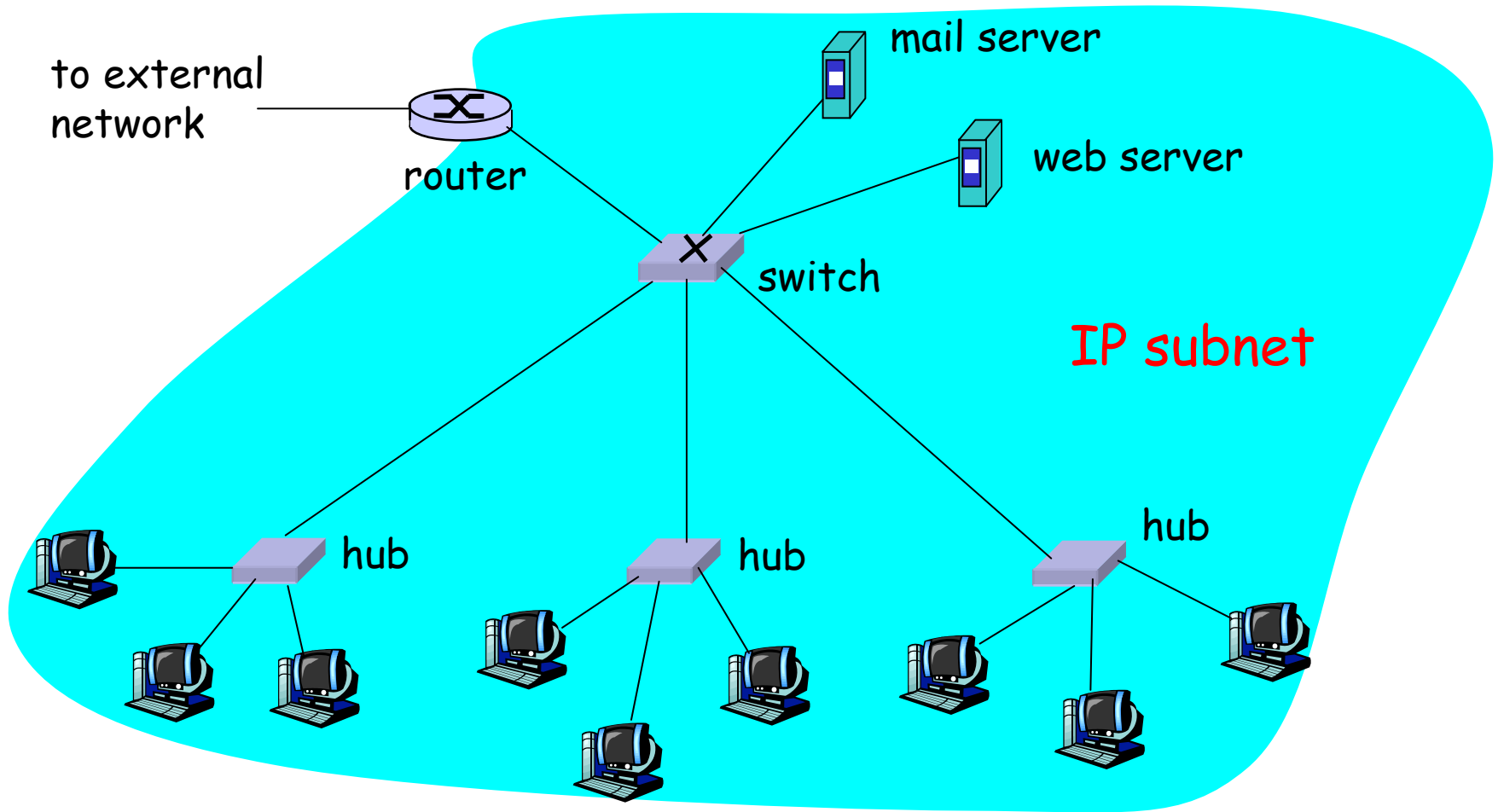- Hosts have direct connection to switch
- No collisions; full duplex

**Switching:** A-to-A' and B-to-B' simultaneously, no collisions

switch

A

C'

B

C

B'

A'

# More on Switches
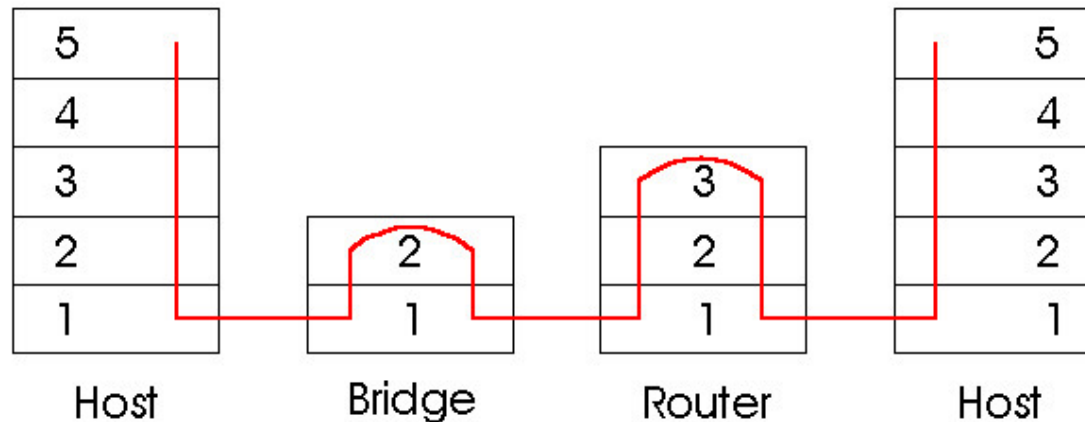
□ cut-through switching: frame forwarded from input to output port without first collecting entire frame

  ○ slight reduction in latency

□ combinations of shared/dedicated, 10/100/1000 Mbps interfaces

# Institutional network



to external network

router

mail server

web server

switch

IP subnet

hub

hub

hub

# Switches vs. Routers

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

# Summary comparison

|                  | hubs | routers | switches |
|------------------|------|---------|----------|
| traffic isolation | no   | yes     | yes      |
| plug & play       | yes  | no      | yes      |
| optimal routing   | no   | yes     | no       |
| cut through       | yes  | no      | yes      |

# 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