



WWW.DAUPHINE.PSL.EU

## Computer Networks

<https://guilhermeir.github.io/L3Networks.html>

Prof. Guilherme lecker Ricardo  
guilherme.iecker-ricardo@dauphine.psl.eu

**Dauphine** | PSL★  
UNIVERSITÉ PARIS



# Computer Networks

## Class 4: Physical Layer

### (Continued)

## Class 2 Review

- Types of Transmissions
  - Guided
  - Wireless
- Fundamentals of Digital Communications
  - Signals
  - Channel Capacity

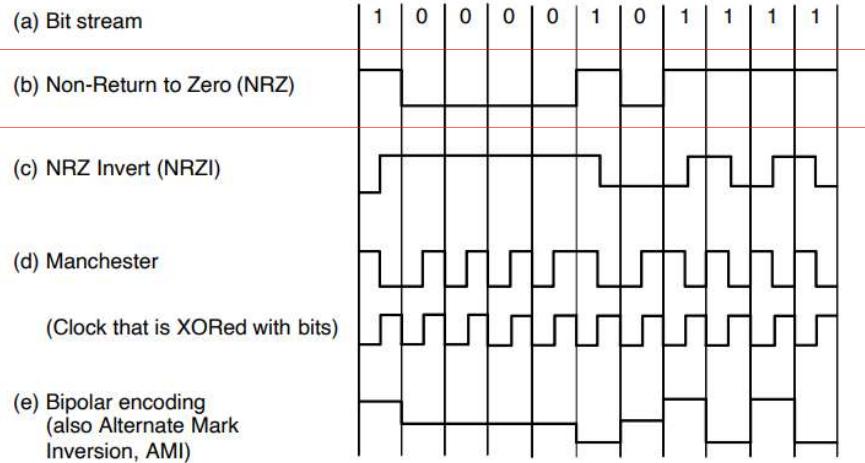
### Physical Layer

## Digital Modulation

### Definition

### Modulation Schemes

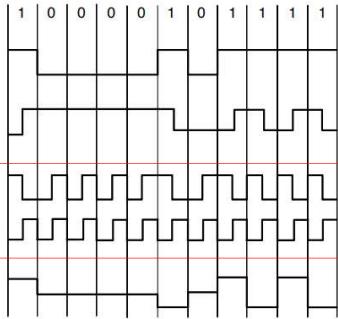
- Baseband Transmission
- Passband Transmission
- Multiplexing

**Physical Layer****Digital Modulation – Baseband Transmission: Non-Return to Zero****Physical Layer****Digital Modulation – Baseband Transmission: Bandwidth Efficiency****Bandwidth Efficiency**

- What is maximum rate for NRZ?
  - $R=2B$  (Nyquist!)
- How to increase max rate?
  - $B$  is fixed in most cases
  - Answer: Increase number of levels
    - 1 symbol carries 2 bits
  - Problem: Attenuation!
- *Symbol rate*: rate at which bits change
- Bit rate = Symbol rate x bits/symbol

**Physical Layer****Digital Modulation – Baseband Transmission: Manchester Encoding**

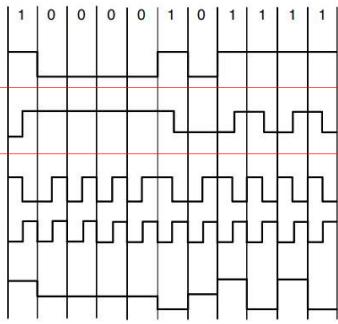
- (a) Bit stream
- (b) Non-Return to Zero (NRZ)
- (c) NRZ Invert (NRZI)
- (d) Manchester**  
(Clock that is XORed with bits)
- (e) Bipolar encoding  
(also Alternate Mark Inversion, AMI)



- NRZ: Hard to detect bounds
- Clock Recovery
- Clock = 2 x Bit Rate
- Encoded Signal = Clock  $\oplus$  Data

**Physical Layer****Digital Modulation – Baseband Transmission: NRZ Inverted Encoding**

- (a) Bit stream
- (b) Non-Return to Zero (NRZ)
- (c) NRZ Invert (NRZI)**
- (d) Manchester
- (Clock that is XORed with bits)
- (e) Bipolar encoding  
(also Alternate Mark Inversion, AMI)



- Manchester: Requires 2x bandwidth
- Transitions are easier to detect
- Encoding:
  - 1-bit: transition
  - 0-bit: no transition

**Physical Layer****Digital Modulation – Baseband Transmission: 4B/5B Encoding**

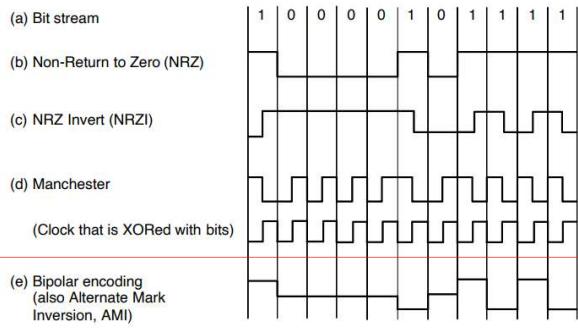
Data (4B)	Codeword (5B)	Data (4B)	Codeword (5B)
0000	11110	1000	10010
0001	01001	1001	10011
0010	10100	1010	10110
0011	10101	1011	10111
0100	01010	1100	11010
0101	01011	1101	11011
0110	01110	1110	11100
0111	01111	1111	11101

- **NRZI: Long streams of zeros**
- 4 bits maps to unique 5 bits
  - No more than 3 zeros in a row
- 25% overhead
  - Still better than Manchester
- 16 combinations are not used for data
  - Control codes:
    - 11111: idle line
    - 11000: start of a frame

**Physical Layer****Digital Modulation – Baseband Transmission: 4B/5B “Scrambling”**

Data (4B)	Codeword (5B)	Data (4B)	Codeword (5B)
0000	11110	1000	10010
0001	01001	1001	10011
0010	10100	1010	10110
0011	10101	1011	10111
0100	01010	1100	11010
0101	01011	1101	11011
0110	01110	1110	11100
0111	01111	1111	11101

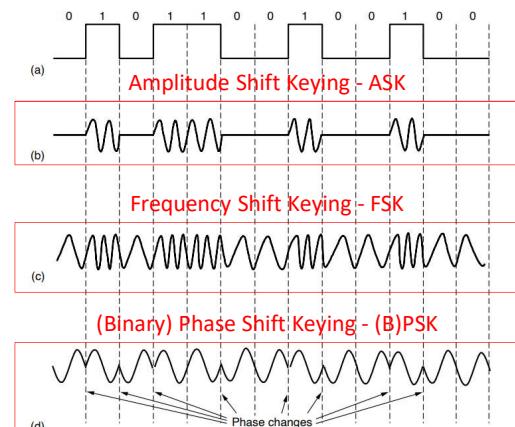
- **4B/5B: Easy to intercept (unsafe)**
- Shared pseudorandom (PR) sequence
- Signal = Data  $\oplus$  PR
- Advantages:
  - No overhead
  - No additional bandwidth needed
- Disadvantage: sequence of zeros!

**Physical Layer****Digital Modulation – Baseband Modulation: Bipolar Encoding**

- Previous coding: unbalanced voltage
- 3-voltage-level Encoding:
  - 0-bit = 0V
  - 1-bit = +1V or -1V (alternating)
- 8B/10B Mapping:
  - 8 bits → 10 bits (80% efficiency)
  - 5 first bits → 6 bits
  - 3 last bits → 4 bits
  - Help balance: at most 2 disparity bits
  - Clock recovery: at most 5 consecutive 1s/0s

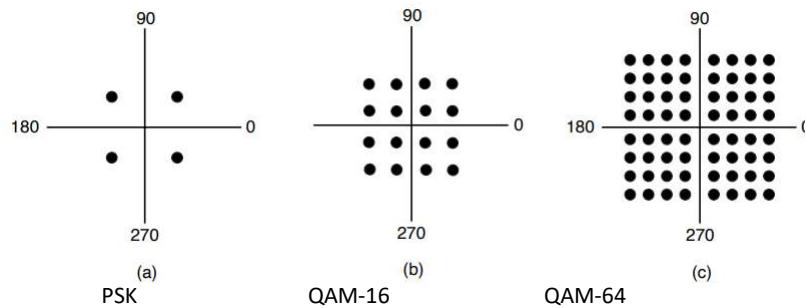
**Physical Layer****Digital Modulation – Passband Modulation**

- Baseband is not convenient for wireless communication
- Wireless transmission often needs to take place in “shifted slices” of the broadband

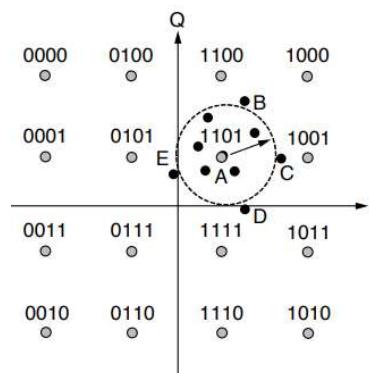
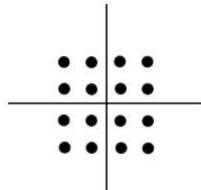


**Physical Layer****Digital Modulation – Passband Modulation: Combined Modulations**

- Quadrature Amplitude Modulation - QAM
  - Amplitude and phase are modulated together
- Constellation Diagram

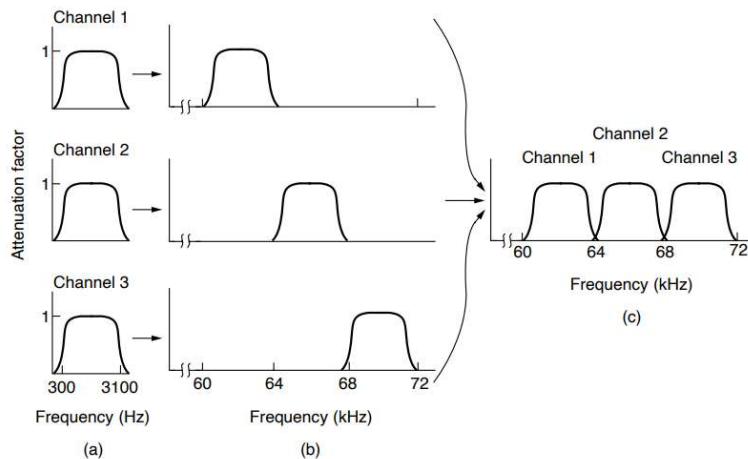
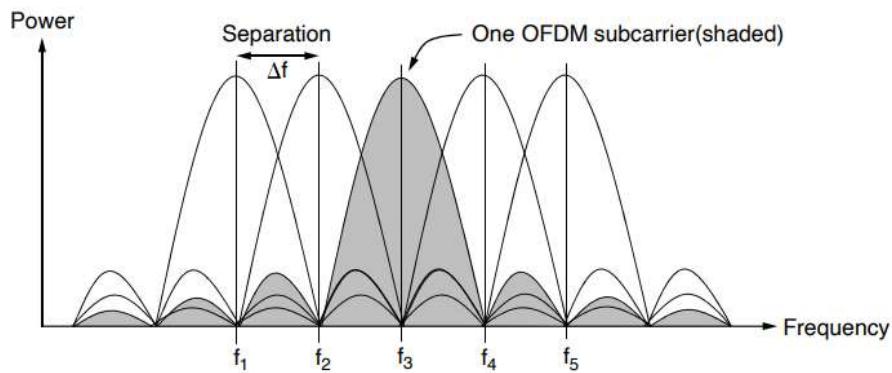
**Physical Layer****Digital Modulation – Passband Modulation: Avoiding Errors**

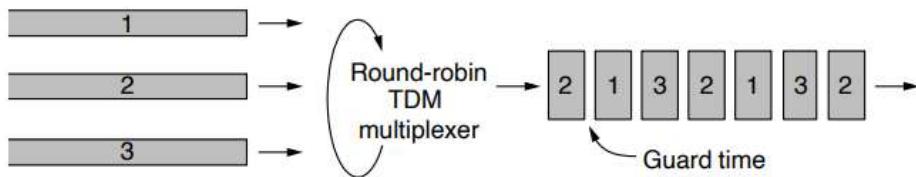
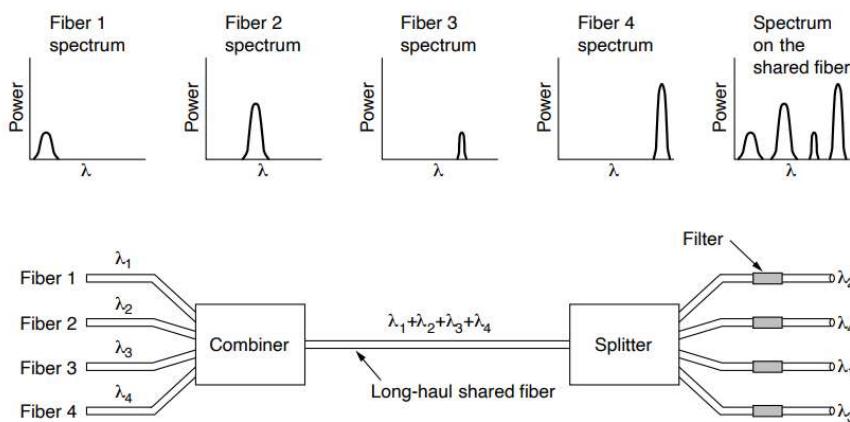
- How bits are mapped to symbols?
- Bit (Translation) Error



When 1101 is sent:

Point	Decodes as	Bit errors
A	1101	0
B	110 <u>0</u>	1
C	1001	1
D	111 <u>1</u>	1
E	0101	1

**Physical Layer****Digital Transmission – Multiplexing: Frequency Division Multiplexing****Physical Layer****Digital Transmission – Multiplexing: Orthogonal FDM**

**Physical Layer****Digital Transmission – Multiplexing: Time Division Multiplexing****Physical Layer****Digital Transmission – Multiplexing: Wavelength Division Multiplexing**

**Physical Layer****Digital Transmission – Multiplexing: Spatial Division Multiplexing****Physical Layer****Digital Transmission – Multiplexing: Code Division Multiplexing**

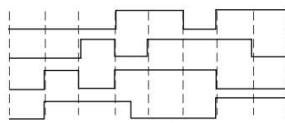
$$A = (-1 -1 -1 +1 +1 -1 +1 +1)$$

$$B = (-1 -1 +1 -1 +1 +1 +1 -1)$$

$$C = (-1 +1 -1 +1 +1 +1 -1 -1)$$

$$D = (-1 +1 -1 -1 -1 -1 +1 -1)$$

(a)



(b)

$$\begin{aligned} S_1 &= C &= (-1 +1 -1 +1 +1 +1 -1 -1) \\ S_2 &= B+C &= (-2 \ 0 \ 0 \ 0 +2 +2 \ 0 -2) \\ S_3 &= A+B &= ( \ 0 \ 0 -2 +2 \ 0 -2 \ 0 +2) \\ S_4 &= A+B+C &= (-1 +1 -3 +3 +1 -1 -1 +1) \\ S_5 &= A+B+C+D &= (-4 \ 0 -2 \ 0 +2 \ 0 +2 -2) \\ S_6 &= A+B+C+D &= (-2 -2 \ 0 -2 \ 0 -2 +4 \ 0) \end{aligned}$$

(c)

$$\begin{aligned} S_1 \cdot C &= [1+1+1+1+1+1+1]/8 = 1 \\ S_2 \cdot C &= [2+0+0+0+2+2+0+2]/8 = 1 \\ S_3 \cdot C &= [0+0+2+2+0-2+0-2]/8 = 0 \\ S_4 \cdot C &= [1+1+3+3+1-1+1-1]/8 = 1 \\ S_5 \cdot C &= [4+0+2+0+2+0-2+2]/8 = 1 \\ S_6 \cdot C &= [2-2+0-2+0-2-4+0]/8 = -1 \end{aligned}$$

(d)

# Computer Networks

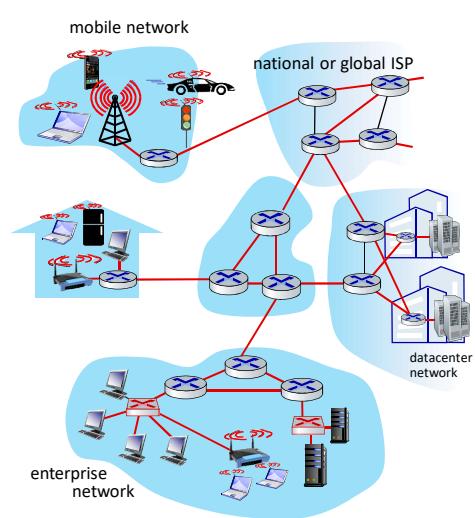
## Class 4: Link Layer

### Link layer: introduction

terminology:

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

*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., WiFi on first link, Ethernet on next link
- each link protocol provides different services
  - e.g., may or may not provide reliable data transfer over link

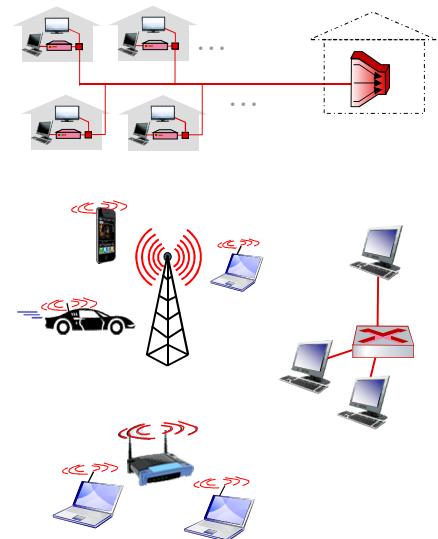
### transportation analogy:

- 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: 6-23

## Link layer: services

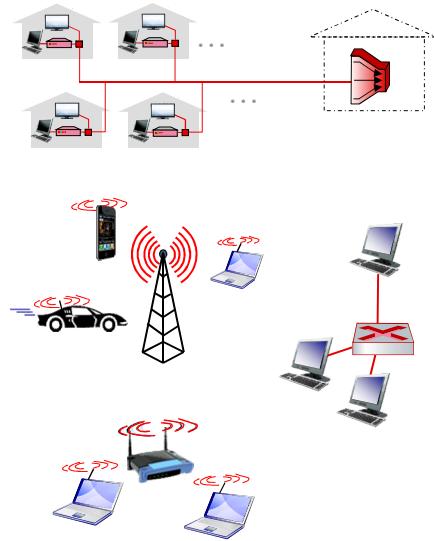
- **framing, link access:**
  - encapsulate datagram into frame, adding header, trailer
  - channel access if shared medium
  - “MAC” addresses in frame headers identify source, destination (different from IP address!)
- **reliable delivery between adjacent nodes**
  - we already know how to do this!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



Link Layer: 6-24

## Link layer: services (more)

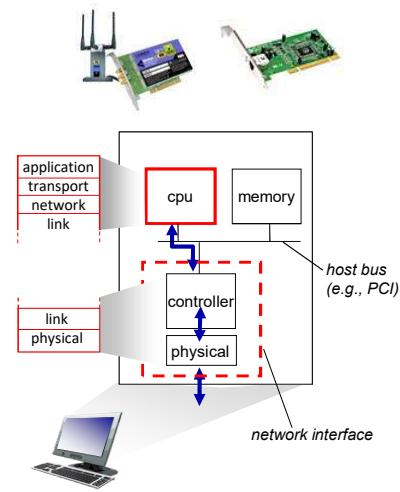
- **flow control:**
  - pacing between adjacent sending and receiving nodes
- **error detection:**
  - errors caused by signal attenuation, noise.
  - receiver detects errors, signals retransmission, or drops frame
- **error correction:**
  - receiver identifies *and corrects* bit error(s) without retransmission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



Link Layer: 6-25

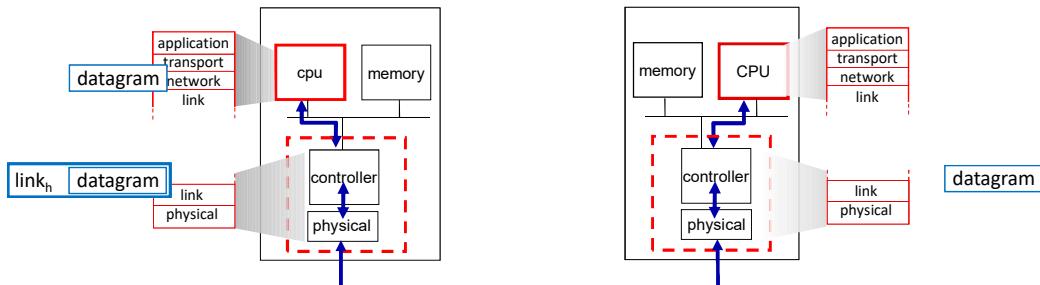
## Where is the link layer implemented?

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



Link Layer: 6-26

## 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: 6-27

## Link layer, LANs: roadmap

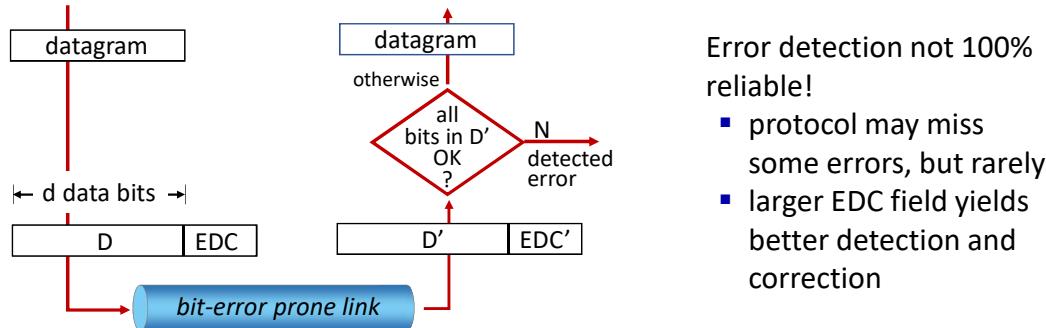
- introduction
- **error detection, correction**
- multiple access protocols
- LANs
  - addressing, ARP
  - Ethernet
  - switches
  - VLANs
- link virtualization: MPLS
- data center networking
- a day in the life of a web request

Link Layer: 6-28

## Error detection

EDC: error detection and correction bits (e.g., redundancy)

D: data protected by error checking, may include header fields

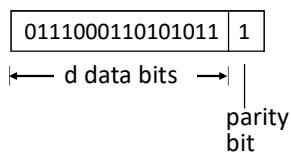


Link Layer: 6-29

## Parity checking

### single bit parity:

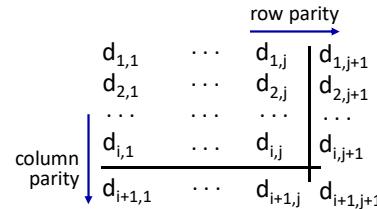
- detect single bit errors



**Even parity:** set parity bit so there is an even number of 1's

### two-dimensional bit parity:

- detect *and correct* single bit errors



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

detected and correctable single-bit error:  
1 0 1 0 1 | 1  
1 0 1 1 0 | 0  
0 1 1 1 0 | 1  
1 0 1 0 1 | 0  
↓  
parity error

\* Check out the online interactive exercises for more examples: [http://gaia.cs.umass.edu/kurose\\_ross/interactive/](http://gaia.cs.umass.edu/kurose_ross/interactive/)

Link Layer: 6-30

## Internet checksum (review)

**Goal:** detect errors (*i.e.*, flipped bits) in transmitted segment

### sender:

- treat contents of UDP segment (including UDP header fields and IP addresses) as sequence of 16-bit integers
- **checksum:** addition (one's complement sum) of segment content
- checksum value put into UDP checksum field

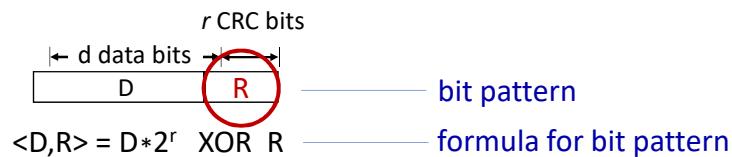
### receiver:

- compute checksum of received segment
- check if computed checksum equals checksum field value:
  - not equal - error detected
  - equal - no error detected. *But maybe errors nonetheless?* More later ....

Transport Layer: 3-31

## Cyclic Redundancy Check (CRC)

- more powerful error-detection coding
- **D:** data bits (given, think of these as a binary number)
- **G:** bit pattern (generator), of  $r+1$  bits (given)



goal: choose  $r$  CRC bits,  $R$ , such that  $\langle D, R \rangle$  exactly divisible by  $G$  ( $\text{mod } 2$ )

- 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)

Link Layer: 6-32

## Cyclic Redundancy Check (CRC): example

We want:

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

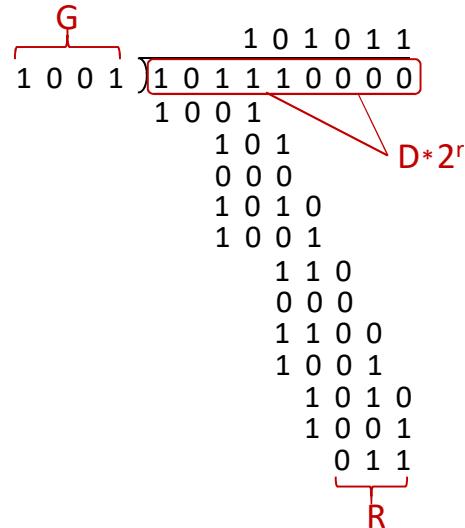
or equivalently:

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

or equivalently:

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

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



## Multiple access links, protocols

two types of “links”:

- **point-to-point**
  - point-to-point link between Ethernet switch, host
  - PPP for dial-up access
- **broadcast (shared wire or medium)**
  - old-fashioned Ethernet
  - upstream HFC in cable-based access network
  - 802.11 wireless LAN, 4G/4G, satellite



Link Layer: 6-35

## 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

Link Layer: 6-36

## An ideal multiple access protocol

*given:* multiple access channel (MAC) 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

Link Layer: 6-37

## MAC protocols: taxonomy

three broad classes:

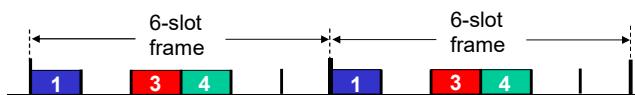
- **channel partitioning**
  - divide channel into smaller “pieces” (time slots, frequency, code)
  - allocate piece to node for exclusive use
- **random access**
  - channel not divided, allow collisions
  - “recover” from collisions
- **“taking turns”**
  - nodes take turns, but nodes with more to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

Link Layer: 6-38

## Channel partitioning MAC protocols: TDMA

### TDMA: time division multiple access

- access to channel in “rounds”
- each station gets fixed length slot (length = packet transmission time) in each round
- unused slots go idle
- example: 6-station LAN, 1,3,4 have packets to send, slots 2,5,6 idle

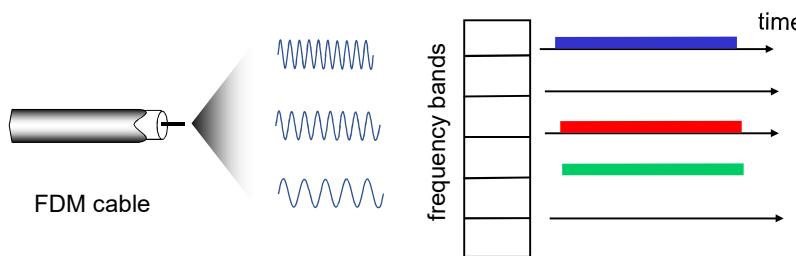


Link Layer: 6-39

## 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 packet to send, frequency bands 2,5,6 idle



Link Layer: 6-40

## 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

Link Layer: 6-41

## 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

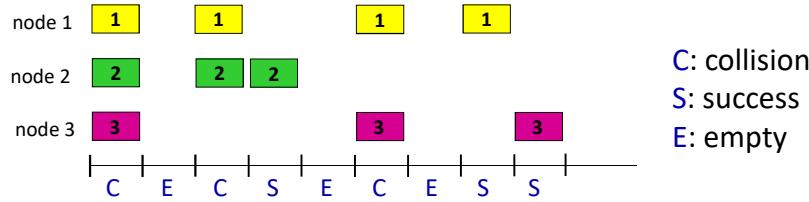
### operation:

- 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 probability  $p$  until success

**randomization – why?**

Link Layer: 6-42

## 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

Link Layer: 6-43

## Slotted ALOHA: efficiency

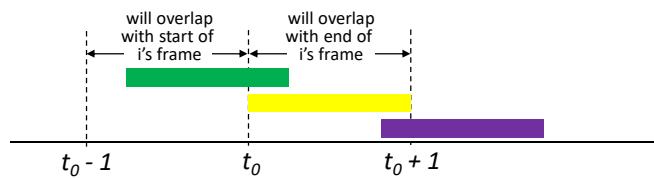
**efficiency:** long-run fraction of successful slots (many nodes, all with many frames to send)

- *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:  
*max efficiency =  $1/e = .37$*
- *at best:* channel used for useful transmissions 37% of time!

Link Layer: 6-44

## Pure ALOHA

- unslotted Aloha: simpler, no synchronization
  - when frame first arrives: transmit immediately
- collision probability increases with no synchronization:
  - frame sent at  $t_0$  collides with other frames sent in  $[t_0-1, t_0+1]$



- pure Aloha efficiency: 18% !

Link Layer: 6-45

## CSMA (carrier sense multiple access)

- simple CSMA: listen before transmit:
- if channel sensed idle: transmit entire frame
  - if channel sensed busy: defer transmission
- human analogy: don't interrupt others!

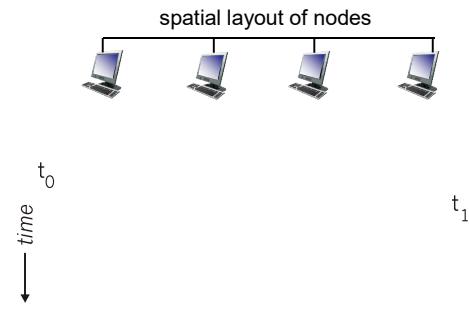
### CSMA/CD: CSMA with *collision detection*

- collisions detected within short time
  - colliding transmissions aborted, reducing channel wastage
  - collision detection easy in wired, difficult with wireless
- human analogy: the polite conversationalist

Link Layer: 6-46

## CSMA: collisions

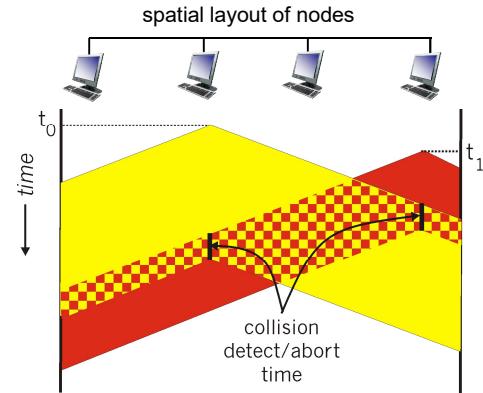
- collisions *can* still occur with carrier sensing:
  - propagation delay means two nodes may not hear each other's just-started transmission
- **collision:** entire packet transmission time wasted
  - distance & propagation delay play role in determining collision probability



Link Layer: 6-47

## CSMA/CD:

- CSMA/CS reduces the amount of time wasted in collisions
  - transmission aborted on collision detection



Link Layer: 6-48

## Ethernet CSMA/CD algorithm

1. NIC receives datagram from network layer, creates frame
2. If NIC senses channel:
  - if **idle**: start frame transmission.
  - if **busy**: wait until channel idle, then transmit
3. If NIC transmits entire frame without collision, NIC is done with frame !
4. If NIC detects another transmission while sending: abort, send jam signal
5. After aborting, NIC enters ***binary (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
  - more collisions: longer backoff interval

Link Layer: 6-49

## CSMA/CD efficiency

- $T_{prop}$  = max prop delay 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
  - as  $t_{trans}$  goes to infinity
- better performance than ALOHA: and simple, cheap, decentralized!

Link Layer: 6-50

## “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

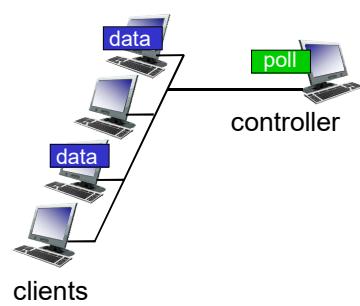
- look for best of both worlds!

Link Layer: 6-51

## “Taking turns” MAC protocols

### polling:

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

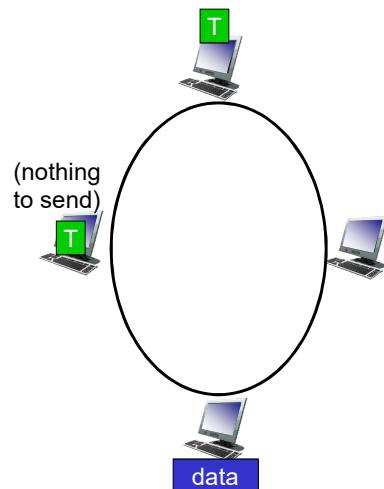


Link Layer: 6-52

## “Taking turns” MAC protocols

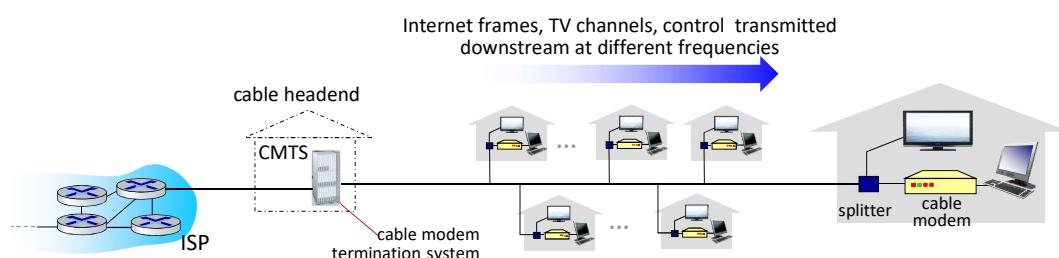
### token passing:

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



Link Layer: 6-53

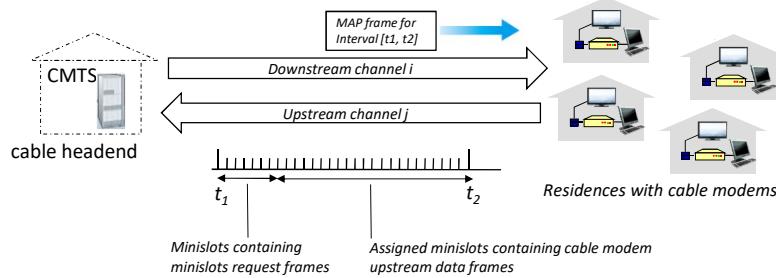
## Cable access network: FDM, TDM and random access!



- **multiple downstream (broadcast) FDM channels:** up to 1.6 Gbps/channel
  - single CMTS transmits into channels
- **multiple upstream channels (up to 1 Gbps/channel)**
  - **multiple access:** all users contend (random access) for certain upstream channel time slots; others assigned TDM

Link Layer: 6-54

## Cable access network:



### DOCSIS: data over cable service interface specification

- 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

Link Layer: 6-55

## 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

Link Layer: 6-56

## Link layer, LANs: roadmap

- introduction
- error detection, correction
- multiple access protocols
- **LANs**
  - addressing, ARP
  - Ethernet
  - switches
  - VLANs
- link virtualization: MPLS
- data center networking
- a day in the life of a web request

Link Layer: 6-57

## MAC addresses

- 32-bit IP address:
  - *network-layer* address for interface
  - used for layer 3 (network layer) forwarding
  - e.g.: 128.119.40.136
- MAC (or LAN or physical or Ethernet) address:
  - function: **used “locally” to get frame from one interface to another physically-connected interface (same subnet, 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

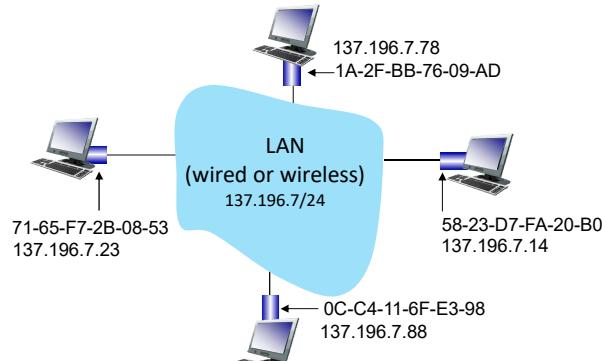
*hexadecimal (base 16) notation  
(each “numeral” represents 4 bits)*

Link Layer: 6-58

## MAC addresses

each interface on LAN

- has unique 48-bit **MAC** address
- has a locally unique 32-bit IP address (as we've seen)



Link Layer: 6-59

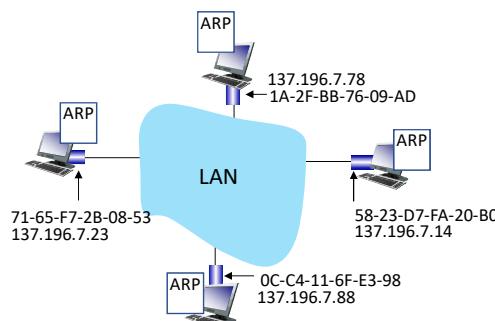
## MAC addresses

- 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 interface from one LAN to another
  - recall IP address *not* portable: depends on IP subnet to which node is attached

Link Layer: 6-60

## ARP: address resolution protocol

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



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

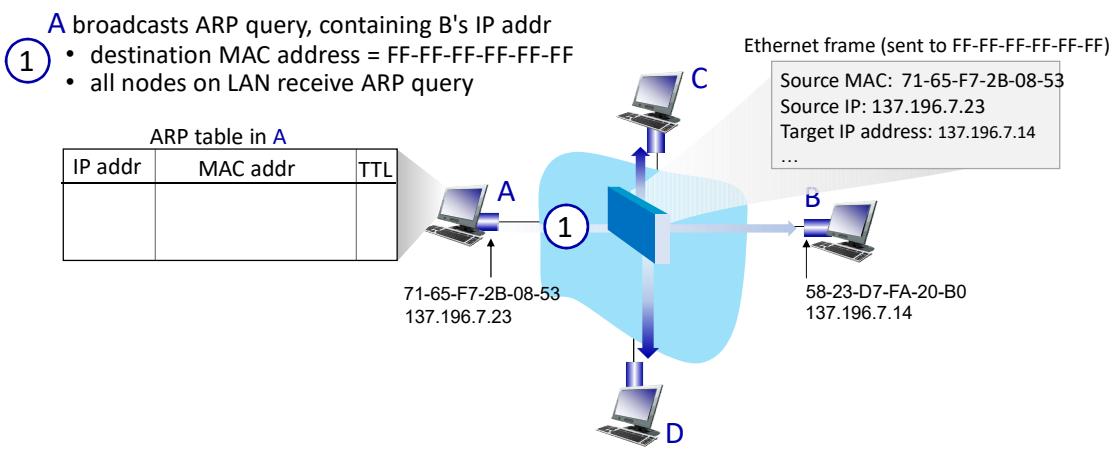
- IP/MAC address mappings for some LAN nodes:  
<IP address; MAC address; TTL>
- TTL (Time To Live): time after which address mapping will be forgotten (typically 20 min)

Link Layer: 6-61

## 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

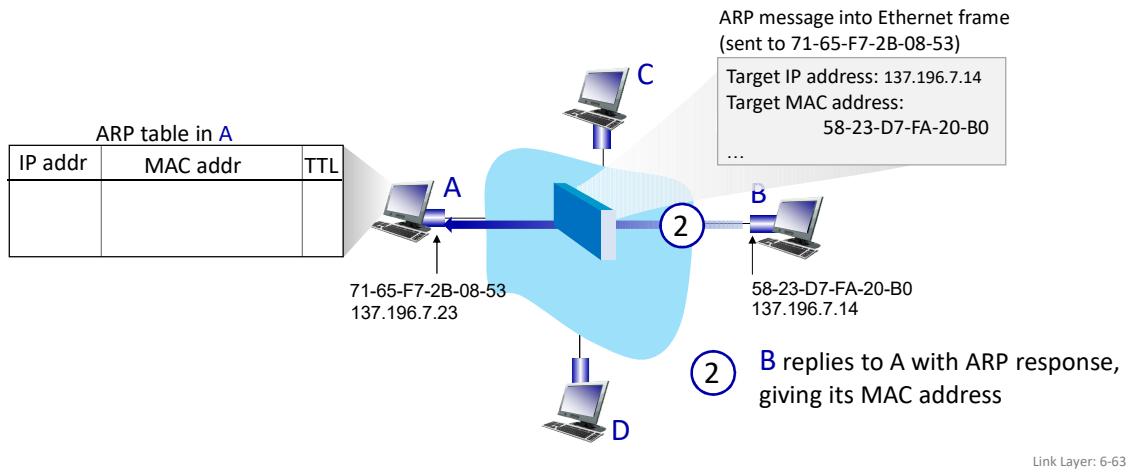


Link Layer: 6-62

## ARP protocol in action

example: A wants to send datagram to B

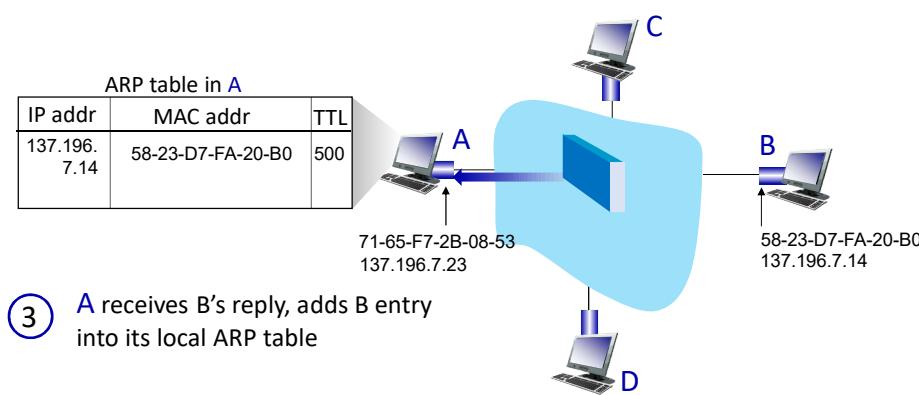
- B's MAC address not in A's ARP table, so A uses ARP to find B's MAC address



## ARP protocol in action

example: A wants to send datagram to B

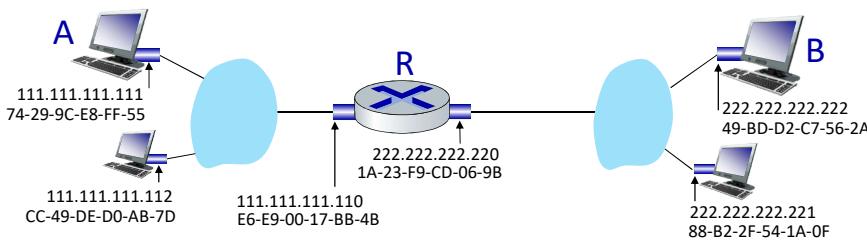
- B's MAC address not in A's ARP table, so A uses ARP to find B's MAC address



## Routing to another subnet: addressing

walkthrough: sending a datagram from A to B via R

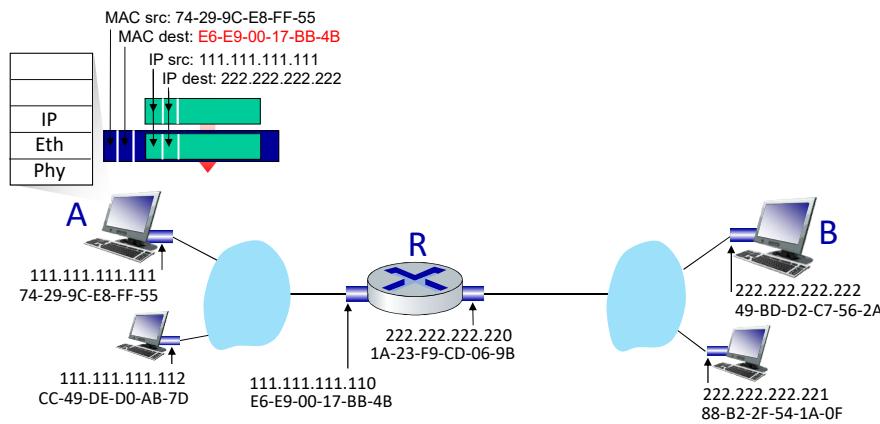
- focus on addressing – at IP (datagram) and MAC layer (frame) levels
- assume that:
  - A knows B's IP address
  - A knows IP address of first hop router, R (how?)
  - A knows R's MAC address (how?)



Link Layer: 6-65

## Routing to another subnet: addressing

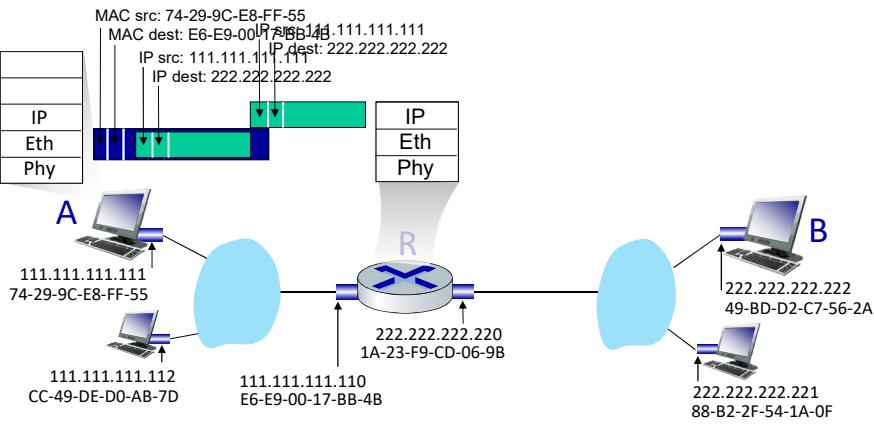
- A creates IP datagram with IP source A, destination B
- A creates link-layer frame containing A-to-B IP datagram
  - R's MAC address is frame's destination



Link Layer: 6-66

## Routing to another subnet: addressing

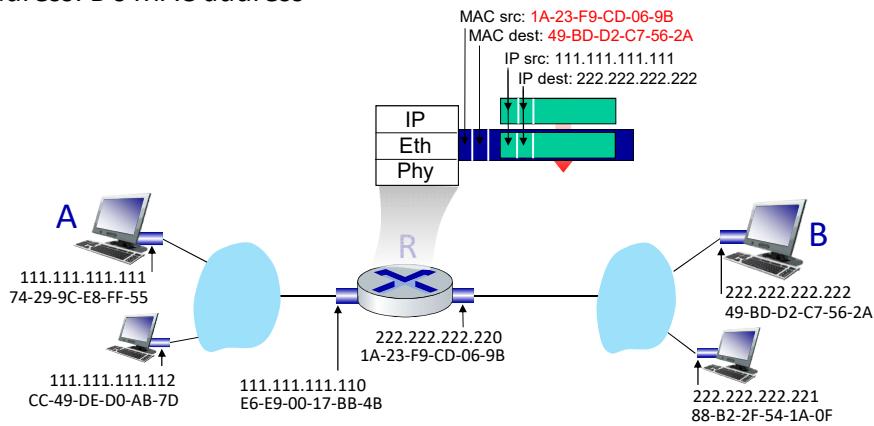
- frame sent from A to R
- frame received at R, datagram removed, passed up to IP



Link Layer: 6-67

## Routing to another subnet: addressing

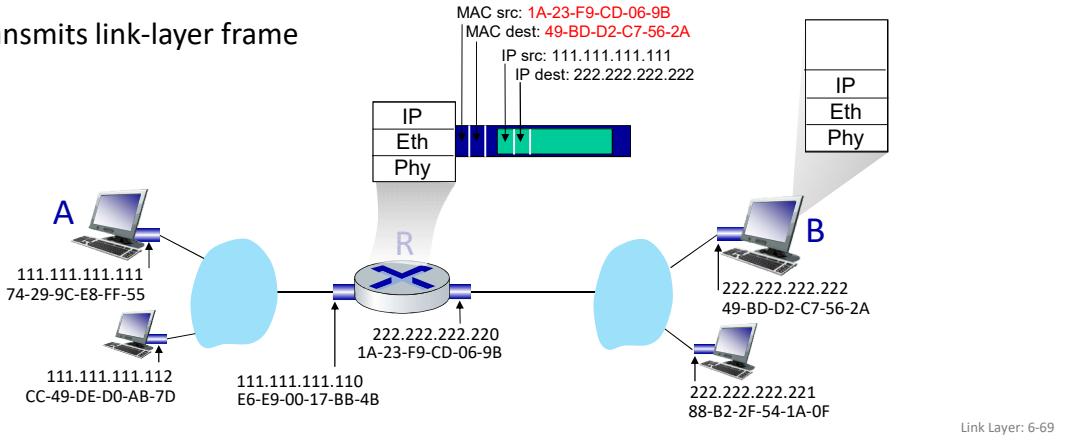
- R determines outgoing interface, passes datagram with IP source A, destination B to link layer
- R creates link-layer frame containing A-to-B IP datagram. Frame destination address: B's MAC address



Link Layer: 6-68

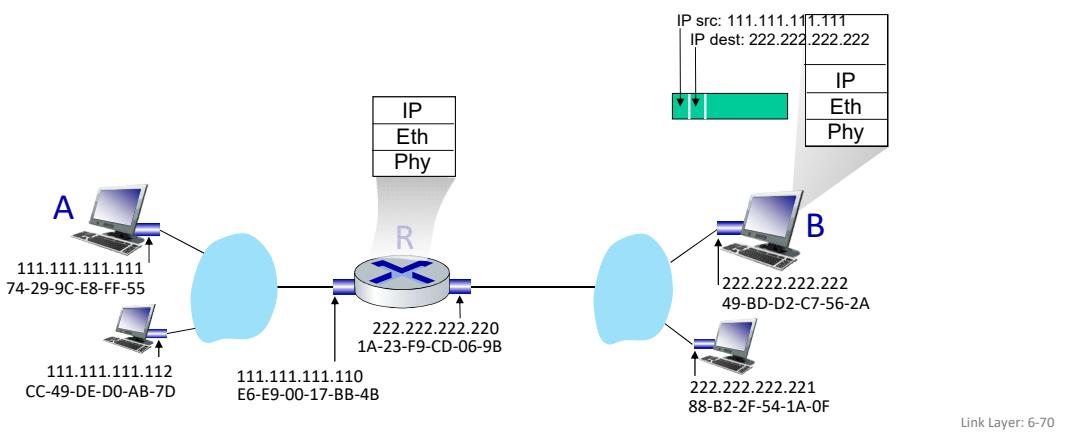
## Routing to another subnet: addressing

- R determines outgoing interface, passes datagram with IP source A, destination B to link layer
- R creates link-layer frame containing A-to-B IP datagram. Frame destination address: B's MAC address
- transmits link-layer frame



## Routing to another subnet: addressing

- B receives frame, extracts IP datagram destination B
- B passes datagram up protocol stack to IP



## Link layer, LANs: roadmap

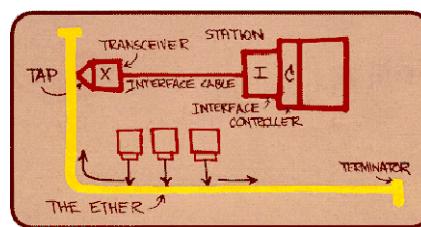
- introduction
- error detection, correction
- multiple access protocols
- LANs**
  - addressing, ARP
  - **Ethernet**
  - switches
  - VLANs
- link virtualization: MPLS
- data center networking
- a day in the life of a web request

Link Layer: 6-71

## Ethernet

“dominant” wired LAN technology:

- first widely used LAN technology
- simpler, cheap
- kept up with speed race: 10 Mbps – 400 Gbps
- single chip, multiple speeds (e.g., Broadcom BCM5761)



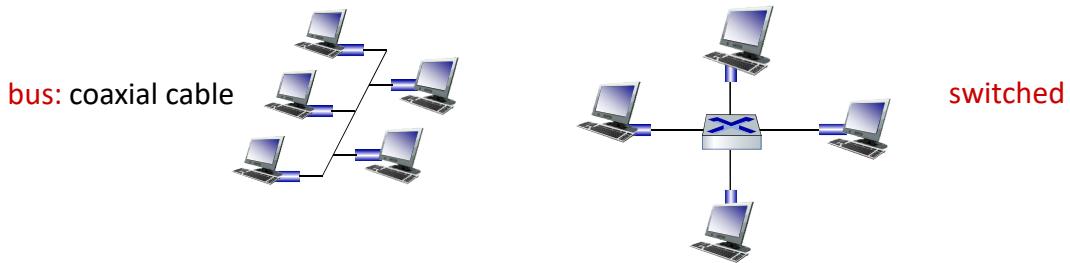
*Metcalfe's Ethernet sketch*

<https://www.uspto.gov/learning-and-resources/journeys-innovation/audio-stories/defying-doubters>

Link Layer: 6-72

## Ethernet: physical topology

- **bus:** popular through mid 90s
  - all nodes in same collision domain (can collide with each other)
- **switched:** prevails today
  - active link-layer 2 *switch* in center
  - each “spoke” runs a (separate) Ethernet protocol (nodes do not collide with each other)



Link Layer: 6-73

## Ethernet frame structure

sending interface encapsulates IP datagram (or other network layer protocol packet) in **Ethernet frame**

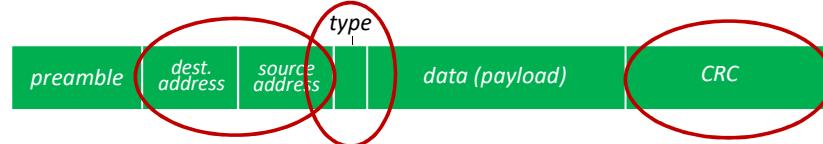
<i>type</i>					
<i>preamble</i>	<i>dest. address</i>	<i>source address</i>		<i>data (payload)</i>	<i>CRC</i>

*preamble:*

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

Link Layer: 6-74

## 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
  - used to demultiplex up at receiver
- **CRC:** cyclic redundancy check at receiver
  - error detected: frame is dropped

Link Layer: 6-75

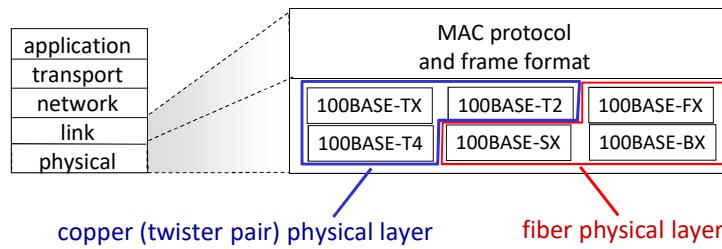
## Ethernet: unreliable, connectionless

- **connectionless:** no handshaking between sending and receiving NICs
- **unreliable:** receiving NIC doesn't send ACKs or NAKs 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 with binary backoff**

Link Layer: 6-76

## 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, 10 Gbps, 40 Gbps
  - different physical layer media: fiber, cable



Link Layer: 6-77

## Link layer, LANs: roadmap

- introduction
- error detection, correction
- multiple access protocols
- **LANs**
  - addressing, ARP
  - Ethernet
  - **switches**
  - VLANs
- link virtualization: MPLS
- data center networking
- a day in the life of a web request

Link Layer: 6-78

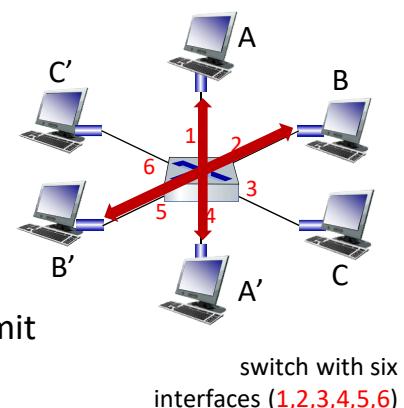
## Ethernet switch

- Switch is a **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 *unaware* of presence of switches
- **plug-and-play, self-learning**
  - switches do not need to be configured

Link Layer: 6-79

## Switch: multiple simultaneous transmissions

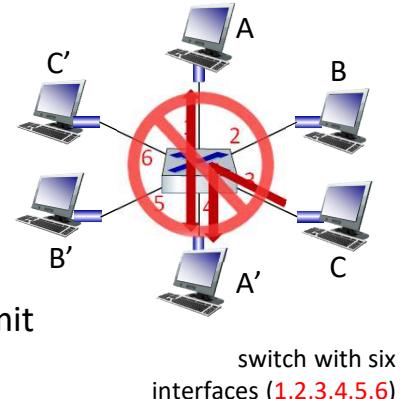
- hosts have dedicated, direct connection to switch
- switches buffer packets
- Ethernet protocol used on *each* incoming link, so:
  - no collisions; full duplex
  - each link is its own collision domain
- **switching**: A-to-A' and B-to-B' can transmit simultaneously, without collisions



Link Layer: 6-80

## Switch: multiple simultaneous transmissions

- hosts have dedicated, direct connection to switch
- switches buffer packets
- Ethernet protocol used on *each* incoming link, so:
  - no collisions; full duplex
  - each link is its own collision domain
- **switching:** A-to-A' and B-to-B' can transmit simultaneously, without collisions
  - but A-to-A' and C to A' can *not* happen simultaneously



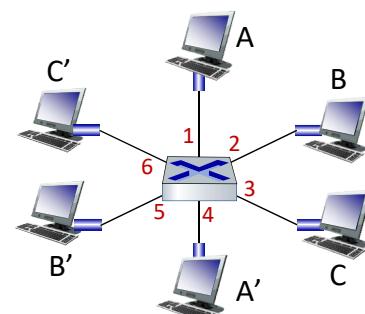
Link Layer: 6-81

## 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!



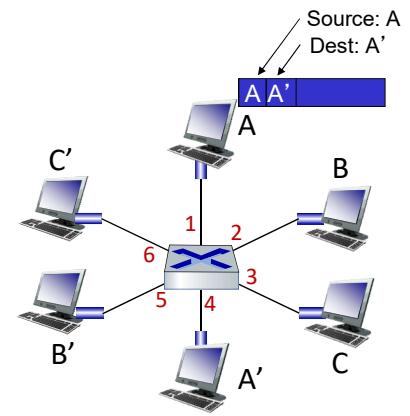
**Q:** how are entries created, maintained in switch table?

- something like a routing protocol?

Link Layer: 6-82

## 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)

Link Layer: 6-83

## Switch: frame filtering/forwarding

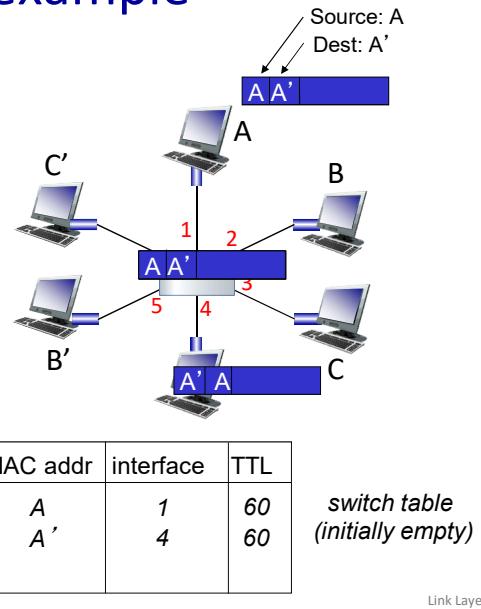
when frame received at switch:

1. record incoming link, MAC address of sending host
2. index switch table using MAC destination address
3. if entry found for destination
  - then {
  - if destination on segment from which frame arrived
    - then drop frame
    - else forward frame on interface indicated by entry
  - }
  - else flood /\* forward on all interfaces except arriving interface \*/

Link Layer: 6-84

## Self-learning, forwarding: example

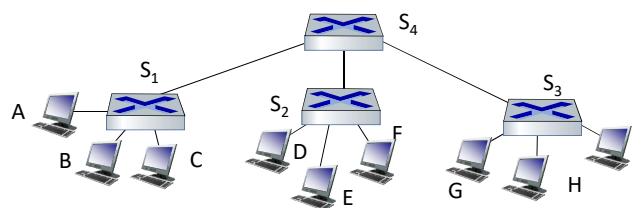
- frame destination, A', location unknown: **flood**
- destination A location known: **selectively send on just one link**



Link Layer: 6-85

## Interconnecting switches

self-learning switches can be connected together:



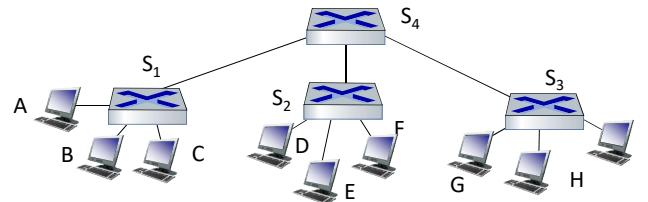
Q: sending from A to G - how does S<sub>1</sub> know to forward frame destined to G via S<sub>4</sub> and S<sub>3</sub>?

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

Link Layer: 6-86

## Self-learning multi-switch example

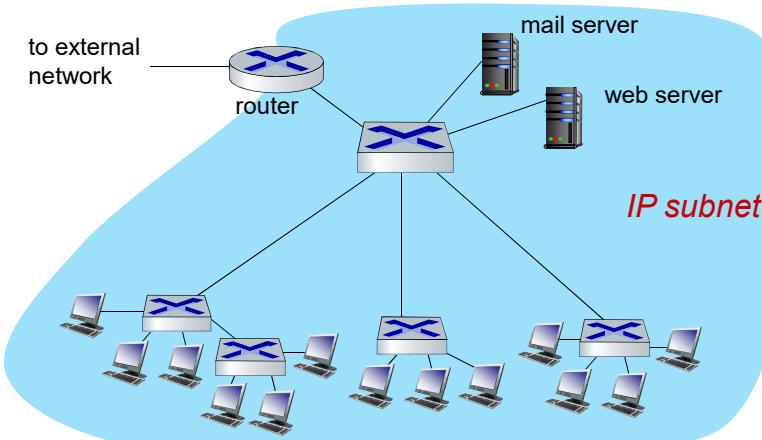
Suppose C sends frame to I, I responds to C



Q: show switch tables and packet forwarding in S<sub>1</sub>, S<sub>2</sub>, S<sub>3</sub>, S<sub>4</sub>

Link Layer: 6-87

## Small institutional network



Link Layer: 6-88

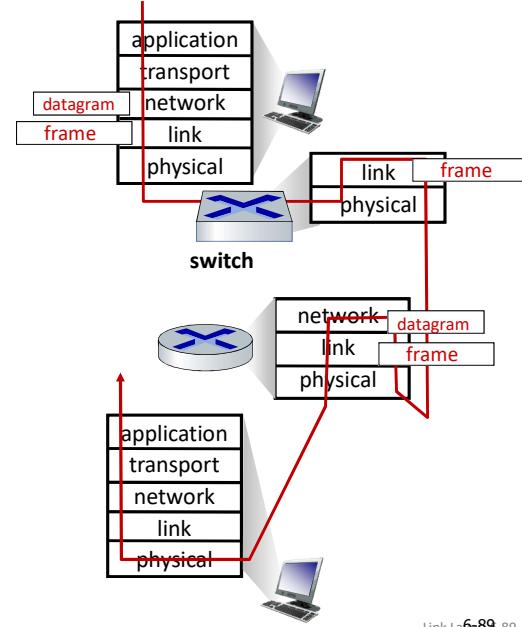
## Switches vs. routers

**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



Link La6-89-89

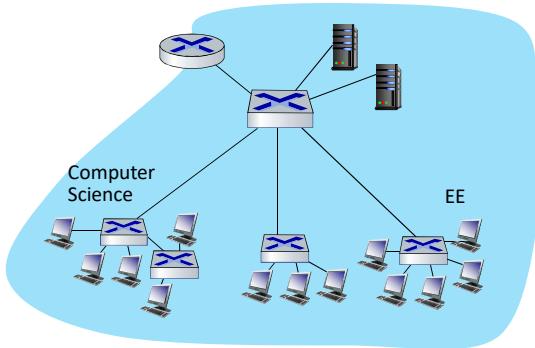
## Link layer, LANs: roadmap

- introduction
- error detection, correction
- multiple access protocols
- **LANs**
  - addressing, ARP
  - Ethernet
  - switches
  - **VLANs**
- link virtualization: MPLS
- data center networking
- a day in the life of a web request

Link Layer: 6-90

## Virtual LANs (VLANs): motivation

**Q:** what happens as LAN sizes scale, users change point of attachment?



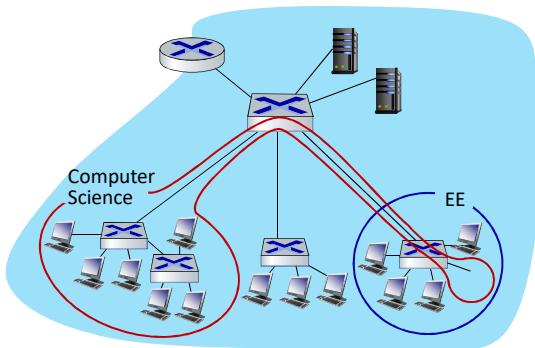
single broadcast domain:

- **scaling:** all layer-2 broadcast traffic (ARP, DHCP, unknown MAC) must cross entire LAN
  - efficiency, security, privacy issues

Link Layer: 6-91

## Virtual LANs (VLANs): motivation

**Q:** what happens as LAN sizes scale, users change point of attachment?



single broadcast domain:

- **scaling:** all layer-2 broadcast traffic (ARP, DHCP, unknown MAC) must cross entire LAN
  - efficiency, security, privacy, efficiency issues

### **administrative issues:**

- CS user moves office to EE - *physically* attached to EE switch, but wants to remain *logically* attached to CS switch

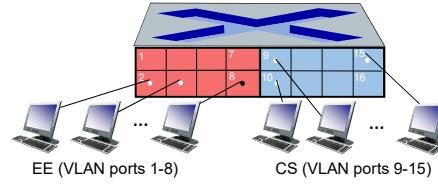
Link Layer: 6-92

## Port-based VLANs

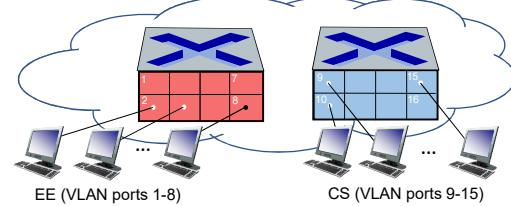
### Virtual Local Area Network (VLAN)

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

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



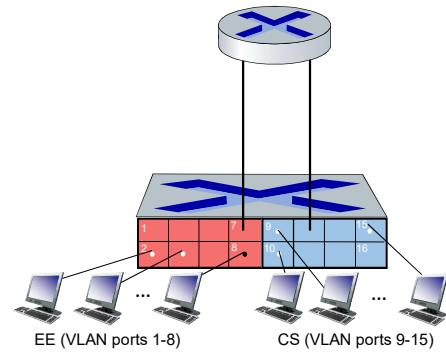
... operates as **multiple virtual switches**



Link Layer: 6-93

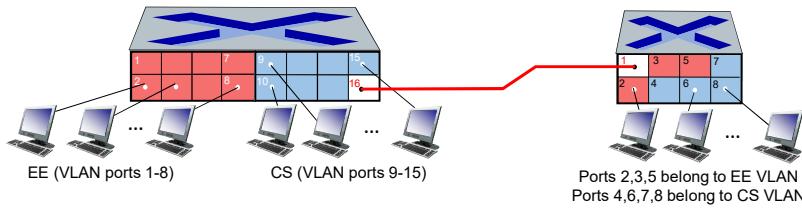
## Port-based VLANs

- **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



Link Layer: 6-94

## VLANs spanning multiple switches

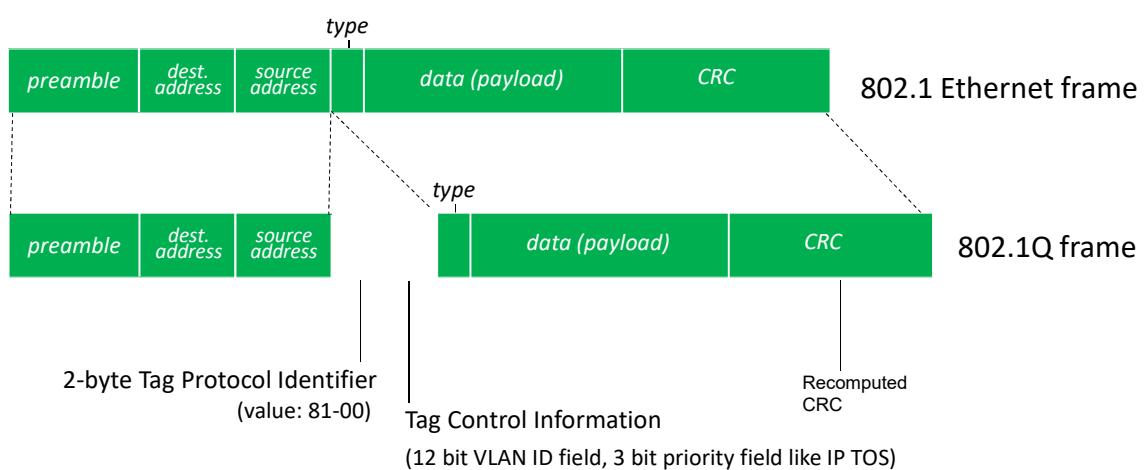


**trunk port:** carries frames between VLANs defined over multiple physical switches

- 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

Link Layer: 6-95

## 802.1Q VLAN frame format



Link Layer: 6-96

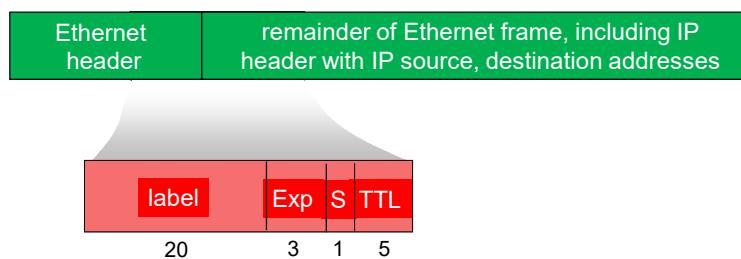
## Link layer, LANs: roadmap

- introduction
- error detection, correction
- multiple access protocols
- LANs
  - addressing, ARP
  - Ethernet
  - switches
  - VLANs
- link virtualization: MPLS
- data center networking
- a day in the life of a web request

Link Layer: 6-97

## Multiprotocol label switching (MPLS)

- **goal:** high-speed IP forwarding among network of MPLS-capable routers, using fixed length label (instead of shortest prefix matching)
  - faster lookup using fixed length identifier
  - borrowing ideas from Virtual Circuit (VC) approach
  - but IP datagram still keeps IP address!



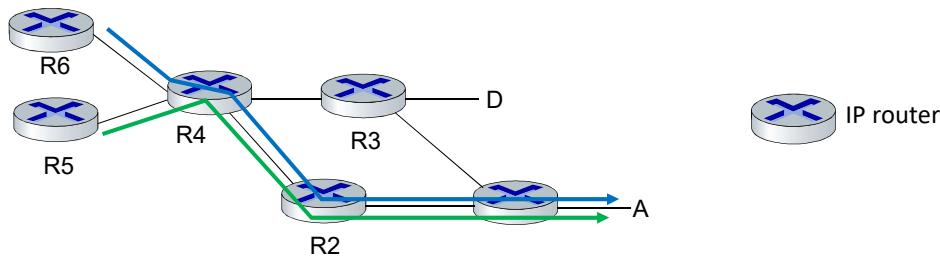
Link Layer: 6-98

## 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

Link Layer: 6-99

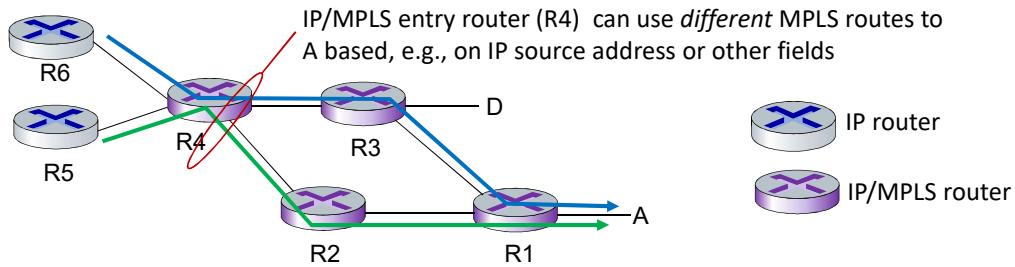
## MPLS versus IP paths



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

Link Layer: 6-100

## MPLS versus IP paths

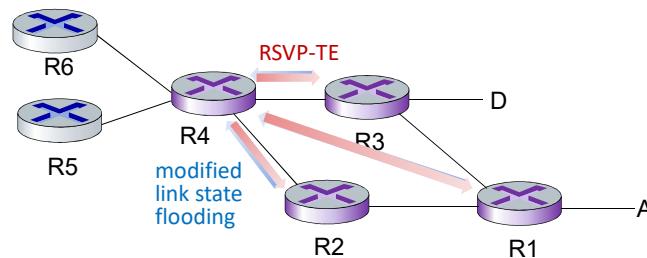


- **IP routing:** path to destination determined by destination address alone
- **MPLS routing:** path to destination can be based on source *and* destination address
  - flavor of generalized forwarding (MPLS 10 years earlier)
  - *fast reroute:* precompute backup routes in case of link failure

Link Layer: 6-101

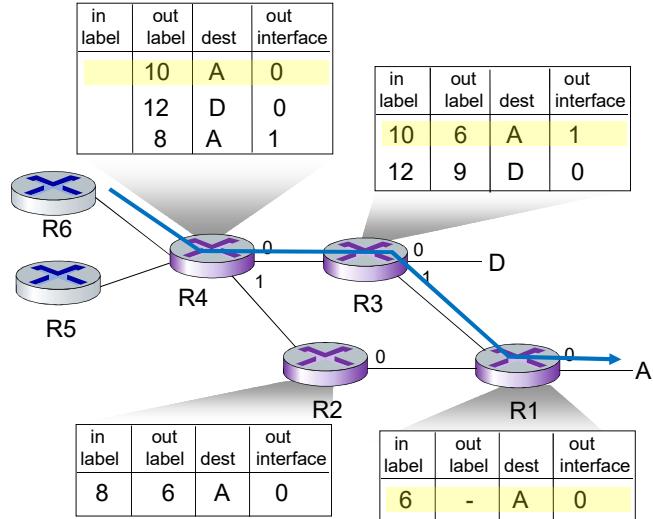
## MPLS signaling

- modify OSPF, IS-IS 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



Link Layer: 6-102

## MPLS forwarding tables



Link Layer: 6-103

## Link layer, LANs: roadmap

- introduction
  - error detection, correction
  - multiple access protocols
  - LANs
    - addressing, ARP
    - Ethernet
    - switches
    - VLANs
  - link virtualization: MPLS
  - data center networking
- a day in the life of a web request

Link Layer: 6-104

## Datacenter 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:

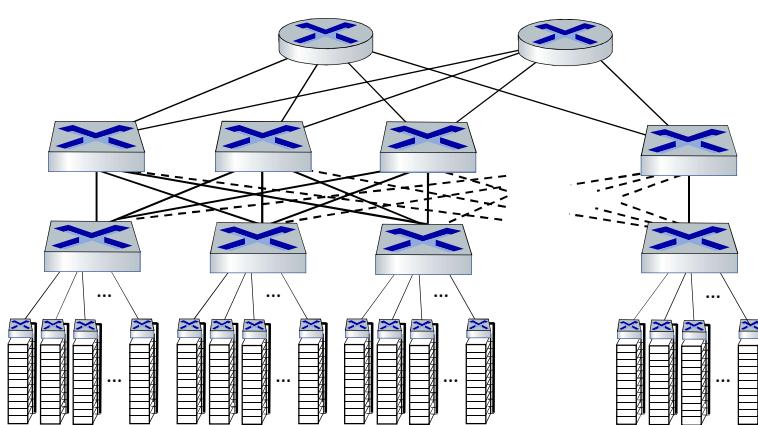
- multiple applications, each serving massive numbers of clients
- reliability
- managing/balancing load, avoiding processing, networking, data bottlenecks



Inside a 40-ft Microsoft container, Chicago data center

Link Layer: 6-105

## 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

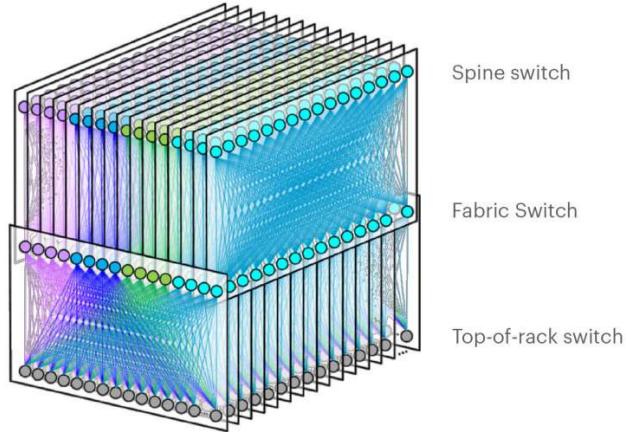
### Server racks

- 20- 40 server blades: hosts

Link Layer: 6-106

## Datacenter networks: network elements

Facebook F16 data center network topology:

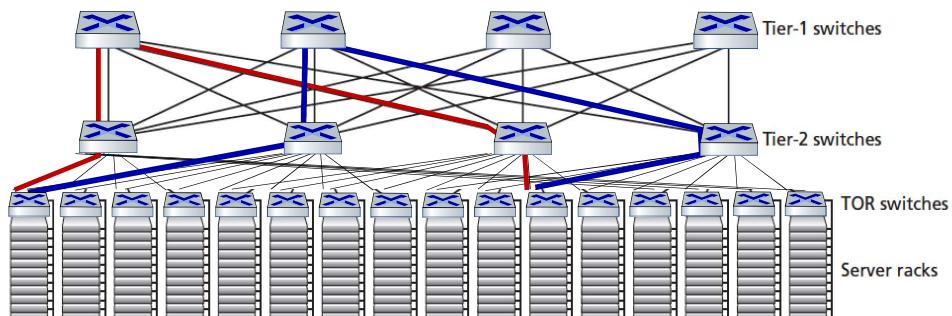


<https://engineering.fb.com/data-center-engineering/f16-minipack/> (posted 3/2019)

Link Layer: 6-107

## Datacenter networks: multipath

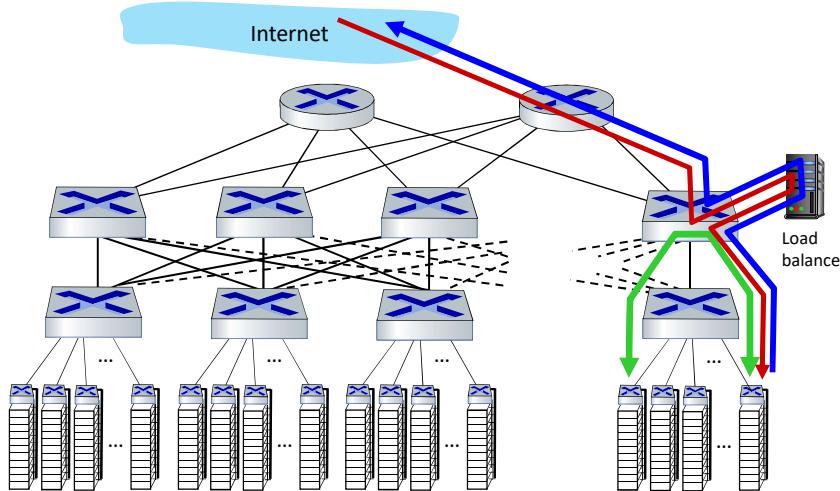
- rich interconnection among switches, racks:
  - increased throughput between racks (multiple routing paths possible)
  - increased reliability via redundancy



two **disjoint** paths highlighted between racks 1 and 11

Link Layer: 6-108

## 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: 6-109

## 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: 6-110

## Link layer, LANs: roadmap

- introduction
- error detection, correction
- multiple access protocols
- LANs
  - addressing, ARP
  - Ethernet
  - switches
  - VLANs
- link virtualization: MPLS
- data center networking
- a day in the life of a web request

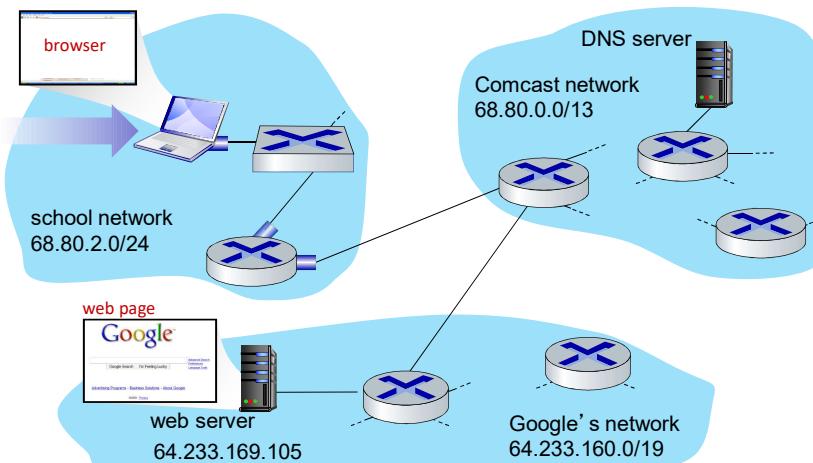
Link Layer: 6-111

## Synthesis: a day in the life of a web request

- our journey down the protocol stack is now 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

Link Layer: 6-112

## A day in the life: scenario



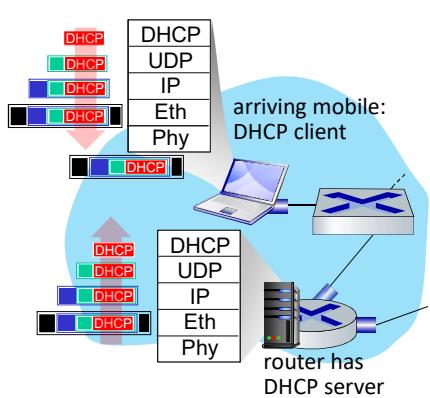
scenario:

- arriving mobile client attaches to network ...
- requests web page: [www.google.com](http://www.google.com)

*Sounds simple!*

Link Layer: 6-113

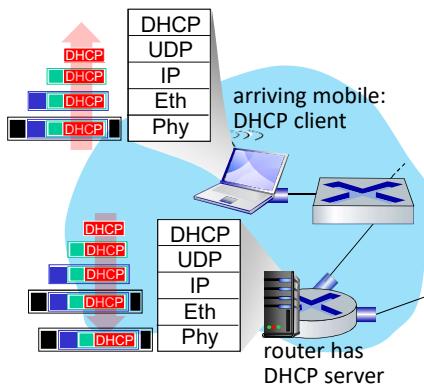
## A day in the life: connecting to the Internet



- 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

Link Layer: 6-114

## A day in the life: connecting to the Internet

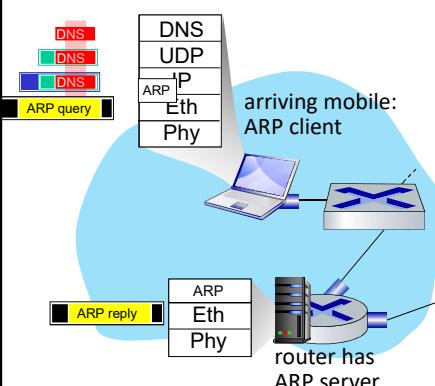


- DHCP server formulates **DHCP ACK** containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server
- encapsulation at DHCP server, frame forwarded (**switch learning**) through LAN, demultiplexing at client
- DHCP client receives DHCP ACK reply

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

Link Layer: 6-115

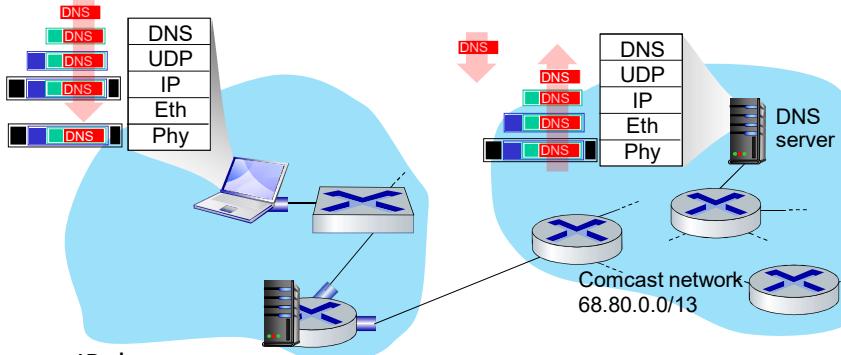
## A day in the life... ARP (before DNS, before HTTP)



- 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

Link Layer: 6-116

## A day in the life... using DNS

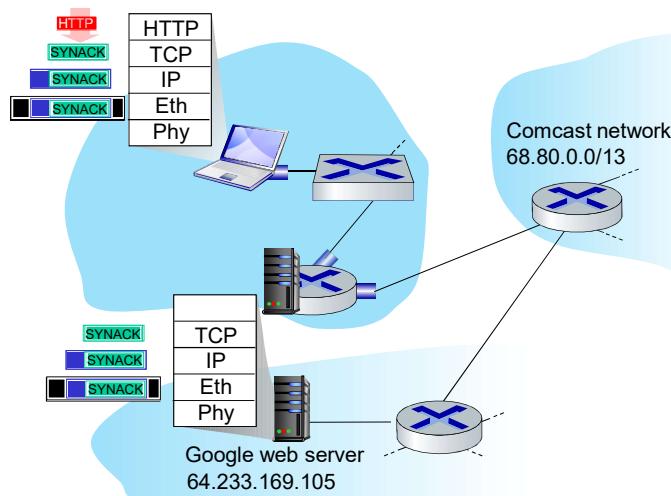


- IP datagram containing DNS query forwarded via LAN switch from client to 1<sup>st</sup> hop router
- IP datagram forwarded from campus network into Comcast network, routed (tables created by **RIP, OSPF, IS-IS** and/or **BGP** routing protocols) to DNS server

- demuxed to DNS
- DNS replies to client with IP address of [www.google.com](http://www.google.com)

Link Layer: 6-117

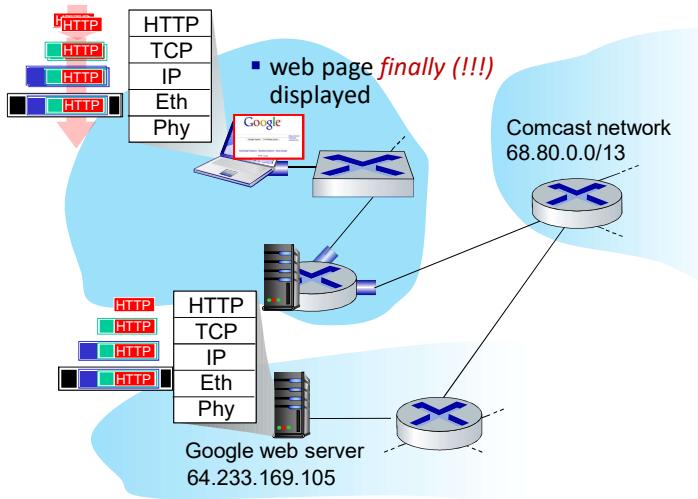
## A day in the life...TCP connection carrying HTTP



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

Link Layer: 6-118

## A day in the life... HTTP request/reply



- **HTTP request** sent into TCP socket
- IP datagram containing HTTP request routed to [www.google.com](http://www.google.com)
- web server responds with **HTTP reply** (containing web page)
- IP datagram containing HTTP reply routed back to client

Link Layer: 6-119

## Chapter 6: Summary

- principles behind data link layer services:
  - error detection, correction
  - sharing a broadcast channel: multiple access
  - link layer addressing
- instantiation, 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

Link Layer: 6-120

## Chapter 6: let's take a breath

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

Link Layer: 6-121

## Additional Chapter 6 slides

Network Layer: 5-122

## Pure ALOHA efficiency

$$\begin{aligned} P(\text{success by given node}) &= P(\text{node transmits}) * \\ &\quad P(\text{no other node transmits in } [t_0-1, t_0] *) * \\ &\quad P(\text{no other node transmits in } [t_0-1, t_0]) \\ &= p \cdot (1-p)^{N-1} \cdot (1-p)^{N-1} \\ &= p \cdot (1-p)^{2(N-1)} \\ \dots \text{ choosing optimum } p \text{ and then letting } n \\ &= 1/(2e) = .18 \rightarrow \infty \end{aligned}$$

even worse than slotted Aloha!