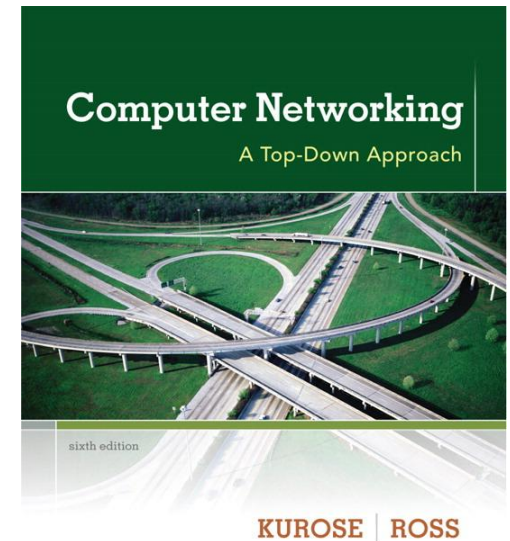# Chapter 5
# Link Layer

**A note on the use of these ppt 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*
*6th edition*
*Jim Kurose, Keith Ross*
*Addison-Wesley*
*March 2012*

# Chapter 5: 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

5.1 introduction, services

5.2 error detection, correction

5.3 multiple access protocols

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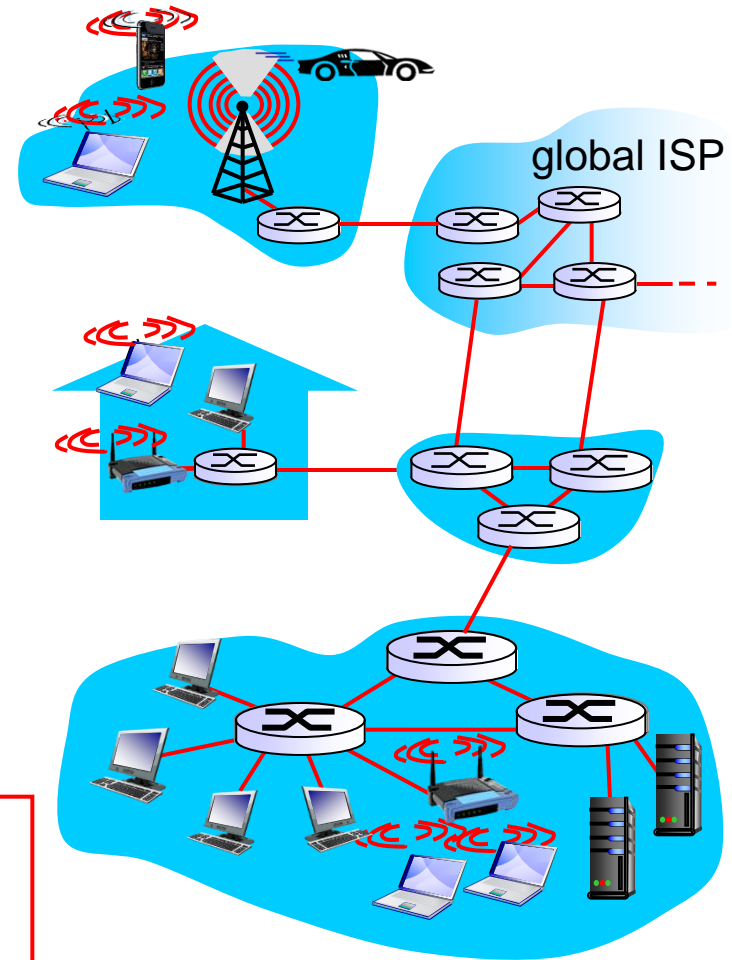- addressing, ARP
- Ethernet
- switches
- VLANS

5.5 Subnetting

5.6 a day in the life of a web request

# Link layer: introduction

*terminology:*

❖ hosts and routers: nodes

❖ communication channels that connect adjacent nodes along communication path: links

  ▪ wired links

  ▪ wireless links

  ▪ LANs

❖ layer-2 packet: frame, encapsulates datagram

global ISP

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

# Link layer: context

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

*transportation analogy:*

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

# Link layer services

❖ *framing, link access:*
- encapsulate datagram into frame, adding header, trailer
- channel access if shared medium
- "MAC" addresses used in frame headers to identify source, dest
  - different from IP address!

❖ *reliable delivery between adjacent nodes*
- we learned how to do this already (chapter 3)!
- seldom used on low bit-error link (fiber, some twisted pair)
- wireless links: high error rates
  - *Q:* why both link-level and end-end reliability?

# Link layer services (more)

❖ *flow control:*
  ▪ pacing between adjacent sending and receiving nodes

❖ *error detection:*
  ▪ errors caused by signal attenuation, noise.
  ▪ receiver detects presence of errors:
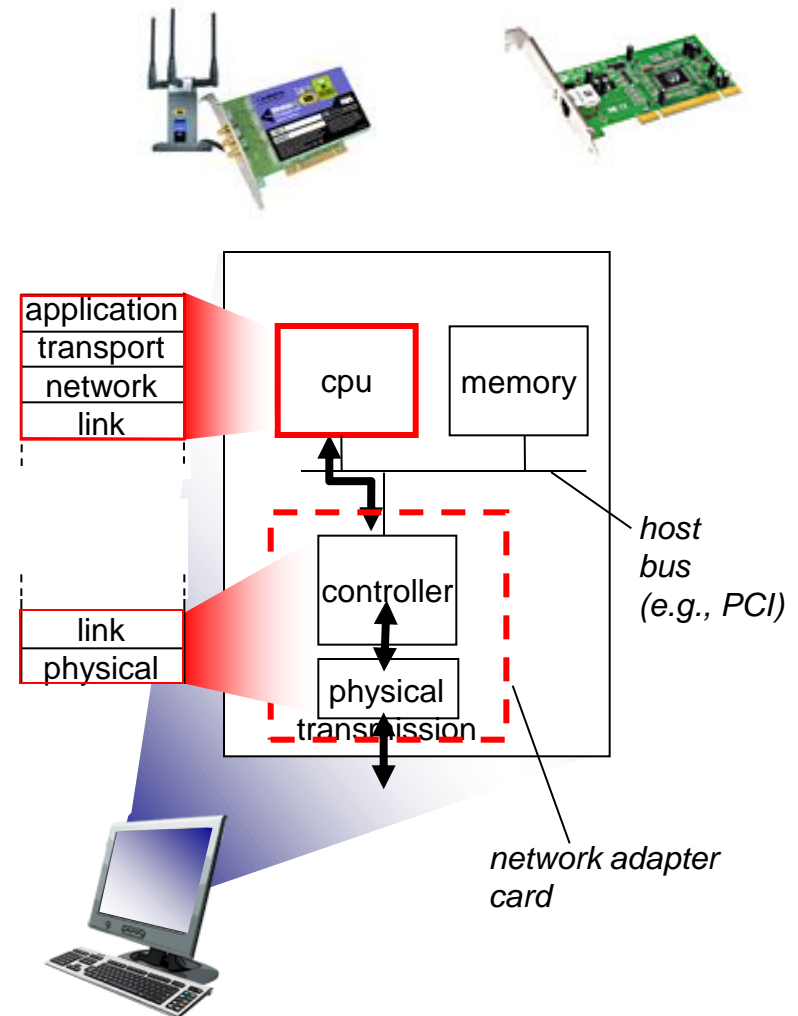    • signals sender for retransmission or drops frame

❖ *error correction:*
  ▪ receiver identifies *and corrects* bit error(s) without resorting to retransmission

❖ *half-duplex and full-duplex*
  ▪ with half duplex, nodes at both ends of link can transmit, but not at same time

# 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

controller

host bus (e.g., PCI)

link
physical

physical

transmission

network adapter card

# Adaptors communicating



*sending host*        *receiving host*

*frame*

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

# Link layer, LANs: outline
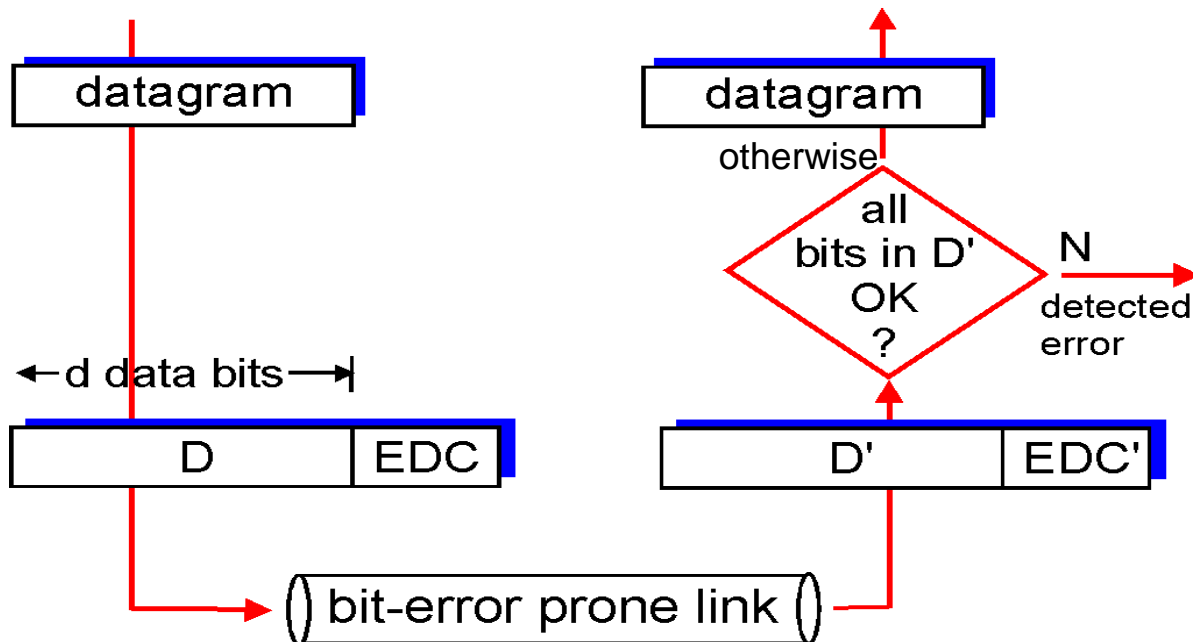
# Error detection

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

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

# Parity checking

*single bit parity:*

❖ *d*etect single bit errors

*two-dimensional bit parity:*

❖ detect and correct single bit errors

←— d data bits —→ | parity bit

| 011100110101011 | 0 |

row parity →

$$\begin{array}{cccc} d_{1,1} & \cdots & d_{1,j} & d_{1,\,j+1} \\ d_{2,1} & \cdots & d_{2,j} & d_{2,j+1} \\ \cdots & \cdots & \cdots & \cdots \\ d_{i,1} & \cdots & d_{i,j} & d_{i,j+1} \\ \hline d_{i+1,1} & \cdots & d_{i+1,j} & d_{i+1,j+1} \end{array}$$

column parity ↓

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

*no errors*

parity error ↓

*correctable single bit error*

# Internet checksum (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?*

# UDP checksum

**Στόχος:** ανίχνευση σφαλμάτων (ανεστραμμένων bits) στο λαμβανόμενο segment

## Αποστολέας:

- χειρίζεται το περιεχόμενο του segment ως ακολουθία ακεραίων των 16 bits
- checksum: συμπλήρωμα ως προς 1 του αθροίσματος των περιεχομένων του segment
- ο αποστολέας τοποθετεί την τιμή του checksum στο πεδίο checksum του UDP segment

## Παραλήπτης:

- υπολογίζει το checksum του λαμβανόμενου segment
- ελέγχει εάν η υπολογισθείσα τιμή του checksum είναι ίση με την τιμή στο πεδίο checksum:
  - Όχι – ανίχνευση σφάλματος
    - απόρριψη segment
    - προώθηση με προειδοποίηση
  - Ναι – δεν έχει ανιχνευθεί κανένα σφάλμα

# UDP checksum

Έστω segment:

0110011001100110 0101010101010101

0000111100001111

## Αποστολέας:

πρόσθεση

0110011001100110

+ 0101010101010101

1011101110111011

0000111100001111

———————————————

1100101011001010

CKS 0011010100110101

## Παραλήπτης:

πρόσθεση

0110011001100110

0101010101010101

+ 1000111100001111

σφάλμα 0011010100110101

———————————————

0111111111111111

# Cyclic redundancy check

❖ more powerful error-detection coding

❖ view data bits, D, as a binary number

❖ choose r+1 bit pattern (generator), G

❖ goal: choose r CRC bits, R, such that
  ▪ <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)

$\longleftarrow$ d bits $\longrightarrow$  $\longleftarrow$ r bits $\longrightarrow$

| D: data bits to be sent | R: CRC bits |

*bit pattern*

$$D * 2^r \quad XOR \quad R$$

*mathematical formula*

# CRC example

want:

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

*equivalently:*

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

*equivalently:*

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

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

```
11010011101100 000 <--- input right padded by 3 bits
1011                <--- divisor
01100011101100 000 <--- result
 1011               <--- divisor ...
00111011101100 000
   1011
00010111101100 000
    1011
0000000110100 000
        1011
0000000110100 000
         1011
00000000011000 000
          1011
00000000001110 000
           1011
0000000000101 000
            101 1
----------------
00000000000000 100 <--- remainder (3 bits)
```

# Hamming Code

❖ In telecommunication, Hamming codes are a family of linear error-correcting codes that generalize the Hamming(7,4)-code invented by Richard Hamming in 1950. Hamming codes can detect up to two-bit errors or correct one-bit errors without detection of uncorrected errors. By contrast, the simple parity code cannot correct errors, and can detect only an odd number of bits in error. Hamming codes are perfect codes, that is, they achieve the highest possible rate for codes with their block length and minimum distance 3.

## Calculating the Hamming Code

❖ The key to the Hamming Code is the use of extra parity bits to allow the identification of a single error. Create the code word as follows:

❖    Mark all bit positions that are powers of two as parity bits. (positions 1, 2, 4, 8, 16, 32, 64, etc.)

❖    All other bit positions are for the data to be encoded. (positions 3, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 17, etc.)

❖    Each parity bit calculates the parity for some of the bits in the code word. The position of the parity bit determines the sequence of bits that it alternately checks and skips.

# Hamming Code

❖ Position 1: check 1 bit, skip 1 bit, check 1 bit, skip 1 bit, etc. (1,3,5,7,9,11,13,15,...)

❖ Position 2: check 2 bits, skip 2 bits, check 2 bits, skip 2 bits, etc. (2,3,6,7,10,11,14,15,...)

❖ Position 4: check 4 bits, skip 4 bits, check 4 bits, skip 4 bits, etc. (4,5,6,7,12,13,14,15,20,21,22,23,...)

❖ Position 8: check 8 bits, skip 8 bits, check 8 bits, skip 8 bits, etc. (8-15,24-31,40-47,...)

❖ Position 16: check 16 bits, skip 16 bits, check 16 bits, skip 16 bits, etc. (16-31,48-63,80-95,...)

❖ Position 32: check 32 bits, skip 32 bits, check 32 bits, skip 32 bits, etc. (32-63,96-127,160-191,...)

❖ etc.

❖ Set a parity bit to 1 if the total number of ones in the positions it checks is odd. Set a parity bit to 0 if the total number of ones in the positions it checks is even.

# Example

❖ A byte of data: 10011010

❖ Create the data word, leaving spaces for the parity bits: _ _ 1 _ 0 0 1 _ 1 0 1 0

❖ Calculate the parity for each parity bit (a ? represents the bit position being set):

❖ Position 1 checks bits 1,3,5,7,9,11:

❖ ? _ 1 _ 0 0 1 _ 1 0 1 0. Even parity so set position 1 to a 0: 0 _ 1 _ 0 0 1 _ 1 0 1 0

❖ Position 2 checks bits 2,3,6,7,10,11:

❖ 0 ? 1 _ 0 0 1 _ 1 0 1 0. Odd parity so set position 2 to a 1: 0 1 1 _ 0 0 1 _ 1 0 1 0

❖ Position 4 checks bits 4,5,6,7,12:

❖ 0 1 1 ? 0 0 1 _ 1 0 1 0. Odd parity so set position 4 to a 1: 0 1 1 1 0 0 1 _ 1 0 1 0

❖ Position 8 checks bits 8,9,10,11,12:

❖ 0 1 1 1 0 0 1 ? 1 0 1 0. Even parity so set position 8 to a 0: 0 1 1 1 0 0 1 0 1 0 1 0

❖ Code word: 011100101010.

# Example

❖ Finding and fixing a bad bit

❖ The above example created a code word of 01110010101010. Suppose the word that was received was 01110010111110 instead. Then the receiver could calculate which bit was wrong and correct it. The method is to verify each check bit. Write down all the incorrect parity bits. Doing so, you will discover that parity bits 2 and 8 are incorrect. It is not an accident that 2 + 8 = 10, and that bit position 10 is the location of the bad bit. In general, check each parity bit, and add the positions that are wrong, this will give you the location of the bad bit.

# Αρχές αξιόπιστης μεταφοράς δεδομένων

# Principles of reliable data transfer

❖ **important in application, transport, link layers**
 ▪ top-10 list of important networking topics!



(a) provided service

❖ **characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)**

# Principles of reliable data transfer

❖ **important in application, transport, link layers**
  ▪ top-10 list of important networking topics!



(a) provided service     (b) service implementation

❖ **characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)**

# Principles of reliable data transfer

❖ **important in application, transport, link layers**
  ▪ top-10 list of important networking topics!

sending process — receiver process

data — data

reliable channel

rdt_send() — data

reliable data transfer protocol (sending side)

udt_send() — packet

data — deliver_data()

reliable data transfer protocol (receiving side)

packet — rdt_rcv()

unreliable channel

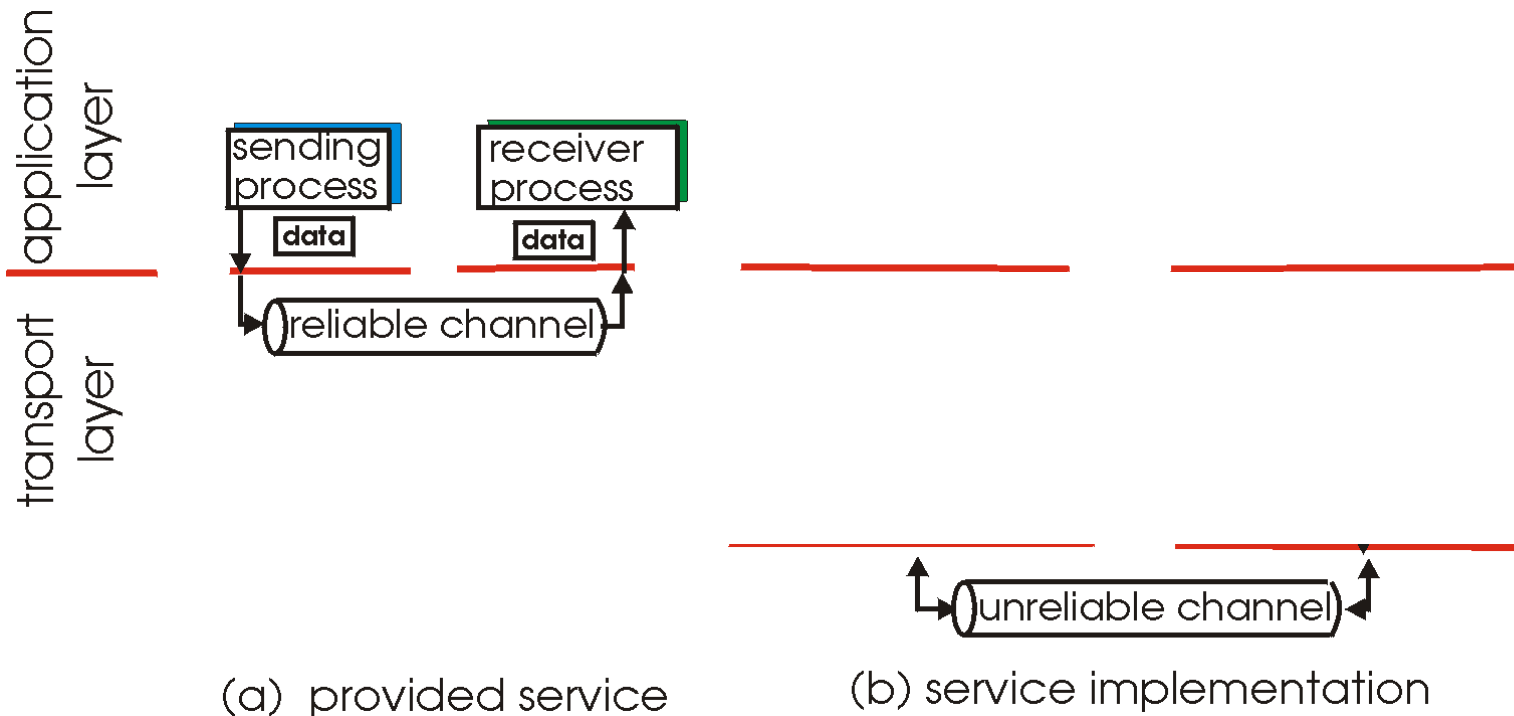(a) provided service                    (b) service implementation

❖ **characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)**

# Reliable data transfer: getting started

**rdt_send():** called from above, (e.g., by app.). Passed data to deliver to receiver upper layer

**deliver_data():** called by **rdt** to deliver data to upper

**rdt_send()** ↓ data

data ↑ **deliver_data()**

send side

reliable data transfer protocol (sending side)

reliable data transfer protocol (receiving side)

receive side

**udt_send()** ↕ packet

packet ↕ **rdt_rcv()**

unreliable channel

**udt_send():** called by rdt, to transfer packet over unreliable channel to receiver

**rdt_rcv():** called when packet arrives on rcv-side of channel

# Reliable data transfer: getting started

we'll:

❖ incrementally develop sender, receiver sides of **r**eliable **d**ata **t**ransfer protocol (rdt)

❖ consider only unidirectional data transfer
  ▪ but control info will flow on both directions!

❖ use finite state machines (FSM)  to specify sender, receiver

state: when in this "state" next state uniquely determined by next event

event causing state transition
actions taken on state transition

state 1

event
actions

state 2

# rdt1.0: reliable transfer over a reliable channel

❖ **underlying channel perfectly reliable**
  - no bit errors
  - no loss of packets
❖ **separate FSMs for sender, receiver:**
  - sender sends data into underlying channel
  - receiver reads data from underlying channel

Wait for
call from
above
$\underline{\text{rdt\_send(data)}}$
packet = make_pkt(data)
udt_send(packet)

Wait for
call from
below
$\underline{\text{rdt\_rcv(packet)}}$
extract (packet,data)
deliver_data(data)

sender                              receiver

# rdt2.0: channel with bit errors

❖ underlying channel may flip bits in packet
  ▪ checksum to detect bit errors
❖ *the* question: how to recover from errors:

*How do humans recover from "errors"
during conversation?*

# rdt2.0: channel with bit errors

❖ **underlying channel may flip bits in packet**
  ▪ checksum to detect bit errors
❖ *the* **question: how to recover from errors:**

  ▪ *acknowledgements (ACKs):* receiver explicitly tells sender that pkt received OK

  ▪ *negative acknowledgements (NAKs):* receiver explicitly tells sender that pkt had errors

  ▪ sender retransmits pkt on receipt of NAK
❖ **new mechanisms in `rdt2.0` (beyond `rdt1.0`):**
  ▪ error detection
  ▪ feedback: control msgs (ACK,NAK) from receiver to sender

# rdt2.0: FSM specification

rdt_send(data)
-------------------
sndpkt = make_pkt(data, checksum)
udt_send(sndpkt)

receiver

Wait for call from above

Wait for ACK or NAK

rdt_rcv(rcvpkt) &&
   isNAK(rcvpkt)
-------------------
udt_send(sndpkt)

rdt_rcv(rcvpkt) && isACK(rcvpkt)
-------------------
$\Lambda$

sender

rdt_rcv(rcvpkt) &&
   corrupt(rcvpkt)
-------------------
udt_send(NAK)

Wait for call from below

rdt_rcv(rcvpkt) &&
   notcorrupt(rcvpkt)
-------------------
extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

# rdt2.0: operation with no errors

rdt_send(data)
_____
snkpkt = make_pkt(data, checksum)
udt_send(sndpkt)

**Wait for call from above**

**Wait for ACK or NAK**

rdt_rcv(rcvpkt) && isNAK(rcvpkt)
_____
udt_send(sndpkt)

rdt_rcv(rcvpkt) && corrupt(rcvpkt)
_____
udt_send(NAK)

rdt_rcv(rcvpkt) && isACK(rcvpkt)
_____
$\Lambda$

**Wait for call from below**

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

# rdt2.0: error scenario

rdt_send(data)
_____
snkpkt = make_pkt(data, checksum)
udt_send(sndpkt)

Wait for call from above

Wait for ACK or NAK

rdt_rcv(rcvpkt) && isNAK(rcvpkt)
_____
udt_send(sndpkt)

rdt_rcv(rcvpkt) && corrupt(rcvpkt)
_____
udt_send(NAK)

rdt_rcv(rcvpkt) && isACK(rcvpkt)
_____
Λ

Wait for call from below

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

# rdt2.0 has a fatal flaw!

## what happens if ACK/NAK corrupted?

- ❖ sender doesn't know what happened at receiver!
- ❖ can't just retransmit: possible duplicate

## handling duplicates:

- ❖ sender retransmits current pkt if ACK/NAK corrupted
- ❖ sender adds *sequence number* to each pkt
- ❖ receiver discards (doesn't deliver up) duplicate pkt

### stop and wait
sender sends one packet, then waits for receiver response

# rdt2.1: sender, handles garbled ACK/NAKs

rdt_send(data)
_____

sndpkt = make_pkt(0, data, checksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isNAK(rcvpkt) )
_____
udt_send(sndpkt)

Wait for
call 0 from
above

Wait for
ACK or
NAK 0

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)
_____
$\Lambda$

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)
_____
$\Lambda$

Wait for
ACK or
NAK 1

Wait for
call 1 from
above

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isNAK(rcvpkt) )
_____
udt_send(sndpkt)

rdt_send(data)
_____

sndpkt = make_pkt(1, data, checksum)
udt_send(sndpkt)

# rdt2.1: receiver, handles garbled ACK/NAKs

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
 && has_seq0(rcvpkt)
---
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) && (corrupt(rcvpkt)
---
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
  not corrupt(rcvpkt) &&
  has_seq1(rcvpkt)
---
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) && (corrupt(rcvpkt)
---
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
  not corrupt(rcvpkt) &&
  has_seq0(rcvpkt)
---
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

**Wait for 0 from below**

**Wait for 1 from below**

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
 && has_seq1(rcvpkt)
---
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

# rdt2.1: discussion

**sender:**

- ❖ seq # added to pkt
- ❖ two seq. #'s (0,1) will suffice. Why?
- ❖ must check if received ACK/NAK corrupted
- ❖ twice as many states
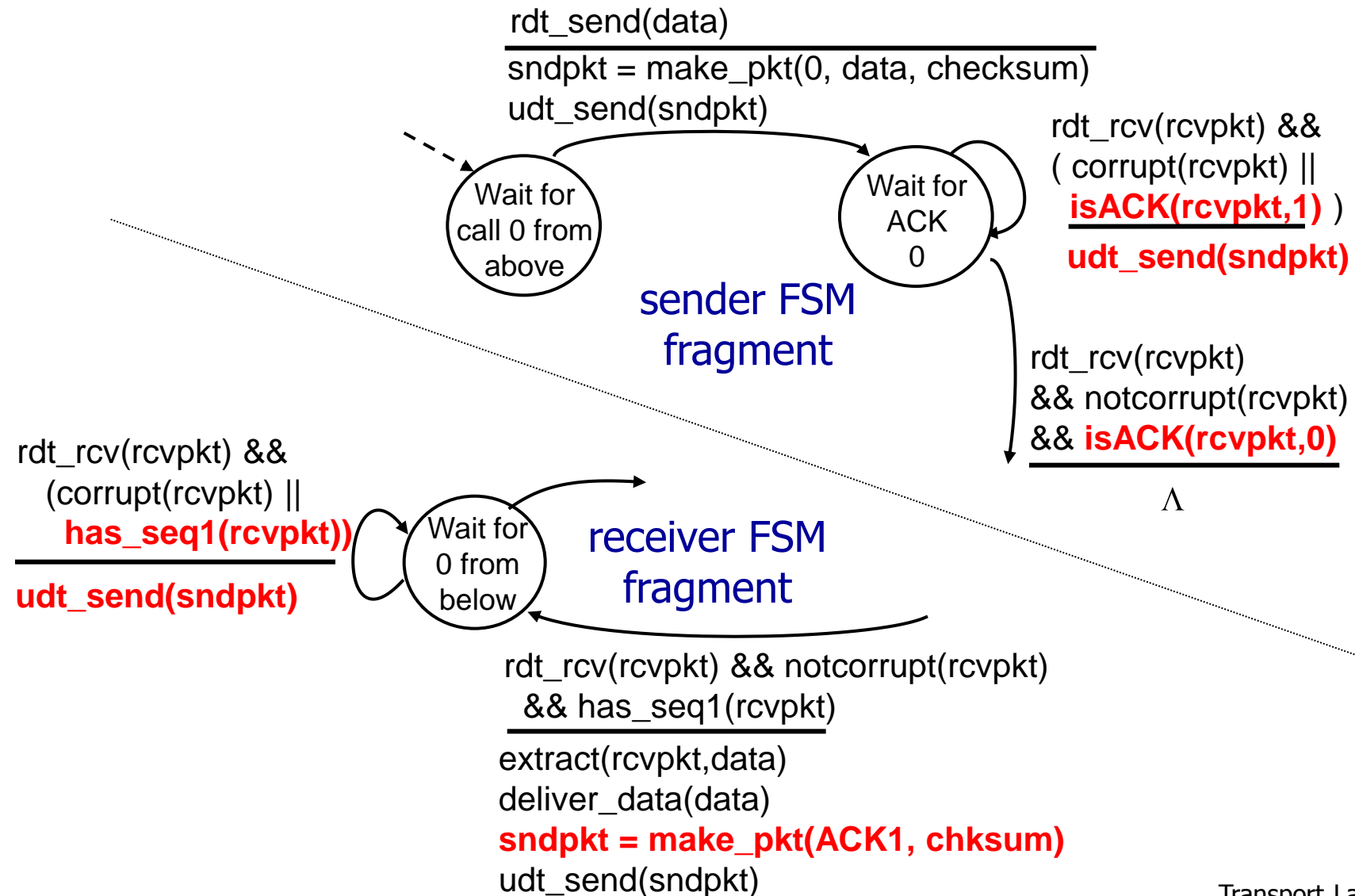  - state must "remember" whether "expected" pkt should have seq # of 0 or 1

**receiver:**

- ❖ must check if received packet is duplicate
  - state indicates whether 0 or 1 is expected pkt seq #
- ❖ note: receiver can *not* know if its last ACK/NAK received OK at sender

# rdt2.2: a NAK-free protocol

❖ same functionality as rdt2.1, using ACKs only

❖ instead of NAK, receiver sends ACK for last pkt received OK
  ▪ receiver must *explicitly* include seq # of pkt being ACKed

❖ duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

# rdt2.2: sender, receiver fragments

rdt_send(data)
_____
sndpkt = make_pkt(0, data, checksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
**isACK(rcvpkt,1)** )
**udt_send(sndpkt)**

Wait for call 0 from above

Wait for ACK 0

sender FSM fragment

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& **isACK(rcvpkt,0)**
_____
$\Lambda$

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
**has_seq1(rcvpkt))**
_____
**udt_send(sndpkt)**

Wait for 0 from below

receiver FSM fragment

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
**sndpkt = make_pkt(ACK1, chksum)**
udt_send(sndpkt)

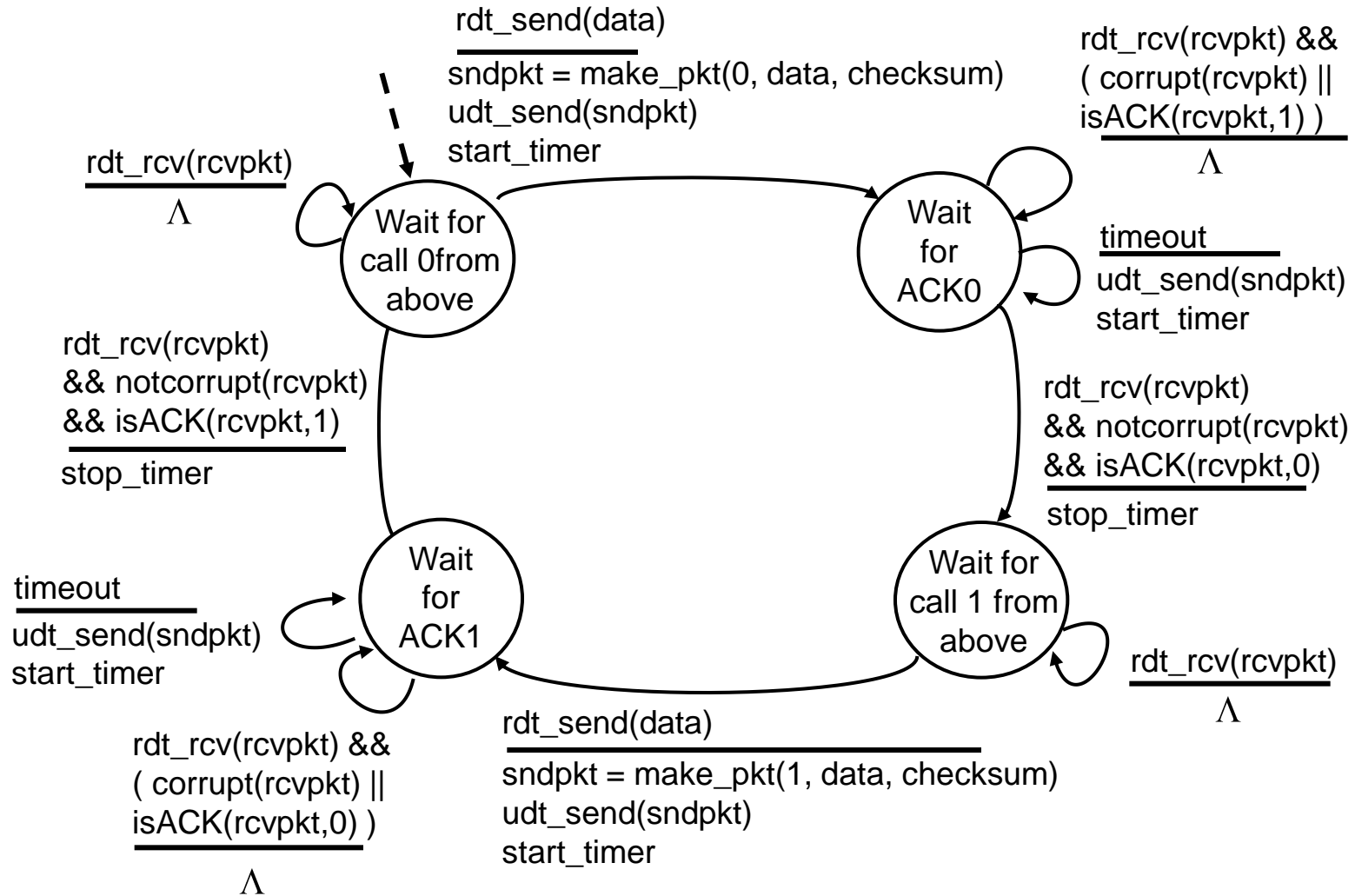# rdt3.0: channels with errors *and* loss

**new assumption:**
underlying channel can also lose packets (data, ACKs)

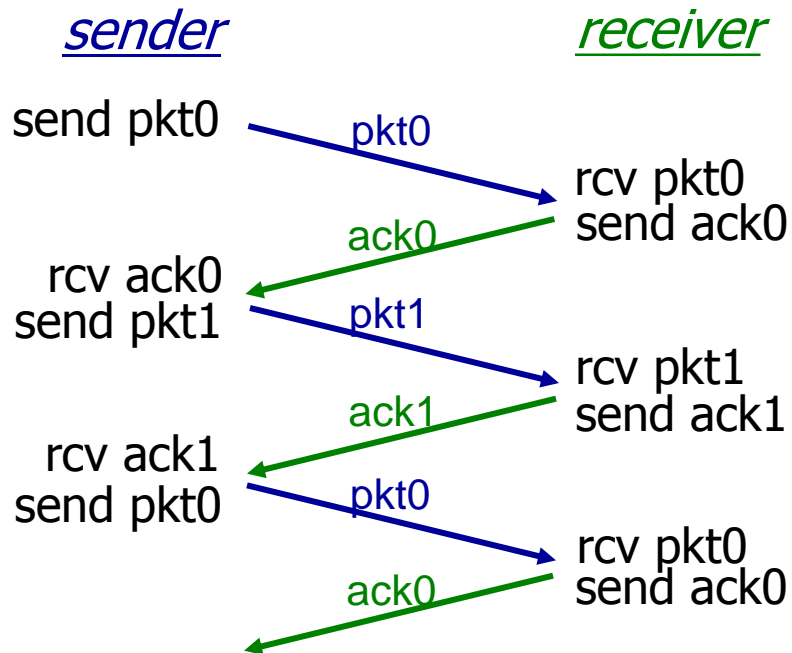- checksum, seq. #, ACKs, retransmissions will be of help … but not enough

**approach:** sender waits "reasonable" amount of time for ACK

- ❖ retransmits if no ACK received in this time
- ❖ if pkt (or ACK) just delayed (not lost):
  - retransmission will be duplicate, but seq. #'s already handles this
  - receiver must specify seq # of pkt being ACKed
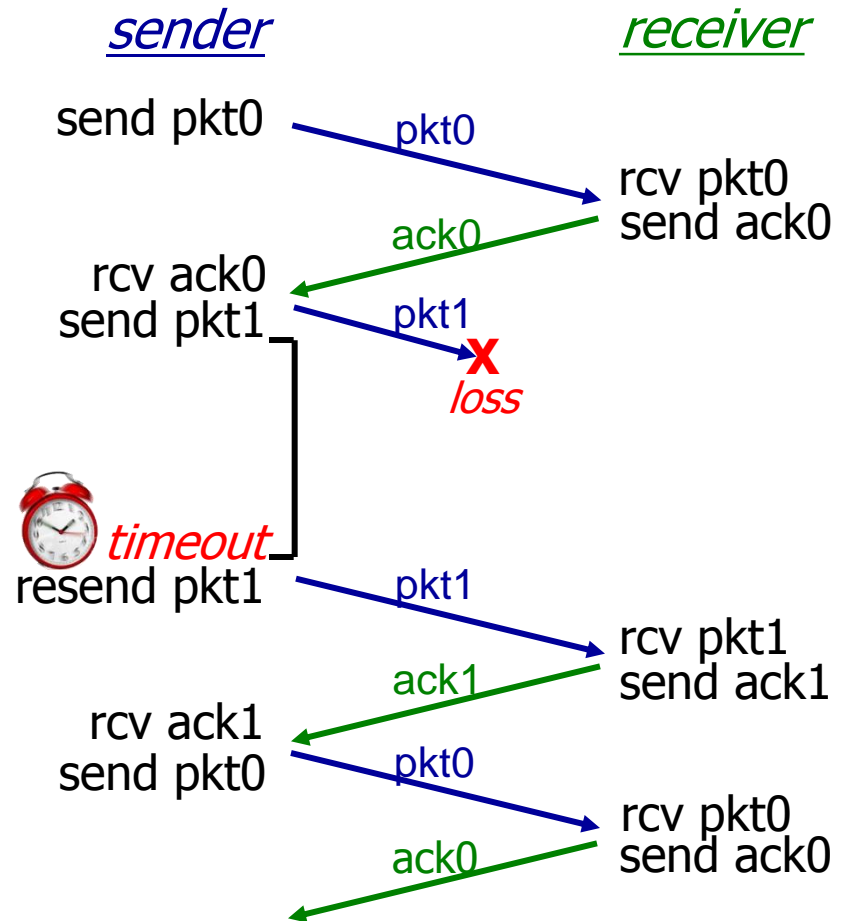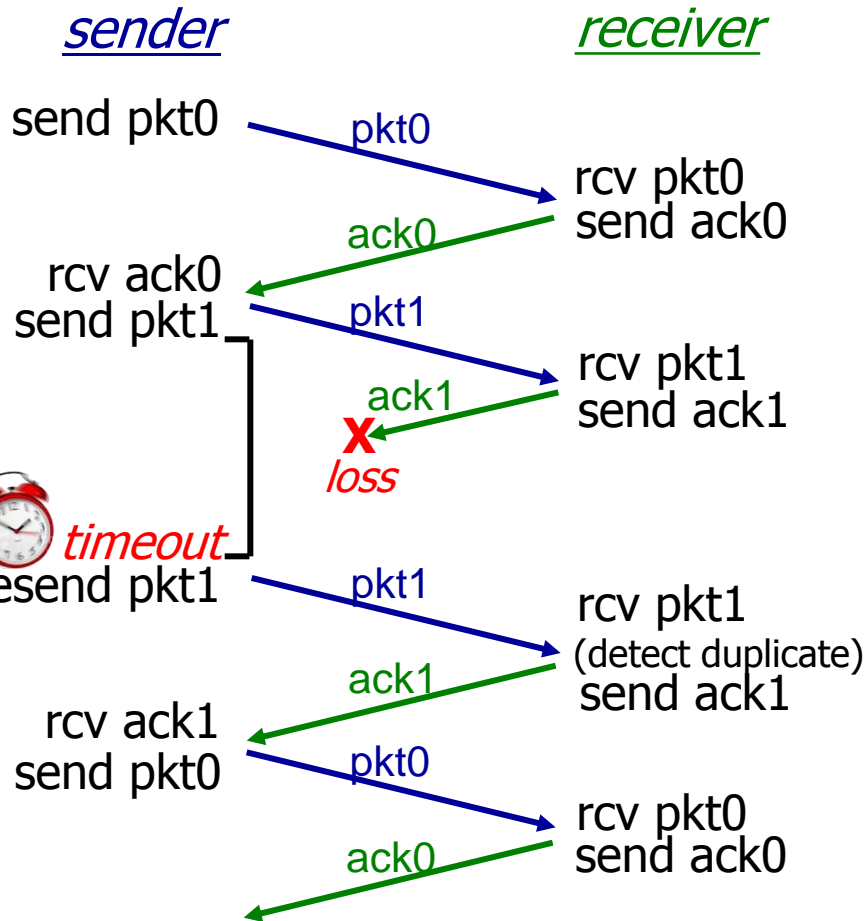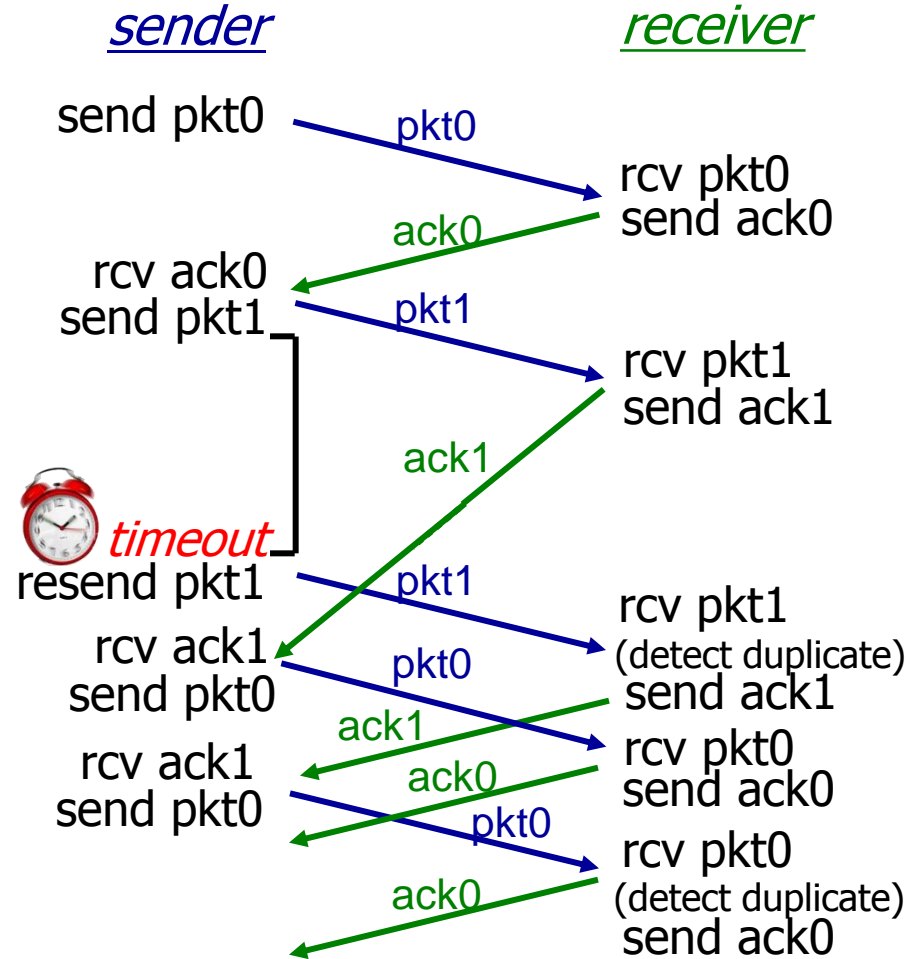- ❖ requires countdown timer

# rdt3.0 sender

rdt_send(data)
$\overline{\text{sndpkt = make\_pkt(0, data, checksum)}}$
udt_send(sndpkt)
start_timer

rdt_rcv(rcvpkt)
$\overline{\Lambda}$

**Wait for call 0 from above**

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,1) )
$\overline{\Lambda}$

**Wait for ACK0**

timeout
$\overline{\text{udt\_send(sndpkt)}}$
start_timer

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,1)
$\overline{\text{stop\_timer}}$

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,0)
$\overline{\text{stop\_timer}}$

**Wait for ACK1**

timeout
$\overline{\text{udt\_send(sndpkt)}}$
start_timer

**Wait for call 1 from above**

rdt_rcv(rcvpkt)
$\overline{\Lambda}$

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,0) )
$\overline{\Lambda}$

rdt_send(data)
$\overline{\text{sndpkt = make\_pkt(1, data, checksum)}}$
udt_send(sndpkt)
start_timer

# rdt3.0 in action

**sender**                    **receiver**

send pkt0 → pkt0
                              rcv pkt0
                              send ack0
rcv ack0 ← ack0
send pkt1 → pkt1
                              rcv pkt1
                              send ack1
rcv ack1 ← ack1
send pkt0 → pkt0
                              rcv pkt0
                              send ack0
          ← ack0

(a) no loss

**sender**                    **receiver**

send pkt0 → pkt0
                              rcv pkt0
                              send ack0
rcv ack0 ← ack0
send pkt1 → pkt1
                              ✗
                              loss

*timeout*
resend pkt1 → pkt1
                              rcv pkt1
                              send ack1
rcv ack1 ← ack1
send pkt0 → pkt0
                              rcv pkt0
                              send ack0
          ← ack0

(b) packet loss

# rdt3.0 in action

**sender**                      **receiver**

send pkt0 → pkt0 → rcv pkt0
send ack0
rcv ack0 ← ack0
send pkt1 → pkt1 → rcv pkt1
send ack1
**X** ack1
loss
**timeout**
resend pkt1 → pkt1 → rcv pkt1
(detect duplicate)
send ack1
rcv ack1 ← ack1
send pkt0 → pkt0 → rcv pkt0
send ack0
← ack0

(c) ACK loss

**sender**                      **receiver**

send pkt0 → pkt0 → rcv pkt0
send ack0
rcv ack0 ← ack0
send pkt1 → pkt1 → rcv pkt1
send ack1
**timeout**
resend pkt1 → pkt1 → rcv pkt1
ack1 (detect duplicate)
rcv ack1 send ack1
send pkt0 → pkt0 → rcv pkt0
ack1 send ack0
rcv ack1 ← ack0
send pkt0 → pkt0 → rcv pkt0
(detect duplicate)
ack0 send ack0

(d) premature timeout/ delayed ACK

# Performance of rdt3.0

❖ rdt3.0 is correct, but performance stinks
❖ e.g.: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$D_{trans} = \frac{L}{R} = \frac{8000 \; bits}{10^9 \; bits/sec} = 8 \; microsecs$$

- U $_{sender}$: *utilization* – fraction of time sender busy sending

$$U_{sender} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

- if RTT=30 msec, 1KB pkt every 30 msec: 33kB/sec throughput over 1 Gbps link

❖ network protocol limits use of physical resources!

# rdt3.0: stop-and-wait operation



sender                                    receiver

first packet bit transmitted, t = 0

last packet bit transmitted, t = L / R

first packet bit arrives

last packet bit arrives, send ACK

RTT

ACK arrives, send next packet, t = RTT + L / R

$$U_{sender} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

# Pipelined protocols

pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged pkts

- range of sequence numbers must be increased
- buffering at sender and/or receiver

data packet

data packets

ACK packets

(a) a stop-and-wait protocol in operation

(b) a pipelined protocol in operation

❖ two generic forms of pipelined protocols: *go-Back-N, selective repeat*

# Pipelining: increased utilization

sender                                    receiver

first packet bit transmitted, t = 0

last bit transmitted, t = L / R

RTT

first packet bit arrives

last packet bit arrives, send ACK

last bit of 2nd packet arrives, send ACK

last bit of 3rd packet arrives, send ACK

ACK arrives, send next
packet, t = RTT + L / R

3-packet pipelining increases
utilization by a factor of 3!

$$U_{sender} = \frac{3L / R}{RTT + L / R} = \frac{.0024}{30.008} = 0.00081$$

# Pipelined protocols: overview

## Go-back-N:

❖ sender can have up to N unacked packets in pipeline

❖ receiver only sends *cumulative ack*

  ▪ doesn't ack packet if there's a gap

❖ sender has timer for oldest unacked packet

  ▪ when timer expires, retransmit *all* unacked packets

## Selective Repeat:

❖ sender can have up to N unack'ed packets in pipeline

❖ rcvr sends *individual ack* for each packet

❖ sender maintains timer for each unacked packet

  ▪ when timer expires, retransmit only that unacked packet

# Go-Back-N: sender

❖ k-bit seq # in pkt header
❖ "window" of up to N, consecutive unack'ed pkts allowed



❖ ACK(n): ACKs all pkts up to, including seq # n - *"cumulative ACK"*
  ▪ may receive duplicate ACKs (see receiver)
❖ timer for oldest in-flight pkt
❖ *timeout(n):* retransmit packet n and all higher seq # pkts in window

# GBN: sender extended FSM

rdt_send(data)
_____

if (nextseqnum < base+N) {
   sndpkt[nextseqnum] = make_pkt(nextseqnum,data,chksum)
   udt_send(sndpkt[nextseqnum])
   if (base == nextseqnum)
     start_timer
   nextseqnum++
   }
else
 refuse_data(data)

$\Lambda$
_____
base=1
nextseqnum=1

( Wait )

rdt_rcv(rcvpkt)
  && corrupt(rcvpkt)
_____

timeout
_____
start_timer
udt_send(sndpkt[base])
udt_send(sndpkt[base+1])
…
udt_send(sndpkt[nextseqnum-1])

rdt_rcv(rcvpkt) &&
  notcorrupt(rcvpkt)
_____
base = getacknum(rcvpkt)+1
If (base == nextseqnum)
  stop_timer
 else
  start_timer

# GBN: receiver extended FSM

default
_____
udt_send(sndpkt)

rdt_rcv(rcvpkt)
  && notcurrupt(rcvpkt)
  && hasseqnum(rcvpkt,expectedseqnum)
_____

Λ
_____
expectedseqnum=1
sndpkt =
  make_pkt(expectedseqnum,ACK,chksum)

Wait

extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(expectedseqnum,ACK,chksum)
udt_send(sndpkt)
expectedseqnum++

ACK-only: always send ACK for correctly-received pkt with highest *in-order* seq #
- may generate duplicate ACKs
- need only remember `expectedseqnum`

❖ out-of-order pkt:
- discard (don't buffer): *no receiver buffering!*
- re-ACK pkt with highest in-order seq #

# GBN in action

_sender window (N=4)_     _sender_     _receiver_

| | | |
|---|---|---|
| 0 1 2 3 4 5 6 7 8 | send pkt0 | |
| 0 1 2 3 4 5 6 7 8 | send pkt1 | |
| 0 1 2 3 4 5 6 7 8 | send pkt2 | **X** _loss_ |
| 0 1 2 3 4 5 6 7 8 | send pkt3 | |
| | (wait) | |

receive pkt0, send ack0
receive pkt1, send ack1

receive pkt3, discard,
        (re)send ack1

| | |
|---|---|
| 0 1 2 3 4 5 6 7 8 | rcv ack0, send pkt4 |
| 0 1 2 3 4 5 6 7 8 | rcv ack1, send pkt5 |

receive pkt4, discard,
        (re)send ack1
receive pkt5, discard,
        (re)send ack1

ignore duplicate ACK

_pkt 2 timeout_

| | |
|---|---|
| 0 1 2 3 4 5 6 7 8 | send pkt2 |
| 0 1 2 3 4 5 6 7 8 | send pkt3 |
| 0 1 2 3 4 5 6 7 8 | send pkt4 |
| 0 1 2 3 4 5 6 7 8 | send pkt5 |

rcv pkt2, deliver, send ack2
rcv pkt3, deliver, send ack3
rcv pkt4, deliver, send ack4
rcv pkt5, deliver, send ack5

# Selective repeat

❖ receiver *individually* acknowledges all correctly received pkts
  ▪ buffers pkts, as needed, for eventual in-order delivery to upper layer

❖ sender only resends pkts for which ACK not received
  ▪ sender timer for each unACKed pkt

❖ sender window
  ▪ *N* consecutive seq #'s
  ▪ limits seq #s of sent, unACKed pkts

# Selective repeat: sender, receiver windows

send_base   nextseqnum

already ack'ed

usable, not yet sent

sent, not yet ack'ed

not usable

window size
N

(a) sender view of sequence numbers

out of order (buffered) but already ack'ed

acceptable (within window)

Expected, not yet received

not usable

window size
N

rcv_base

(b) receiver view of sequence numbers

# Selective repeat

## sender

### data from above:

- ❖ if next available seq # in window, send pkt

### timeout(n):

- ❖ resend pkt n, restart timer

### ACK(n) in [sendbase,sendbase+N]:

- ❖ mark pkt n as received
- ❖ if n smallest unACKed pkt, advance window base to next unACKed seq #

## receiver

### pkt n in [rcvbase, rcvbase+N-1]

- ❖ send ACK(n)
- ❖ out-of-order: buffer
- ❖ in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

### pkt n in [rcvbase-N,rcvbase-1]

- ❖ ACK(n)

### otherwise:

- ❖ ignore

# Selective repeat in action

sender window (N=4)      sender              receiver

`0 1 2 3`4 5 6 7 8     send  pkt0

`0 1 2 3`4 5 6 7 8     send  pkt1

`0 1 2 3`4 5 6 7 8     send  pkt2            receive pkt0, send ack0

`0 1 2 3`4 5 6 7 8     send  pkt3    **X** *loss*    receive pkt1, send ack1

                        (wait)

                                 receive pkt3, buffer,
                                           send ack3

0 `1 2 3 4`5 6 7 8    rcv ack0, send pkt4

0 1 `2 3 4 5`6 7 8    rcv ack1, send pkt5      receive pkt4, buffer,
                                           send ack4

       record ack3 arrived         receive pkt5, buffer,
                                           send ack5

          *pkt 2 timeout*

0 1 `2 3 4 5`6 7 8      send  pkt2

0 1 `2 3 4 5`6 7 8   record ack4 arrived

0 1 `2 3 4 5`6 7 8   record ack4 arrived       rcv pkt2; deliver pkt2,

0 1 `2 3 4 5`6 7 8                         pkt3, pkt4, pkt5; send ack2

       *Q: what happens when ack2 arrives?*

# Selective repeat: dilemma

example:
- ❖ seq #'s: 0, 1, 2, 3
- ❖ window size=3

- ❖ receiver sees no difference in two scenarios!
- ❖ duplicate data accepted as new in (b)

Q: what relationship between seq # size and window size to avoid problem in (b)?



sender window (after receipt)

receiver window (after receipt)

pkt0
pkt1
pkt2
pkt3
pkt0

(a) no problem

*receiver can't see sender side.*
*receiver behavior identical in both cases!*
*something's (very) wrong!*

will accept packet with seq number 0

pkt0
pkt1
pkt2

timeout retransmit pkt0
pkt0

(b) oops!

will accept packet with seq number 0

# Link layer, LANs: outline

# 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)

# Multiple access protocols

❖ single shared broadcast channel

❖ two or more simultaneous transmissions by nodes: interference
  ▪ *collision* if node receives two or more signals at the same time

*multiple access protocol*

❖ distributed algorithm that determines how nodes share channel, i.e., determine when node can transmit

❖ communication about channel sharing must use channel itself!
  ▪ no out-of-band channel for coordination

# An ideal multiple access protocol

*given:* broadcast channel of rate R bps

*desiderata:*

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

# MAC protocols: taxonomy

three broad classes:

❖ *channel partitioning*
 ▪ divide channel into smaller "pieces" (time slots, frequency, code)
 ▪ allocate piece to node for exclusive use

❖ *random access*
 ▪ channel not divided, allow collisions
 ▪ "recover" from collisions

❖ *"taking turns"*
 ▪ nodes take turns, but nodes with more to send can take longer turns

# Channel partitioning MAC protocols: TDMA

## TDMA: time division multiple access

❖ access to channel in "rounds"

❖ each station gets fixed length slot (length = pkt trans time) in each round

❖ unused slots go idle

❖ example: 6-station LAN, 1,3,4 have pkt, slots 2,5,6 idle

# Channel partitioning MAC protocols: FDMA

## FDMA: frequency division multiple access

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

FDM cable

frequency bands

time

# Code Division Multiple Access (CDMA)

❖ unique "code" assigned to each user; i.e., code set partitioning
- all users share same frequency, but each user has own "chipping" sequence (i.e., code) to encode data
- allows multiple users to "coexist" and transmit simultaneously with minimal interference (if codes are "orthogonal")

❖ *encoded signal* = (original data) X (chipping sequence)

❖ *decoding:* inner-product of encoded signal and chipping sequence

# CDMA encode/decode



$Z_{i,m} = d_i \cdot c_m$

channel output $Z_{i,m}$

sender

data bits: $d_1 = -1$, $d_0 = 1$

code

slot 1   slot 0

slot 1 channel output   slot 0 channel output

$$D_i = \frac{\sum_{m=1}^{M} Z_{i,m} \cdot c_m}{M}$$

received input

code

receiver

slot 1   slot 0

$d_1 = -1$, $d_0 = 1$

slot 1 channel output   slot 0 channel output

# CDMA: two-sender interference

senders



*channel sums together transmissions by sender 1 and 2*

*using same code as sender 1, receiver recovers sender 1's original data from summed channel data!*

# 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",
❖ <span style="color:red">random access MAC protocol</span> specifies:
  ▪ how to detect collisions
  ▪ how to recover from collisions (e.g., via delayed retransmissions)
❖ examples of random access MAC protocols:
  ▪ slotted ALOHA
  ▪ ALOHA
  ▪ CSMA, CSMA/CD, CSMA/CA

# Slotted ALOHA

*assumptions:*

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

*operation:*

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

# Slotted ALOHA

| 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

# 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!

!

# Pure (unslotted) ALOHA

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

will overlap
with start of
← i's frame →

will overlap
with end of
← i's frame →

node i frame

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

# Pure ALOHA efficiency

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

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

P(no other node transmits in $[t_0-1,t_0]$

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

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

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

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

even *worse* than slotted Aloha!

# CSMA (carrier sense multiple access)

*CSMA:* listen before transmit:

if channel sensed idle: transmit entire frame

❖ if channel sensed busy, defer transmission

❖ human analogy: don't interrupt others!

# CSMA collisions

spatial layout of nodes



❖ **collisions** *can* **still occur:** propagation delay means two nodes may not hear each other's transmission

❖ **collision:** entire packet transmission time wasted

- distance & propagation delay play role in in determining collision probability

$t_0$

*time*

$t_1$

# CSMA/CD (collision detection)

*CSMA/CD:* carrier sensing, deferral as in CSMA
- collisions *detected* within short time
- colliding transmissions aborted, reducing channel wastage

❖ collision detection:
- easy in wired LANs: measure signal strengths, compare transmitted, received signals
- difficult in wireless LANs: received signal strength overwhelmed by local transmission strength

❖ human analogy: the polite conversationalist

# CSMA/CD (collision detection)

spatial layout of nodes

time

$t_0$

$t_1$

collision
detect/abort
time

# 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 $\{0,1,2, \ldots, 2^m\text{-}1\}$. NIC waits $K \cdot 512$ bit times, returns to Step 2
   - longer backoff interval with more collisions

# 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
  - ▪ as $t_{prop}$ goes to 0
  - ▪ as $t_{trans}$ goes to infinity
- ❖ better performance than ALOHA: and simple, cheap, decentralized!

# Wireless Link Characteristics

*important* differences from wired link ….

- *decreased signal strength:* radio signal attenuates as it propagates through matter (path loss)
- *interference from other sources:* standardized wireless network frequencies (e.g., 2.4 GHz) shared by other devices (e.g., phone); devices (motors) interfere as well
- *multipath propagation:* radio signal reflects off objects ground, arriving ad destination at slightly different times

…. make communication across (even a point to point) wireless link much more "difficult"

# Wireless network characteristics

Multiple wireless senders and receivers create additional problems (beyond multiple access):



*Hidden terminal problem*

❖ B, A hear each other

❖ B, C hear each other

❖ A, C can not hear each other means A, C unaware of their interference at B

*Signal attenuation:*

❖ B, A hear each other

❖ B, C hear each other

❖ A, C can not hear each other interfering at B

# IEEE 802.11: multiple access

❖ avoid collisions: 2⁺ nodes transmitting at same time

❖ 802.11: CSMA - sense before transmitting
  ▪ don't collide with ongoing transmission by other node

❖ 802.11: *no* collision detection!
  ▪ difficult to receive (sense collisions) when transmitting due to weak received signals (fading)
  ▪ can't sense all collisions in any case: hidden terminal, fading
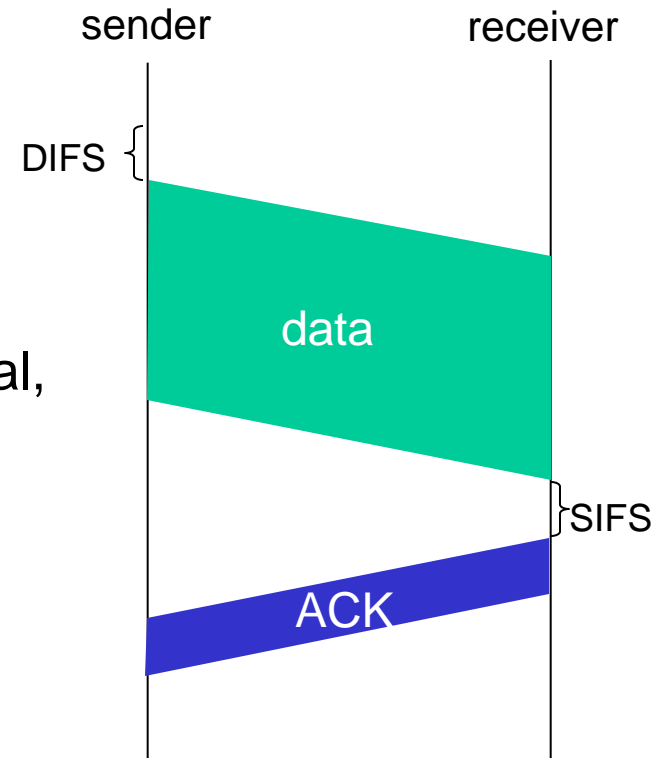  ▪ goal: *avoid collisions:* CSMA/C(ollision)A(voidance)



A's signal strength

C's signal strength

space

# IEEE 802.11 MAC Protocol: CSMA/CA

*802.11 sender*

**1** if sense channel idle for **DIFS**  then

   transmit entire frame (no CD)

2 if sense channel busy then

   start random backoff time

   timer counts down while channel idle

   transmit when timer expires

   if no ACK, increase random backoff interval, repeat 2

*802.11 receiver*

- if frame received OK

   return ACK after **SIFS** (ACK needed due to hidden terminal problem)

# Avoiding collisions (more)

*idea:* allow sender to "reserve" channel rather than random access of data frames: avoid collisions of long data frames

❖ sender first transmits *small* request-to-send (RTS) packets to BS using CSMA
  ▪ RTSs may still collide with each other (but they're short)
❖ BS broadcasts clear-to-send CTS in response to RTS
❖ CTS heard by all nodes
  ▪ sender transmits data frame
  ▪ other stations defer transmissions

*avoid data frame collisions completely using small reservation packets!*

# Collision Avoidance: RTS-CTS exchange

A          AP          B

RTS(A)          RTS(B)

reservation collision

RTS(A)

CTS(A)          CTS(A)

DATA (A)          defer

time

ACK(A)          ACK(A)

# "Taking turns" MAC protocols

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

random access MAC protocols
- efficient at low load: single node can fully utilize channel
- high load: collision overhead

"taking turns" protocols
look for best of both worlds!
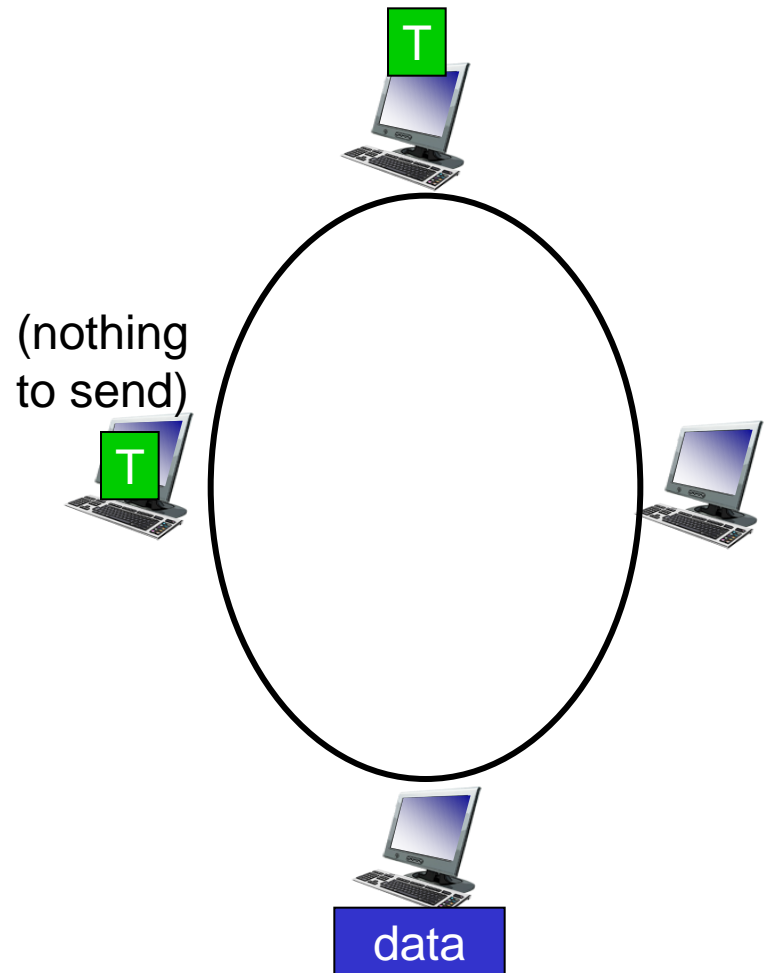
# "Taking turns" MAC protocols

## polling:

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

❖ typically used with "dumb" slave devices

❖ concerns:
  - polling overhead
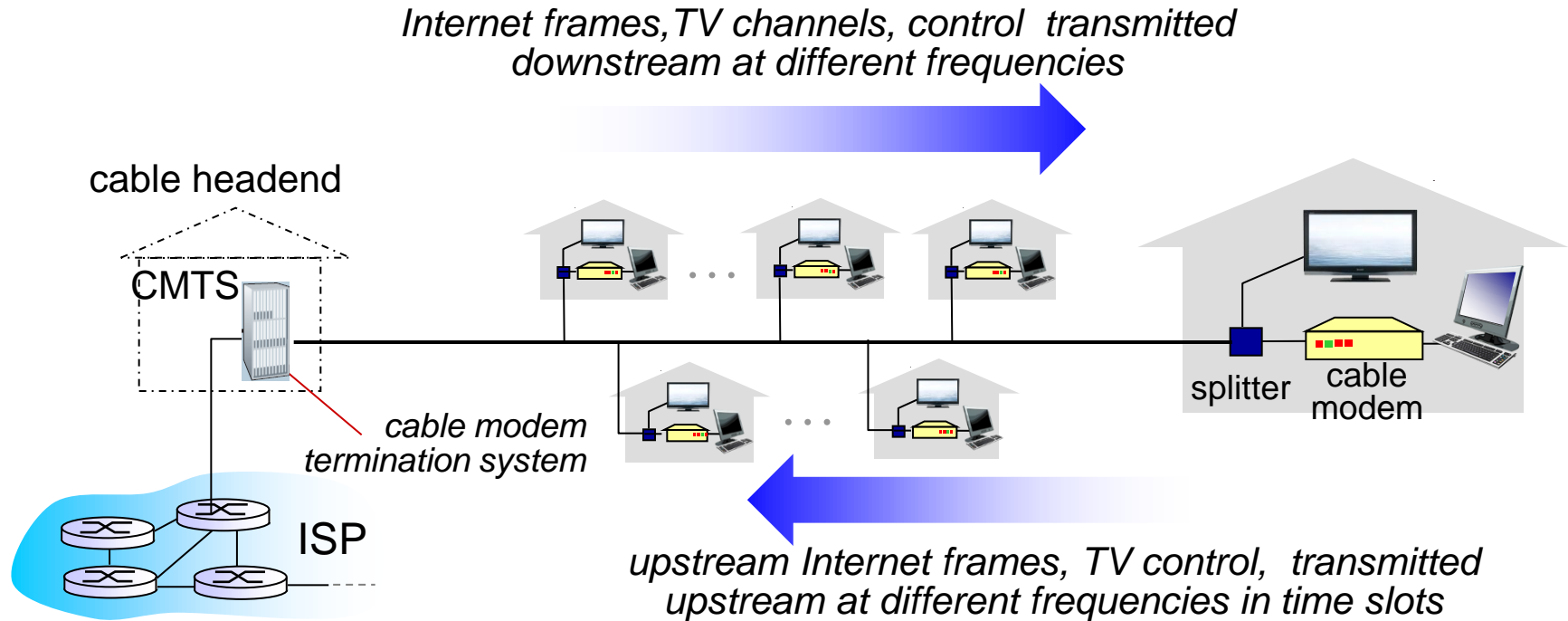  - latency
  - single point of failure (master)



slaves

master

# "Taking turns" MAC protocols

*token passing:*

❖ control *token* passed from one node to next sequentially.

❖ token message

❖ concerns:
  ▪ token overhead
  ▪ latency
  ▪ single point of failure (token)
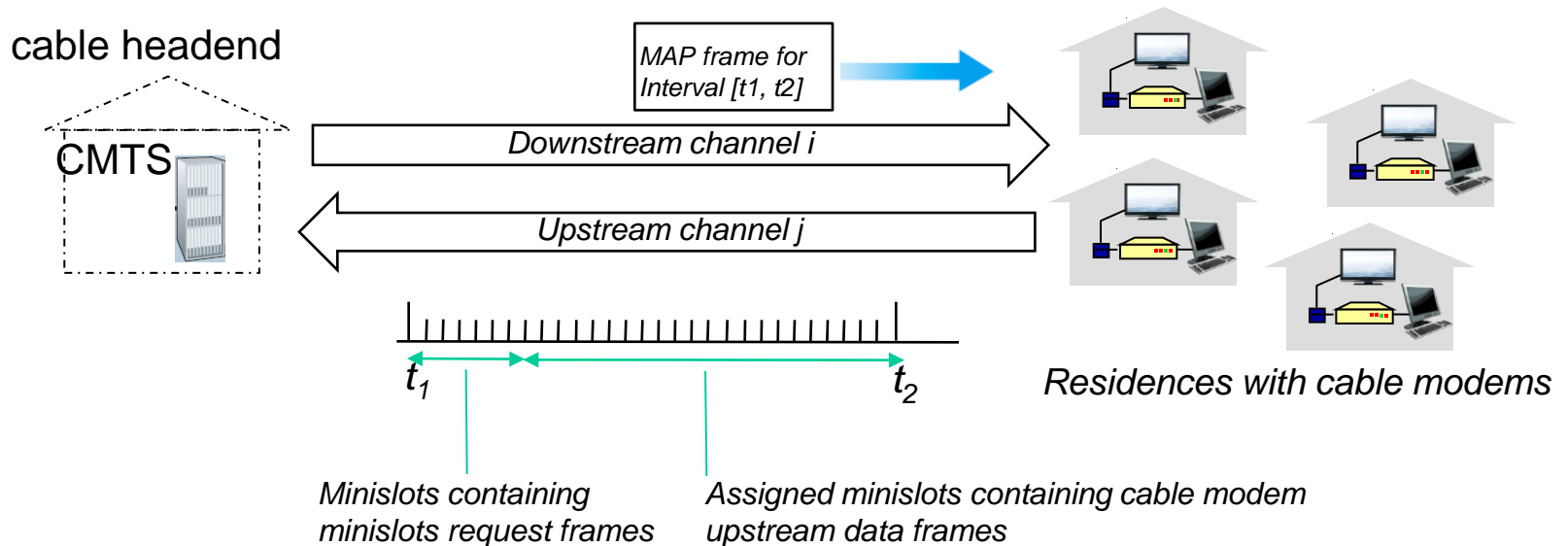


(nothing to send)

T

data

# Cable access network



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

cable headend

CMTS

cable modem termination system

ISP

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

splitter   cable modem

- ❖ *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)

# Cable access network



cable headend

MAP frame for Interval [t1, t2]

Downstream channel i

Upstream channel j

CMTS

$t_1$        $t_2$

Residences with cable modems

Minislots containing minislots request frames

Assigned minislots containing cable modem upstream data frames

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

# 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, 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 Subnetting

5.6 a day in the life of a web request

# MAC addresses and ARP

❖ 32-bit IP address:
  ▪ *network-layer* address for interface
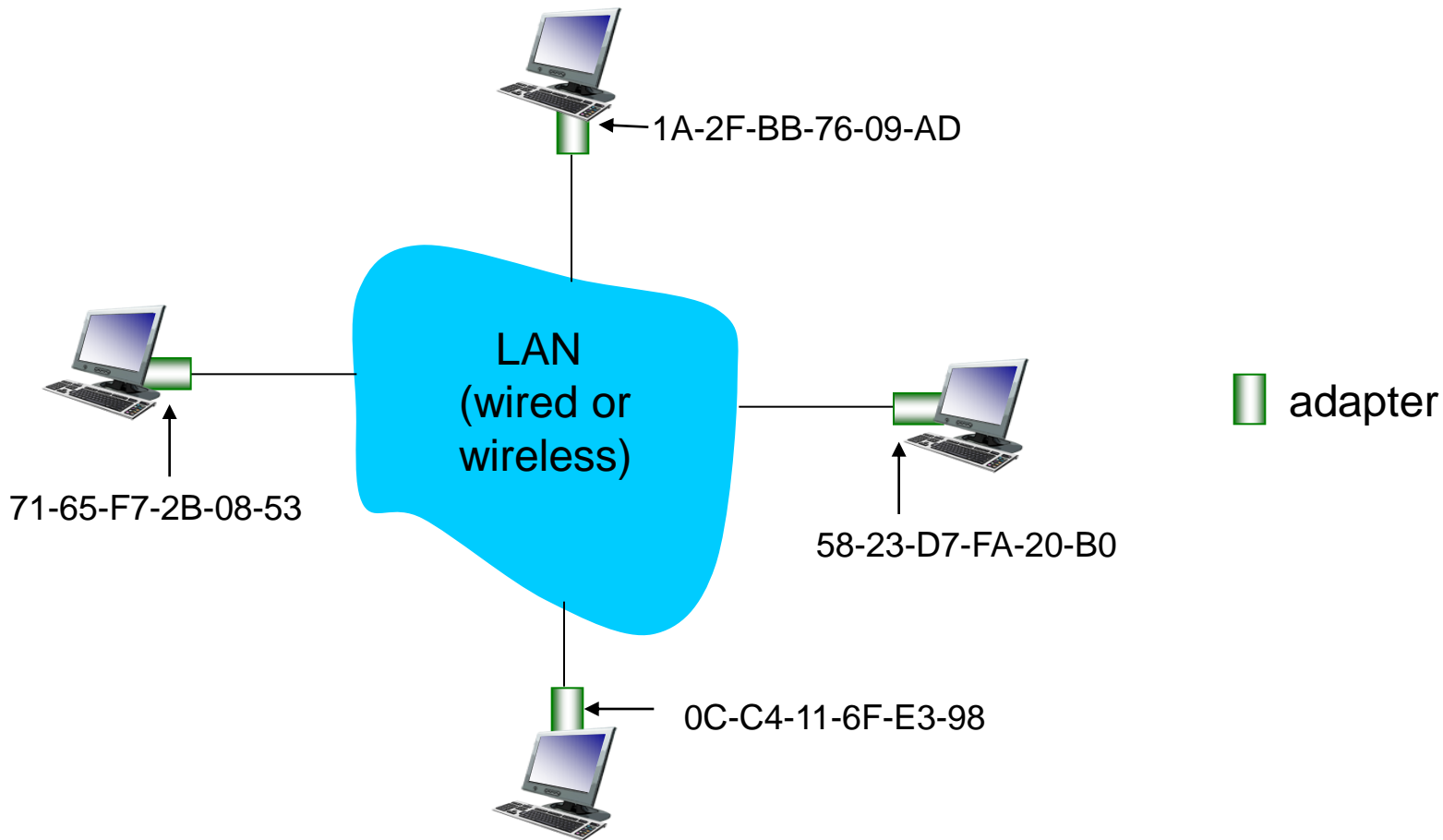  ▪ used for layer 3 (network layer) forwarding

❖ 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)

# LAN addresses and ARP

each adapter on LAN has unique *LAN* address



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

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

LAN
(wired or
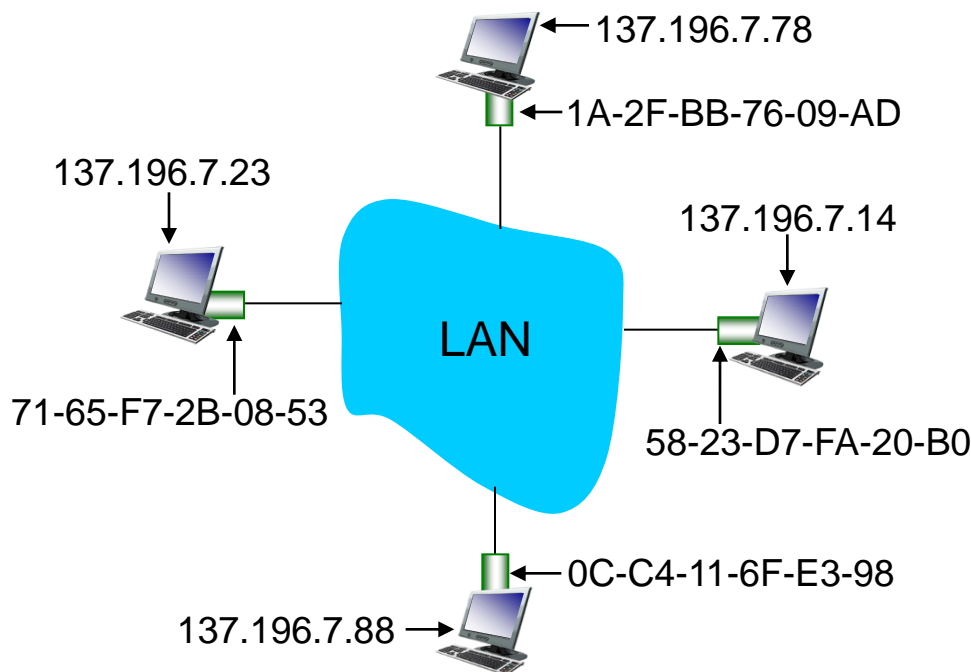wireless)

58-23-D7-FA-20-B0

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

adapter

# LAN addresses (more)

❖ MAC address allocation administered by IEEE

❖ manufacturer buys portion of MAC address space (to assure uniqueness)

❖ analogy:
  - MAC address: like Social Security Number
  - IP address: like postal address

❖ MAC flat address ➜ portability
  - can move LAN card from one LAN to another

❖ IP hierarchical address *not* portable
  - address depends on IP subnet to which node is attached

# ARP: address resolution protocol

*Question:* how to determine 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>

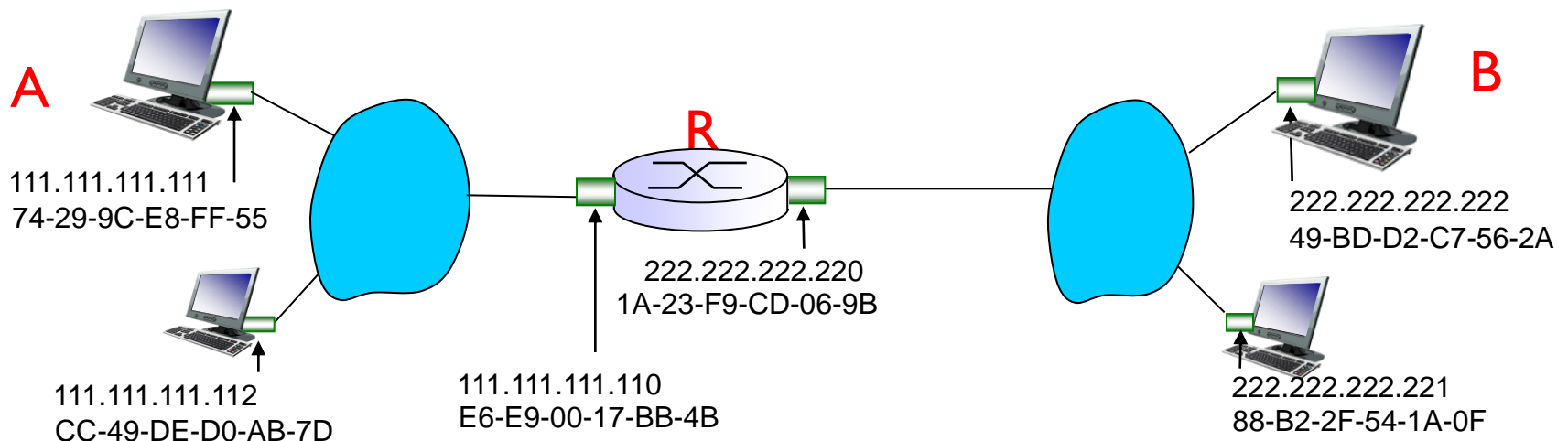- TTL (Time To Live): time after which address mapping will be forgotten (typically 20 min)

137.196.7.78

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

137.196.7.23

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

LAN

137.196.7.14

58-23-D7-FA-20-B0

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

137.196.7.88

# 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*
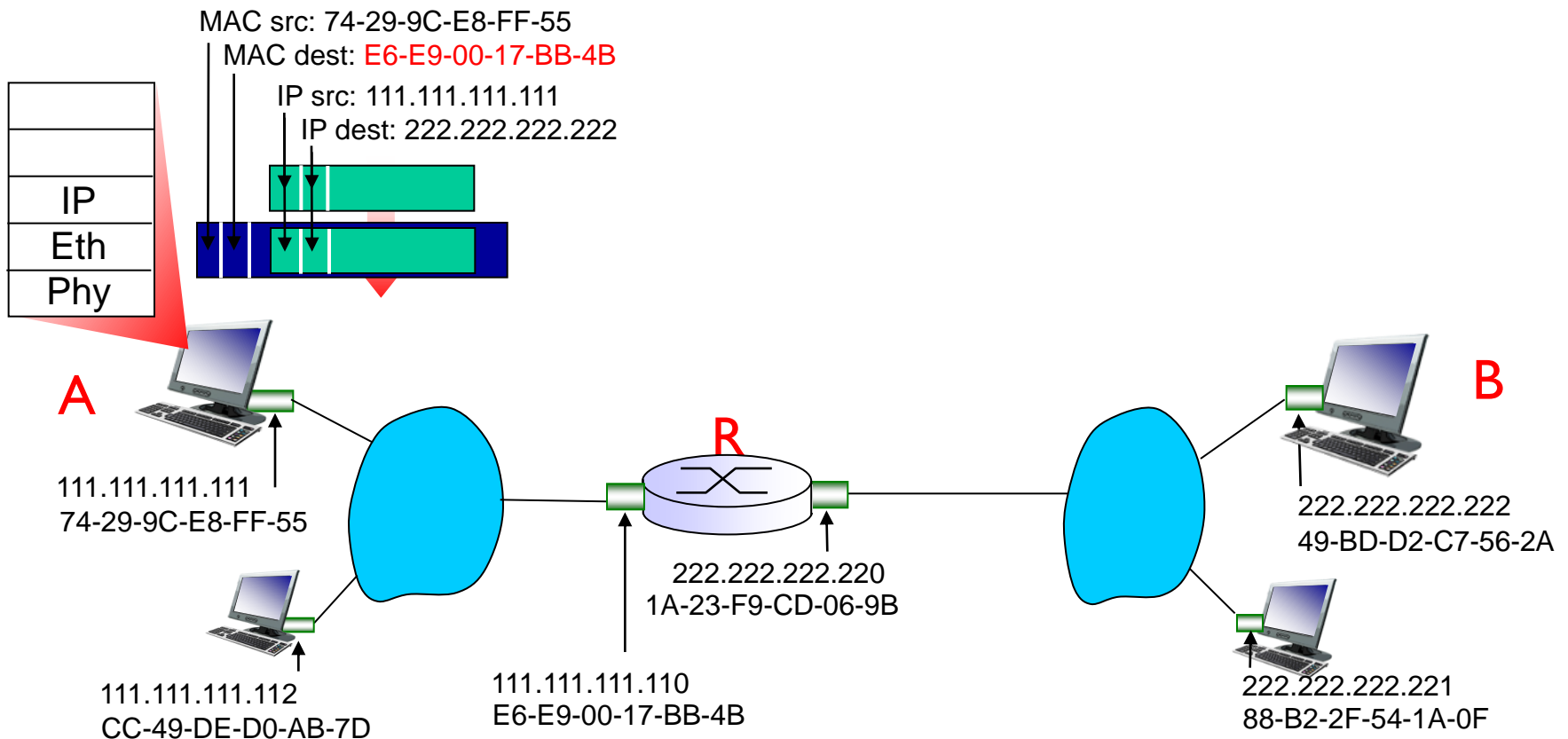
# Addressing: routing to another LAN

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?)
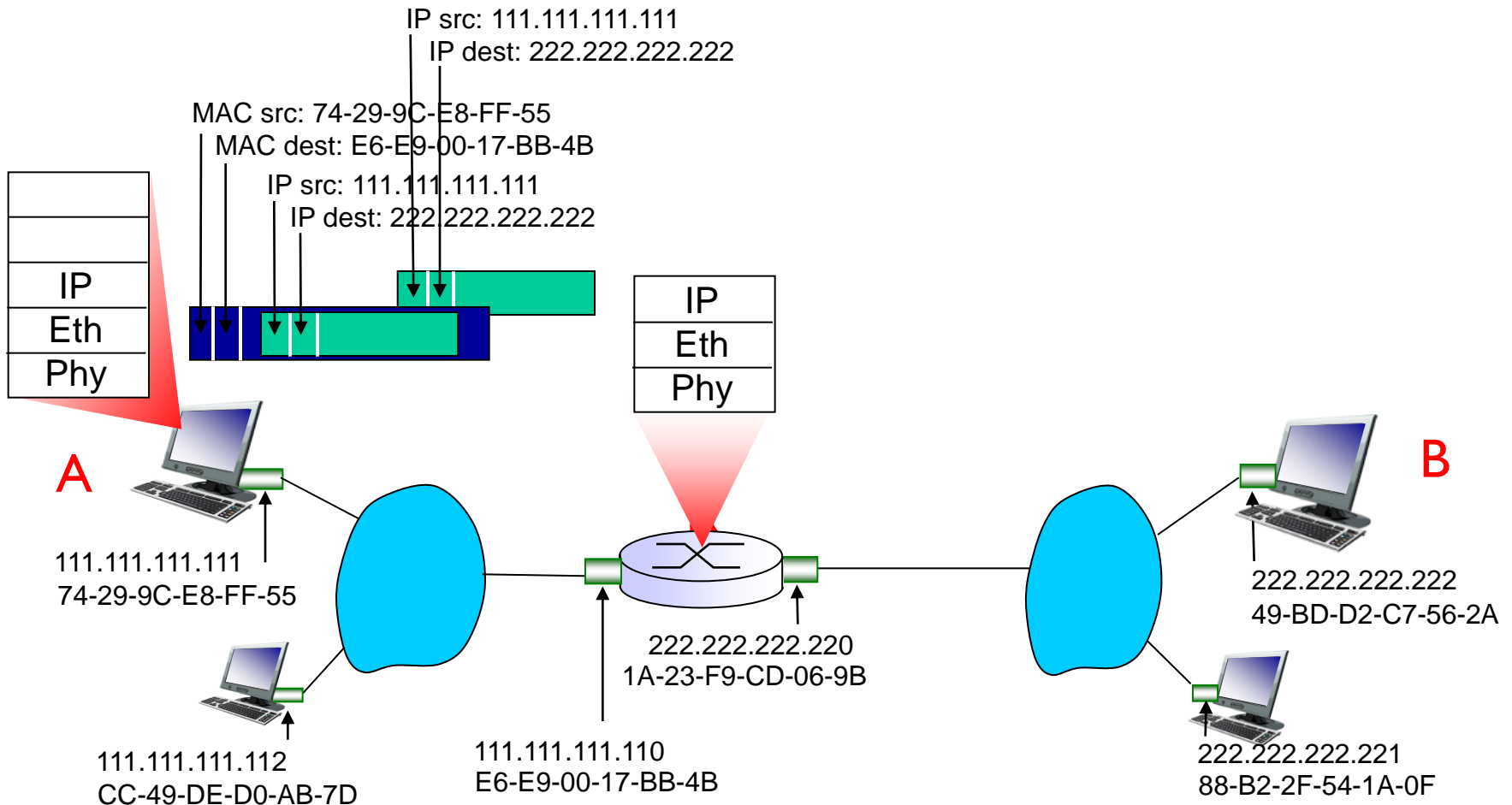- assume A knows R's MAC address (how?)



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

# Addressing: routing to another LAN

❖ A creates IP datagram with IP source A, destination B

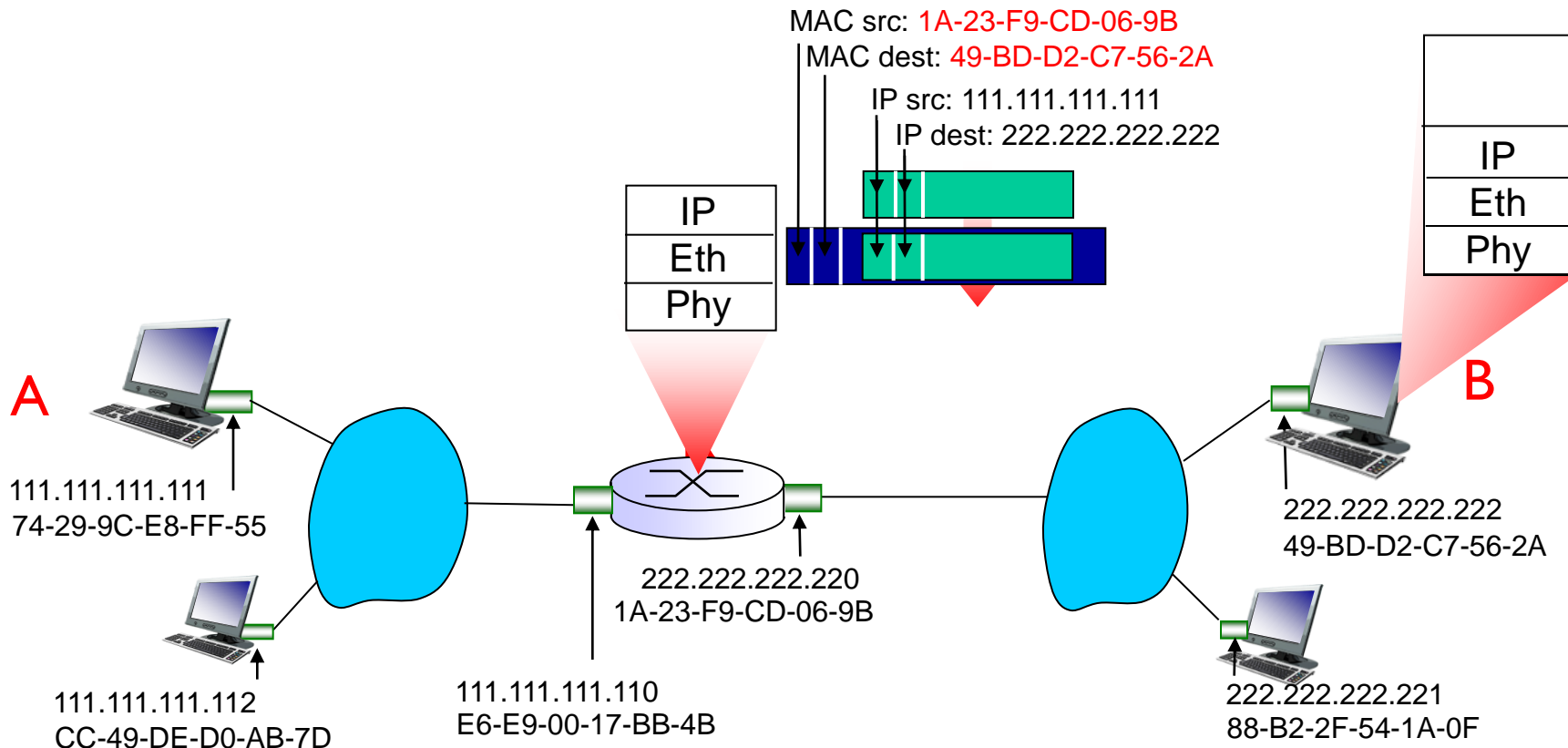❖ A creates link-layer frame with R's MAC address as dest, frame contains A-to-B IP datagram

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

B

R

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

# Addressing: routing to another LAN

❖ frame sent from A to R

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

IP src: 111.111.111.111
IP dest: 222.222.222.222

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

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

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

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

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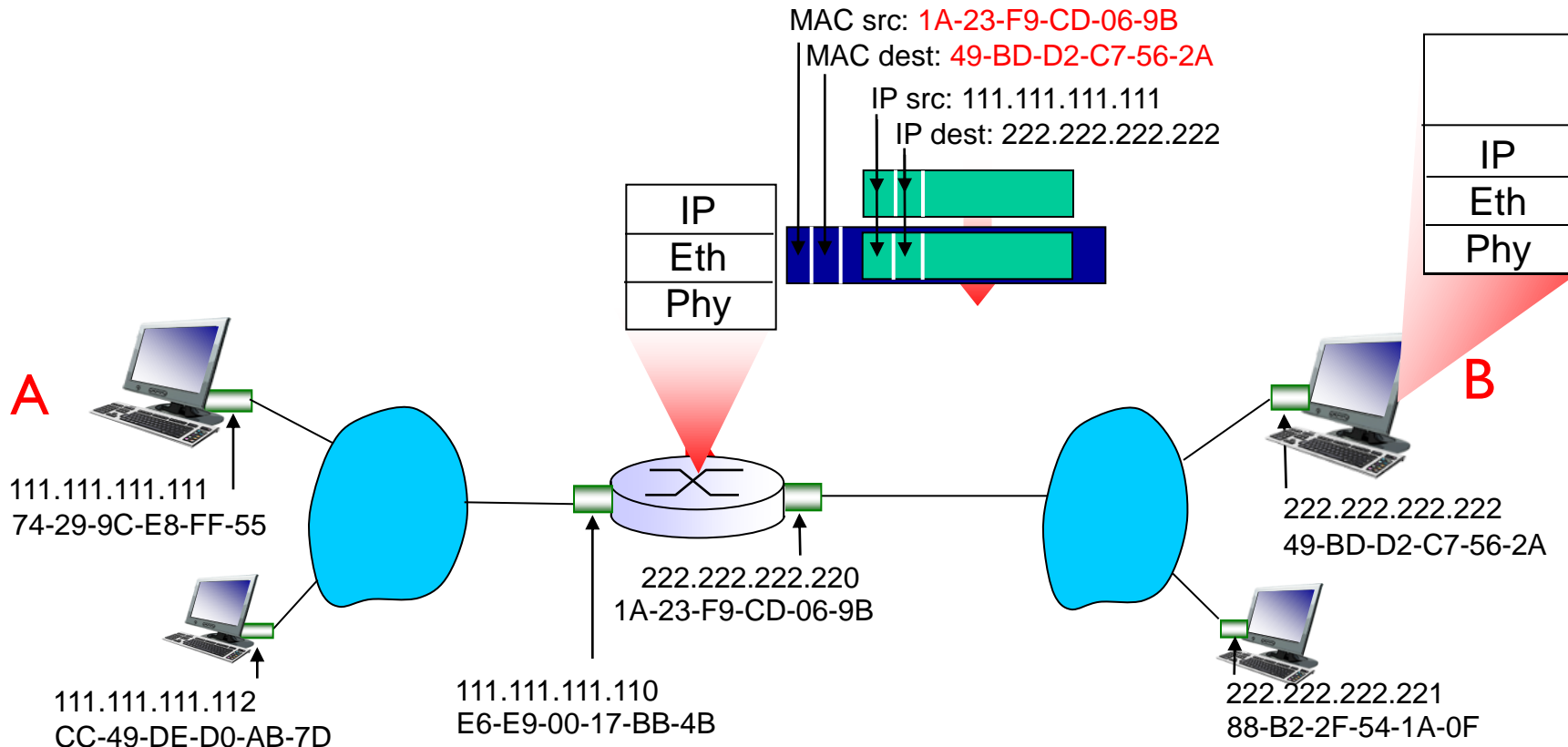E6-E9-00-17-BB-4B

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

# Addressing: routing to another LAN

❖ R forwards datagram with IP source A, destination B

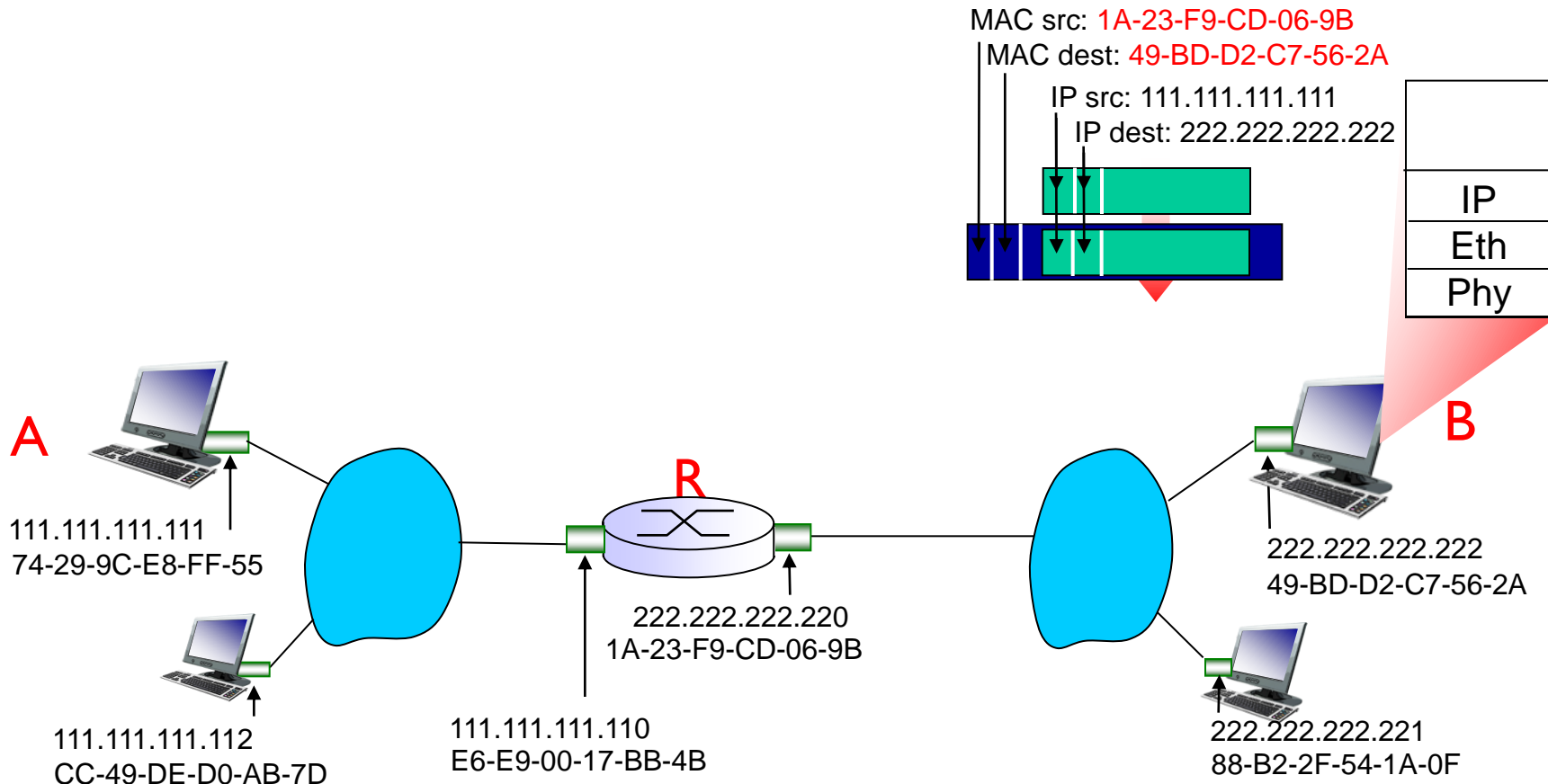❖ R creates link-layer frame with B's MAC address as dest, frame contains A-to-B IP datagram

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

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

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

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

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

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

# 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

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

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

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

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

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

# 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

A

R

B

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

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

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

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

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

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

# 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 Subnetting

5.6 a day in the life of a web request

# Ethernet

"dominant" wired LAN technology:

❖ cheap $20 for NIC

❖ first widely used LAN technology

❖ simpler, cheaper than token LANs and ATM
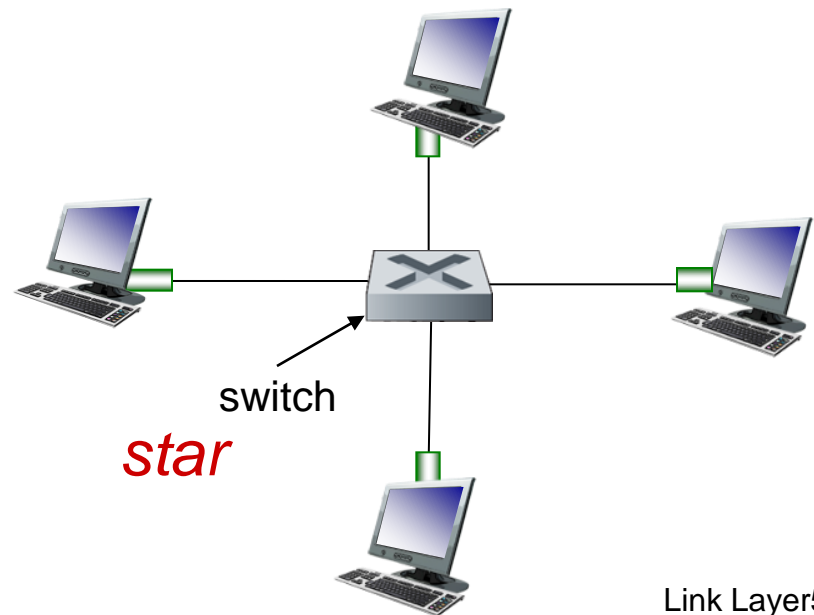
❖ kept up with speed race: 10 Mbps – 10 Gbps

*Metcalfe's Ethernet sketch*

# 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*

# Ethernet frame structure

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

| preamble | dest. address | source address | *type* | data (payload) | CRC |

*preamble:*

❖ 7 bytes with pattern 10101010 followed by one byte with pattern 10101011

❖ used to synchronize receiver, sender clock rates

# Ethernet frame structure (more)

❖ *addresses:* 6 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
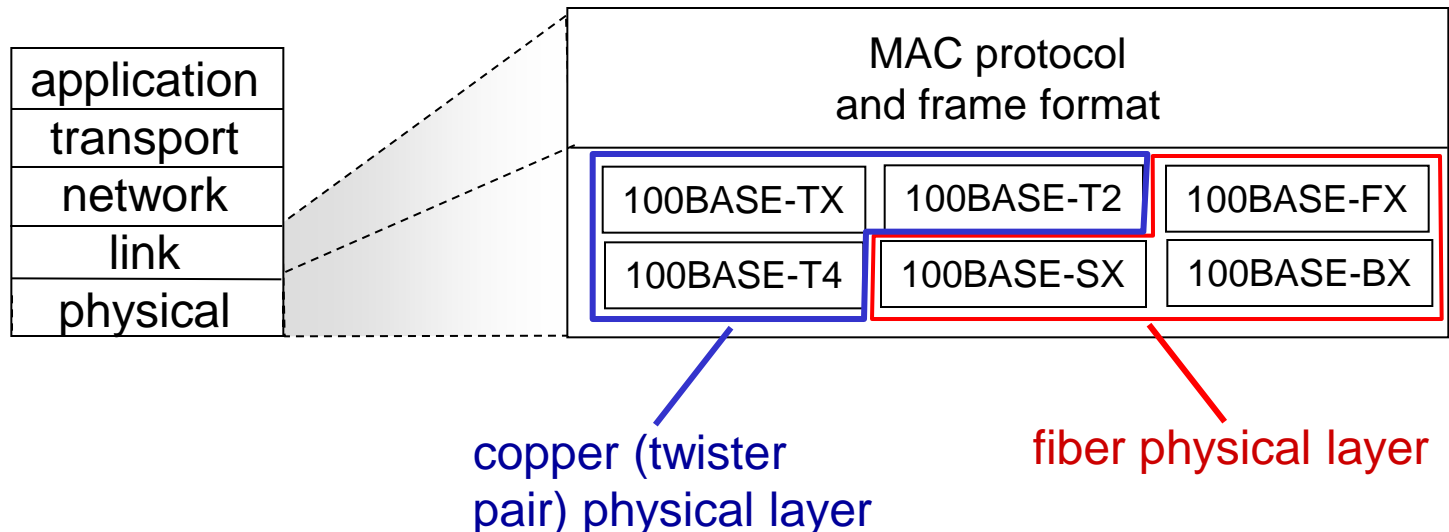  - ▪ error detected: frame is dropped

*type*

| preamble | dest. address | source address | | data (payload) | CRC |
|---|---|---|---|---|---|

# Ethernet: unreliable, connectionless

❖ *connectionless:* no handshaking between sending and receiving NICs

❖ *unreliable:* receiving NIC doesnt send acks 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*

# 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<br>and frame format | | |
|---|---|---|
| 100BASE-TX | 100BASE-T2 | 100BASE-FX |
| 100BASE-T4 | 100BASE-SX | 100BASE-BX |

copper (twister pair) physical layer

fiber physical layer

# 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 Subnetting

5.6 a day in the life of a web request

# Ethernet switch

❖ link-layer device: takes an *active* role
  - store, forward Ethernet frames
  - examine incoming frame's MAC address, selectively forward  frame to one-or-more outgoing links when frame is to be forwarded on segment, uses CSMA/CD to access segment

❖ *transparent*
  - hosts are unaware of presence of switches

❖ *plug-and-play, self-learning*
  - switches do not need to be configured

# Switch: *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*
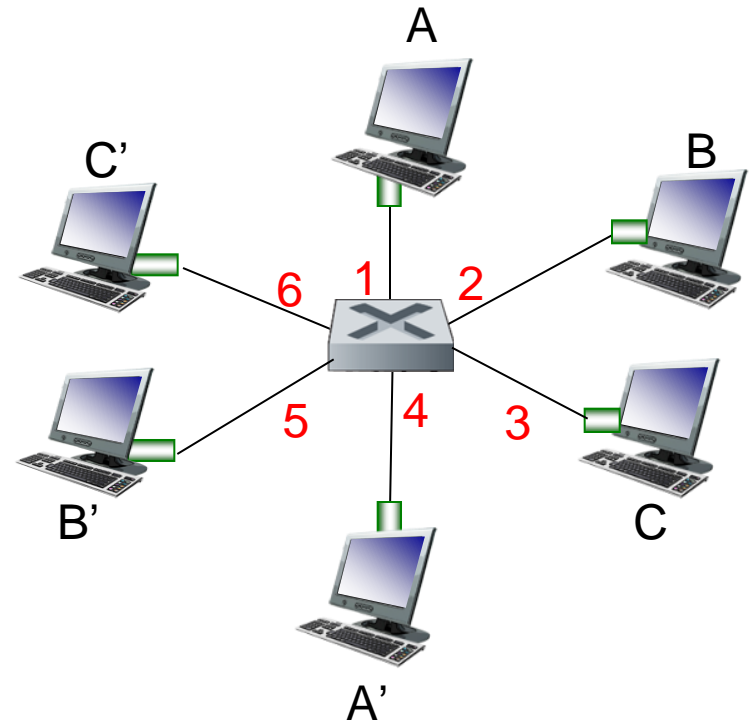*(1,2,3,4,5,6)*

# 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!*

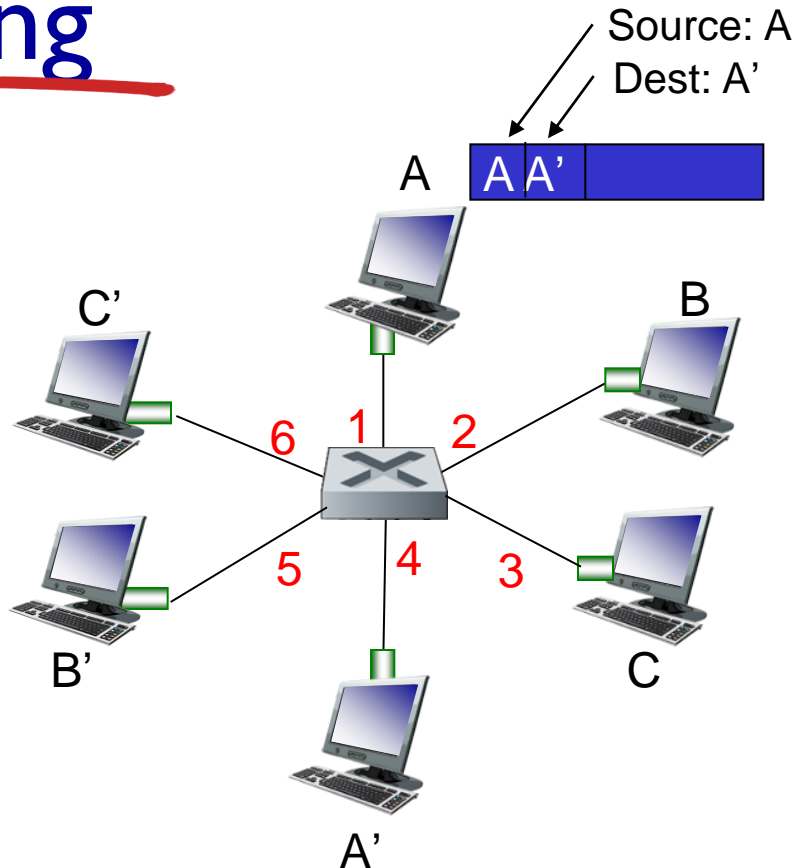*Q: how are entries created, maintained in switch table?*

  ▪ *something like a routing protocol?*



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

# Switch: self-learning

Source: A
Dest: A'

A  A A'

❖ switch *learns* which hosts can be reached through which interfaces

- ■ when frame received, switch "learns" location of sender: incoming LAN segment
- ■ records sender/location pair in switch table

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

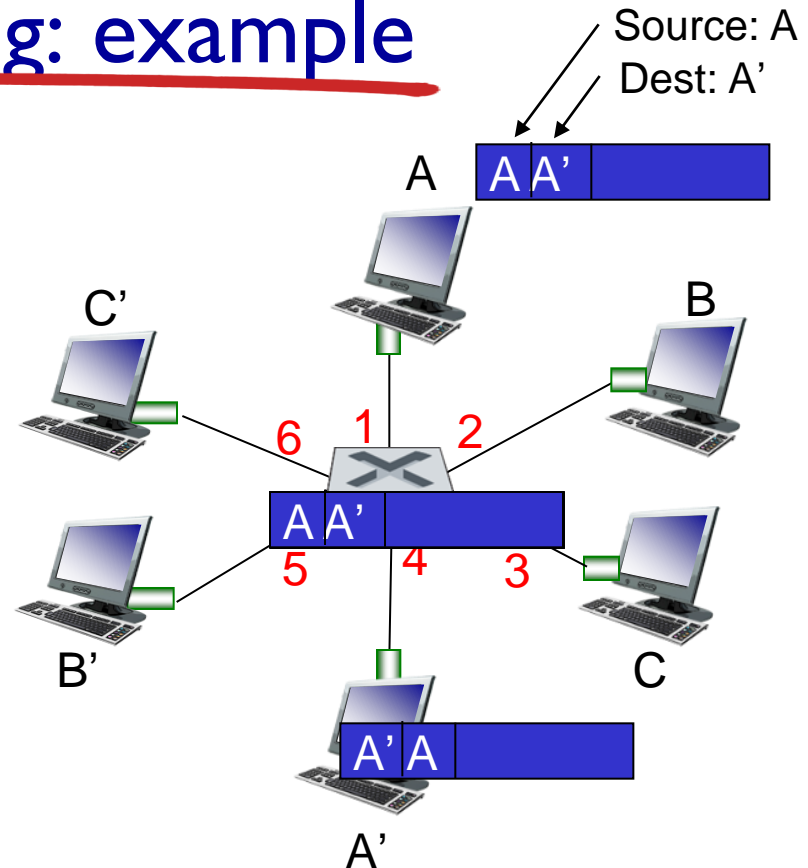*Switch table (initially empty)*

# Switch: frame filtering/forwarding

when  frame received 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 */

# Self-learning, forwarding: example

Source: A
Dest: A'

* frame destination, A', locaton unknown: *flood*

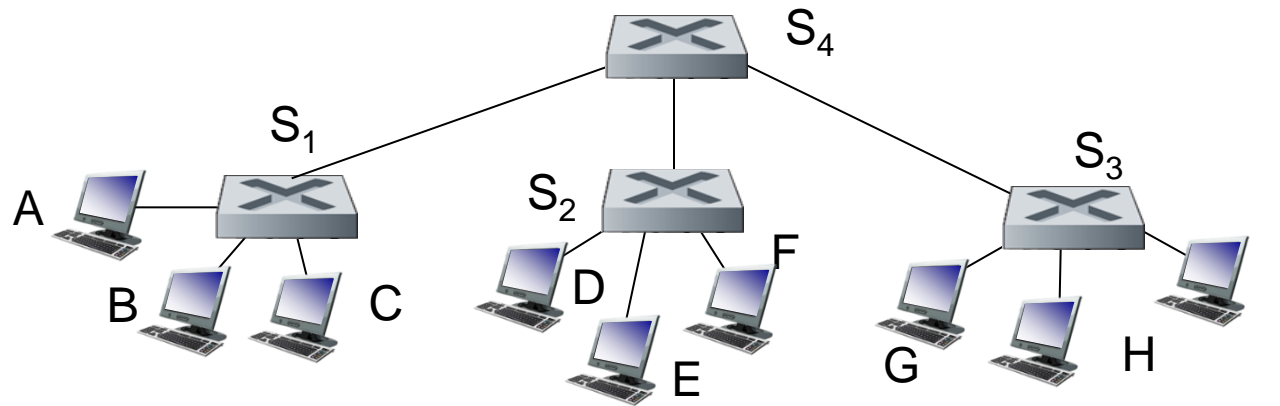* destination A location known: *selectively send on just one link*

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

*switch table (initially empty)*

# Interconnecting switches
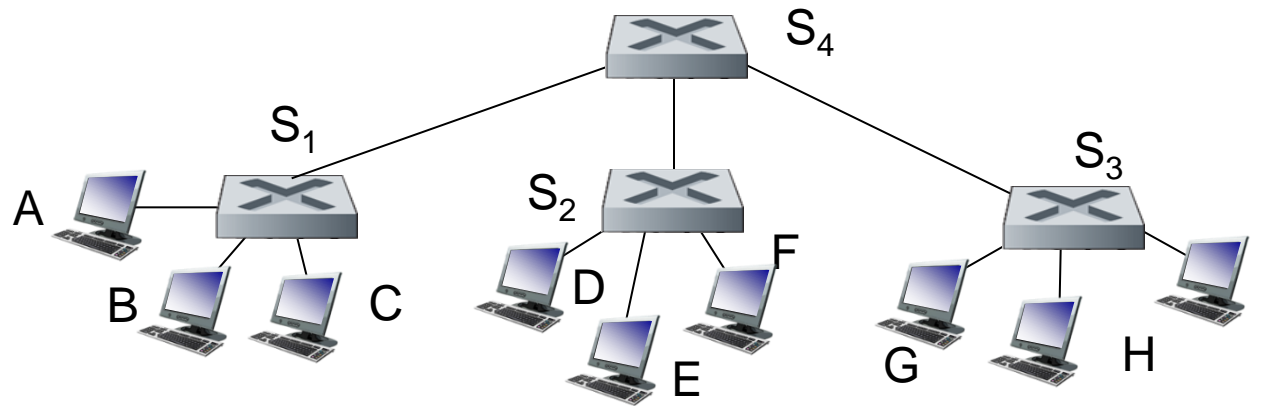
❖ switches can be connected together



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

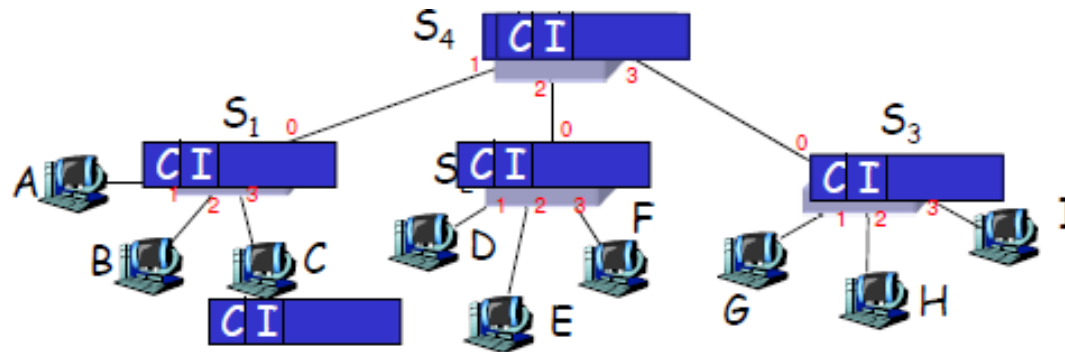❖ *A:* self learning! (works *exactly* the same as in single-switch case!)

# Self-learning multi-switch example

Suppose C sends frame to I, I responds to C



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

# Self-learning multi-switch example



|     | MAC addr | interface | TTL |
| --- | --- | --- | --- |
| $S_1$ | C | 3 | 60 |
|     |     |     |     |
|     |     |     |     |

|     | MAC addr | interface | TTL |
| --- | --- | --- | --- |
| $S_3$ | C | 0 | 60 |
|     |     |     |     |
|     |     |     |     |

|     | MAC addr | interface | TTL |
| --- | --- | --- | --- |
| $S_2$ | C | 0 | 60 |
|     |     |     |     |
|     |     |     |     |

|     | MAC addr | interface | TTL |
| --- | --- | --- | --- |
| $S_4$ | C | 1 | 60 |
|     |     |     |     |
|     |     |     |     |

# Self-learning multi-switch example



| | MAC addr | interface | TTL |
|---|---|---|---|
| $S_1$ | C | 3 | 60 |
| | I | 0 | 60 |

| | MAC addr | interface | TTL |
|---|---|---|---|
| $S_3$ | C | 0 | 60 |
| | I | 3 | 60 |

| | MAC addr | interface | TTL |
|---|---|---|---|
| $S_2$ | C | 0 | 60 |

| | MAC addr | interface | TTL |
|---|---|---|---|
| $S_4$ | C | 1 | 60 |
| | I | 3 | 60 |

# Institutional network



to external network

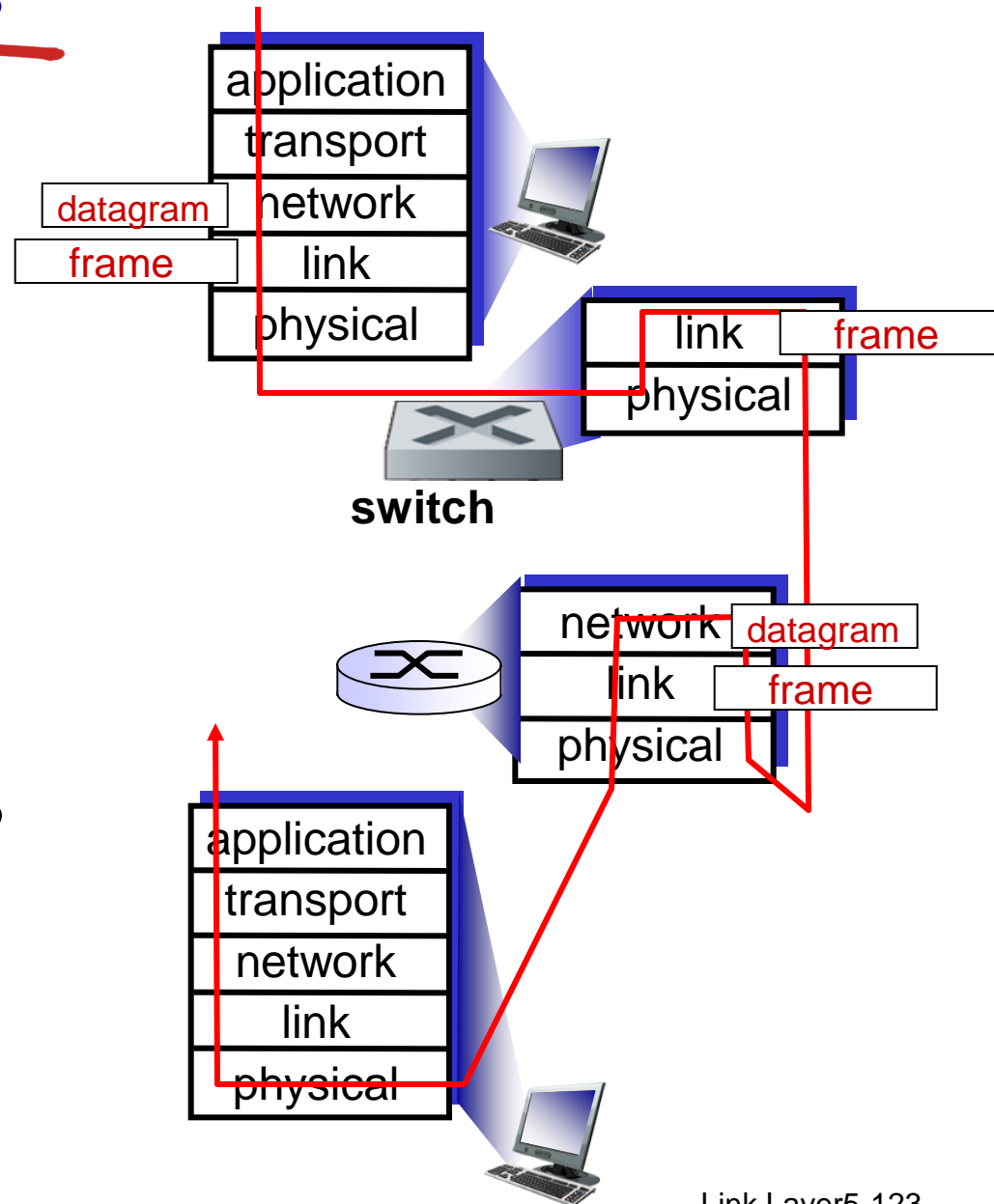router

mail server

web server

IP subnet
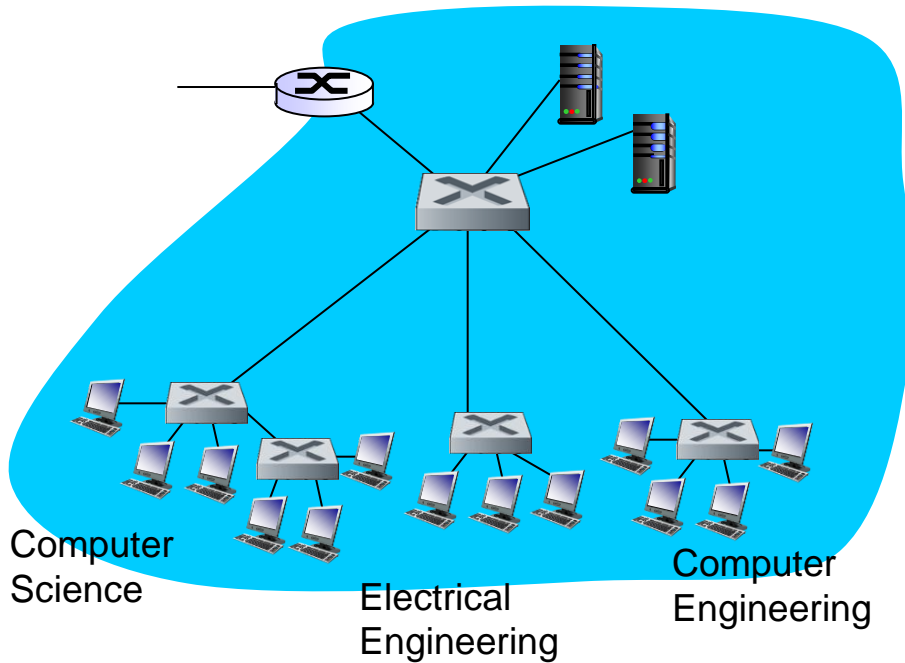
# Switches vs. routers

**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
physical
frame

**switch**

network
link
physical

datagram
frame

application
transport
network
link
physical

# VLANs: motivation



Computer
Science

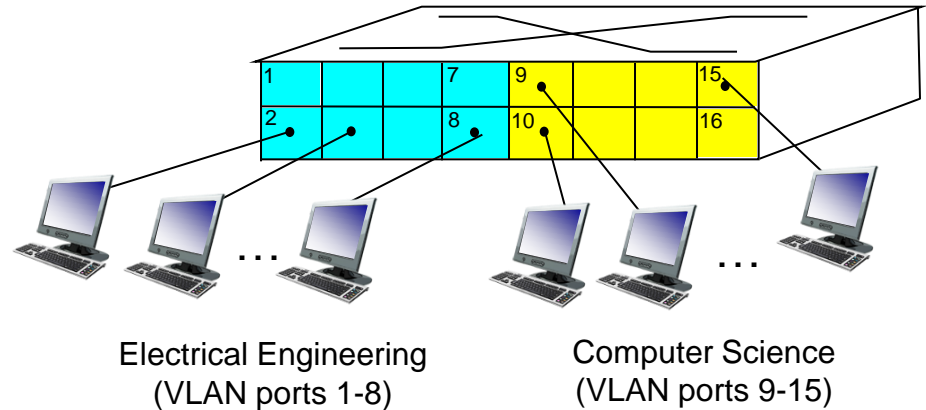Electrical
Engineering

Computer
Engineering

*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

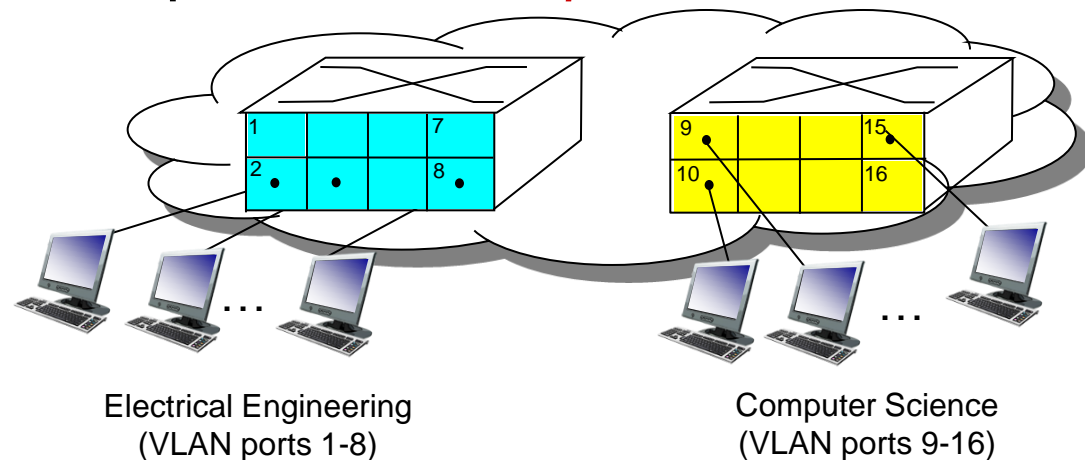  ▪ security/privacy, efficiency issues

# VLANs

**_Virtual Local Area Network_**

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 ……
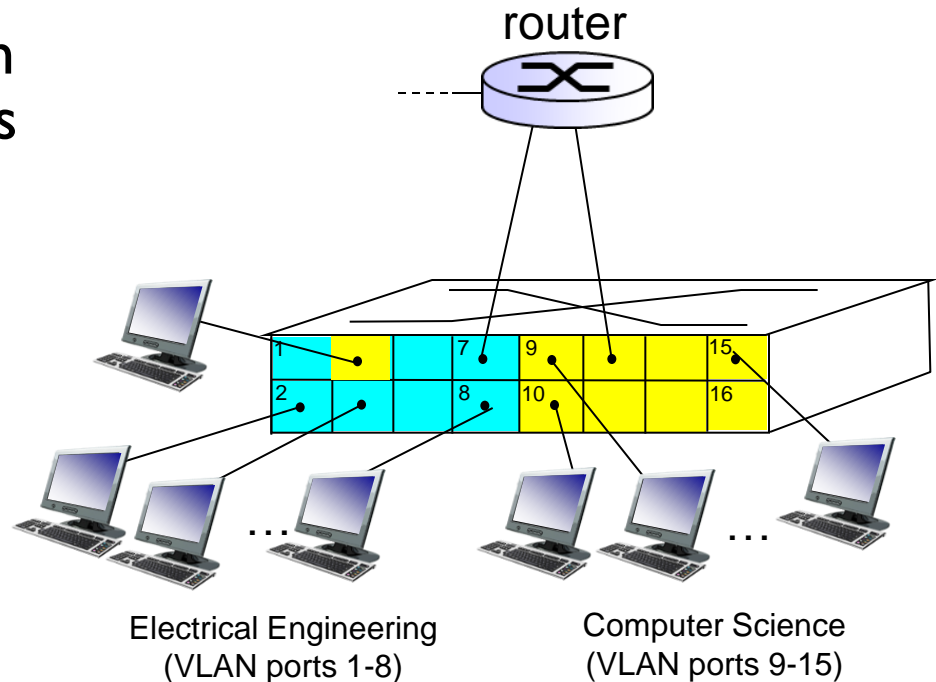


Electrical Engineering
(VLAN ports 1-8)

Computer Science
(VLAN ports 9-15)

… operates as *multiple* virtual switches



Electrical Engineering
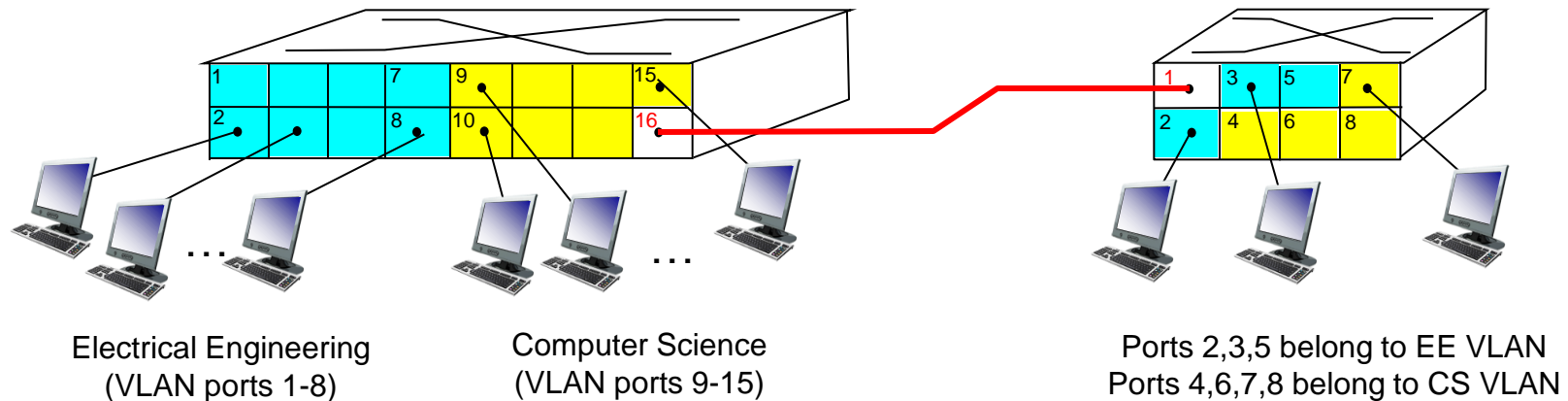(VLAN ports 1-8)

Computer Science
(VLAN ports 9-16)

# 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)

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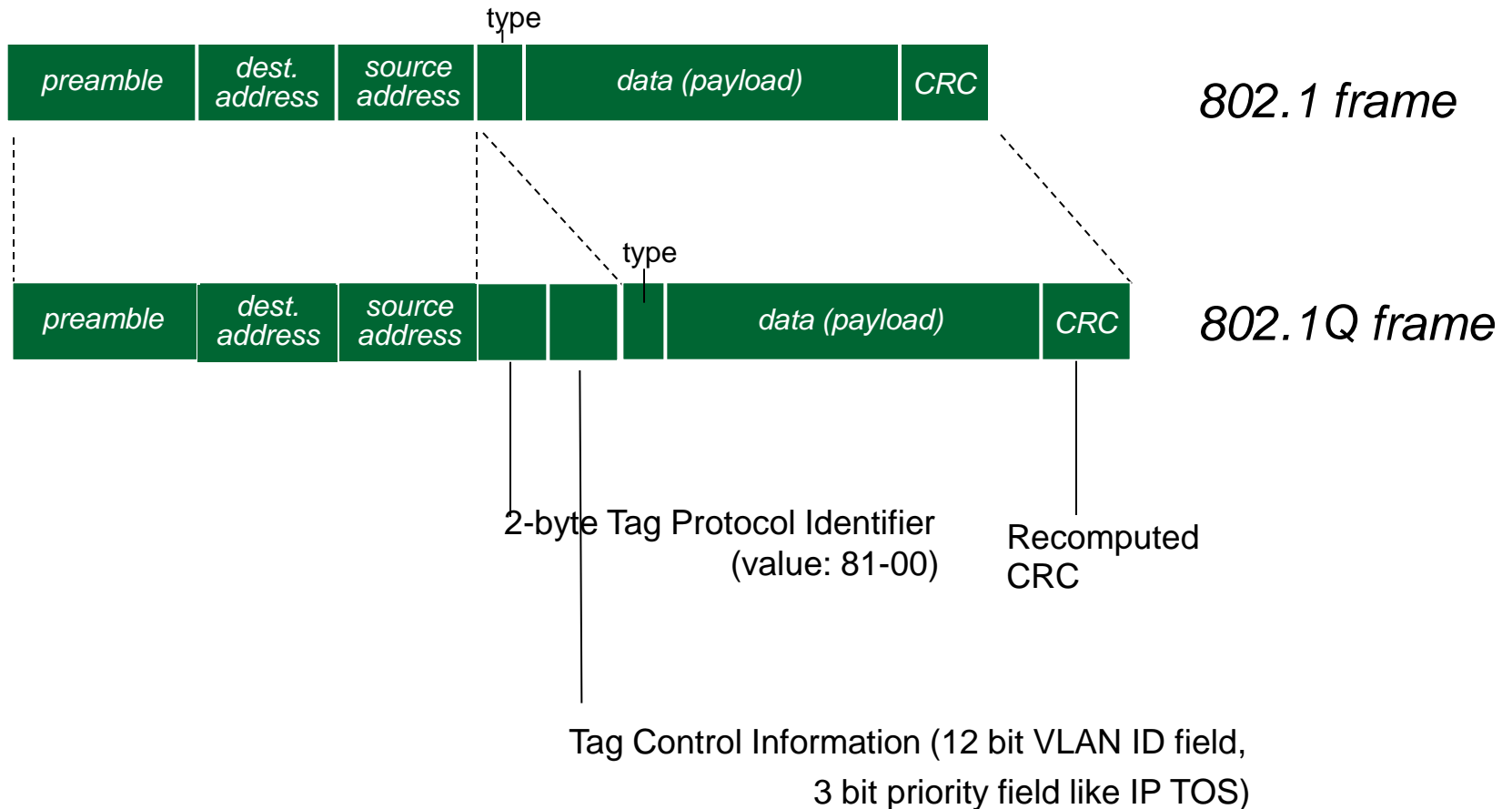
❖ *trunk port:* carries frames between VLANS defined over multiple physical switches
- frames forwarded within VLAN between switches can't be vanilla 802.1 frames (must carry VLAN ID info)
- 802.1q protocol adds/removed additional header fields for frames forwarded between trunk ports

# 802.1Q VLAN frame format

type

| preamble | dest. address | source address | | data (payload) | CRC |

*802.1 frame*

type

| preamble | dest. address | source address | | | | data (payload) | CRC |

*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)

# Link layer, LANs: outline

# Datagram networks

❖ no call setup at network layer
❖ routers: no state about end-to-end connections
  ▪ no network-level concept of "connection"
❖ packets forwarded using destination host address



1. send datagrams

2. receive datagrams

# Datagram forwarding table



| local forwarding table | |
| --- | --- |
| dest address | output link |
| address-range 1 | 3 |
| address-range 2 | 2 |
| address-range 3 | 2 |
| address-range 4 | 1 |

routing algorithm

4 billion IP addresses, so rather than list individual destination address list *range* of addresses (aggregate table entries)

IP destination address in arriving packet's header

# Datagram forwarding  table

| Destination Address Range | Link Interface |
|---|---|
| 11001000 00010111 00010000 00000000<br>through<br>11001000 00010111 00010111 11111111 | 0 |
| 11001000 00010111 00011000 00000000<br>through<br>11001000 00010111 00011000 11111111 | 1 |
| 11001000 00010111 00011001 00000000<br>through<br>11001000 00010111 00011111 11111111 | 2 |
| otherwise | 3 |

*Q:* but what happens if ranges don't divide up so nicely?

# Longest prefix matching

*longest prefix matching*
when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

| Destination Address Range | Link interface |
|---|---|
| `11001000 00010111 00010*** ********` | 0 |
| `11001000 00010111 00011000 ********` | 1 |
| `11001000 00010111 00011*** ********` | 2 |
| otherwise | 3 |

examples:

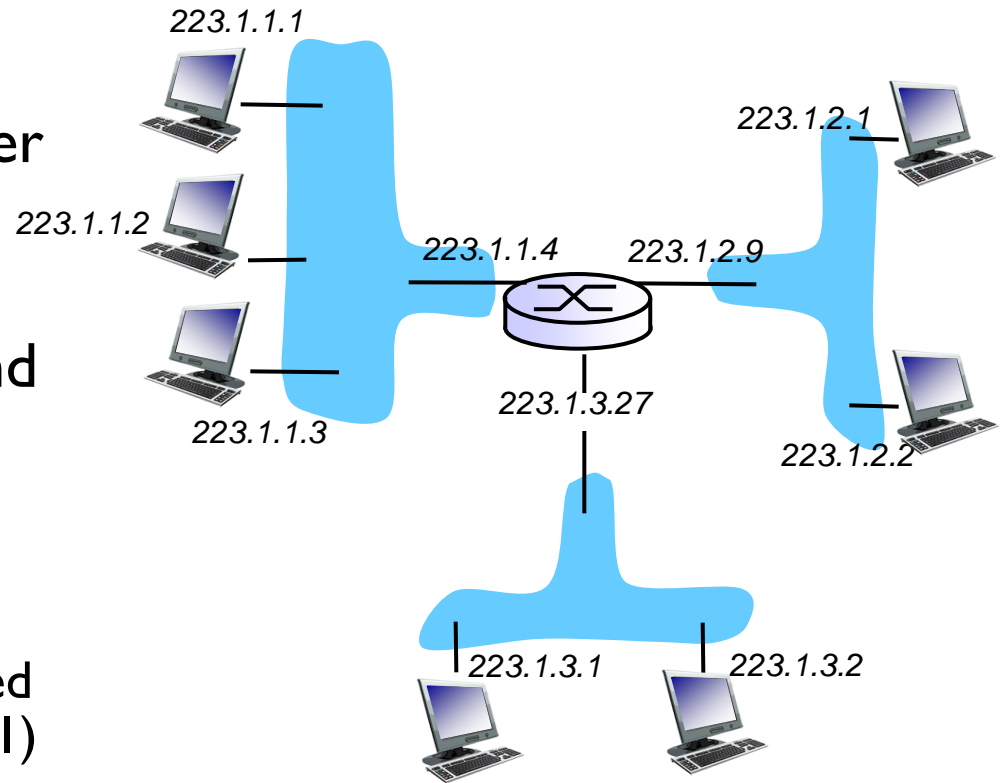DA: 11001000  00010111  00010110  10100001      which interface?  0

DA: 11001000  00010111  00011000  10101010      which interface?  1

# IP addressing: introduction

❖ *IP address:* 32-bit identifier for host, router *interface*

❖ *interface:* connection between host/router and physical link
  ▪ router's typically have multiple interfaces
  ▪ host typically has one or two interfaces (e.g., wired Ethernet, wireless 802.11)

❖ *IP addresses associated with each interface*

223.1.1.1

223.1.1.2

223.1.1.4    223.1.2.9

223.1.1.3

223.1.3.27

223.1.2.1

223.1.2.2

223.1.3.1    223.1.3.2

*223.1.1.1 = 11011111 00000001 00000001 00000001*

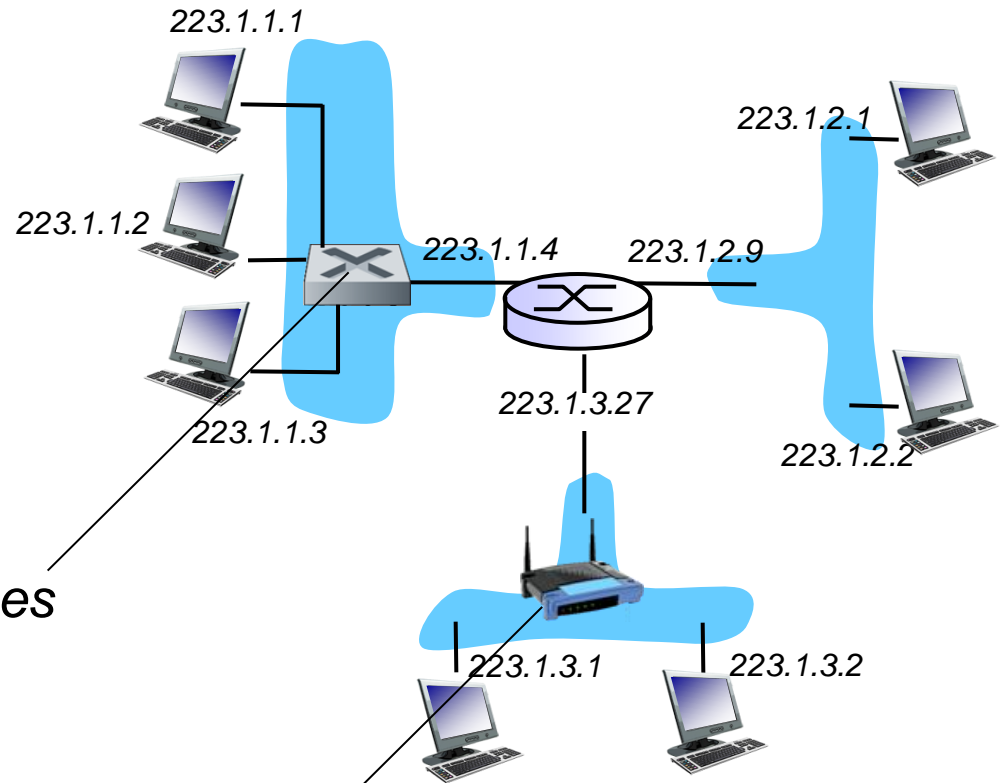*223          1          1          1*

# IP addressing: introduction

*Q: how are interfaces actually connected?*

*A: we'll learn about that in chapter 5, 6.*

A: *wired Ethernet interfaces connected by Ethernet switches*

*For now: don't need to worry about how one interface is connected to another (with no intervening router)*

223.1.1.1

223.1.2.1

223.1.1.2

223.1.1.4    223.1.2.9

223.1.1.3

223.1.3.27

223.1.2.2

223.1.3.1    223.1.3.2

A: *wireless WiFi interfaces connected by WiFi base station*

# Subnets
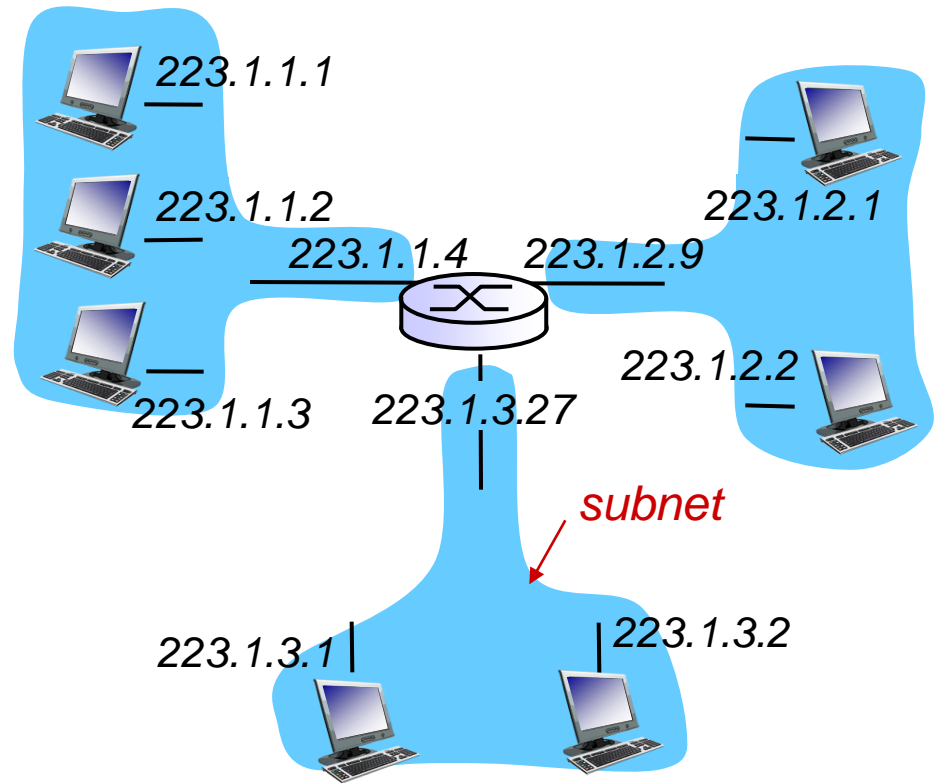
❖ IP address:

- subnet part - high order bits
- host part - low order bits

❖ *what's a subnet ?*

- device interfaces with same subnet part of IP address
- can physically reach each other *without intervening router*



*223.1.1.1*

*223.1.1.2*
*223.1.1.4*  *223.1.2.9*

*223.1.1.3*  *223.1.3.27*

*223.1.2.1*

*223.1.2.2*

*subnet*

*223.1.3.1*  *223.1.3.2*

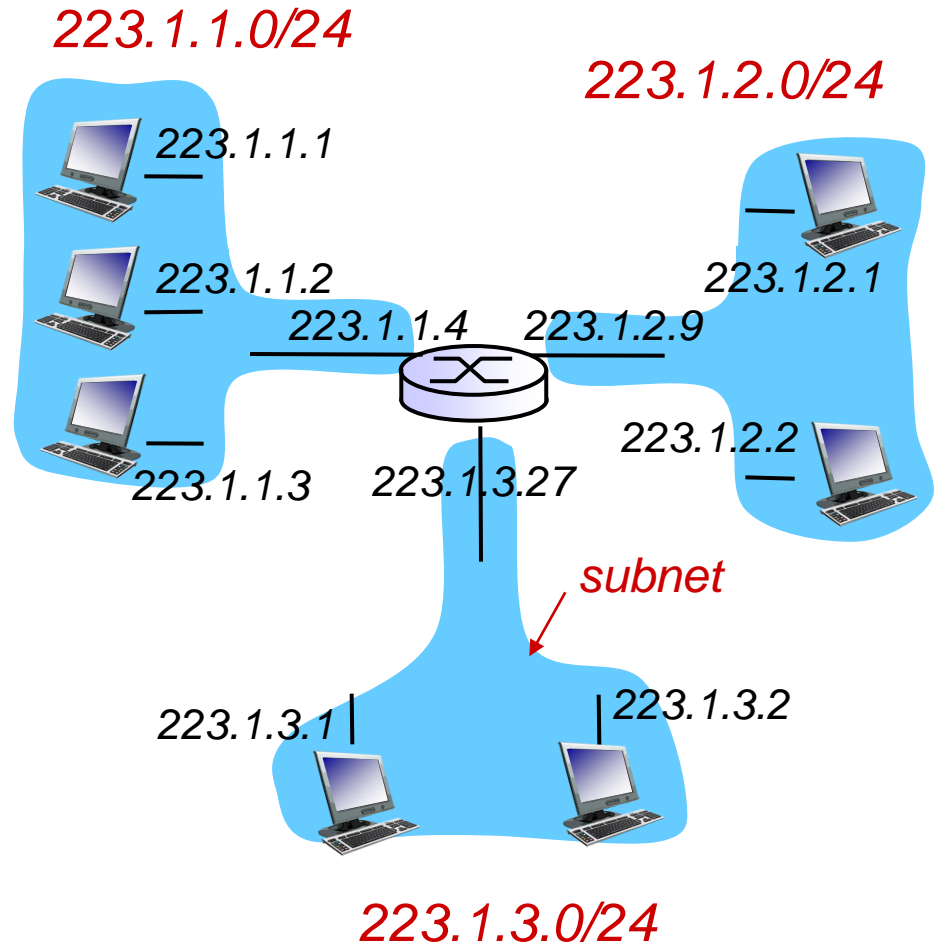*network consisting of 3 subnets*

# Subnets

*recipe*

❖ to determine the subnets, detach each interface from its host or router, creating islands of isolated networks

❖ each isolated network is called a *subnet*

223.1.1.0/24

223.1.2.0/24

223.1.1.1

223.1.1.2

223.1.1.4      223.1.2.9

223.1.2.1

223.1.2.2

223.1.1.3     223.1.3.27

*subnet*

223.1.3.1

223.1.3.2

223.1.3.0/24

*subnet mask: /24*

# Subnets

how many?



223.1.1.2

223.1.1.1

223.1.1.4

223.1.1.3

223.1.9.2

223.1.7.0

223.1.9.1

223.1.7.1

223.1.8.1

223.1.8.0

223.1.2.6

223.1.3.27

223.1.2.1

223.1.2.2

223.1.3.1

223.1.3.2

# IP addressing: CIDR

CIDR: Classless InterDomain Routing
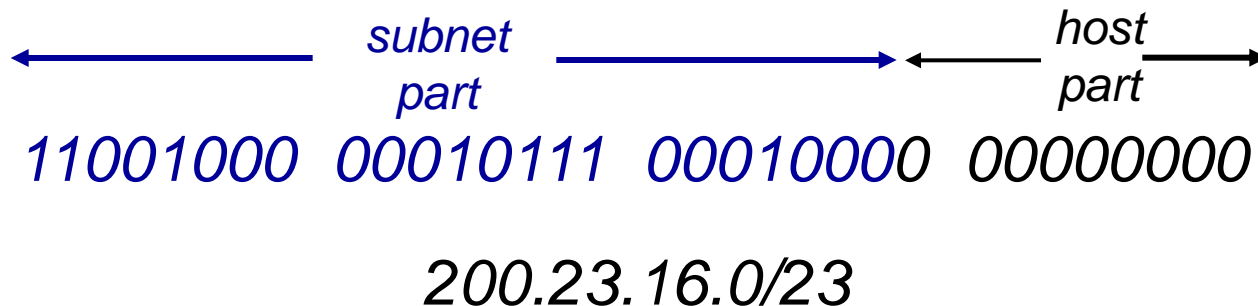- subnet portion of address of arbitrary length
- address format: a.b.c.d/x, where x is # bits in subnet portion of address



*subnet part* ←→ *host part*

*11001000  00010111  00010000  00000000*

*200.23.16.0/23*

# IP addresses: how to get one?

Q: How does a *host* get IP address?

❖ hard-coded by system admin in a file
  ▪ Windows: control-panel->network->configuration->tcp/ip->properties
  ▪ UNIX: /etc/rc.config
❖ DHCP: Dynamic Host Configuration Protocol: dynamically get address from as server
  ▪ "plug-and-play"

# DHCP: Dynamic Host Configuration Protocol

*goal:* allow host to *dynamically* obtain its IP address from network server when it joins network

- can renew its lease on address in use
- allows reuse of addresses (only hold address while connected/"on")
- support for mobile users who want to join network (more shortly)

*DHCP overview:*

- host broadcasts "DHCP discover" msg [optional]
- DHCP server responds with "DHCP offer" msg [optional]
- host requests IP address: "DHCP request" msg
- DHCP server sends address: "DHCP ack" msg

# DHCP client-server scenario



*223.1.1.0/24*

223.1.1.1

223.1.1.2

223.1.1.4    223.1.2.9

223.1.1.3    223.1.3.27

DHCP server

223.1.2.1

223.1.2.2

*223.1.2.0/24*

*arriving DHCP client needs address in this network*

223.1.3.1    223.1.3.2

*223.1.3.0/24*

# DHCP client-server scenario

**DHCP server: 223.1.2.5**

*arriving client*

**DHCP discover**

> src : 0.0.0.0, 68
> dest.: 255.255.255.255,67
> yiaddr:   0.0.0.0
> transaction ID: 654

**DHCP offer**

> src: 223.1.2.5, 67
> dest:  255.255.255.255, 68
> yiaddrr: 223.1.2.4
> transaction ID: 654
> lifetime: 3600 secs

**DHCP request**

> src:  0.0.0.0, 68
> dest:: 255.255.255.255, 67
> yiaddrr: 223.1.2.4
> transaction ID: 655
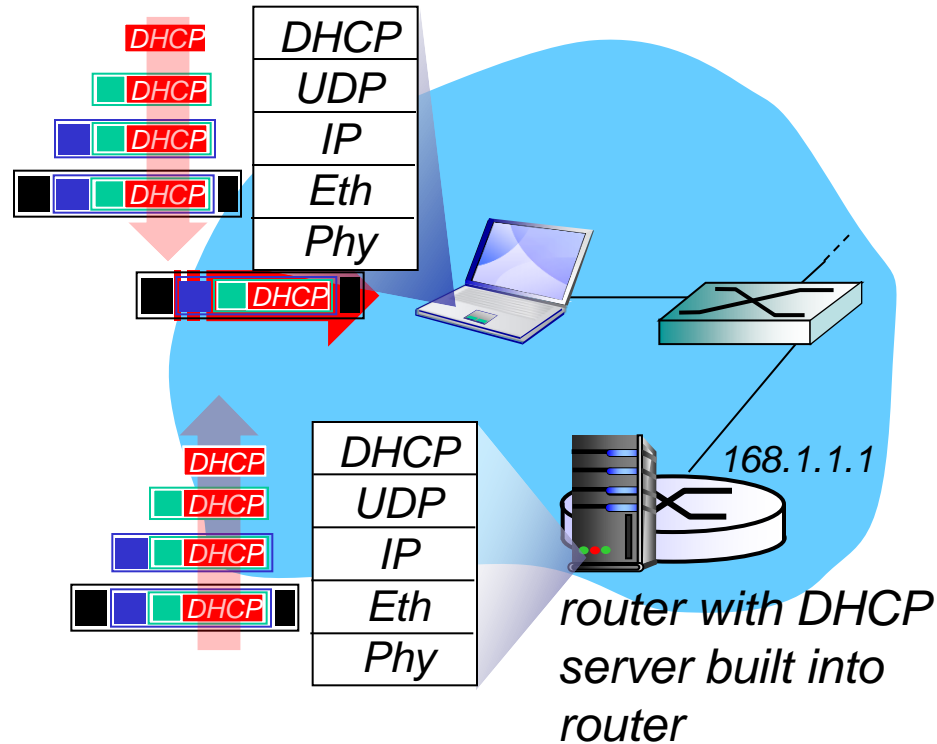> lifetime: 3600 secs

**DHCP ACK**

> src: 223.1.2.5, 67
> dest:  255.255.255.255, 68
> yiaddrr: 223.1.2.4
> transaction ID: 655
> lifetime: 3600 secs

# DHCP: more than IP addresses

DHCP can return more than just allocated IP address on subnet:

- address of first-hop router for client
- name and IP address of DNS sever
- network mask (indicating network versus host portion of address)

# DHCP: example



DHCP · UDP · IP · Eth · Phy

DHCP · UDP · IP · Eth · Phy

*router with DHCP server built into router*

168.1.1.1

- ❖ connecting laptop needs its IP address, addr of first-hop router, addr of DNS server: use DHCP

- ❖ *DHCP request encapsulated in UDP, encapsulated in IP, encapsulated in 802.1 Ethernet*

- ❖ *Ethernet frame broadcast (dest: FFFFFFFFFFFF) on LAN, received at router running DHCP server*

- ❖ *Ethernet demuxed to IP demuxed, UDP demuxed to DHCP*

# DHCP: example



router with DHCP
server built into
router

❖ 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 of DHCP
server, frame forwarded to
client, demuxing up to
DHCP at client*

❖ *client now knows its IP
address, name and IP
address of DSN server, IP
address of its first-hop
router*

# DHCP: Wireshark output (home LAN)

*request*

Message type: **Boot Request (1)**
Hardware type: Ethernet
Hardware address length: 6
Hops: 0
**Transaction ID: 0x6b3a11b7**
Seconds elapsed: 0
Bootp flags: 0x0000 (Unicast)
Client IP address: 0.0.0.0 (0.0.0.0)
Your (client) IP address: 0.0.0.0 (0.0.0.0)
Next server IP address: 0.0.0.0 (0.0.0.0)
Relay agent IP address: 0.0.0.0 (0.0.0.0)
**Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)**
Server host name not given
Boot file name not given
Magic cookie: (OK)
Option: (t=53,l=1) **DHCP Message Type = DHCP Request**
Option: (61) Client identifier
   Length: 7; Value: 010016D323688A;
   Hardware type: Ethernet
   Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)
Option: (t=50,l=4) Requested IP Address = 192.168.1.101
Option: (t=12,l=5) Host Name = "nomad"
**Option: (55) Parameter Request List**
   Length: 11; Value: 010F03062C2E2F1F21F92B
   **1 = Subnet Mask; 15 = Domain Name**
   **3 = Router; 6 = Domain Name Server**
   44 = NetBIOS over TCP/IP Name Server
   ……

*reply*

Message type: **Boot Reply (2)**
Hardware type: Ethernet
Hardware address length: 6
Hops: 0
**Transaction ID: 0x6b3a11b7**
Seconds elapsed: 0
Bootp flags: 0x0000 (Unicast)
**Client IP address: 192.168.1.101 (192.168.1.101)**
Your (client) IP address: 0.0.0.0 (0.0.0.0)
**Next server IP address: 192.168.1.1 (192.168.1.1)**
Relay agent IP address: 0.0.0.0 (0.0.0.0)
Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)
Server host name not given
Boot file name not given
Magic cookie: (OK)
**Option: (t=53,l=1) DHCP Message Type = DHCP ACK**
**Option: (t=54,l=4) Server Identifier = 192.168.1.1**
**Option: (t=1,l=4) Subnet Mask = 255.255.255.0**
**Option: (t=3,l=4) Router = 192.168.1.1**
**Option: (6) Domain Name Server**
   **Length: 12; Value: 445747E2445749F244574092;**
   **IP Address: 68.87.71.226;**
   **IP Address: 68.87.73.242;**
   **IP Address: 68.87.64.146**
**Option: (t=15,l=20) Domain Name = "hsd1.ma.comcast.net."**

# IP addresses: how to get one?

*Q:* how does *network* get subnet part of IP addr?

*A:* gets allocated portion of its provider ISP's address space

| | | | | |
|---|---|---|---|---|
| *ISP's block* | *11001000  00010111  00010000* | *00000000* | *200.23.16.0/20* | |
| | | | | |
| *Organization 0* | *11001000  00010111  00010000* | *00000000* | *200.23.16.0/23* | |
| *Organization 1* | *11001000  00010111  00010010* | *00000000* | *200.23.18.0/23* | |
| *Organization 2* | *11001000  00010111  00010100* | *00000000* | *200.23.20.0/23* | |
| *...* | *…..* | *….* | *….* | |
| *Organization 7* | *11001000  00010111  00011110* | *00000000* | *200.23.30.0/23* | |

# Hierarchical addressing: route aggregation

*hierarchical addressing allows efficient advertisement of routing information:*

*Organization 0*

200.23.16.0/23

*Organization 1*

200.23.18.0/23

*Organization 2*

200.23.20.0/23

*Organization 7*

200.23.30.0/23

*Fly-By-Night-ISP*

*ISPs-R-Us*

*"Send me anything with addresses beginning 200.23.16.0/20"*

*"Send me anything with addresses beginning 199.31.0.0/16"*

*Internet*

# Hierarchical addressing: more specific routes

*ISPs-R-Us has a more specific route to Organization 1*

Organization 0
200.23.16.0/23

Organization 2
200.23.20.0/23

Organization 7
200.23.30.0/23

Fly-By-Night-ISP

"Send me anything
with addresses
beginning
200.23.16.0/20"

Internet

ISPs-R-Us

Organization 1
200.23.18.0/23

"Send me anything
with addresses
beginning 199.31.0.0/16
or 200.23.18.0/23"

# IP addressing: the last word...
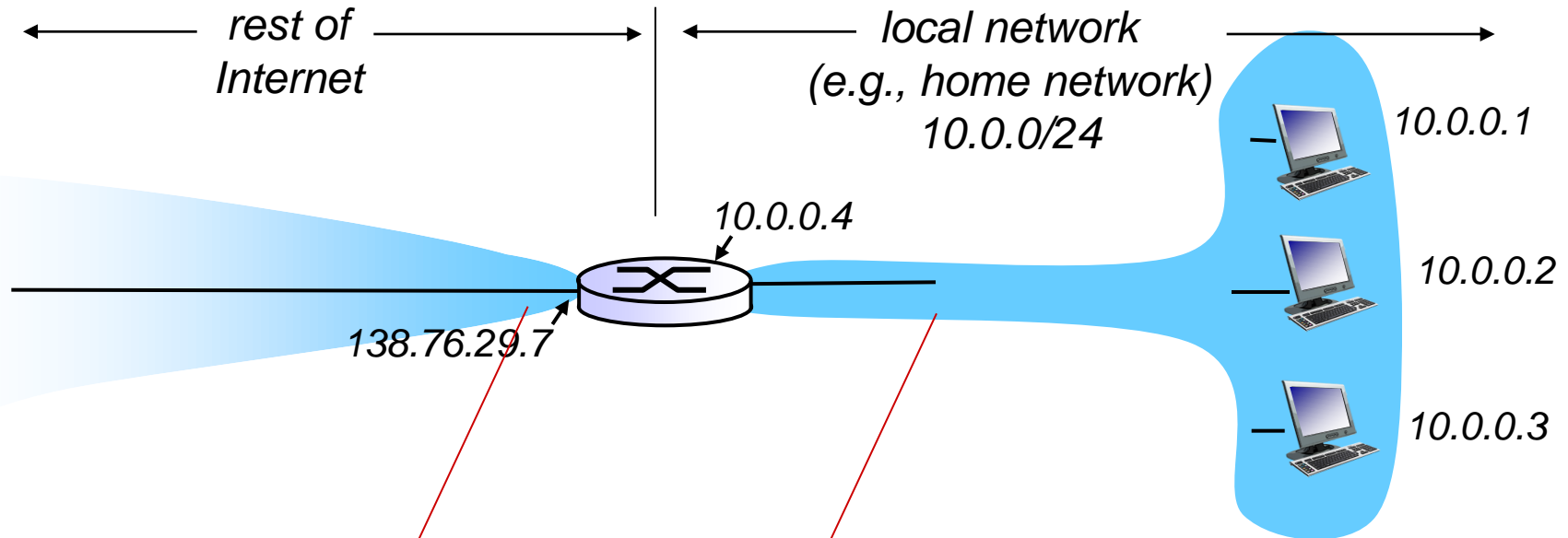
*Q:* how does an ISP get block of addresses?

*A:* ICANN: Internet Corporation for Assigned
   Names and Numbers http://www.icann.org/
   - allocates addresses
   - manages DNS
   - assigns domain names, resolves disputes

# NAT: network address translation

rest of Internet

local network (e.g., home network) 10.0.0/24

10.0.0.1

10.0.0.4

10.0.0.2

138.76.29.7

10.0.0.3

*all datagrams leaving local network have same single source NAT IP address: 138.76.29.7, different source port numbers*

*datagrams with source or destination in this network have 10.0.0/24 address for source, destination (as usual) (or any address belonging in 10.0.0.0/8, 172.16.0.0/12, 192.168.0.0/16)*

# NAT: network address translation

*motivation:* local network uses just one IP address as far as outside world is concerned:
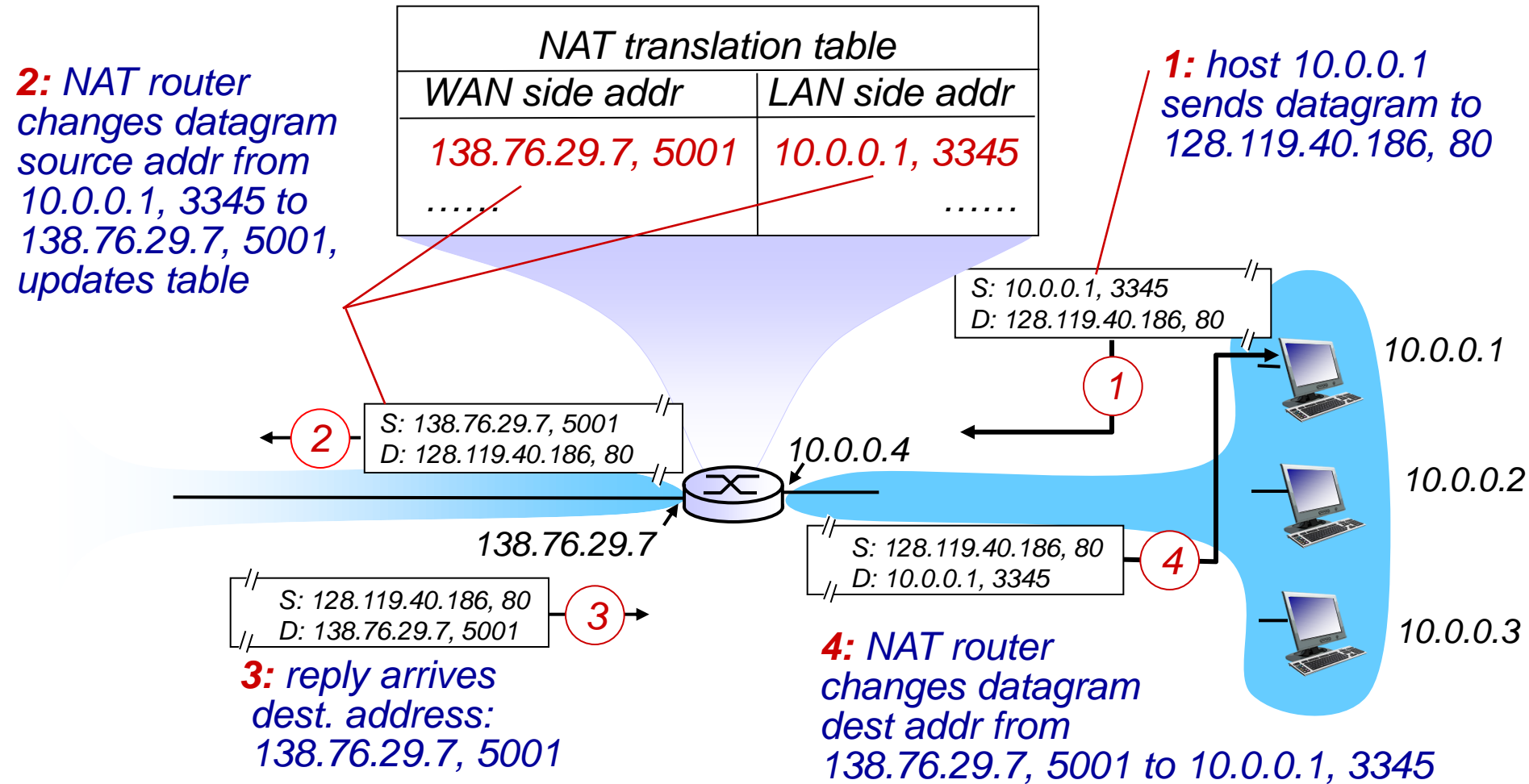
- range of addresses not needed from ISP:  just one IP address for all devices
- can change addresses of devices in local network without notifying outside world
- can change ISP without changing addresses of devices in local network
- devices inside local net not explicitly addressable, visible by outside world (a security plus)

# NAT: network address translation

*implementation*: NAT router must:

- *outgoing datagrams: replace* (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #)
    - . . . remote clients/servers will respond using (NAT IP address, new port #) as destination addr

- *remember (in NAT translation table)* every (source IP address, port #)  to (NAT IP address, new port #) translation pair

- *incoming datagrams: replace* (NAT IP address, new port #) in dest fields of every incoming datagram with corresponding (source IP address, port #) stored in NAT table

# NAT: network address translation

**2:** *NAT router changes datagram source addr from 10.0.0.1, 3345 to 138.76.29.7, 5001, updates table*

| NAT translation table | |
|---|---|
| WAN side addr | LAN side addr |
| 138.76.29.7, 5001 | 10.0.0.1, 3345 |
| …… | …… |

**1:** *host 10.0.0.1 sends datagram to 128.119.40.186, 80*

S: 10.0.0.1, 3345
D: 128.119.40.186, 80

① 

10.0.0.1

2   S: 138.76.29.7, 5001
    D: 128.119.40.186, 80

10.0.0.4

138.76.29.7

S: 128.119.40.186, 80
D: 10.0.0.1, 3345    ④

10.0.0.2

3   S: 128.119.40.186, 80
    D: 138.76.29.7, 5001

**3:** *reply arrives dest. address: 138.76.29.7, 5001*

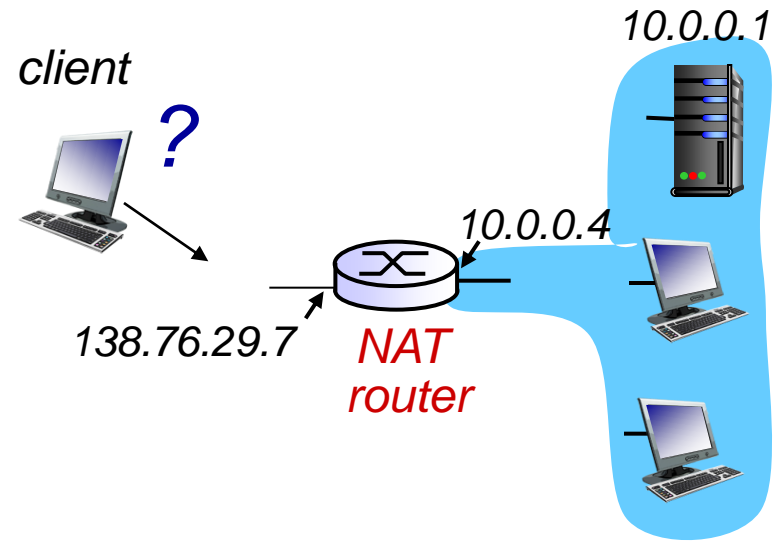**4:** *NAT router changes datagram dest addr from 138.76.29.7, 5001 to 10.0.0.1, 3345*

10.0.0.3

# NAT: network address translation

❖ 16-bit port-number field:
- 60,000 simultaneous connections with a single LAN-side address!

❖ NAT is controversial:
- routers should only process up to layer 3
- violates end-to-end argument
  - NAT possibility must be taken into account by app designers, e.g., P2P applications
- address shortage should instead be solved by IPv6

# NAT traversal problem

❖ **client wants to connect to server with address 10.0.0.1**

  ▪ server address 10.0.0.1 local to LAN (client can't use it as destination addr)

  ▪ only one externally visible NATed address: 138.76.29.7

❖ *solution1:* statically configure NAT to forward incoming connection requests at given port to server

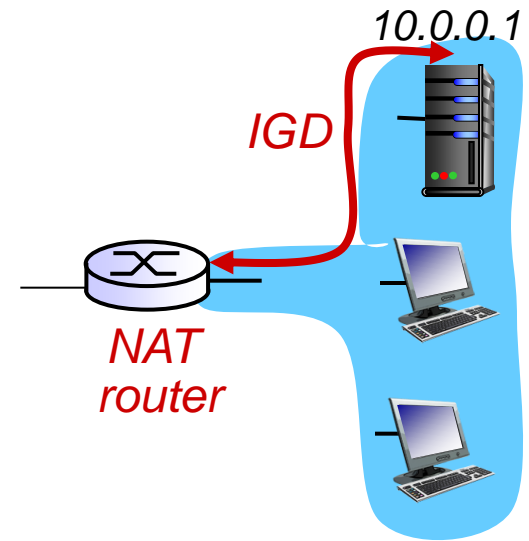  ▪ e.g., (123.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000

*client*

?

10.0.0.1

10.0.0.4

138.76.29.7

*NAT router*

# NAT traversal problem

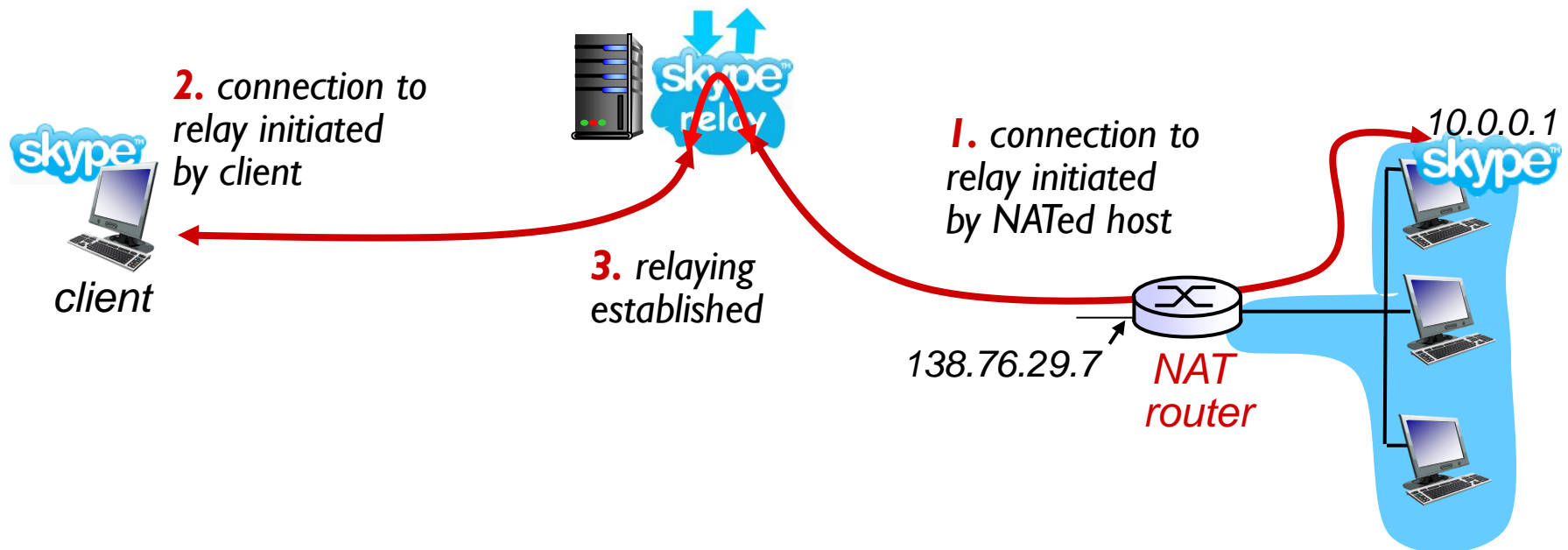❖ *solution 2:* Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol.  Allows NATed host to:
  ❖ learn public IP address (138.76.29.7)
  ❖ add/remove port mappings (with lease times)

  i.e., automate static NAT port map configuration

10.0.0.1

*IGD*

*NAT router*

# NAT traversal problem

❖ *solution 3:* relaying (used in Skype)
- ▪ NATed client establishes connection to relay
- ▪ external client connects to relay
- ▪ relay bridges packets between to connections



**2.** *connection to relay initiated by client*

*client*

**3.** *relaying established*

**1.** *connection to relay initiated by NATed host*

*10.0.0.1*

*138.76.29.7*

*NAT router*

# DNS: Domain Name System

Για την αναγνώριση των hosts, routers στο Διαδίκτυο χρησιμοποιούνται:

- διεύθυνση IP (32 bits) – χρησιμοποιείται για διευθυνσιοδότηση των datagrams π.χ. 147.27.18.11
- "όνομα", π.χ. kronos.telecom.tuc.gr – διευκολύνει ανθρώπους

## Domain Name System:

- υπηρεσία καταλόγου που μεταφράζει hostname ⇔ διεύθυνση IP
- κατανεμημένη βάση δεδομένων που υλοποιείται από μία ιεραρχία από πολλούς name servers

- πρωτόκολλο επιπέδου εφαρμογής hosts και name servers επικοινωνούν για την ανάλυση ονομάτων σε διευθύνσεις IP και αντίστροφα

  - σημείωση: λειτουργία του πυρήνα του Διαδικτύου που υλοποιείται ως πρωτόκολλο του επιπέδου εφαρμογής
  - πολυπλοκότητα στο άκρο του δικτύου

- χρησιμοποιεί UDP, port 53
- μοντέλο client/server

- DNS χρησιμοποιείται από άλλα πρωτόκολλα εφαρμογών π.χ. HTTP, SMTP, FTP
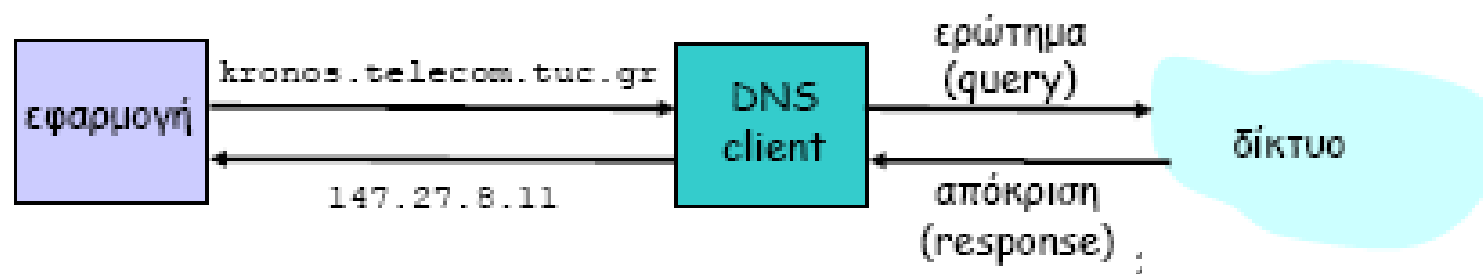
# DNS: Domain Name System

## Υπηρεσίες DNS

- Μετάφραση ονόματος host σε διεύθυνση IP
- Ψευδώνυμο host
  - κανονικό (canonical) όνομα και ψευδώνυμα
- Ψευδώνυμο mail server
- Κατανομή φορτίου
  - επαναληπτικοί web servers: σε ένα κανονικό όνομα αντιστοιχούν περισσότερες από μία διευθύνσεις IP

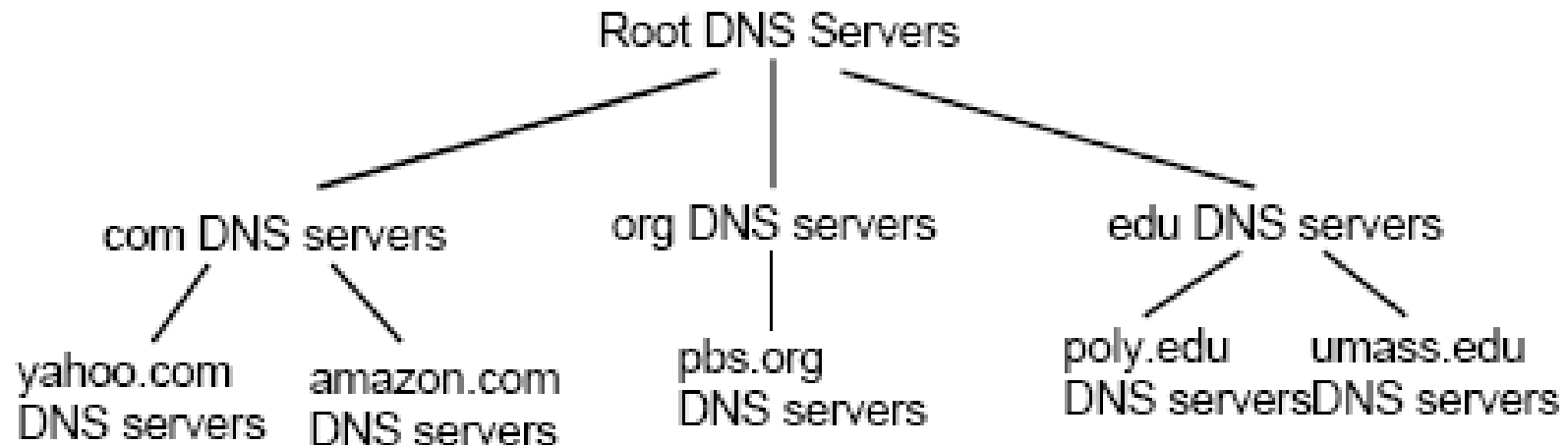## Γιατί όχι κεντρικοποιημένο DNS;

- μοναδικό σημείο αστοχίας
- όγκος κίνησης
- κεντρικοποιημένη βάση δεδομένων σε μεγάλη απόσταση
- συντήρηση τεράστιας βάσης δεδομένων με συχνές ενημερώσεις

## δεν έχει *scalability*

# Κατανεμημένη, ιεραρχική βάση δεδομένων



**Ένας client επιθυμεί τη διεύθυνση IP του host www.amazon.com (χονδρική προσέγγιση):**

- ❑ Ο client στέλνει ερώτημα σε έναν root server αναζητώντας τον DNS server του domain com

- ❑ Ο client στέλνει ερώτημα σε έναν DNS server του domain com αναζητώντας τον DNS server του subdomain amazon.com

- ❑ Ο client στέλνει ερώτημα στον DNS server του amazon.com αναζητώντας τη διεύθυνση IP του host www.amazon.com

# DNS name servers

□ κανένας server δεν γνωρίζει όλα τα ζεύγη hostname-διεύθυνση IP

**local name server:**

- ○ κάθε ISP, εταιρεία έχει ένα *local (default) name server*
- ○ Σε μία DNS query ο host απευθύνεται πρώτα στον local name server

**authoritative name server (ενός host):**

- ○ έχει αποθηκευμένα τη διεύθυνση IP και το όνομα του host
- ○ μπορεί να μεταφράσει το όνομα αυτού του host στην IP διεύθυνσή του

# DNS: Root name servers

□ σ' αυτούς απευθύνονται οι τοπικοί (local) name servers που δεν μπορούν να μεταφράσουν ένα hostname



a Verisign, Dulles, VA
c Cogent, Herndon, VA (also Los Angeles)
d U Maryland College Park, MD
g US DoD Vienna, VA
h ARL Aberdeen, MD
j Verisign, ( 11 locations)

k RIPE London (also Amsterdam, Frankfurt)
l Autonomica, Stockholm (plus 3 other locations)
m WIDE Tokyo

e NASA Mt View, CA
f Internet Software C. Palo Alto, CA (and 17 other locations)

b USC-ISI Marina del Rey, CA
l ICANN Los Angeles, CA

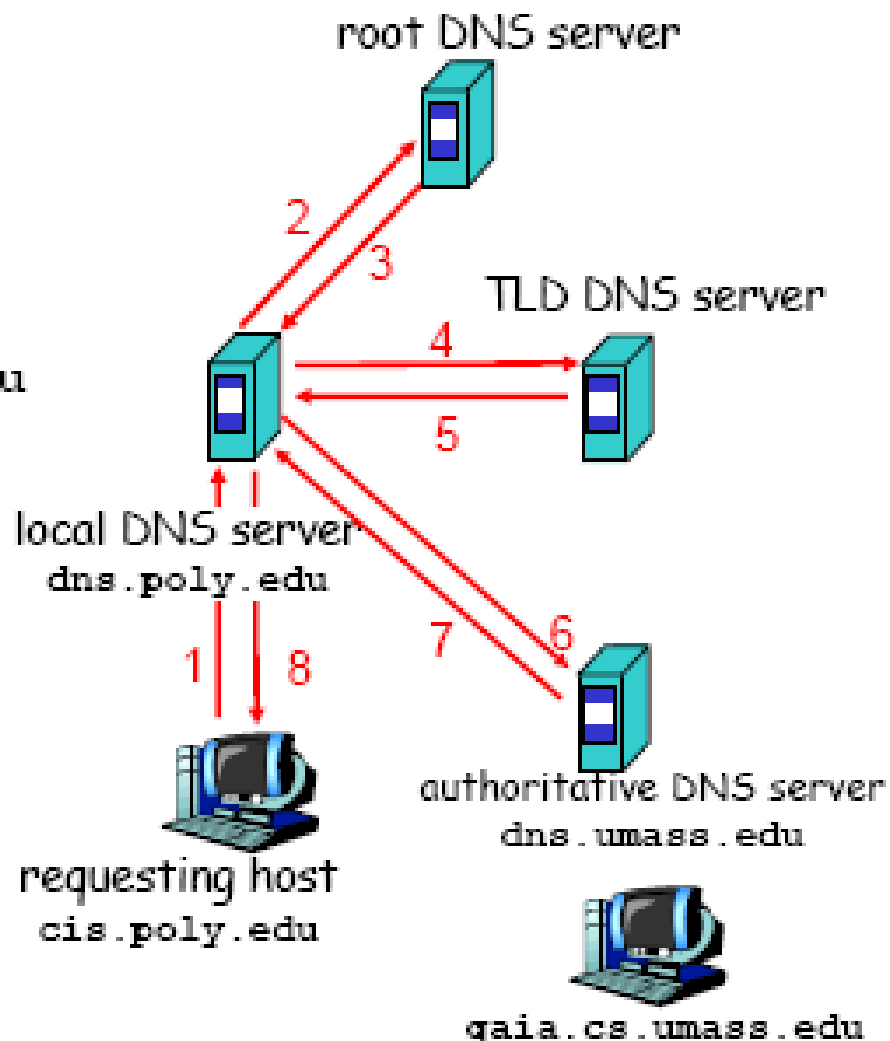13 root name servers παγκοσμίως

# TLD Servers και Authoritative Servers

□ **Top-level domain (TLD) servers:** υπεύθυνοι για τα domains com, org, net, edu, κλπ. καθώς και για όλα τα top-level domains των κρατών: π.χ. gr, uk, fr, ca, jp

- o Η εταιρεία Network solutions διαχειρίζεται τους servers του TLD com
- o Η εταιρεία Educause διαχειρίζεται τους servers του TLD edu

□ **Authoritative DNS servers:** οι DNS servers ενός οργανισμού, που παρέχουν τις αντιστοιχίες hostname σε διεύθυνση IP για τους servers του οργανισμού (π.χ. Web και mail servers).

- o Μπορεί να διατηρούνται από τον οργανισμό τον ίδιο ή από service provider

# Local Name Server

- Δεν ανήκει αυστηρά στην ιεραρχία των DNS servers

- Κάθε ISP (οικιακός ISP, εταιρεία, πανεπιστήμιο) έχει ένα τοπικό name server
  - Καλείται επίσης "default name server"

- Όταν ένας host κάνει ένα ερώτημα (query) DNS, το ερώτημα στέλνεται στον τοπικό του DNS server
  - Ο τοπικός name server προωθεί το ερώτημα στην ιεραρχία

# Παράδειγμα

□ Ο host `cis.poly.edu` θέλει τη διεύθυνση IP του `gaia.cs.umass.edu`

root DNS server

2

3

TLD DNS server

4

5

local DNS server
`dns.poly.edu`

1

8

7

6

requesting host
`cis.poly.edu`

authoritative DNS server
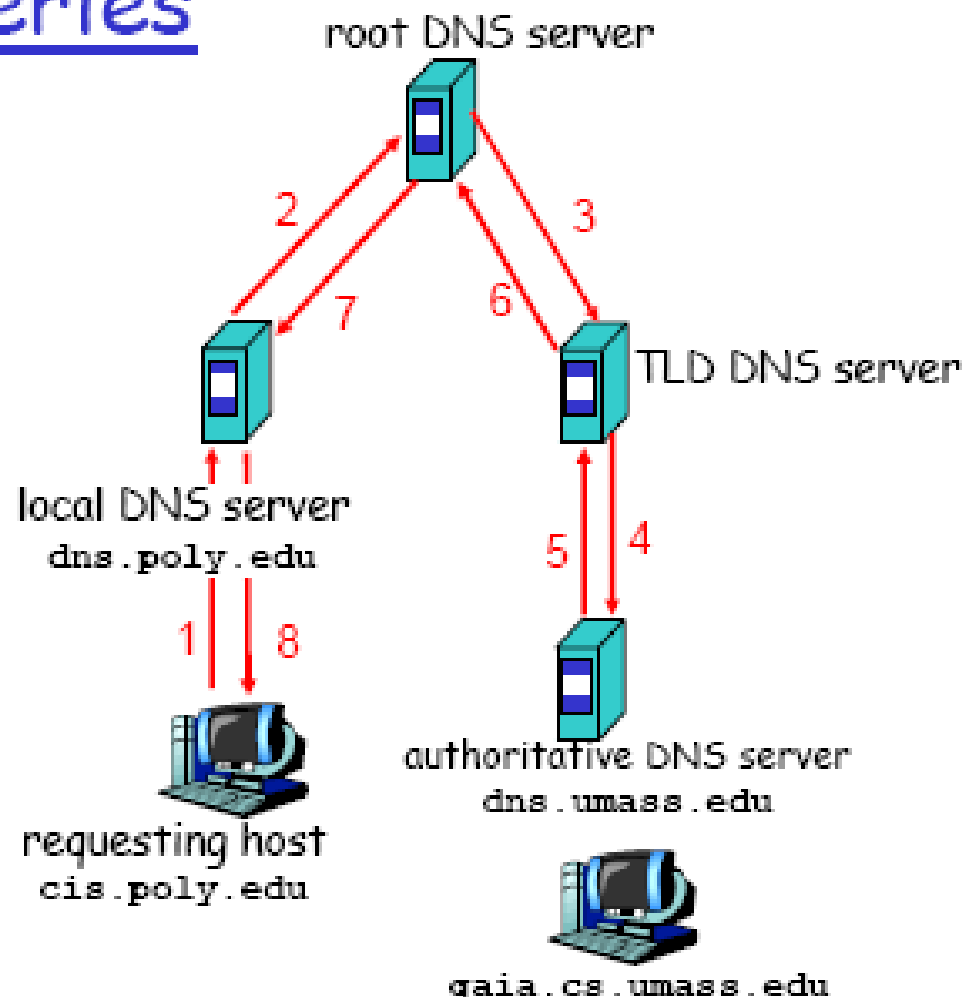`dns.umass.edu`

`gaia.cs.umass.edu`

# Recursive queries

## recursive query (αναδρομικό ερώτημα):

- αναθέτει τη μετάφραση του ονόματος στον name server στον οποίο απευθύνεται
- υψηλό φορτίο;

## iterated query (επαναληπτικό ερώτημα):

- ο ερωτηθείς name server παραπέμπει σε άλλο name server
- "Δεν γνωρίζω αυτό το όνομα, ρώτα αυτόν τον server"



root DNS server

2  7

local DNS server
dns.poly.edu

1  8

requesting host
cis.poly.edu

3  6

TLD DNS server

5  4

authoritative DNS server
dns.umass.edu

gaia.cs.umass.edu

# DNS: caching και ενημέρωση records

Για τη βελτίωση των καθυστερήσεων και του αριθμού DNS queries στο δίκτυο:

☐ όταν (οποιοσδήποτε) name server λάβει ένα record, καταχωρεί αντίγραφο σε μία *cache*

- ○ η καταχώρηση της cache εξαφανίζεται (timeout) μετά από κάποιο χρόνο
- ○ καταχωρήσεις για τους TLD servers υπάρχουν συνήθως στην cache του τοπικού name server (οι root name servers δεν δέχονται συχνές επισκέψεις)

Η εισαγωγή δεδομένων στη βάση γινόταν μέχρι πρόσφατα στατικά από το διαχειριστή

☐ Υπό σχεδίαση (IETF) βρίσκονται δυναμικοί μηχανισμοί ενημέρωσης που χρησιμοποιούν μηνύματα DNS:

- ○ RFC 2136
- ○ http://www.ietf.org/html.charters/dnsind-charter.html

# DNS records

<u>DNS:</u> κατανεμημένη βάση δεδομένων αποθηκεύει <span style="color:red">resource records (RR)</span>

> **RR format**: `(name, value, type, ttl)`

- **Type=A**
  - `name`: hostname
  - `value`: IP address

- **Type=NS**
  - `name`: domain (π.χ. `tuc.gr`)
  - `value`: διεύθυνση IP του authoritative name server για αυτό το domain (π.χ. `147.27.18.1` – `dedalus.noc.tuc.gr`

- **Type=CNAME**
  - `name`: ψευδώνυμο (alias name) για κάποιο "canonical" name (πραγματικό όνομα)
  
    `www.ibm.com` στην πραγματικότητα `servereast.backup2.ibm.com`
  - `value`: canonical name
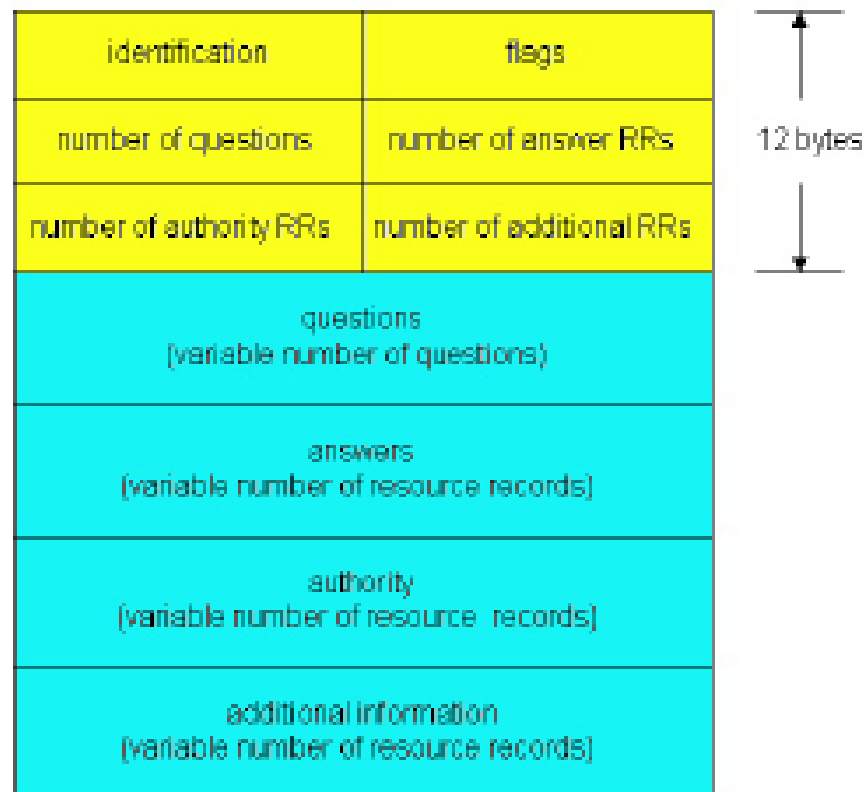
- **Type=MX**
  - `name`: ψευδώνυμο mail server
  - `value`: "canonical" name mail server
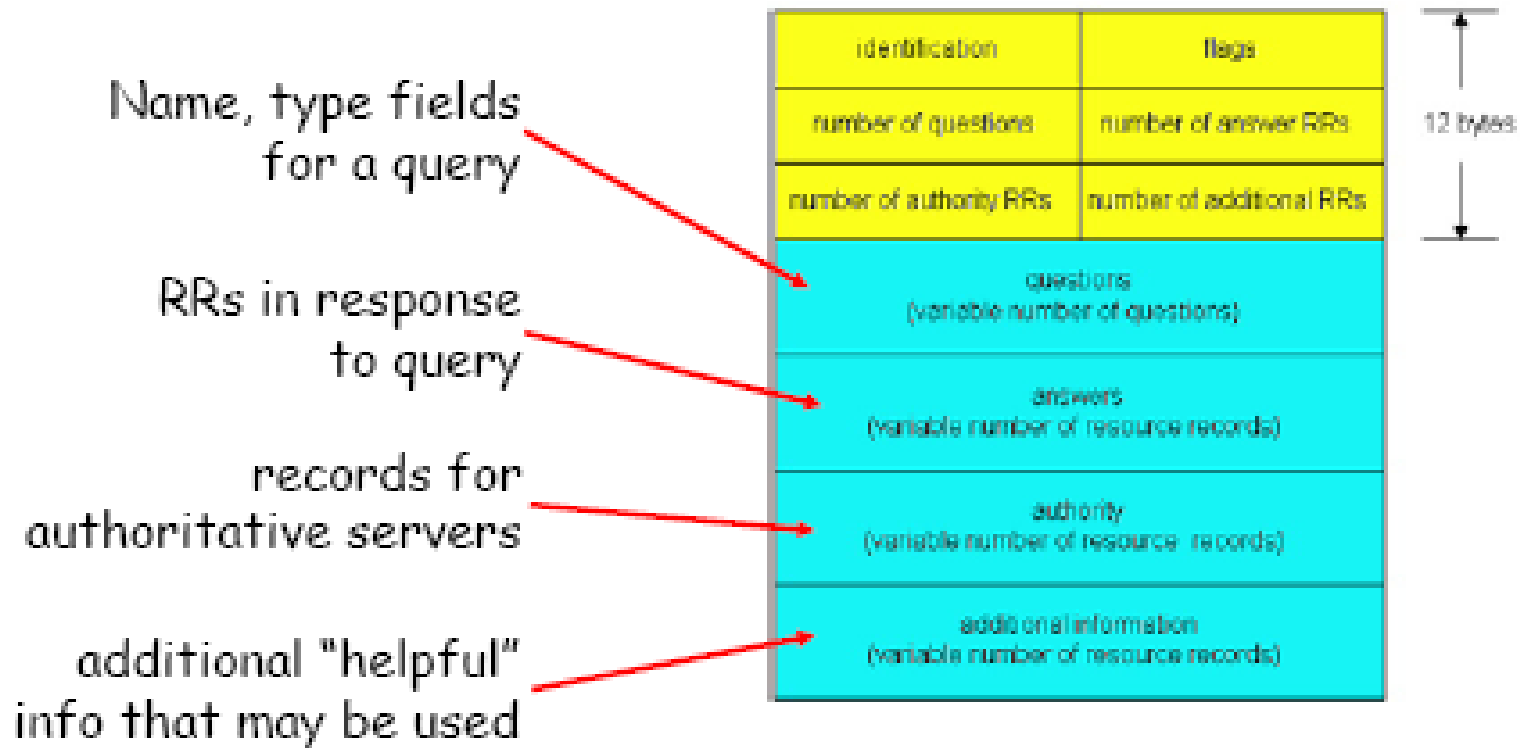
# Πρωτόκολλο DNS, μηνύματα

**Πρωτόκολλο DNS:** μηνύματα *query* and *reply* έχουν το ίδιο *format*

**επικεφαλίδα μηνύματος**

□ identification: αριθμός από 16 bits που προσδιορίζει το ερώτημα (query), η απόκριση σε ένα ερώτημα χρησιμοποιεί τον ίδιο αριθμό

□ flags:

    ○ query ή reply

    ○ recursion desired

    ○ recursion available

    ○ reply is authoritative

| identification | flags |
|---|---|
| number of questions | number of answer RRs |
| number of authority RRs | number of additional RRs |

12 bytes

| questions<br>(variable number of questions) |
|---|
| answers<br>(variable number of resource records) |
| authority<br>(variable number of resource records) |
| additional information<br>(variable number of resource records) |

# Πρωτόκολλο DNS, μηνύματα

Name, type fields
for a query

RRs in response
to query

records for
authoritative servers

additional "helpful"
info that may be used

| identification | flags |
|---|---|
| number of questions | number of answer RRs |
| number of authority RRs | number of additional RRs |

12 bytes

questions
(variable number of questions)

answers
(variable number of resource records)

authority
(variable number of resource records)

additional information
(variable number of resource records)

# Link layer, LANs: outline

5.1 introduction, services

5.2 error detection, correction

5.3 multiple access protocols

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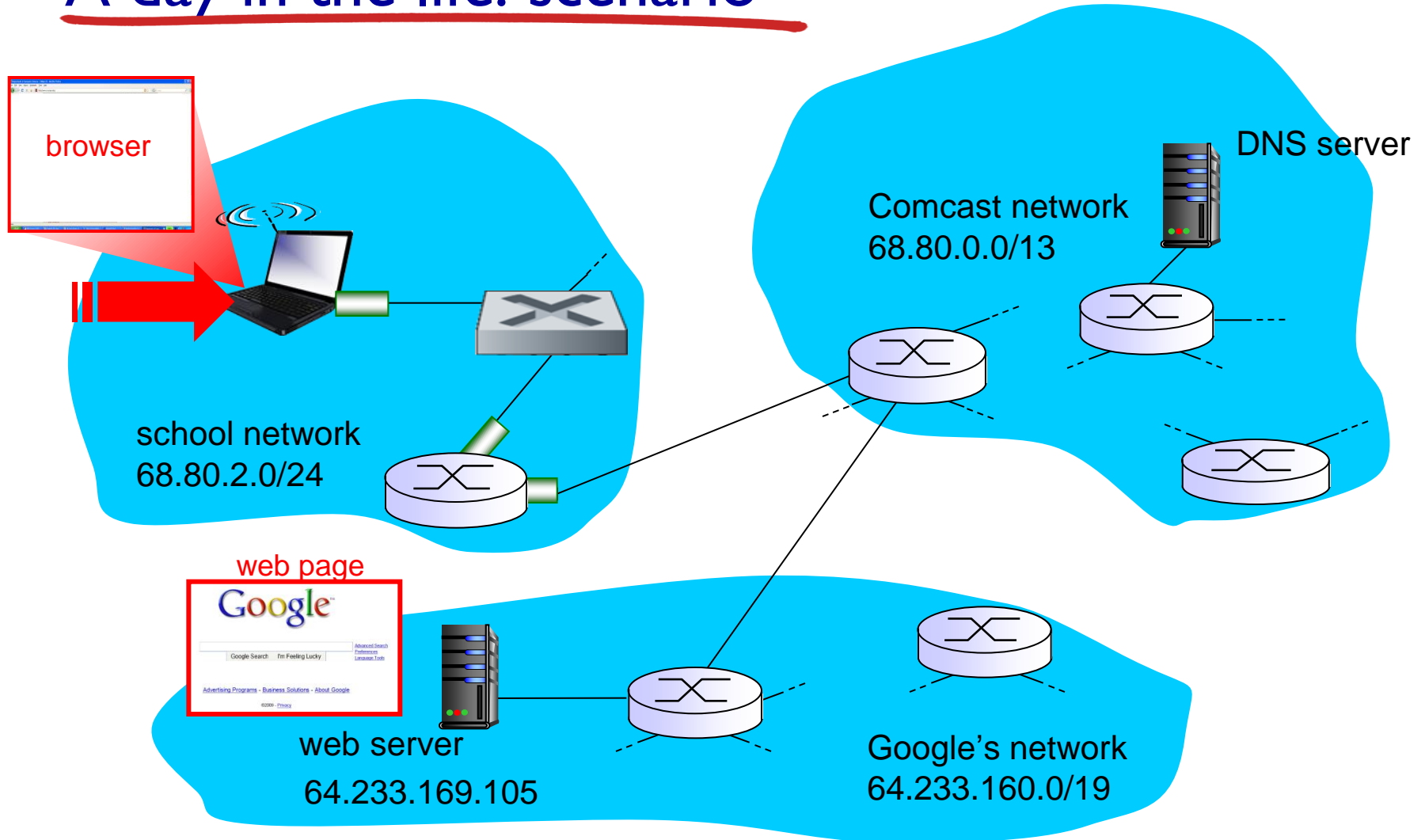- addressing, ARP
- Ethernet
- switches
- VLANS

5.5 Subnetting

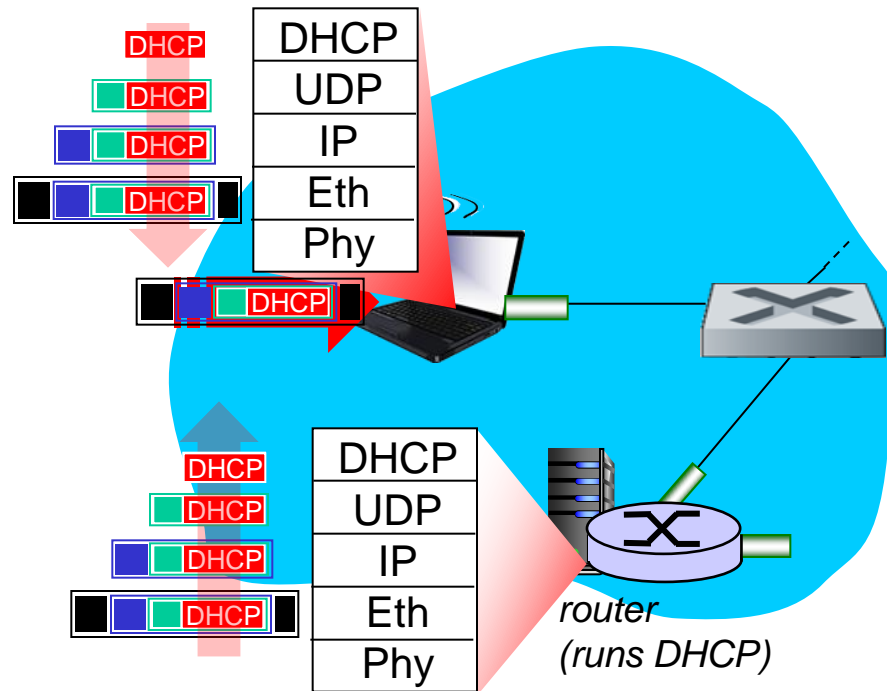5.6 a day in the life of a web request

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

❖ journey down protocol stack complete!
- application, transport, network, link

❖ putting-it-all-together: synthesis!
- *goal:* identify, review, understand protocols (at all layers) involved in seemingly simple scenario: requesting www page
- *scenario:* student attaches laptop to campus network, requests/receives www.google.com
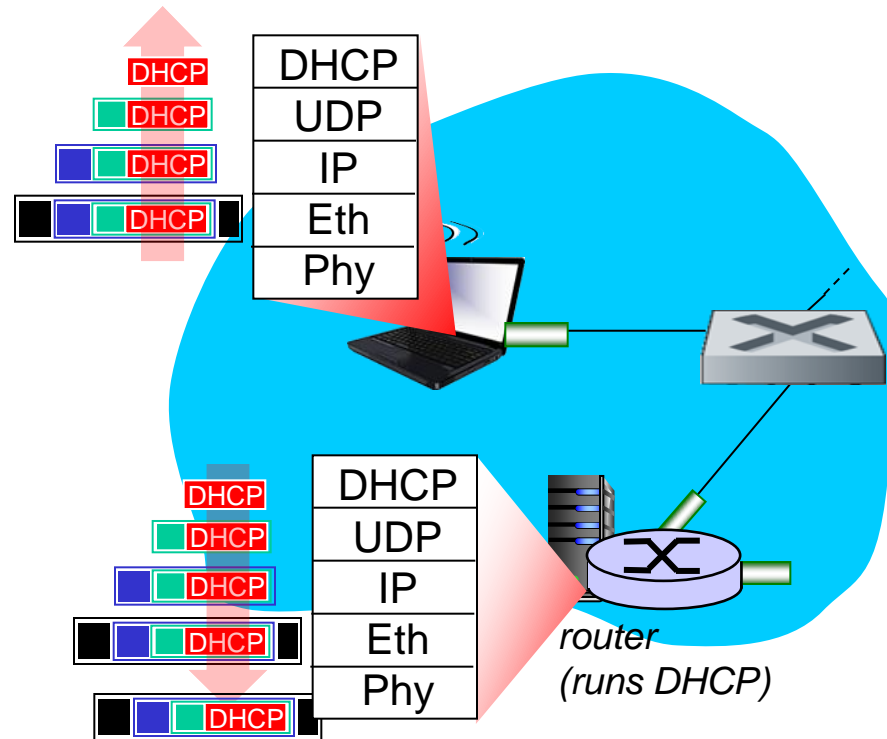
# A day in the life: scenario



browser

DNS server

Comcast network
68.80.0.0/13

school network
68.80.2.0/24

web page

Google

web server
64.233.169.105

Google's network
64.233.160.0/19

# A day in the life… connecting to the Internet



*router (runs DHCP)*

- ❖ 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

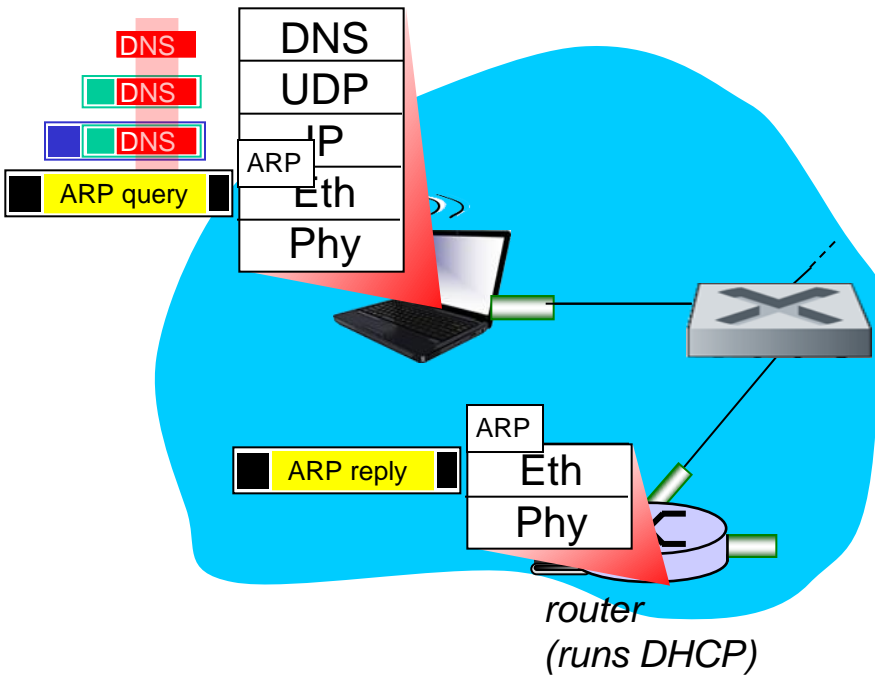- ❖ Ethernet *demuxed* to IP demuxed, UDP demuxed to DHCP

# 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
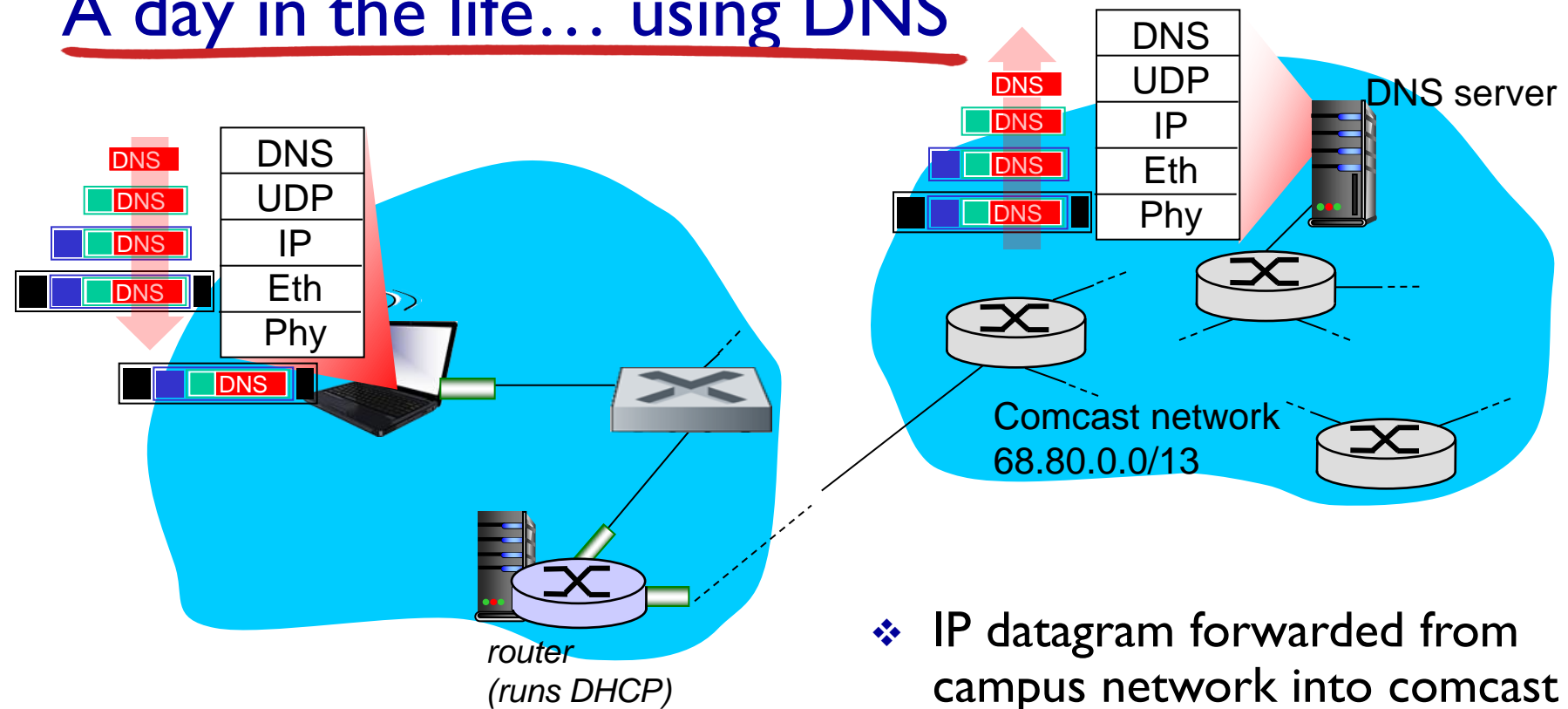
❖ DHCP client receives DHCP ACK reply

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

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



DNS
UDP
IP
ARP
Eth
Phy

ARP query

ARP reply
ARP
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

# A day in the life… using DNS



router
(runs DHCP)

Comcast network
68.80.0.0/13
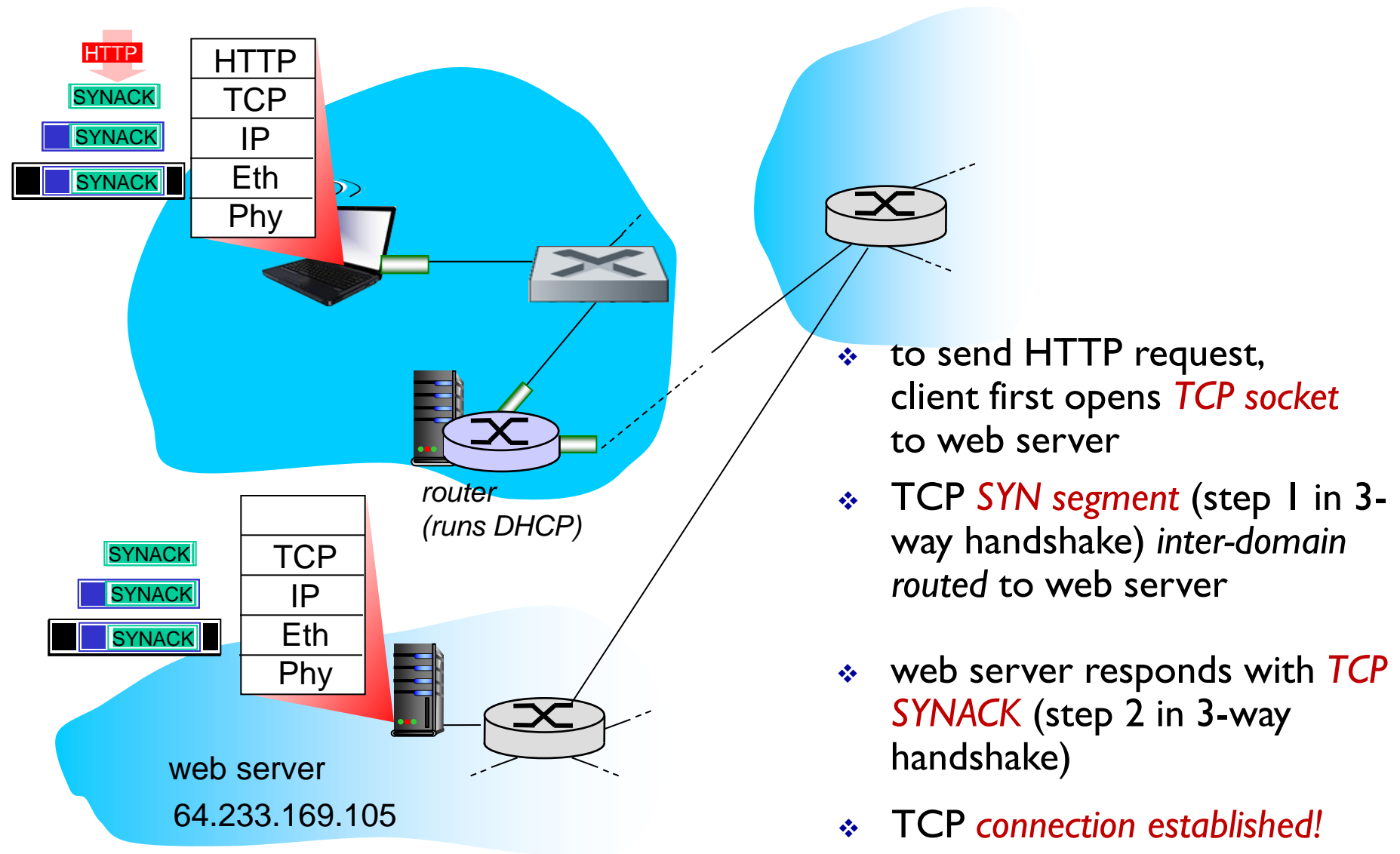
DNS server

❖ 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

# 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 3-way handshake) *inter-domain routed* to web server

- ❖ web server responds with *TCP SYNACK* (step 2 in 3-way handshake)

- ❖ TCP *connection established!*

# A day in the life… HTTP request/reply



**Google**

❖ web page *finally (!!!)* displayed

*router*
*(runs DHCP)*

web server
64.233.169.105

❖ *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