# Computer Networking

谢 逸

中山大学·计算机学院

**2023. Fall**

---

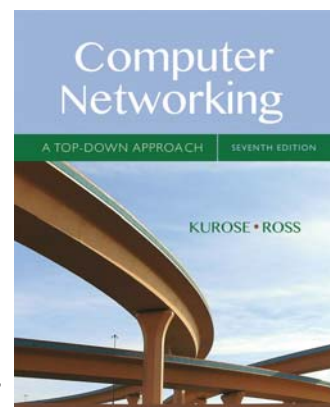# Chapter 6
# The Link Layer
# and LANs

A note on the use of these Powerpoint slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you see the animations; and can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- If you use these slides (e.g., in a class) that you mention their source (after all, we'd like people to use our book!)
- If you post any slides on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy! JFK/KWR

*Computer Networking: A Top Down Approach*

7th edition
Jim Kurose, Keith Ross
Pearson/Addison Wesley
April 2016

# Homework

- Ch6 (ver7), 6, 11, 13, 14, 15, 17, 18, 21, 23, 29, 31,32

- Keywords: VLAN, CSMA/CD/CA, CRC, TDM, CDM, ALOHA, ARP, self-learning, ATM, MPLS, PPP, Ethernet, MAC protocols

# Chapter 6: 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

**6.1 introduction, services**

**6.2** error detection, correction

**6.3** multiple access protocols

**6.4 LANs**

- addressing, ARP
- Ethernet
- switches
- VLANS

**6.5** link virtualization: MPLS

**6.6** data center networking
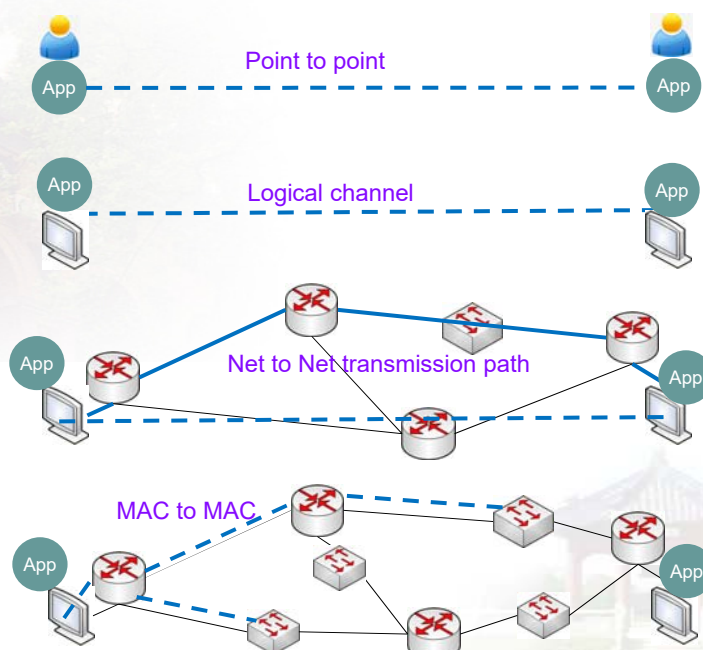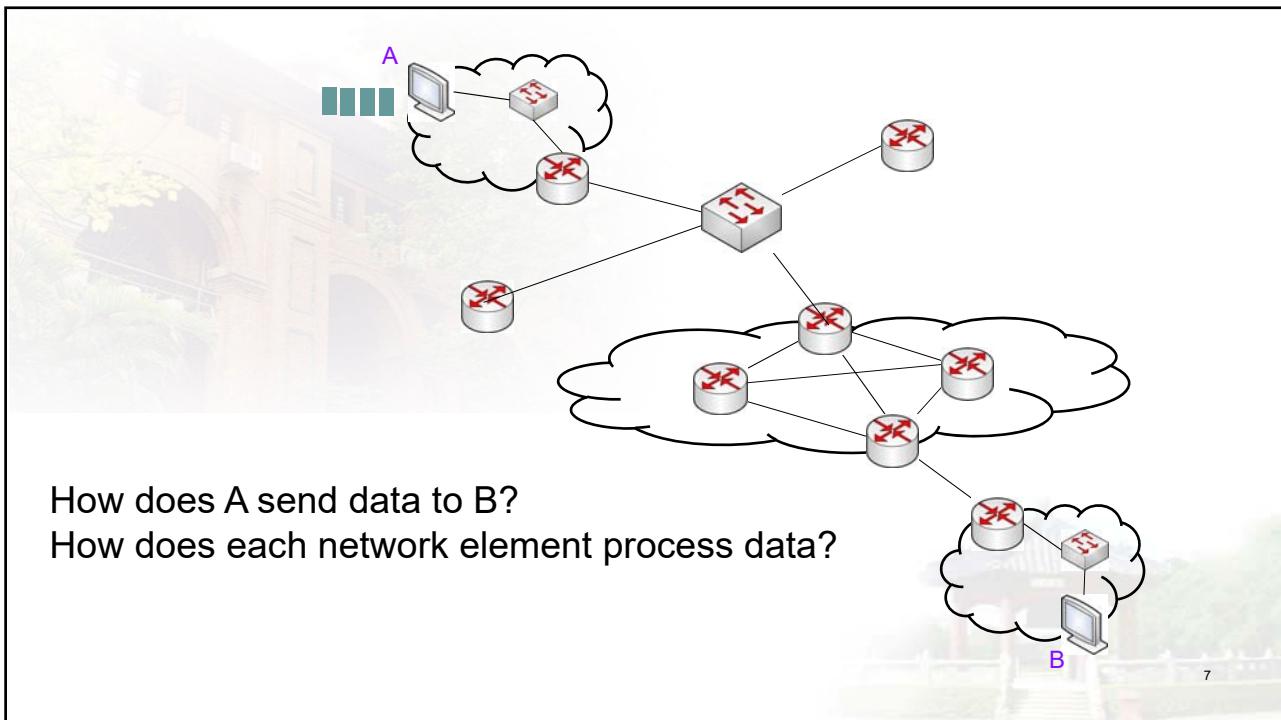
**6.7** a day in the life of a web request

5-5

---

应用层网络：　App - - - - - Point to point - - - - - App

传输层网络：　App - - - - - Logical channel - - - - - App

网络层网络：　App　　Net to Net transmission path　App

链路层网络：　MAC to MAC　App

How does A send data to B?
How does each network element process data?
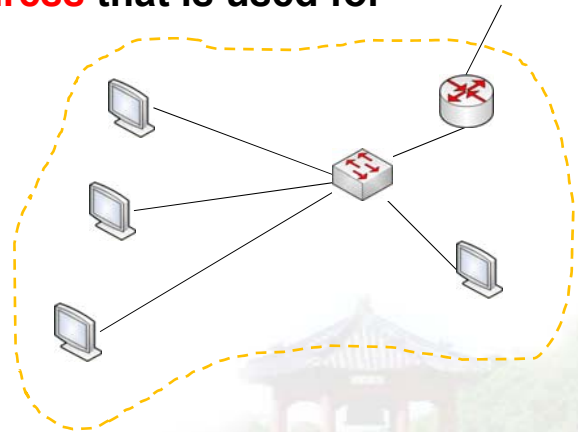
# Link layer: introduction

- **Each network interface card has an address named physical address**
  - **TV-net, telephone-net, data-net, …**
  - **Physical address is different from IP address**
  - **Physical addresses do NOT have a uniform format and standard. It may be 12bits, 48bits, or others.**

# Link layer: introduction

- **In a same subnet**, the nodes are considered to have **a uniform standard physical address** that is used for addressing.

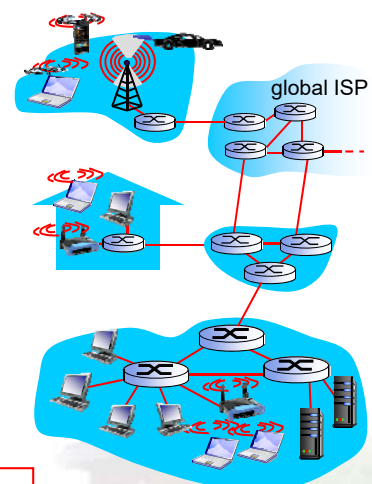**MAC address is used for addressing in the same subnet**



9

# Link layer: introduction

*terminology:*

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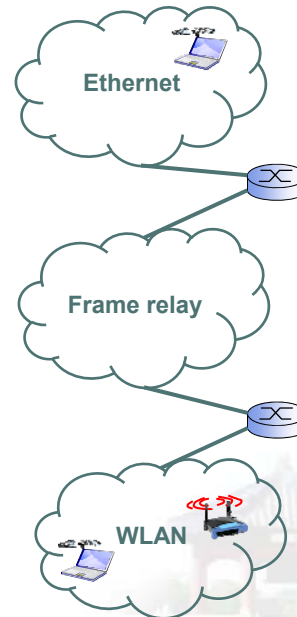- layer-2 packet: **frame,** encapsulates datagram



global ISP

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

5-10

# 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**

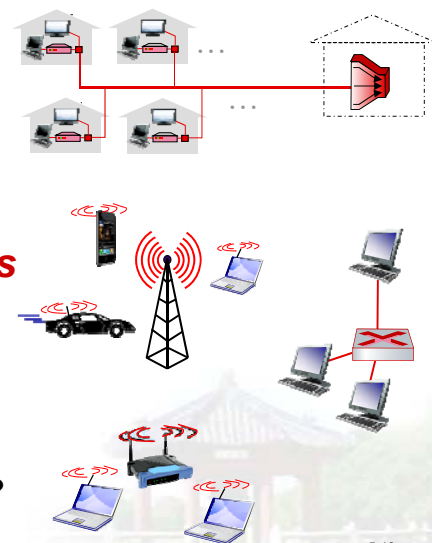  ▪ e.g., may or may not provide rdt over link

Ethernet

Frame relay

WLAN

5-11

# Link layer services

| IP | MAC Hdr |

● *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?

5-12

# 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

  **Why no congestion control?**

5-13

# Where is the link layer implemented?

- in each and every host
- link layer implemented in "adaptor" (aka *network interface card* NIC) or on a chip
  - Ethernet card, 802.11 card; Ethernet chipset
  - implements link, physical layer
- attaches into host's system buses
- combination of hardware, software, firmware

application
transport
network
link

cpu        memory

link
physical

controller

physical
transmission

host
bus
(e.g., PCI)

network adapter
card

5-14

## Interfaces communicating



sending side:
- encapsulates datagram in frame
- adds error checking bits, reliable data transfer, flow control, etc.

receiving side:
- looks for errors, reliable data transfer, flow control, etc.
- extracts datagram, passes to upper layer at receiving side

Link Layer: 6-15

---



目标：**C1→C2**
- **C1**获取**C2**的**IP**，填入IP报头
- **C1**根据$IP_{C2}$查路由表得到下一跳路由器**R1**的地址$IP_{R1}$
- **C1**广播查询$IP_{R1}$对应的物理地址$\textbf{MAC}_{R1}$
- C1把$\textbf{MAC}_{R1}$填入链路帧头；
- C1把链路帧发给**S1**；
- **S1**根据帧头的$\textbf{MAC}_{R1}$转发给**R1**；
- **R1**抽取IP，根据$IP_{C2}$查路由表，决定下一跳转发给$IP_{R2}$；
- **R1**广播查询$IP_{R2}$对应的物理地址$\textbf{MAC}_{R2}$
- **R1**把$\textbf{MAC}_{R2}$填入帧头，转发给**R2**
- **R2** ？

16

# Link layer, LANs: outline

**6.1** introduction, services
**6.2** error detection, correction
**6.3** multiple access protocols
**6.4** LANs
  - addressing, ARP
  - Ethernet
  - switches
  - VLANS

**6.5** link virtualization: MPLS
**6.6** data center networking
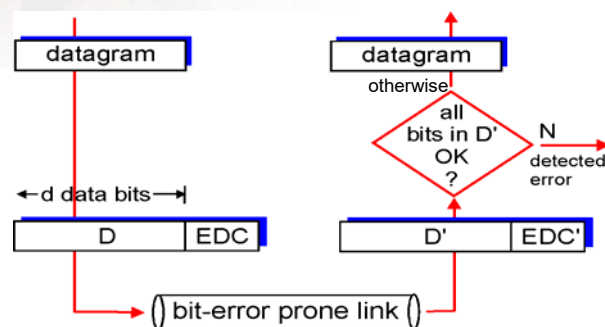**6.7** a day in the life of a web request

5-17

# 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



5-18

# Parity checking

*single bit parity:*
❖ detect single bit errors

*two-dimensional bit parity:*
❖ detect and correct single bit errors



5-19

# 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?***

5-20

# Cyclic redundancy check

❖ **more powerful error-detection coding**
❖ **view data bits, D, as a binary number**   $R=f(D)$
❖ **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 (Ethernet, 802.11 WiFi, ATM)**

```
←———— d bits ————→←— r bits —→
┌─────────────────────┬──────────────┐
│ D data bits to be sent │ R CRC bits │
└─────────────────────┴──────────────┘
```
*bit pattern*

$$D * 2^r \ \ XOR \ \ R$$
*mathematical formula*

5-21

---

# Cyclic redundancy check

● **generator**

| $2^6$ | $2^5$ | $2^4$ | $2^3$ | $2^2$ | $2^1$ | $2^0$ |
|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 1 | 1 | 1 |

$$\Rightarrow x^6 \cdot 1 + x^5 \cdot 0 + x^4 \cdot 1 + x^3 \cdot 0 + x^2 \cdot 1 + x^1 \cdot 1 + x^0 \cdot 1$$
$$= x^6 + x^4 + x^2 + x + 1$$
$$= G(x)$$

$$G(x) = x^4 + x^3 + x + 1$$
$$\Rightarrow ?$$

22

Here:



# 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:

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

**Sender:**

余数的位数一定要是比除数位数只能少一位，哪怕前面位是0，甚至是全为0（整除时）也都不能省略

---

# Link layer, LANs: outline

**5.1** introduction, services

**5.2** error detection, correction

**5.3 multiple access protocols**

**5.4** LANs
- addressing, ARP
- Ethernet
- switches
- VLANS

**5.5** link virtualization: MPLS

**5.6** data center networking

**5.7** a day in the life of a web request

# Multiple access links, protocols

two types of "links" :

- **point-to-point**
  - PPP for dial-up access
  - point-to-point link between **Ethernet switch**, host

- *broadcast (shared wire or medium)*
  - old-fashioned Ethernet
  - upstream HFC
  - 802.11 wireless LAN

shared wire (e.g., cabled Ethernet)

shared RF (e.g., 802.11 WiFi)

shared RF (satellite)

humans at a cocktail party (shared air, acoustical)

5-25

# 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

5-26

# An ideal multiple access protocol

*given:* broadcast channel of rate R bps

*desiderata:*

1. when one node wants to transmit, it can send at **full 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

5-27

# 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

5-28

# 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**



5-29

# 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**



5-30

# Random access protocols

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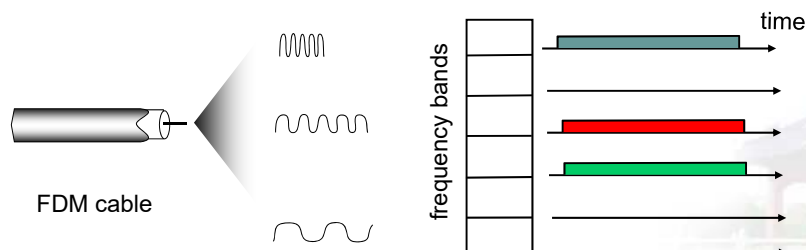- **two or more transmitting nodes ➜ "collision",**
- **random access MAC protocol specifies:**
  - **how to detect collisions**
  - **how to recover from collisions (e.g., via delayed retransmissions)**
- **examples of random access MAC protocols:**
  - **ALOHA**
  - **slotted ALOHA**
  - **CSMA, CSMA/CD,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

5-31

# 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 →

**Collides may
happen at any time**

node i frame

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

5-32

# Pure ALOHA efficiency

Pr[success by given node] = Pr[node transmits] ·

Pr[no other node transmits in $[t_0-1,t_0]$ ] ·

Pr[no other node transmits in $[t_0,t_0+1]$]

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

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

Pr[a successful slots] = $Np(1-p)^{2N-1}$

… choosing optimum p and then letting N $\longrightarrow \infty$

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

*worse* than the following slotted Aloha!

# Slotted ALOHA

*assumptions:*

❖ all frames <u>same size</u>
❖ time divided into <u>equal size slots</u> (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

**Collides only happen at the beginning**

node 1 | 1 | 1 | 1 | 1
node 2 | 2 | 2 | 2
node 3 | 3 | 3 | 3

C E C S E C E S S

*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

5-35

---

# 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 of a **successful slots** $= Np(1-p)^{N-1}$

- 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:

*max efficiency = 1/e = .37*

*at best:* channel used for useful transmissions 37% of time!

**!**

5-36

# 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!**

5-37

# CSMA collisions

- **collisions *can* still occur:** propagation delay means two nodes may not hear each other's transmission
- **Nodes don't detect collision during transmission**
- **collision: entire packet transmission time wasted**
  - **distance & propagation delay play role in in determining collision probability**



spatial layout of nodes

space

time

$t_0$

$t_1$

5-38

# 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**

5-39

# CSMA/CD (collision detection)



spatial layout of nodes

5-40

# 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 jam signal**

5. **After aborting, NIC enters *binary (exponential) backoff*:**
   - **after $m^{th}$ collision, NIC chooses $K$ at random from *{0,1,2, …, $2^m$- 1}*. NIC waits K·512 bit times, returns to Step 2**
   - **longer backoff interval with more collisions**

5-41

# CSMA/CD efficiency

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

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

- **efficiency goes to 1 (i.e., NO collision)**
  - **as $t_{prop}$ goes to 0**   **→ No time difference**
  - **as $t_{trans}$ goes to infinity**   **→ Always occupying the channel**
- **better performance than ALOHA: and simple, cheap, decentralized!**

5-42

# "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!

5-43

# "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)

data

poll

master

data

slaves

5-44

# "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)

(nothing to send)

T

data

5-45

# A Case: Cable access network

Internet frames, TV channels, control transmitted downstream at **different frequencies**

cable headend

CMTS

cable modem termination system

ISP

splitter

cable modem

upstream Internet frames, TV control, transmitted upstream **at different frequencies in time slots**

- ❖ multiple 40Mbps downstream (broadcast) channels
  - ▪ single CMTS transmits into channels
- ❖ multiple 30 Mbps upstream channels
  - ▪ multiple access: all users contend for certain upstream channel time slots (others assigned)

46

# Cable access network



cable headend

CMTS

MAP frame for Interval [t1, t2]

Downstream channel i

Upstream channel j

$t_1$   $t_2$

Minislots containing minislots request frames

Assigned minislots containing cable modem upstream data frames

Residences with cable modems

DOCSIS: data over cable service interface spec

❖ FDM over upstream, downstream frequency channels
❖ TDM upstream: some slots assigned, some have contention
  ▪ downstream MAP frame: assigns upstream slots
  ▪ request for upstream slots (and data) transmitted random access (binary backoff) in selected slots

5-47

---

# TDM+FDM+CDM

## CDMA
## Code Division Multiple Access



Code

Frequency

Time

48

# 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,  token ring

A question: why do we use CSMA/CD rather than TDMA/FDMA?

5-49

# Link layer, LANs: outline

6.1 introduction, services

6.2 error detection, correction

6.3 multiple access protocols

6.4 LANs
- addressing, ARP
- Ethernet
- switches
- VLANS

6.5 link virtualization: MPLS

6.6 data center networking

6.7 a day in the life of a web request

5-50

# MAC addresses and ARP

- **32-bit IP address:**
  - *network-layer* address for interface
  - used for layer 3 (network layer) forwarding
- **MAC (or LAN or physical or Ethernet) address:**
  - **function:** *used 'locally" to 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
(each "number" represents 4 bits)

5-51

---

# LAN addresses and ARP
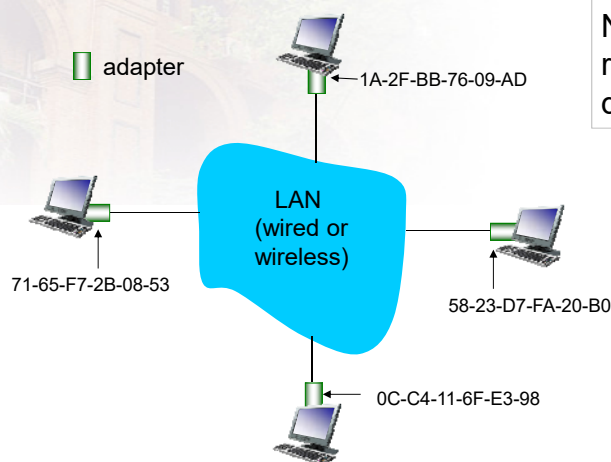
each adapter on LAN has unique *LAN* address

NIC decides the frame is received or not by its destination MAC address.

adapter

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

LAN
(wired or wireless)

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

58-23-D7-FA-20-B0

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

5-52

# LAN addresses (more)

- **MAC address allocation administered by IEEE**
- **manufacturer buys portion of MAC address space (to assure uniqueness)**
- **analogy:**
  - **MAC address: like identity card number (fixed)**
  - **IP address: like postal address (changeable)**
- **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**

| IP | MAC |
|------|------|
| 逻辑地址 | 物理地址 |
| 全局有效 | 局部有效 |
| 层次化 | 平面型 |
| 网间寻址 | 网内寻址 |

5-53

# ARP: address resolution protocol

Question: how to determine interface's MAC address, knowing its IP address?



137.196.7.78
1A-2F-BB-76-09-AD
137.196.7.23
137.196.7.14
LAN
71-65-F7-2B-08-53
58-23-D7-FA-20-B0
0C-C4-11-6F-E3-98
137.196.7.88

*ARP table:* **each IP node (host, router) on LAN has 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)**

5-54

# 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 nodes 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*

5-55

## ARP protocol in action

example: A wants to send datagram to B

- B's MAC address not in A's ARP table, so A uses ARP to find B's MAC address

A broadcasts ARP query, containing B's IP addr

① • destination MAC address = FF-FF-FF-FF-FF-FF
  • all nodes on LAN receive ARP query



Ethernet frame (sent to FF-FF-FF-FF-FF-FF)
Source MAC: 71-65-F7-2B-08-53
Source IP: 137.196.7.23
Target IP address: 137.196.7.14
...

ARP table in A
| IP addr | MAC addr | TTL |
|---------|----------|-----|
|         |          |     |
|         |          |     |

71-65-F7-2B-08-53
137.196.7.23

58-23-D7-FA-20-B0
137.196.7.14

Link Layer: 6-56

# ARP protocol in action

## example: A wants to send datagram to B

- B's MAC address not in A's ARP table, so A uses ARP to find B's MAC address

ARP message into Ethernet frame
(sent to 71-65-F7-2B-08-53)

Target IP address: 137.196.7.14
Target MAC address:
58-23-D7-FA-20-B0
...

ARP table in A

| IP addr | MAC addr | TTL |
|---------|----------|-----|
|         |          |     |

C

A
71-65-F7-2B-08-53
137.196.7.23

B
58-23-D7-FA-20-B0
137.196.7.14

②

D

② B replies to A with ARP response, giving its MAC address

Link Layer: 6-57

---

# ARP protocol in action

## example: A wants to send datagram to B

- B's MAC address not in A's ARP table, so A uses ARP to find B's MAC address

C

ARP table in A

| IP addr | MAC addr | TTL |
|---------|----------|-----|
| 137.196.7.14 | 58-23-D7-FA-20-B0 | 500 |

A
71-65-F7-2B-08-53
137.196.7.23

B
58-23-D7-FA-20-B0
137.196.7.14

③ A receives B's reply, adds B entry into its local ARP table

D

Link Layer: 6-58

2023/11/19

# Addressing: routing to another LAN

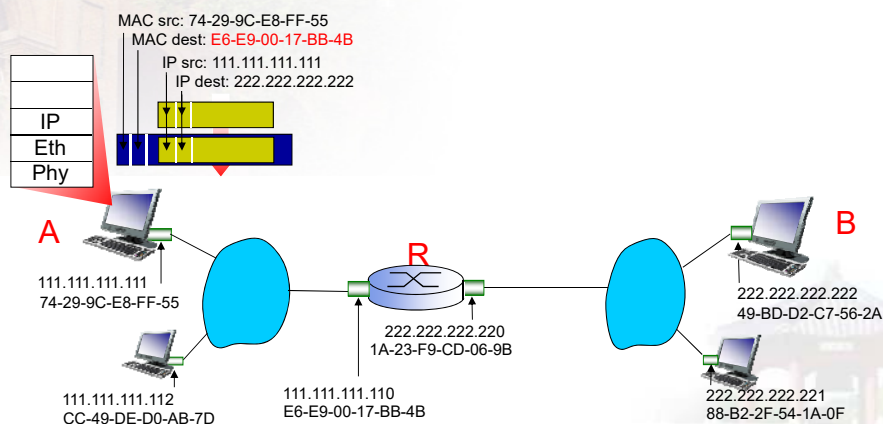**walkthrough: send datagram from A to B via R**

- **focus on addressing – at 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
111.111.111.111
74-29-9C-E8-FF-55

111.111.111.112
CC-49-DE-D0-AB-7D

R
222.222.222.220
1A-23-F9-CD-06-9B

111.111.111.110
E6-E9-00-17-BB-4B

B
222.222.222.222
49-BD-D2-C7-56-2A

222.222.222.221
88-B2-2F-54-1A-0F

5-59

---

# Addressing: routing to another LAN

- ❖ **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**

MAC src: 74-29-9C-E8-FF-55
MAC dest: E6-E9-00-17-BB-4B
IP src: 111.111.111.111
IP dest: 222.222.222.222

IP
Eth
Phy

A
111.111.111.111
74-29-9C-E8-FF-55

111.111.111.112
CC-49-DE-D0-AB-7D

R
222.222.222.220
1A-23-F9-CD-06-9B

111.111.111.110
E6-E9-00-17-BB-4B

B
222.222.222.222
49-BD-D2-C7-56-2A

222.222.222.221
88-B2-2F-54-1A-0F

5-60

# Addressing: routing to another LAN

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

MAC src: 74-29-9C-E8-FF-55
MAC dest: E6-E9-00-17-BB-4B
IP src: 111.111.111.111
IP dest: 222.222.222.222

```
IP
Eth
Phy
```

```
IP
Eth
Phy
```

A

B

111.111.111.111
74-29-9C-E8-FF-55

111.111.111.112
CC-49-DE-D0-AB-7D

222.222.222.220
1A-23-F9-CD-06-9B

111.111.111.110
E6-E9-00-17-BB-4B

222.222.222.222
49-BD-D2-C7-56-2A

222.222.222.221
88-B2-2F-54-1A-0F

5-61

# Addressing: routing to another LAN

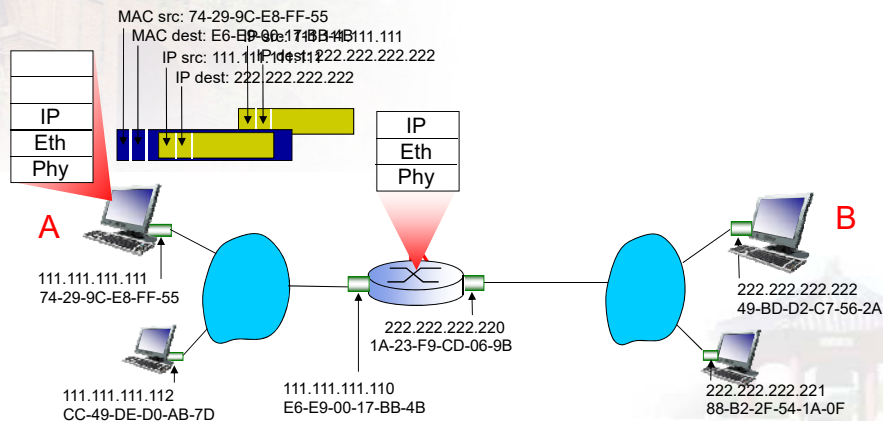❖ R forwards datagram with IP source A, destination B
❖ R creates link-layer frame with B's MAC address as dest, frame contains A-to-B IP datagram

MAC src: 1A-23-F9-CD-06-9B
MAC dest: 49-BD-D2-C7-56-2A
IP src: 111.111.111.111
IP dest: 222.222.222.222

```
IP
Eth
Phy
```

```
IP
Eth
Phy
```

A

B

111.111.111.111
74-29-9C-E8-FF-55

111.111.111.112
CC-49-DE-D0-AB-7D

222.222.222.220
1A-23-F9-CD-06-9B

111.111.111.110
E6-E9-00-17-BB-4B

222.222.222.222
49-BD-D2-C7-56-2A

222.222.222.221
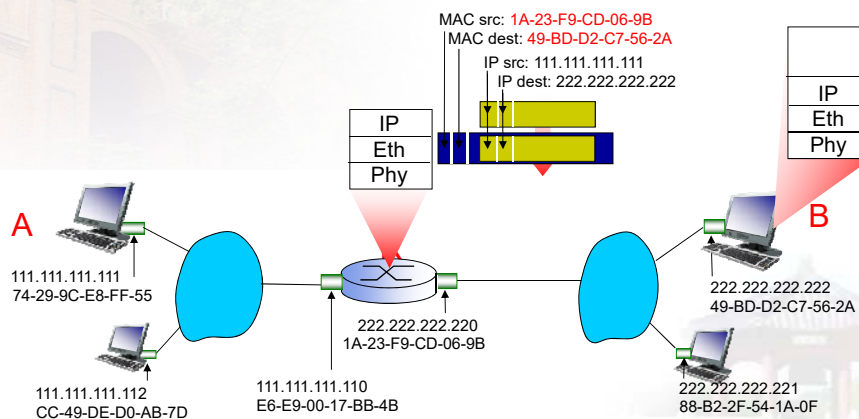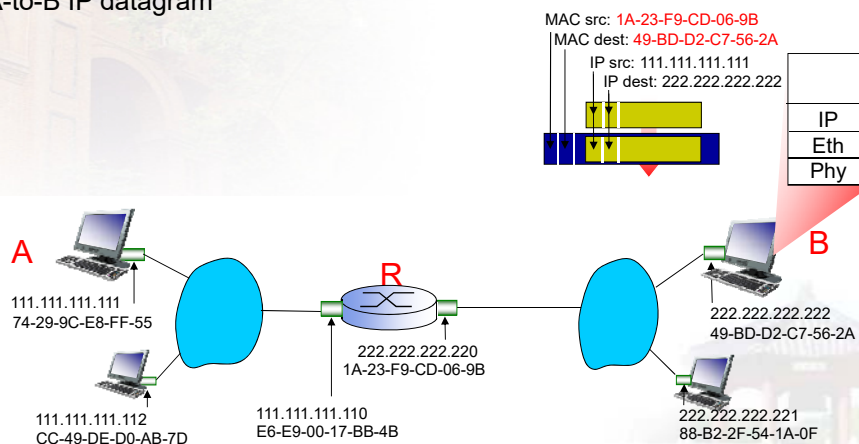88-B2-2F-54-1A-0F

5-62

# Addressing: routing to another LAN

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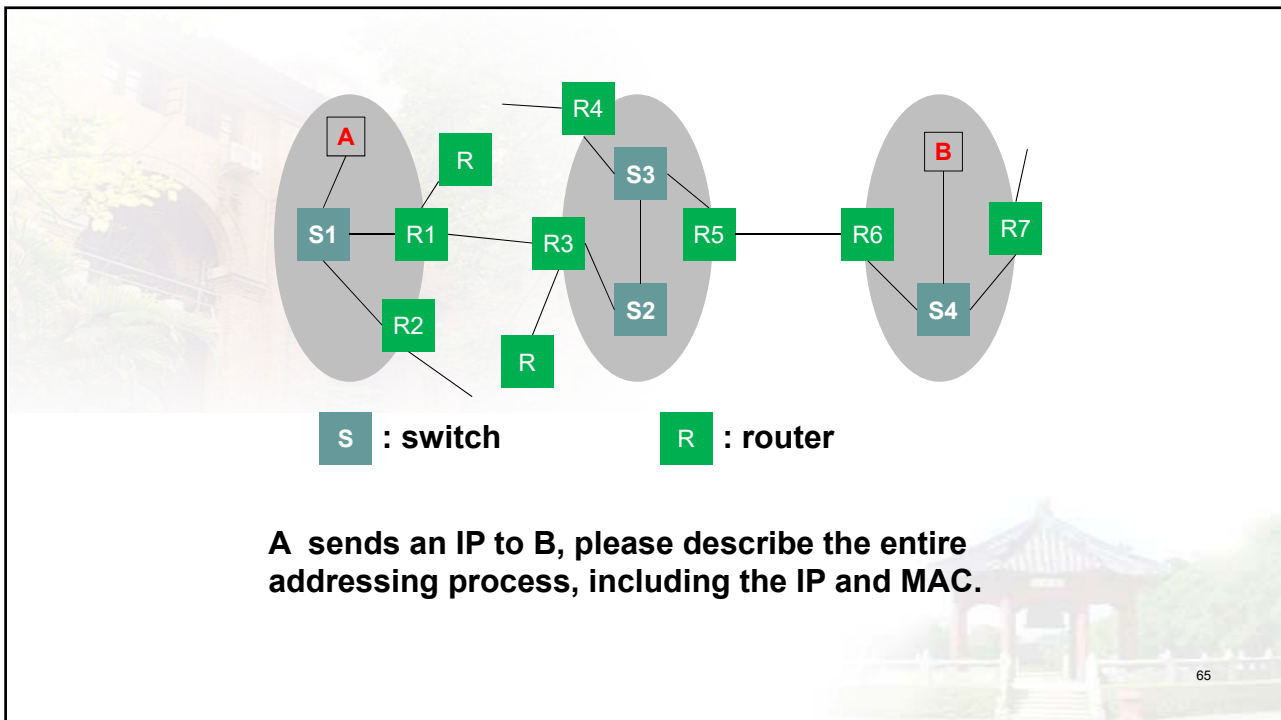❖ R creates link-layer frame with B's MAC address as dest, frame contains A-to-B IP datagram

MAC src: 1A-23-F9-CD-06-9B
MAC dest: 49-BD-D2-C7-56-2A
IP src: 111.111.111.111
IP dest: 222.222.222.222

IP
Eth
Phy

IP
Eth
Phy

A
111.111.111.111
74-29-9C-E8-FF-55

111.111.111.112
CC-49-DE-D0-AB-7D

222.222.222.220
1A-23-F9-CD-06-9B

111.111.111.110
E6-E9-00-17-BB-4B

B
222.222.222.222
49-BD-D2-C7-56-2A

222.222.222.221
88-B2-2F-54-1A-0F

5-63

---

**A sends an IP to B, please describe the entire addressing process, including the IP and MAC.**

65

---

# Link layer, LANs: outline

**6.1** introduction, services

**6.2** error detection, correction

**6.3** multiple access protocols

**6.4 LANs**

- addressing, ARP
- Ethernet
- switches
- VLANS

**6.5** link virtualization: MPLS
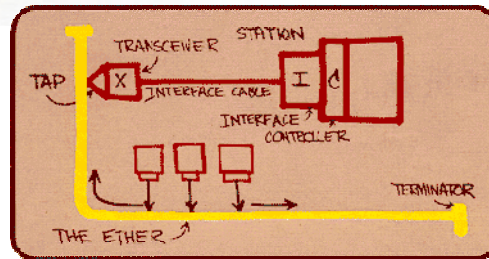
**6.6** data center networking

**6.7** a day in the life of a web request

5-66

# Ethernet

"dominant" wired LAN technology:

- **cheap $20 for NIC**
- **first widely used LAN technology**
- **simpler, cheaper than token LANs and ATM**
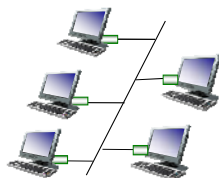- **kept up with speed race: 10 Mbps – 10 Gbps**
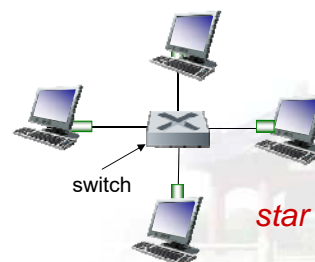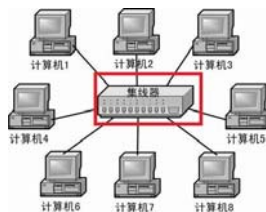
Metcalfe's Ethernet sketch

5-67

# Ethernet: physical topology

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

*bus:* coaxial cable

switch

*star*

5-68

# 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**

5-69

# Ethernet frame structure (more)

❖ *addresses:* **6 byte source, destination MAC addresses**

- **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:* **cyclic redundancy check at receiver**

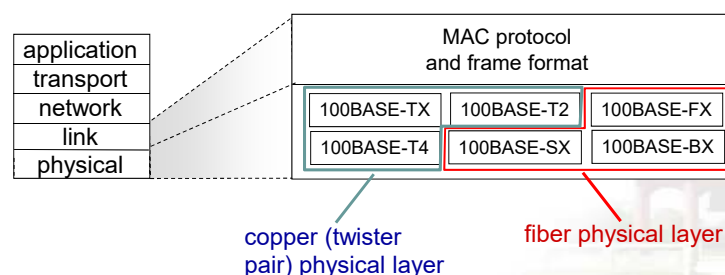- **error detected: frame is dropped**



5-70

2023/11/19

# Ethernet: unreliable, connectionless

- *connectionless:* no handshaking between sending and receiving NICs
- *unreliable:* receiving NIC doesn't send **Ack**s or **Nacks** to sending NIC
  - data in dropped frames recovered only if initial sender uses higher layer rdt (e.g., TCP), otherwise dropped data lost
- Ethernet's MAC protocol: unslotted *CSMA/CD wth binary backoff*

5-71

---

# 802.3 Ethernet standards: link & physical layers

- *many* different Ethernet standards
  - common MAC protocol and frame format
  - different speeds: 2 Mbps, 10 Mbps, 100 Mbps, 1Gbps, 10G bps
  - different physical layer media: fiber, cable

| application |
| transport |
| network |
| link |
| physical |

MAC protocol
and frame format

| 100BASE-TX | 100BASE-T2 | 100BASE-FX |
| 100BASE-T4 | 100BASE-SX | 100BASE-BX |

copper (twister pair) physical layer

fiber physical layer

5-72

# Link layer, LANs: outline

**6.1** introduction, services

**6.2** error detection, correction

**6.3** multiple access protocols

**6.4 LANs**

- addressing, ARP
- Ethernet
- **switches**
- VLANS

**6.5** link virtualization: MPLS

**6.6** data center networking

**6.7** a day in the life of a web request

5-74

# 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 (No IP and MAC address)*
  - **hosts are unaware of presence of switches**
- *plug-and-play, self-learning*
  - **switches do not need to be configured**

5-75

# Switch: *multiple* simultaneous transmissions

- **hosts have dedicated, direct connection to switch**
- **switches buffer packets**
- **Ethernet protocol used on *each* incoming link, but no collisions; full duplex**
  - **each link is its own collision domain**
- *switching:* **A-to-A' and B-to-B' can transmit simultaneously, without collisions**

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

5-76

# Switch forwarding table
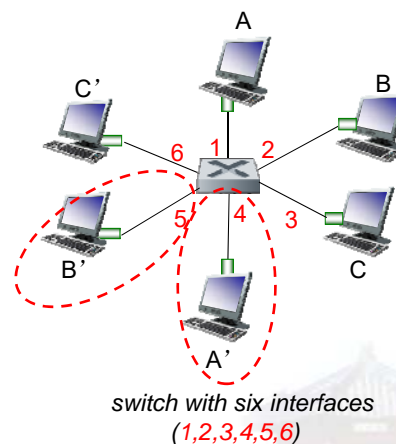
*Q:* **how does switch know A' reachable via interface 4, B' reachable via interface 5?**

- ❖ A: each switch has a switch table, each entry:
  - (MAC address of host, interface to reach host, time stamp)
  - looks like a routing table!

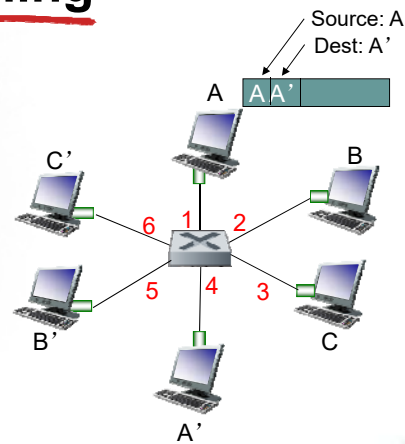Q: how are entries created, maintained in switch table?
  - something like a routing protocol?

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

5-77

# 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**



Source: A
Dest: A'

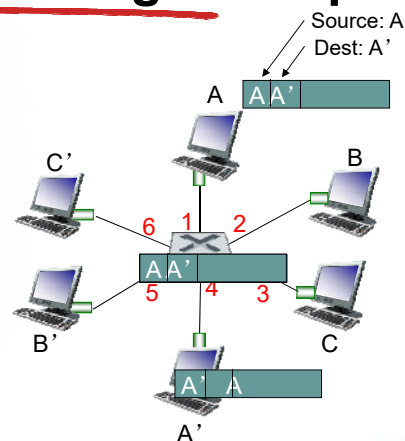| MAC addr | interface | TTL |
|----------|-----------|-----|
| A | 1 | *60* |
|  |  |  |

*Switch table (initially empty)*

5-78

# Self-learning, forwarding: example

- **frame destination, A', location unknown: *flood***

  ❖ destination A location known: selectively send on just one link



Source: A
Dest: A'

| MAC addr | interface | TTL |
|----------|-----------|-----|
| A | 1 | *60* |
| A' | 4 | *60* |

*switch table (initially empty)*

5-80

# Interconnecting switches

❖ **switches can be connected together**



Q: sending from **A to G** - how does $S_1$ know to forward frame destined to F via $S_4$ and $S_3$?

❖ A: self learning! (works exactly the same as in single-switch case!)

# Self-learning multi-switch example

**Suppose C sends frame to I, I responds to C**



❖ Q: show switch tables and packet forwarding in $S_1$, $S_2$, $S_3$, $S_4$

# Institutional network



to external network

router

mail server

web server

*IP subnet*

5-83

# Switches vs. routers

**both are store-and-forward:**

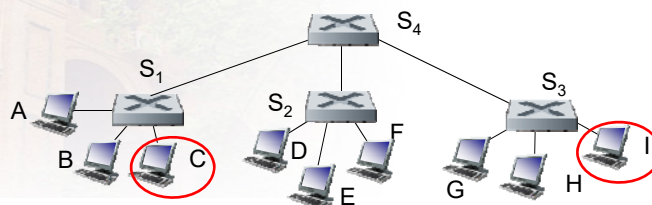▪*routers:* **network-layer devices (examine network-layer headers)**

▪*switches:* **link-layer devices (examine link-layer headers)**

**both have forwarding tables:**

▪*routers:* **compute tables using routing algorithms, IP addresses**

▪*switches:* **learn forwarding table using flooding, learning, MAC addresses**



application
transport
network
link
physical

datagram

frame

link    frame

physical

**switch**

network    datagram
link    frame
physical

application
transport
network
link
physical

5-84

# Link layer, LANs: outline

**6.1 introduction, services**

**6.2 error detection, correction**

**6.3 multiple access protocols**

**6.4 LANs**
- addressing, ARP
- Ethernet
- switches
- VLANS

**6.5 link virtualization: MPLS**

**6.6 data center networking**

**6.7 a day in the life of a web request**

5-85

# VLANs: motivation



Computer Science

Electrical Engineering

Computer Engineering

What is the difference between VLAN and the general LAN formed by router?

*consider:*

- **CS user moves office to EE, but wants connect to CS switch?**

- **single broadcast domain:**
- **all layer-2 broadcast traffic (ARP, DHCP, unknown location of destination MAC address) must cross entire LAN**
- **security/privacy, efficiency issues**

5-86

# VLANs

**port-based VLAN:** switch ports grouped (by switch management software) so that *single* physical switch ......

VLAN1:Port1~8
VLAN2:Port9~16

**Virtual Local Area Network**

switch(es) supporting VLAN capabilities can be configured to define multiple *virtual* LANS over single physical LAN infrastructure.

Electrical Engineering
(VLAN ports 1-8)

Computer Science
(VLAN ports 9-15)
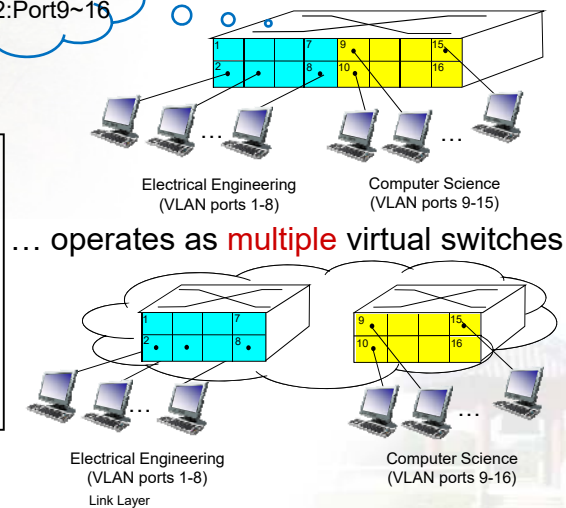
... operates as multiple virtual switches

Electrical Engineering
(VLAN ports 1-8)
Link Layer

Computer Science
(VLAN ports 9-16)

5-87

---

# Port-based VLAN

❖ *traffic isolation:* **frames to/from ports 1-8 can *only* reach ports 1-8**
  ▪ **can also define VLAN based on MAC addresses of endpoints, rather than switch port**

❖ dynamic membership: ports can be dynamically assigned among VLANs

❖ forwarding between VLANS: done via **routing** (just as with separate switches)
  ▪ in practice vendors sell combined switches plus routers

router

Electrical Engineering
(VLAN ports 1-8)

Computer Science
(VLAN ports 9-15)

5-88

# VLANS spanning multiple switches



Electrical Engineering
(VLAN ports 1-8)

Computer Science
(VLAN ports 9-15)

Ports 2,3,5 belong to EE VLAN
Ports 4,6,7,8 belong to CS VLAN

- *trunk port:* **carries frames between VLANS defined over multiple physical switches**
    - **frames forwarded within VLAN between switches can't be vanilla 802.1 frames (must carry VLAN ID info)**
    - **802.1q protocol adds/removed additional header fields for frames forwarded between trunk ports**

5-89

# 802.1Q VLAN frame format



802.1 frame

802.1Q frame

2-byte Tag Protocol Identifier
(value: 81-00)

Recomputed
CRC

Tag Control Information (**12 bit VLAN ID field**,
3 bit priority field like IP TOS)

router

Electrical Engineering
(VLAN ports 1-8)

Computer Science
(VLAN ports 9-15)

5-90

802.1Q TRUNK

vlan 10
vlan 20

Switch1 Dot1.q    Switch2 Dot1.q

1  2  3  4    A  B  C  D

同一个交换机中，端口/MAC与VLAN是**捆绑关系**，因此，主机发送帧**不需要**带上VLAN ID。

在跨越不同交换机时，下一个交换机无法通过普通帧判断它所属的VLAN ID，所以，从Trunk口发出去前，必须插入VLAN ID。

1).1主机发送普通的数据帧；

2).switch1收到此帧首先需要对其解封装，查看二层帧头部帧目的MAC地址；

3).从表中查找其**目的MAC地址**对应的VLAN ID与**接收该帧的接口**对应的VLAN ID 是否相同，如果相同则找到对应的出接口，如果不同则丢弃该帧；

4).找到出接口后，**打上对应的VLAN 标签**，封装成802.1Q的帧，从Trunk接口发送出去；

5).到达switch2后，解封装查看帧头部的目的MAC地址；

6).从表中查找其**目的MAC地址**对应的VLAN ID与**接收该帧头部**的VLAN ID是否匹配，如果匹配，则查找对应的出接口，如果不同则丢弃该帧；

7).找到出接口后，封装成原始的帧，从相应端口转发出去。

91

# HUB、2-layer switch、 3-layer switch



传统 HUB 的工作过程

HUB

哦，这封信是给地瓜的。

土豆，土豆信已收到

是我的信，收下

不是给我的，丢弃

不是给我的，丢弃

92

**HUB、2-layer switch、 3-layer switch**



传统 Switch 的已知地址报文转发过程

**HUB、2-layer switch、 3-layer switch**



传统 Switch 的地址未知报文广播过程

**HUB、2-layer switch、 <mark>3-layer switch</mark>**

- 三层交换技术：二层交换技术**+**三层转发技术。解决<span style="color:red">局域网中</span>网段划分之后，网段中<span style="color:red">子网必须依赖路由器进行管理</span>的局面，解决了传统路由器低速、复杂所造成的网络瓶颈问题。



三层交换机功能模型

ETH0:10.110.0.254/24　ETH2:10.110.2.254/24
ETH1:10.110.1.254/24

10.110.0.113/24　10.110.1.69/24　10.110.1.88/24　10.110.2.200/24
G:10.110.0.254　G:10.110.1.254　G:10.110.1.254　G:10.110.2.254

第三层交换、第四层交换、多层交换、多层数据包分类和路由交换机

95

---

**HUB、2-layer switch、 <mark>3-layer switch</mark>**

## 三层交换机中的路由和二层交换



- 二层交换引擎：实现同一网段内的快速二层转发
- 三层路由引擎：实现跨网段的三层路由转发

96

## HUB、2-layer switch、 ==3-layer switch==

### 报文到报文的三层交换技术



- 传统三层技术对每个报文进行处理，并基于第三层地址转发报文。这一方法称为报文到报文（Px P）。

97

## HUB、2-layer switch、 ==3-layer switch==

### 基于流交换的三层交换技术



- 不在三层处理所有报文的的方法称之为流交换（FS）。
  - —— 第一个报文
  - ········ 后续报文

98

**HUB、2-layer switch、<mark>3-layer switch</mark>**

假设两个使用**IP**协议的站点**A**、**B**通过第三层交换机进行通信

- 发送站点**A**在开始发送时，把自己的**IP**地址与**B**站的**IP**地址比较，判断**B**站是否与自己在同一子网内。
- 若目的站**B**与发送站**A**在同一子网内，则进行二层的转发。
- 若两个站点不在同一子网内，发送站**A**要向"缺省网关"发出**ARP**(地址解析)封包，而"缺省网关"的**IP**地址其实是三层交换机的三层交换模块。

A —— 三层交换 —— B

100

---

**HUB、2-layer switch、<mark>3-layer switch</mark>**

- 当发送站**A**向"缺省网关"的**IP**发送**ARP**请求时，如果三层交换模块在以前的通信过程中已经知道**B**站的**MAC**地址，则<span style="color:red">向发送站**A**回复**B**的**MAC**地址</span>。否则三层交换模块根据路由信息向**B**站广播一个**ARP**请求;
- **B**站得到此**ARP**请求后向三层交换模块回复其**MAC**地址;
- 三层交换模块保存此地址并<span style="color:red">回复给发送站**A**</span>,同时将**B**站的**MAC**地址发送到二层交换引擎的**MAC**地址表中。

101

## HUB、2-layer switch、 <mark>3-layer switch</mark>

● 从这以后，当**A**向**B**发送的数据包便全部交给<span style="color:red">二层交换</span>处理，信息得以高速交换。由于仅仅在路由过程中才需要三层处理，绝大部分数据都通过二层交换转发，因此三层交换机的速度很快，接近二层交换机的速度

方法二: 交换机直接把三层的目的**IP**映射到二层目的主机的**MAC**所对应的端口,三层交换机的交换表:

| 目的IP地址 | 目的主机MAC | 输出端口 |
|---|---|---|
| | | |

102

---

# Link layer, LANs: outline

<mark>15周:2022.12.9</mark>

**6.1** introduction, services

**6.2** error detection, correction

**6.3** multiple access protocols

**6.4** LANs
- addressing, ARP
- Ethernet
- switches
- VLANS

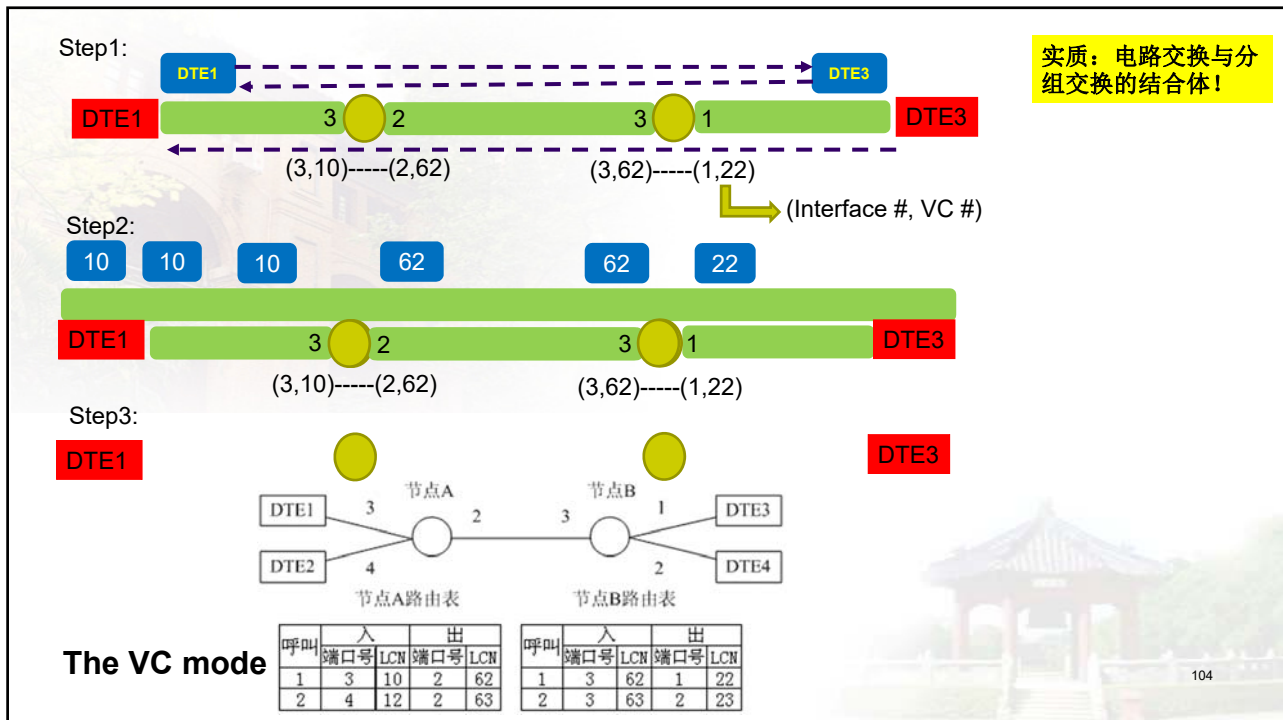<span style="color:red">**6.5** link virtualization: MPLS</span>
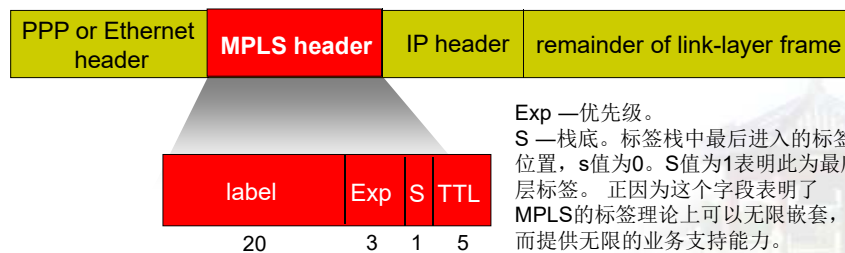
**6.6** data center networking

**6.7** a day in the life of a web request

5-103

Step1:

DTE1 → DTE3

DTE1 — 3 ⬤ 2 — 3 ⬤ 1 — DTE3

(3,10)-----(2,62)    (3,62)-----(1,22)

→ (Interface #, VC #)

实质：电路交换与分组交换的结合体！

Step2:

10  10  10  62  62  22

DTE1 — 3 ⬤ 2 — 3 ⬤ 1 — DTE3

(3,10)-----(2,62)    (3,62)-----(1,22)

Step3:

DTE1  ⬤  ⬤  DTE3

节点A         节点B

| DTE1 | 3 | | 2 | 1 | | DTE3 |
| DTE2 | 4 | | | 2 | | DTE4 |

节点A路由表          节点B路由表

| 呼叫 | 入 | | 出 | | | 呼叫 | 入 | | 出 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 端口号 | LCN | 端口号 | LCN | | | 端口号 | LCN | 端口号 | LCN |
| 1 | 3 | 10 | 2 | 62 | | 1 | 3 | 62 | 1 | 22 |
| 2 | 4 | 12 | 2 | 63 | | 2 | 3 | 63 | 2 | 23 |

**The VC mode**

104

---

# Multiprotocol label switching (MPLS)

- **initial goal: high-speed IP forwarding using fixed length label (instead of IP address)**
  - **fast lookup using fixed length identifier (rather than shortest prefix matching)**
  - **borrowing ideas from Virtual Circuit (VC) approach**
  - **but IP datagram still keeps IP address!**

| PPP or Ethernet header | MPLS header | IP header | remainder of link-layer frame |
|---|---|---|---|

| label | Exp | S | TTL |
|---|---|---|---|
| 20 | 3 | 1 | 5 |

Exp —优先级。
S —栈底。标签栈中最后进入的标签位置，s值为0。S值为1表明此为最底层标签。 正因为这个字段表明了MPLS的标签理论上可以无限嵌套，从而提供无限的业务支持能力。
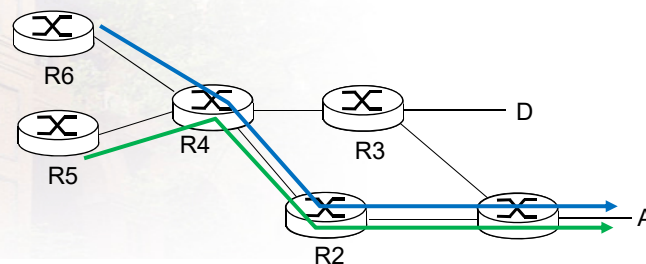
5-105

# MPLS capable routers

- **a.k.a. label-switched router**
- **forward packets to outgoing interface based only on label value (*don't inspect IP address*)**
  - **MPLS forwarding table distinct from IP forwarding tables**
- ***flexibility:*  MPLS forwarding decisions can *differ* from those of IP**
  - **use destination *and* source addresses to route flows to same destination differently (traffic engineering)**
  - **re-route flows quickly if link fails: pre-computed backup paths (useful for VoIP)**
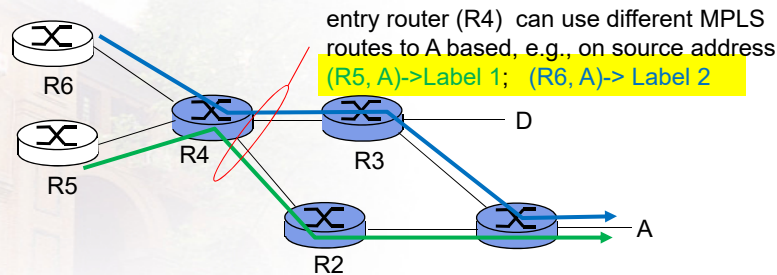
5-106

# MPLS versus IP paths



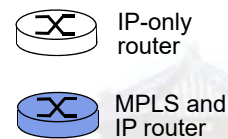- ❖ **IP routing**: path to destination determined by destination address alone

IP router

5-107

2023/11/19

# MPLS versus IP paths

entry router (R4) can use different MPLS routes to A based, e.g., on source address (R5, A)->Label 1;  (R6, A)-> Label 2

R6

R5

R4

R3 — D

R2

A

- ❖ **IP routing:** path to destination determined **by destination address alone**
- ❖ **MPLS routing:** path to destination can be based on **source and dest**. address
  - ▪ **fast reroute:** precompute backup routes in case of link failure

IP-only router

MPLS and IP router

5-108

---

# MPLS signaling

- ● **modify OSPF link-state flooding protocols to carry info used by MPLS routing,**
  - ▪ **e.g., link bandwidth, amount of "reserved" link bandwidth**
    - ❖ *entry MPLS router uses RSVP-TE signaling protocol to set up MPLS forwarding at downstream routers*

R6

R5

R4

RSVP-TE

D

modified link state flooding

A

5-109

# MPLS forwarding tables

R4 table:

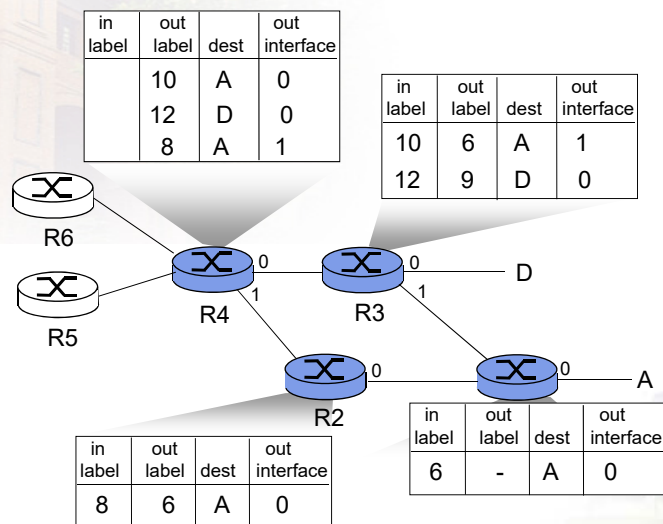| in label | out label | dest | out interface |
|---|---|---|---|
| 10 | | A | 0 |
| 12 | | D | 0 |
| 8 | | A | 1 |

R3 table:

| in label | out label | dest | out interface |
|---|---|---|---|
| 10 | 6 | A | 1 |
| 12 | 9 | D | 0 |

R2 table:

| in label | out label | dest | out interface |
|---|---|---|---|
| 8 | 6 | A | 0 |

R1 table:

| in label | out label | dest | out interface |
|---|---|---|---|
| 6 | - | A | 0 |

R6, R5, R4 (0,1), R3 (0,1), R2 (0), — D, A

5-110

---

匹配查找   匹配查找   匹配查找   匹配查找   匹配查找   匹配查找

IP分组 → 路由器 → 路由器 → 路由器 → 路由器 → 路由器 → 路由器 → IP分组

IP分组   IP分组   IP分组   IP分组   IP分组
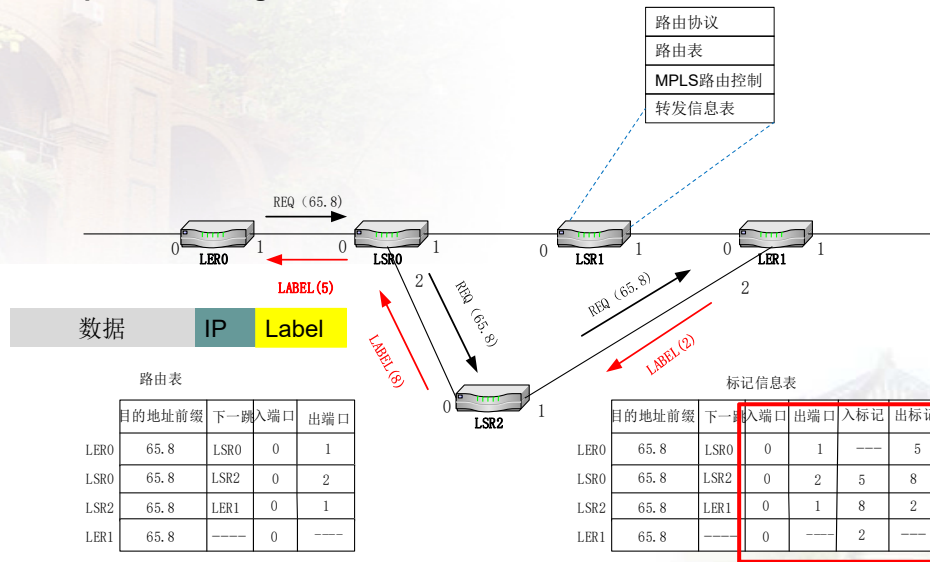
匹配查找

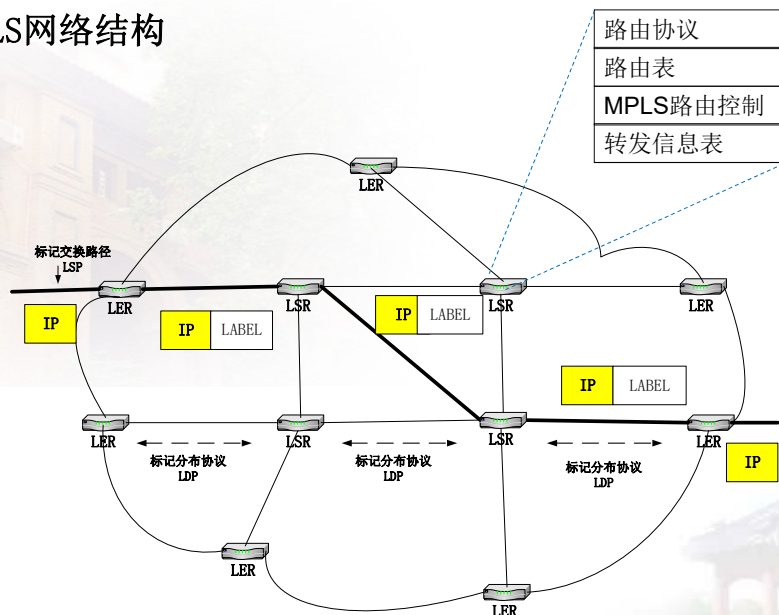IP分组 → 路由器 → IP分组

下一跳中继

**IP网络的逐跳式分组转发**

**Step1: routing for the path;**
**Step2: labeling each link;**
**Step3: forwarding based on label**

Replace the **global** routing table with the
**local** label switching table

路由协议
路由表
MPLS路由控制
转发信息表

REQ（65.8）

0 LER0 1　0 LSR0 1　0 LSR1 1　0 LER1 1

LABEL(5)

| 数据 | IP | Label |
|---|---|---|

2　REQ（65.8）　REQ（65.8）　2

LABEL(8)　LABEL(2)

0 LSR2 1

路由表

| | 目的地址前缀 | 下一跳 | 入端口 | 出端口 |
|---|---|---|---|---|
| LER0 | 65.8 | LSR0 | 0 | 1 |
| LSR0 | 65.8 | LSR2 | 0 | 2 |
| LSR2 | 65.8 | LER1 | 0 | 1 |
| LER1 | 65.8 | ---- | 0 | ---- |

标记信息表

| | 目的地址前缀 | 下一跳 | 入端口 | 出端口 | 入标记 | 出标记 |
|---|---|---|---|---|---|---|
| LER0 | 65.8 | LSR0 | 0 | 1 | --- | 5 |
| LSR0 | 65.8 | LSR2 | 0 | 2 | 5 | 8 |
| LSR2 | 65.8 | LER1 | 0 | 1 | 8 | 2 |
| LER1 | 65.8 | ---- | 0 | ---- | 2 | ---- |

112

---

# MPLS网络结构

路由协议
路由表
MPLS路由控制
转发信息表

LER

标记交换路径
LSP

| IP |

LER　LSR　IP LABEL　LSR　LER

| IP | LABEL |

| IP | LABEL |

LER　LSR　LSR　LER

| IP |

标记分布协议 LDP　标记分布协议 LDP　标记分布协议 LDP

LER

LER

113

# Link layer, LANs: outline

**6.1** introduction, services

**6.2** error detection, correction

**6.3** multiple access protocols

**6.4** LANs
- addressing, ARP
- Ethernet
- switches
- VLANS

**6.5** link virtualization: MPLS

**6.6** data center networking

**6.7** a day in the life of a web request

# Data center networks

- **10's to 100's of thousands of hosts, often closely coupled, in close proximity:**
  - **e-business** (e.g. Amazon)
  - **content-servers** (e.g., YouTube, Akamai, Apple, Microsoft)
  - **search engines**, data mining (e.g., Google)
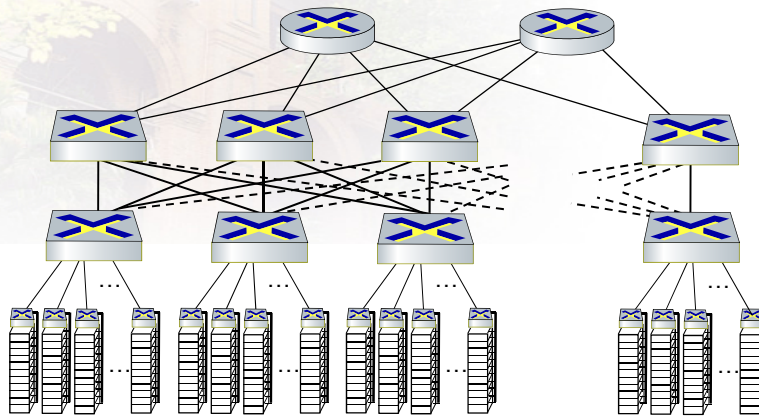
  - ❖ challenges:
    - multiple applications, each serving massive numbers of clients
    - managing/balancing load, avoiding processing, networking, data bottlenecks

Inside a 40-ft Microsoft container, Chicago data center

# Datacenter networks: network elements

**Border routers**
- connections outside datacenter

**Tier-1 switches**
- connecting to ~16 T-2s below

**Tier-2 switches**
- connecting to ~16 TORs below

**Top of Rack (TOR) switch**
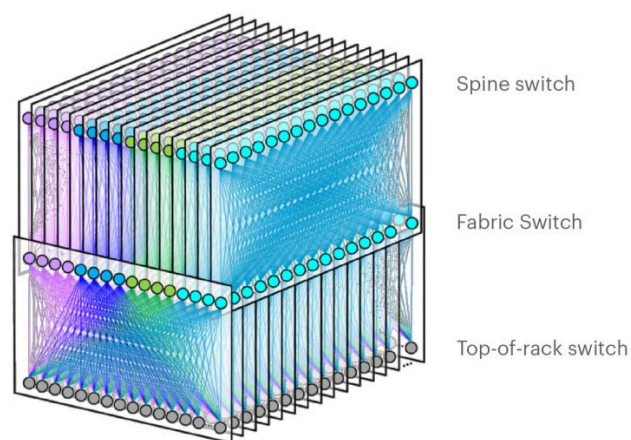- one per rack
- 40-100Gbps Ethernet to blades

**Server racks**
- 20- 40 server blades: hosts

Link Layer: 6-116

---

# Datacenter networks: network elements

Facebook F16 data center network topology:
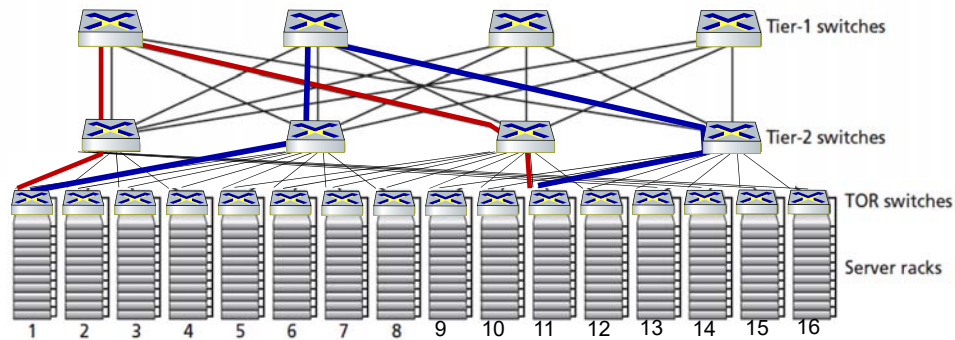
Spine switch

Fabric Switch

Top-of-rack switch

https://engineering.fb.com/data-center-engineering/f16-minipack/     (posted 3/2019)

Link Layer: 6-117
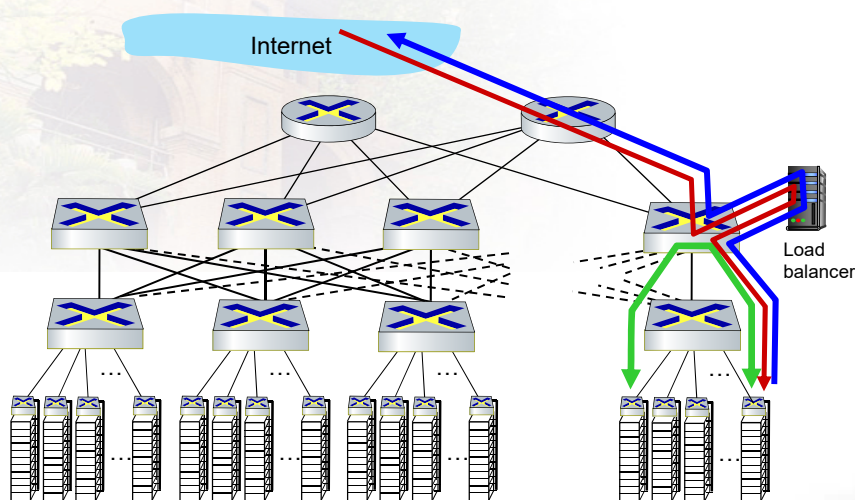
# Datacenter networks: multipath

- rich interconnection among switches, racks:
  - increased throughput between racks (multiple routing paths possible)
  - increased reliability via redundancy



two disjoint paths highlighted between racks 1 and 11

Link Layer: 6-118

# Datacenter networks: application-layer routing



load balancer: application-layer routing

- receives external client requests
- directs workload within data center
- returns results to external client (hiding data center internals from client)

Link Layer: 6-119

## Datacenter networks: protocol innovations

- **link layer:**
  - RoCE: remote DMA (RDMA) over Converged Ethernet
- **transport layer:**
  - ECN (explicit congestion notification) used in transport-layer congestion control (DCTCP, DCQCN)
  - experimentation with hop-by-hop (backpressure) congestion control
- **routing, management:**
  - SDN widely used within/among organizations' datacenters
  - place related services, data as close as possible (e.g., in same rack or nearby rack) to minimize tier-2, tier-1 communication

Link Layer: 6-120

## Link layer, LANs: outline

**6.1** introduction, services

**6.2** error detection, correction

**6.3** multiple access protocols

**6.4** LANs
- addressing, ARP
- Ethernet
- switches
- VLANS

**6.5** link virtualization: MPLS

**6.6** data center networking

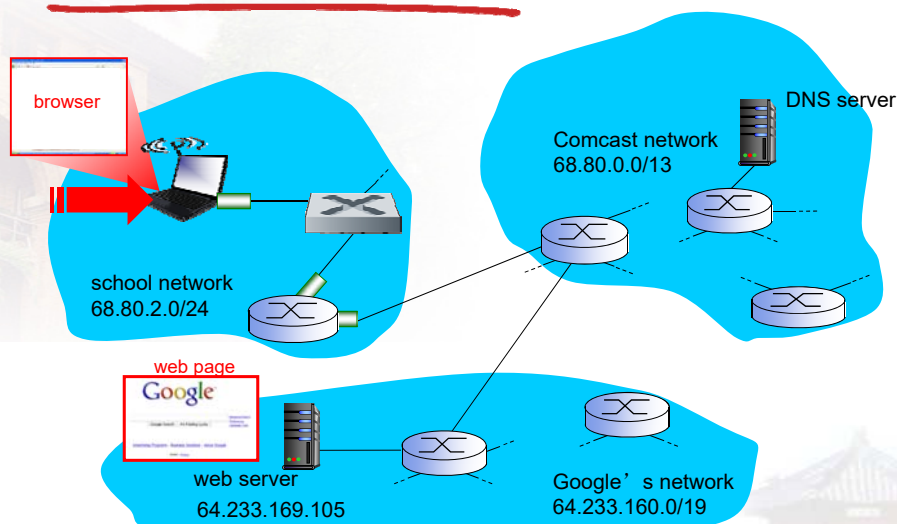**6.7 a day in the life of a web request**

5-121

## *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**
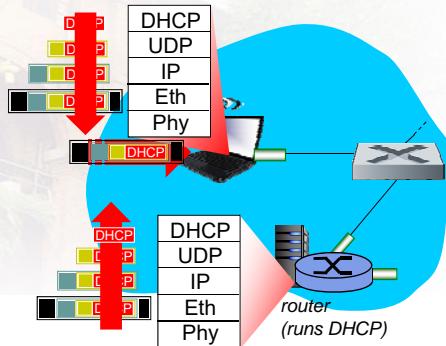
5-122

## A day in the life: scenario



browser

Comcast network
68.80.0.0/13

DNS server

school network
68.80.2.0/24

web page

Google

web server
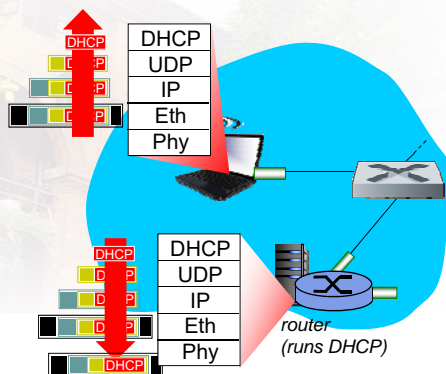64.233.169.105

Google's network
64.233.160.0/19

5-123

# 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: FFFFFFFFFFFF) on LAN, received at router running DHCP server

❖ Ethernet demuxed to IP demuxed, UDP demuxed to DHCP

5-124

# A day in the life… connecting to the Internet



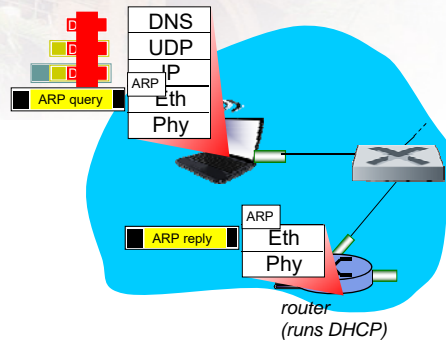● **DHCP server formulates *DHCP ACK* containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server**

❖ encapsulation at DHCP server, frame forwarded (switch learning) through LAN, demultiplexing at client

❖ DHCP client receives DHCP ACK reply

*Client now has IP address, knows name & addr of DNS server, IP address of its first-hop router*
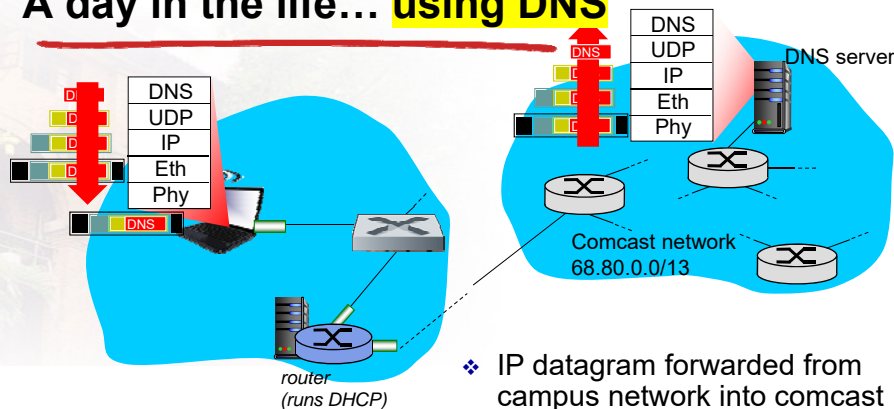
5-125

## A day in the life… ARP (before DNS, before HTTP)



DNS
UDP
IP
ARP  Eth
Phy

ARP query

ARP
ARP reply  Eth
Phy

*router
(runs DHCP)*

❖ **before sending *HTTP* request, need IP address of www.google.com: *DNS***

❖ DNS query created, encapsulated in UDP, encapsulated in IP, encapsulated in Eth.  To send frame to router, need MAC address of router interface: ARP

❖ ARP query broadcast, received by router, which replies with ARP reply giving MAC address of router interface

❖ client now knows MAC address of first hop router, so can now send frame containing DNS query

5-126

## A day in the life… using DNS



DNS
UDP
IP
Eth
Phy

DNS

DNS
UDP
IP
Eth
Phy

DNS

DNS server

Comcast network
68.80.0.0/13

*router
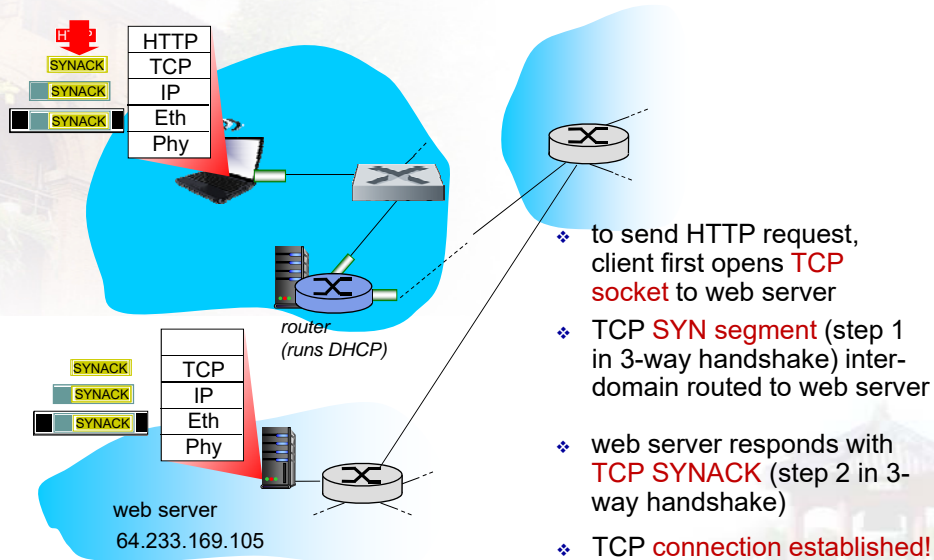(runs DHCP)*

❖ IP datagram containing DNS query forwarded via LAN switch from client to 1st hop router

❖ IP datagram forwarded from campus network into comcast network, routed (tables created by RIP, OSPF, IS-IS and/or BGP routing protocols) to DNS server

❖ Demux'ed to DNS server

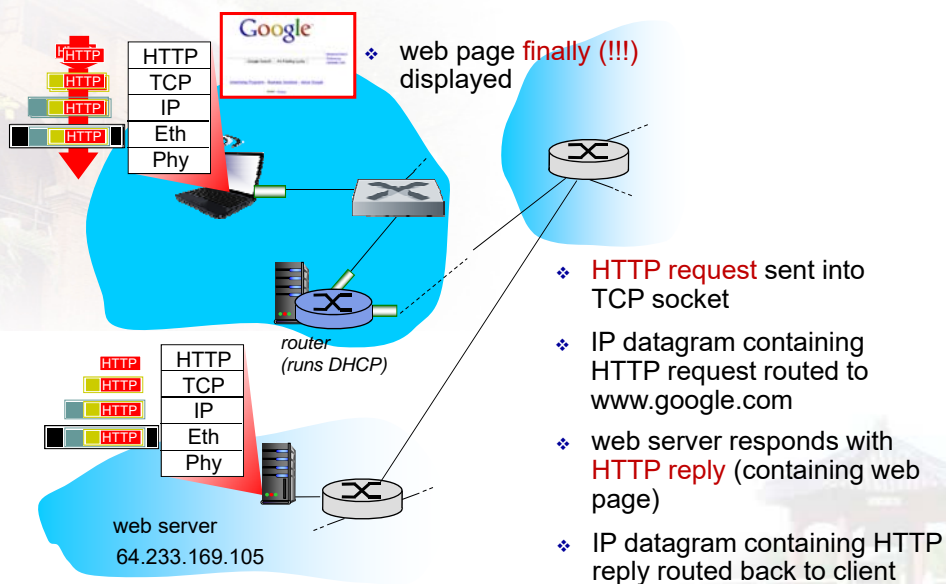❖ DNS server replies to client with IP address of www.google.com

5-127

# A day in the life…TCP connection carrying HTTP

```
HTTP
SYNACK   TCP
SYNACK   IP
SYNACK   Eth
         Phy
```

*router (runs DHCP)*

```
SYNACK   TCP
SYNACK   IP
SYNACK   Eth
         Phy
```

web server
64.233.169.105

- to send HTTP request, client first opens TCP socket to web server
- TCP SYN segment (step 1 in 3-way handshake) inter-domain routed to web server
- web server responds with TCP SYNACK (step 2 in 3-way handshake)
- TCP connection established!

5-128

# A day in the life… HTTP request/reply

```
HTTP   HTTP
HTTP   TCP
HTTP   IP
HTTP   Eth
       Phy
```

Google

*router (runs DHCP)*

```
HTTP   HTTP
HTTP   TCP
HTTP   IP
HTTP   Eth
       Phy
```

web server
64.233.169.105

- web page finally (!!!) displayed
- HTTP request sent into TCP socket
- IP datagram containing HTTP request routed to www.google.com
- web server responds with HTTP reply (containing web page)
- IP datagram containing HTTP reply routed back to client

5-129

# Chapter 6: 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, VLANs**
  - ■ **virtualized networks as a link layer: MPLS**
- **synthesis: a day in the life of a web request**

5-130

# Chapter 6: 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**

**The End of Chapter 6**  5-131

# Thanks

Q & A

**Email: xieyi5@mail.sysu.edu.cn**
**https://cse.sysu.edu.cn/content/2462**