# Δίκτυα Καθοριζόμενα από Λογισμικό

**Ενότητα 1:** 1.2 Επιστροφή στα βασικά: λίγα απαραίτητα

Ξενοφώντας Δημητρόπουλος

Τμήμα Επιστήμης Υπολογιστών

# Back to the Basics: Some Necessary 101

Xenofontas Dimitropoulos
29/9/2014

Credits: Kurose, James F., and Keith W. Ross. Computer Networking: A top-down approach featuring the Internet.

# Agenda

❐ Link Layer
  ❍ Link-Layer Addressing and ARP
  ❍ Switches
  ❍ VLANs
❐ Network Layer
  ❍ IPv4 Addressing
  ❍ Routing algorithms
  ❍ Routing in the Internet
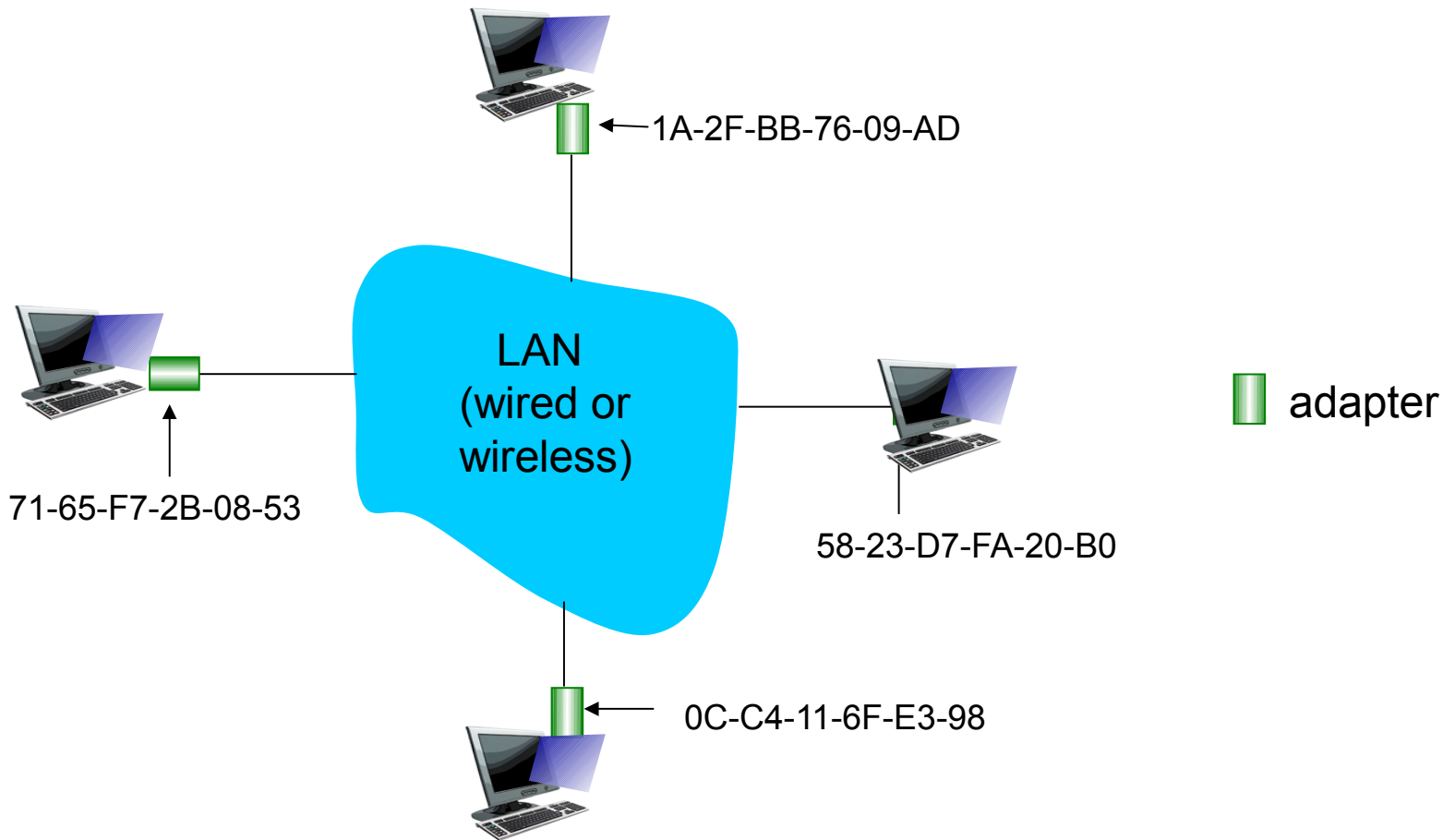❐ A Day in the Life of a Web Request

# 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
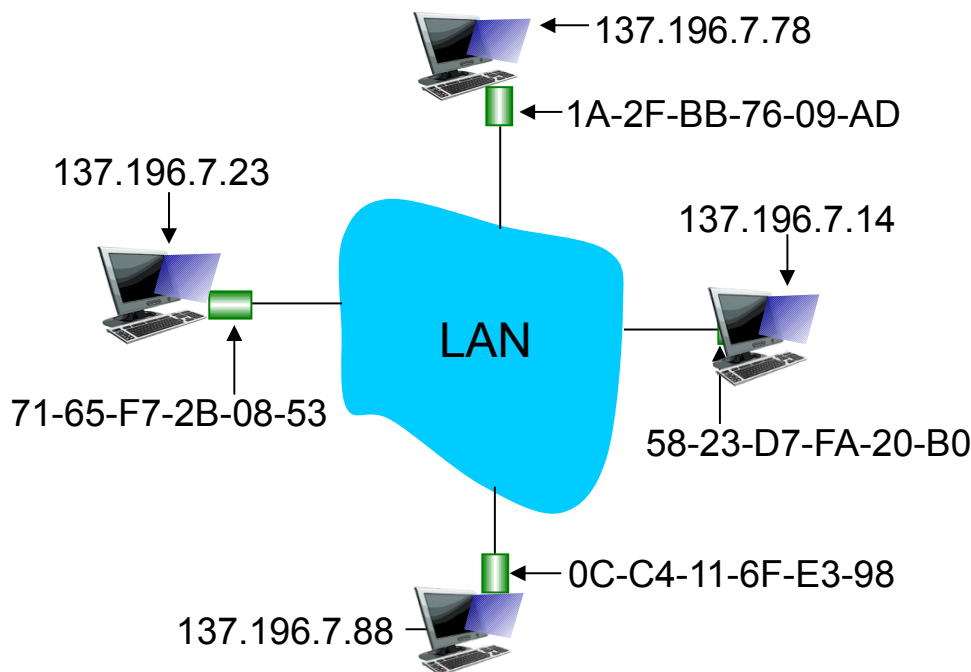wireless)

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>

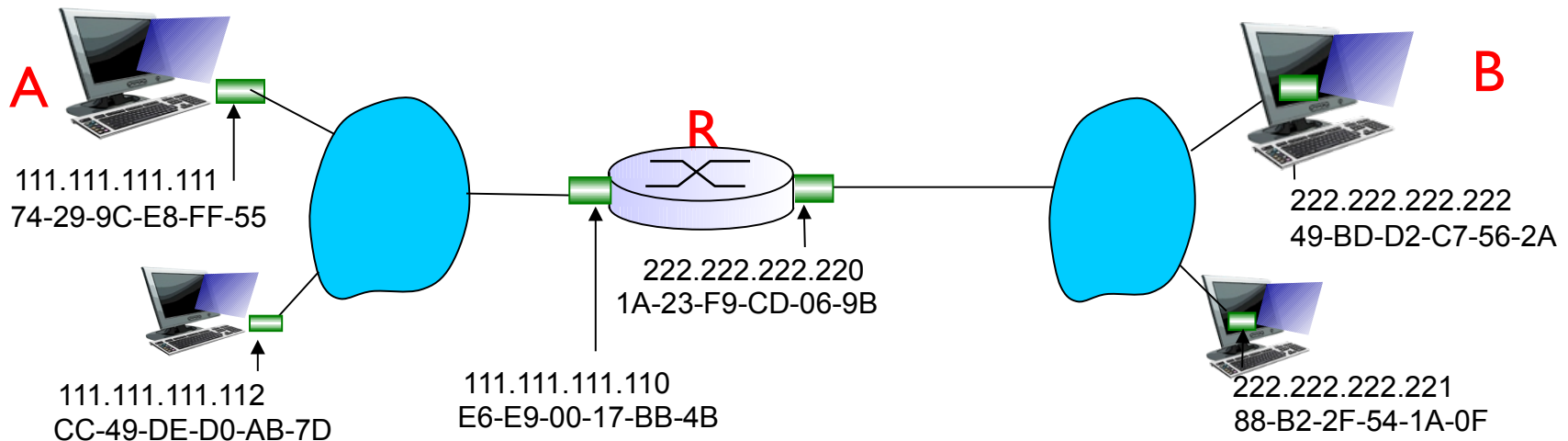- TTL (Time To Live): time after which address mapping will be forgotten (typically 20 min)

137.196.7.78

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

137.196.7.23

137.196.7.14

LAN

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

58-23-D7-FA-20-B0

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

137.196.7.88

# ARP 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*
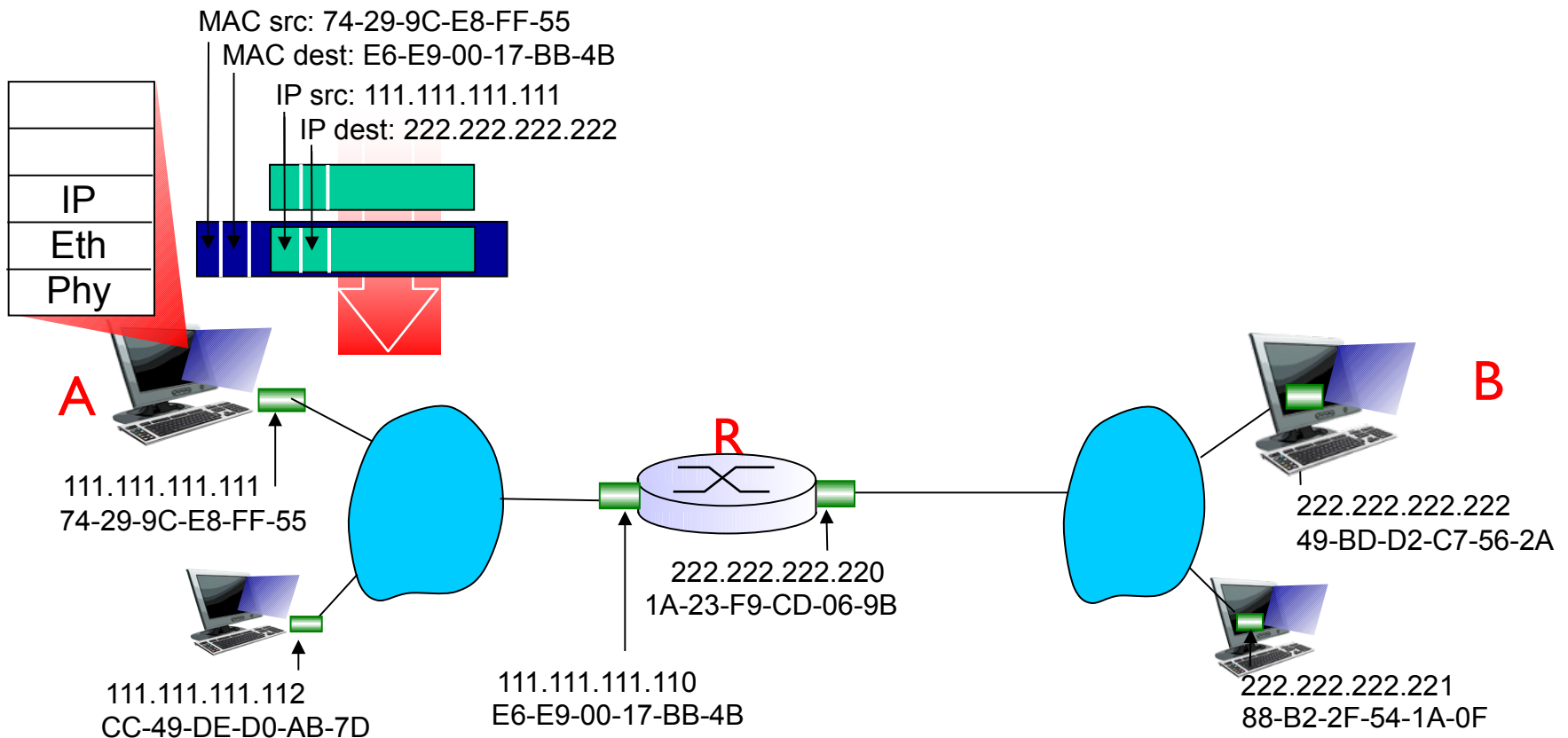
# Addressing: routing to another LAN

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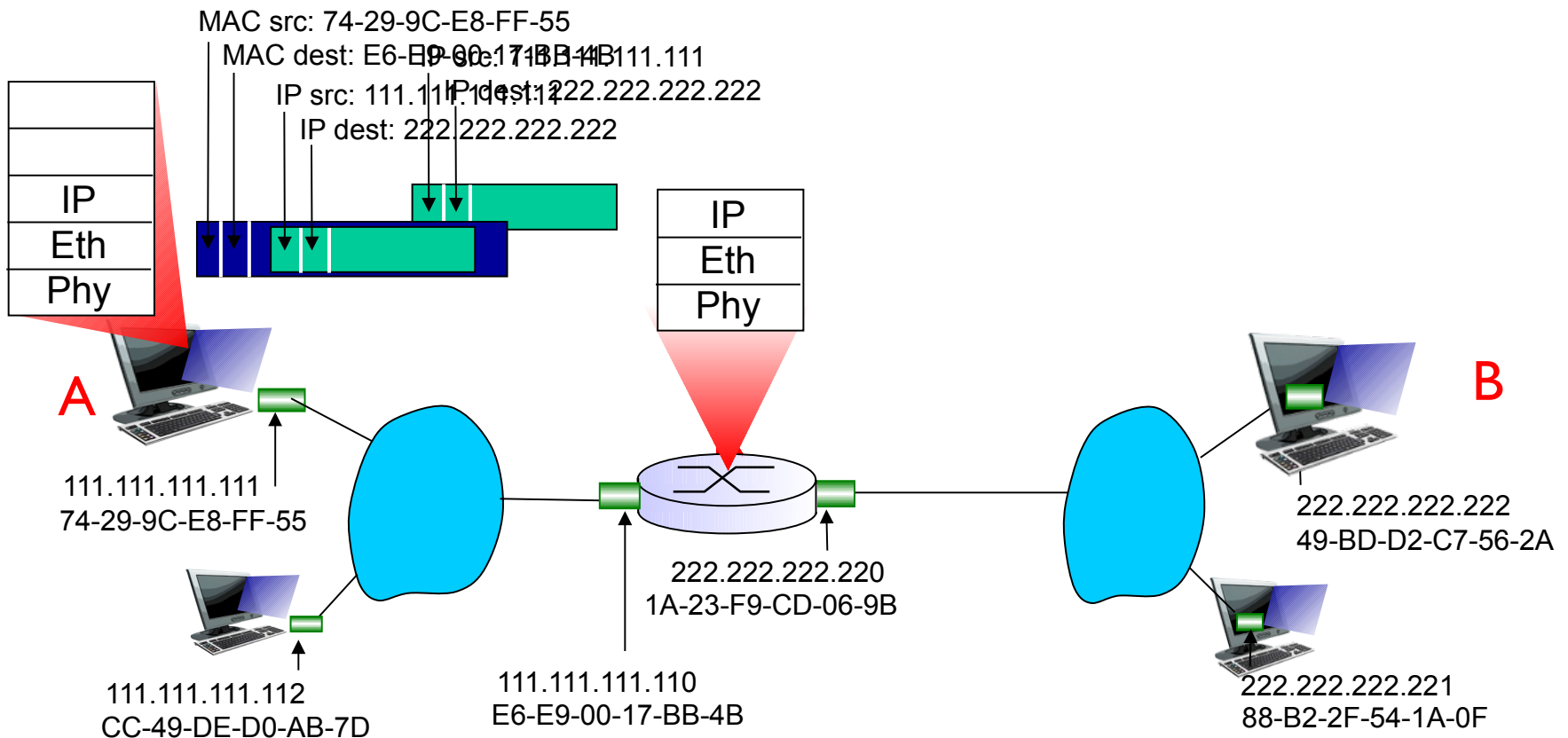- assume A knows R's MAC address (how?)

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

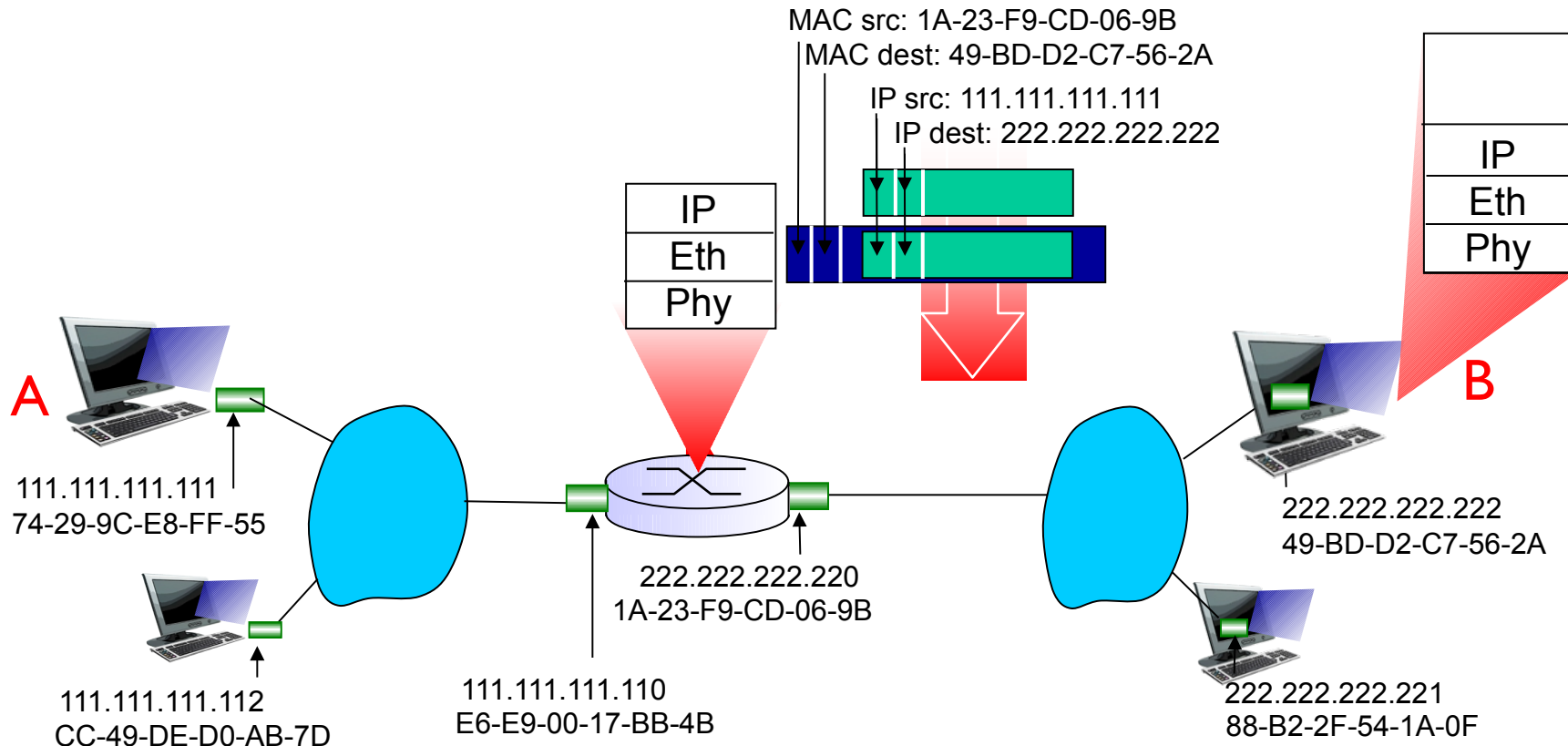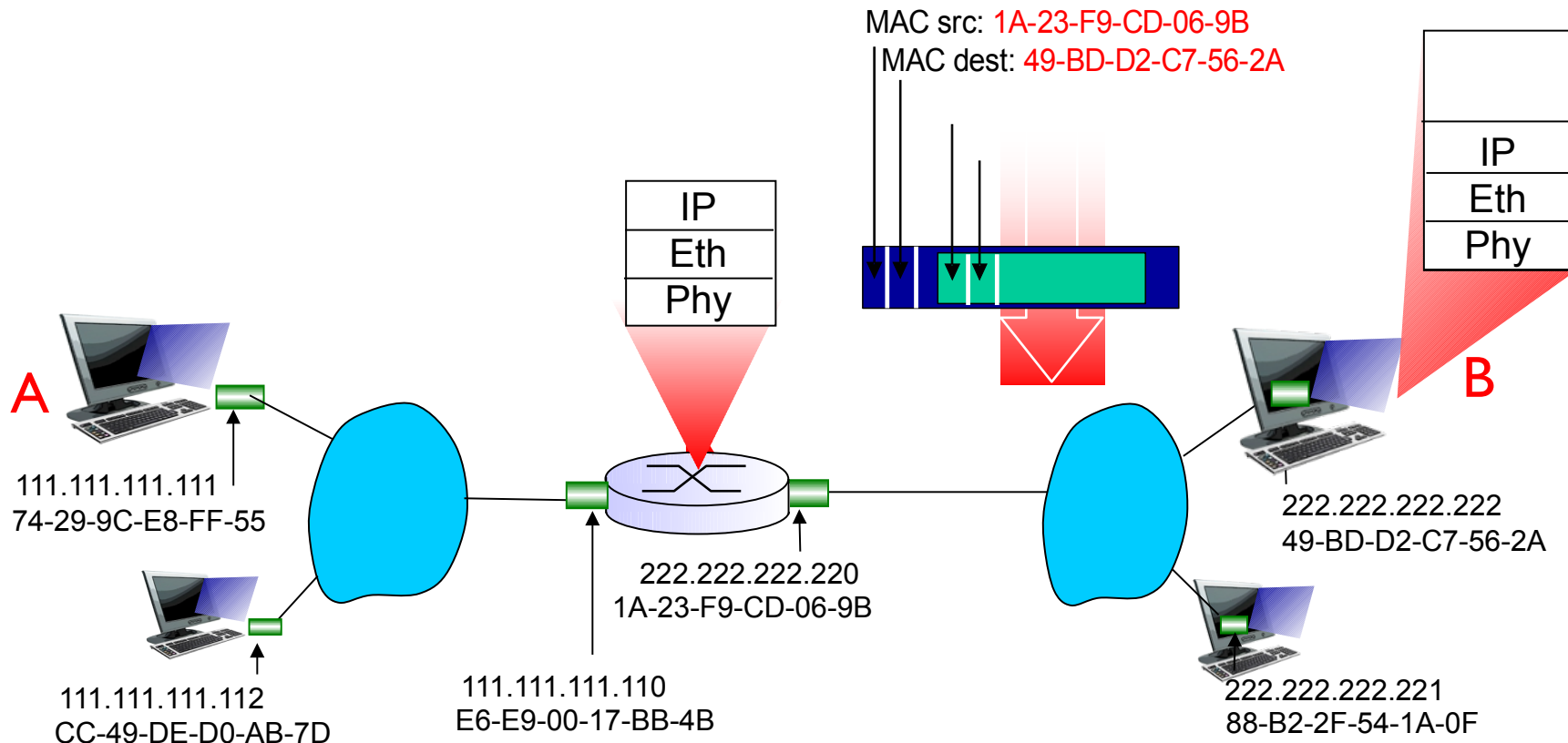❖ A creates link-layer frame with R's MAC address as dest, frame contains A-to-B IP datagram

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

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

❖ frame sent from A to R

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

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

IP src: 111.111.111.111
IP dest: 222.222.222.222

IP src: 111.111.111.111
IP dest: 222.222.222.222

IP dest: 222.222.222.222

| IP  |
|-----|
| Eth |
| Phy |

| IP  |
|-----|
| Eth |
| Phy |

A

B

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

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

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

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

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

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

# Addressing: routing to another LAN

❖ R forwards datagram with IP source A, destination B

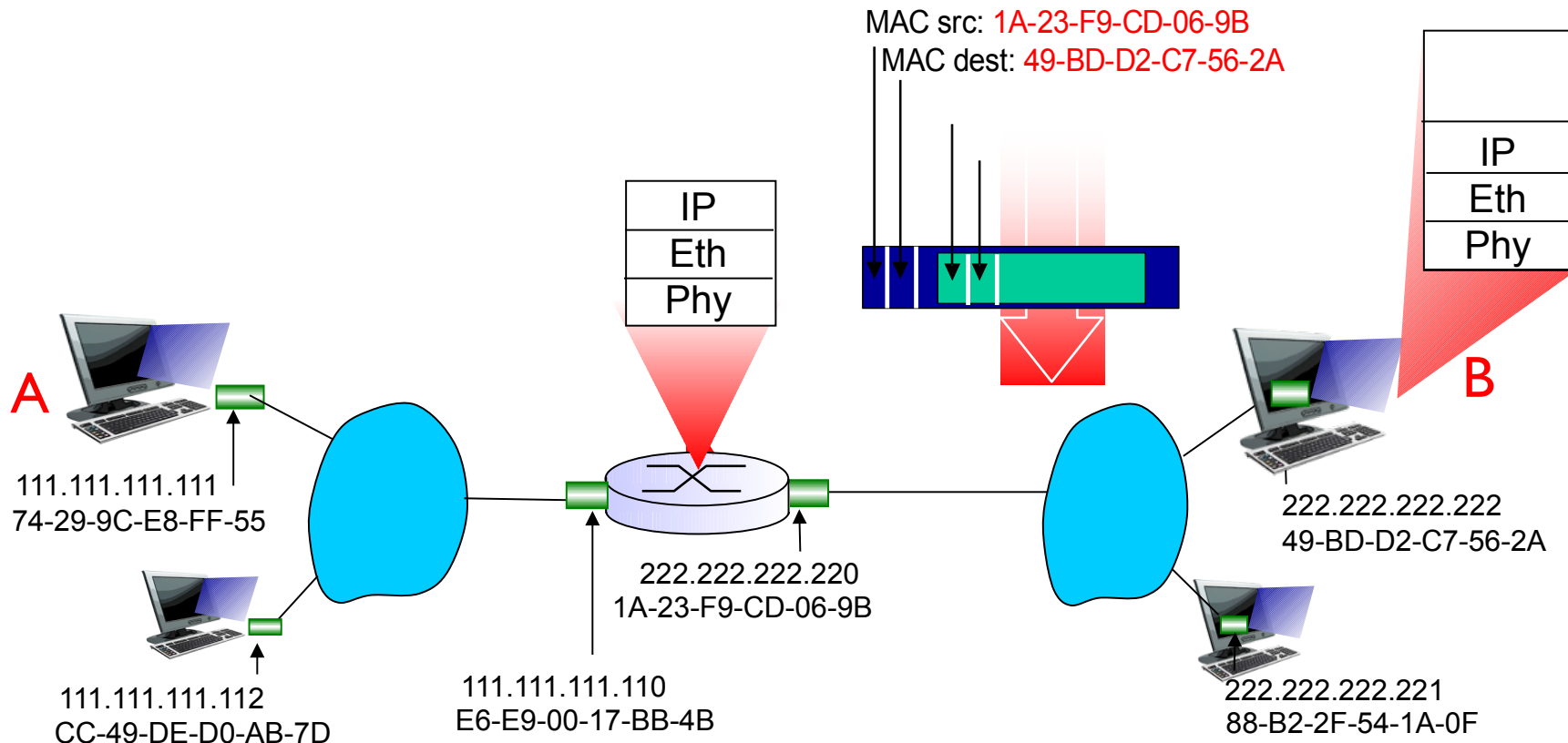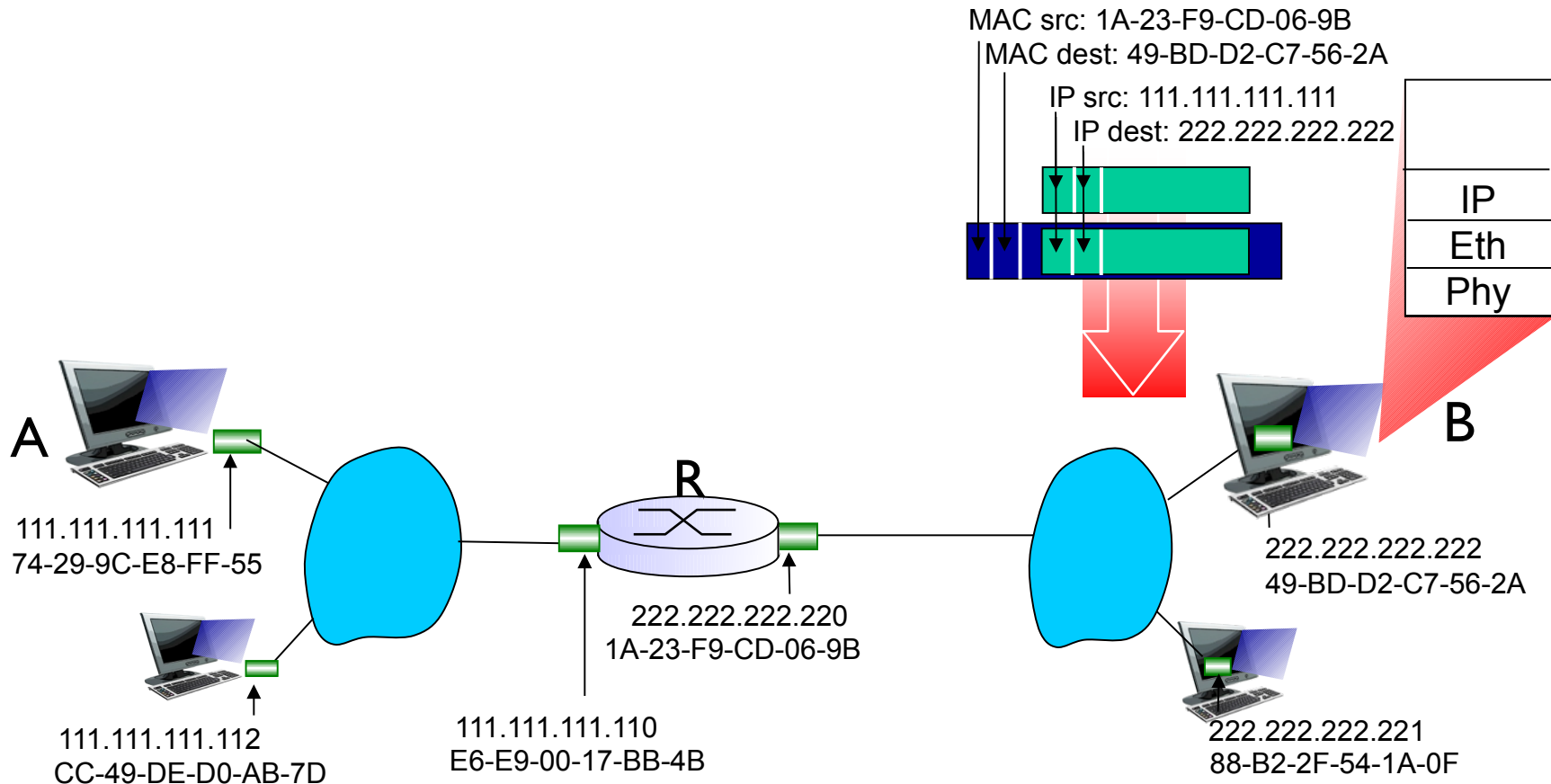❖ R creates link-layer frame with B's MAC address as dest, frame contains A-to-B IP datagram

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

IP src: 111.111.111.111
IP dest: 222.222.222.222

| IP |
|----|
| Eth |
| Phy |

| IP |
|----|
| Eth |
| Phy |

A

B

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

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

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

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

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

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

# Addressing: routing to another LAN

❖ R forwards datagram with IP source A, destination B

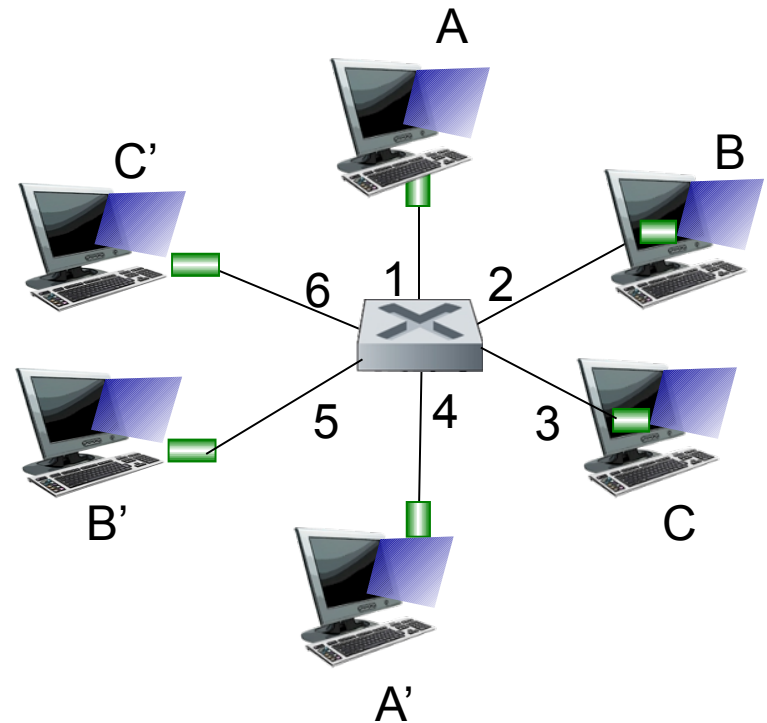❖ R creates link-layer frame with B's MAC address as dest, frame contains A-to-B IP datagram

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

| IP |
|-----|
| Eth |
| Phy |

| IP |
|-----|
| Eth |
| Phy |

A

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

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

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

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

B

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

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

# 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
Eth
Phy

IP
Eth
Phy

A

B

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

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

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

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

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

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

# 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

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

# Agenda

❖ Link Layer
- Link-Layer Addressing and ARP
- Switches
- VLANs

❖ Network Layer
- IPv4 Addressing
- Routing algorithms
- Routing in the Internet

❖ 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

❖ 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 S1 know to forward frame destined to F via S4 and S3?

❖ _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 S1, S2, S3, S4

# 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
frame
physical

**switch**

network
datagram
link
frame
physical

application
transport
network
link
physical

# Hubs

Hubs are essentially physical-layer repeaters:

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

twisted pair

hub

# Agenda

- ❖ Link Layer
    - ▪ Link-Layer Addressing and ARP
    - ▪ Switches
    - ▪ VLANs
- ❖ Network Layer
    - ▪ IPv4 Addressing
    - ▪ Routing algorithms
    - ▪ Routing in the Internet
- ❖ A Day in the Life of a Web Request

# 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

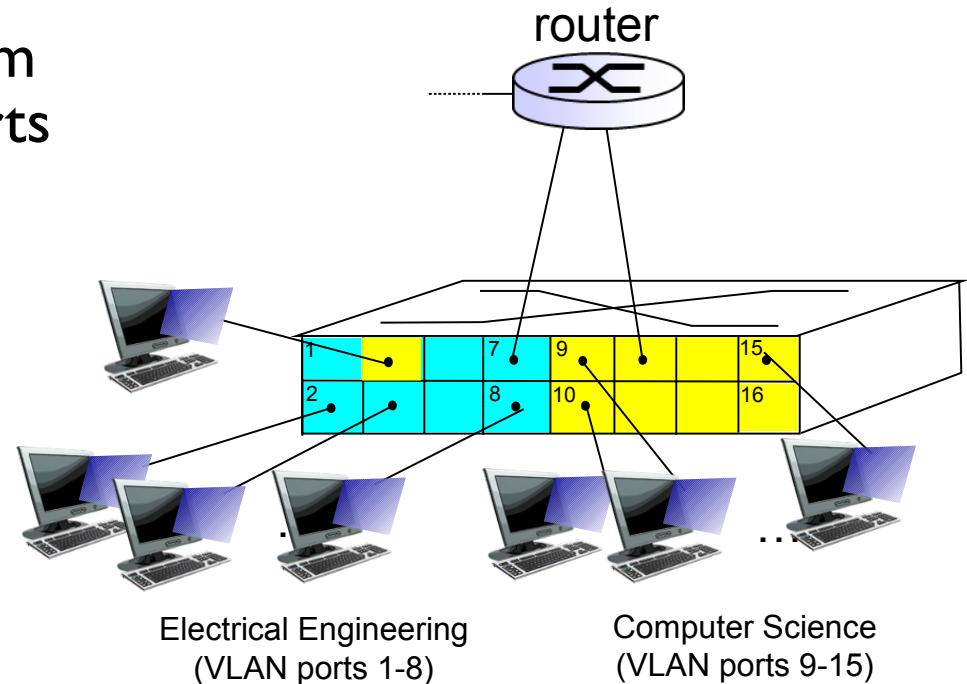**port-based VLAN:** switch ports grouped (by switch management software) so that *single* physical switch ......

**Virtual Local Area Network**

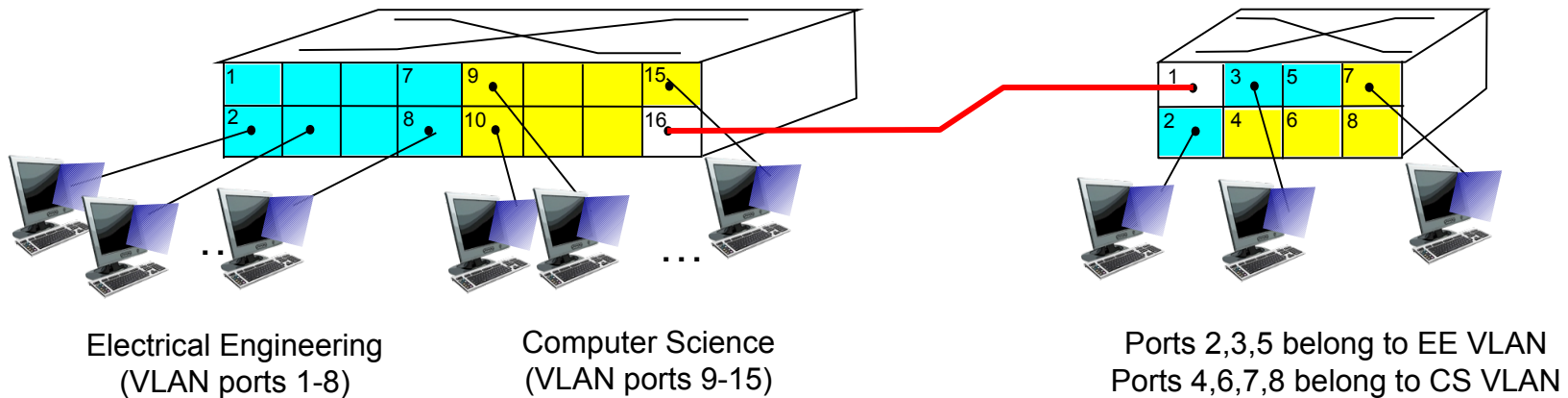switch(es) supporting VLAN capabilities can be configured to define multiple **_virtual_** LANS over single physical LAN infrastructure.



Electrical Engineering
(VLAN ports 1-8)

Computer Science
(VLAN ports 9-15)

… operates as *multiple* virtual switches



Electrical Engineering
(VLAN ports 1-8)

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)

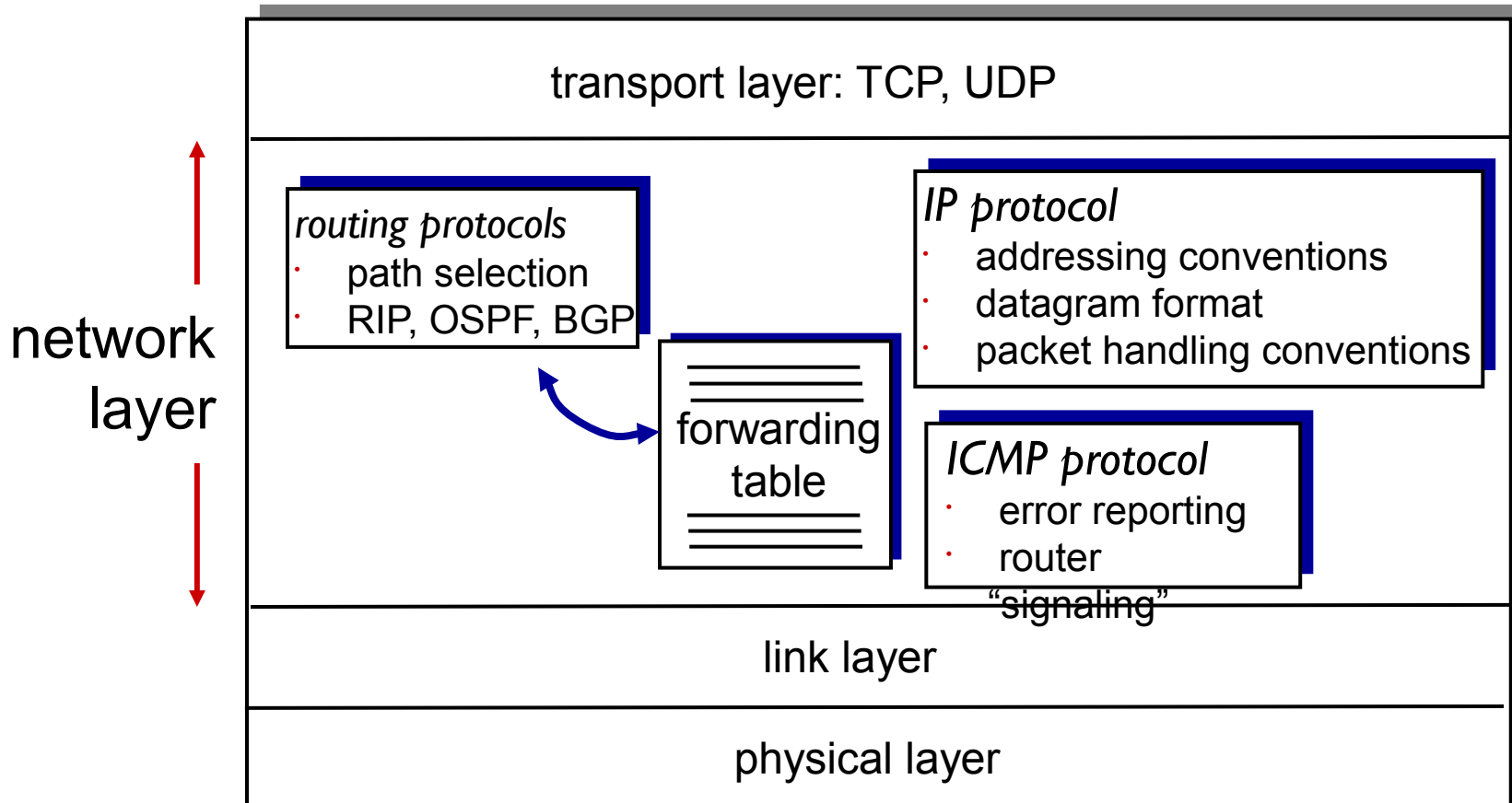  ▪ in practice vendors sell combined switches plus routers



router

| 1 | | 7 | 9 | | | 15 |
| 2 | | 8 | 10 | | | 16 |

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

# Agenda

❖ Link Layer
  ▪ Link-Layer Addressing and ARP
  ▪ Switches
  ▪ VLANs
❖ Network Layer
  ▪ IPv4 Addressing
  ▪ Routing algorithms
  ▪ Routing in the Internet
❖ A Day in the Life of a Web Request

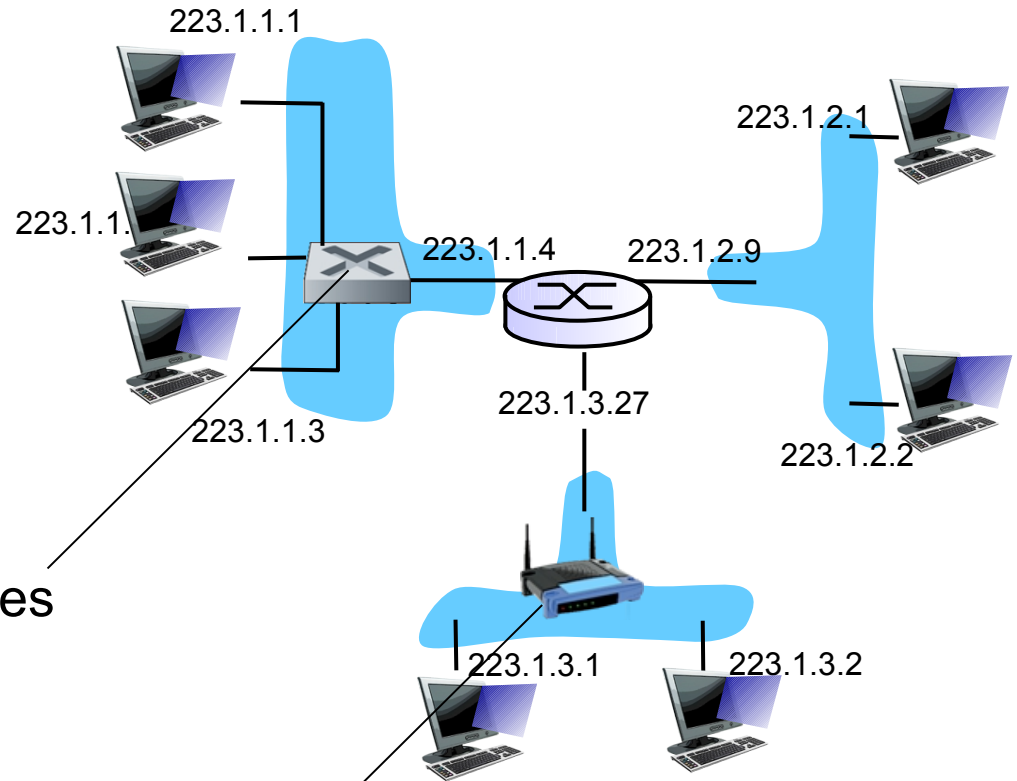# The Internet network layer

host, router network layer functions:



Inside the diagram:

transport layer: TCP, UDP

network layer

*routing protocols*
· path selection
· RIP, OSPF, BGP

forwarding table

*IP protocol*
· addressing conventions
· datagram format
· packet handling conventions

*ICMP protocol*
· error reporting
· router "signaling"

link layer

physical layer

# IP addressing: introduction

*Q: how are interfaces actually connected?*

223.1.1.1

223.1.2.1

223.1.1.

223.1.1.4    223.1.2.9

223.1.1.3

223.1.3.27

223.1.2.2

*A:* wired Ethernet interfaces connected by Ethernet switches

223.1.3.1    223.1.3.2

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

*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

223.1.2.1

223.1.2.2

223.1.1.3    223.1.3.27

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

223.1.1.3    223.1.3.27

223.1.2.2

subnet

223.1.3.1    223.1.3.2

*223.1.3.0/24*

subnet mask: /24

# NAT: network address translation



rest of Internet

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

10.0.0.1

10.0.0.2

10.0.0.3

10.0.0.4

138.76.29.7

*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)

# 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



Network Layer        4-41

# 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

# ICMP: internet control message protocol

- used by hosts & routers to communicate network-level information
  - error reporting: unreachable host, network, port, protocol
  - echo request/reply (used by ping)
- network-layer "above" IP:
  - ICMP msgs carried in IP datagrams
- ICMP message: type, code plus first 8 bytes of IP datagram causing error

| Type | Code | description |
|------|------|-------------|
| 0 | 0 | echo reply (ping) |
| 3 | 0 | dest. network unreachable |
| 3 | 1 | dest host unreachable |
| 3 | 2 | dest protocol unreachable |
| 3 | 3 | dest port unreachable |
| 3 | 6 | dest network unknown |
| 3 | 7 | dest host unknown |
| 4 | 0 | source quench (congestion control - not used) |
| 8 | 0 | echo request (ping) |
| 9 | 0 | route advertisement |
| 10 | 0 | router discovery |
| 11 | 0 | TTL expired |
| 12 | 0 | bad IP header |

# Traceroute and ICMP

- source sends series of UDP segments to dest
  - first set has TTL =1
  - second set has TTL=2, etc.
  - unlikely port number
- when *n*th set of datagrams arrives to nth router:
  - router discards datagrams
  - and sends source ICMP messages (type 11, code 0)
  - ICMP messages includes name of router & IP address

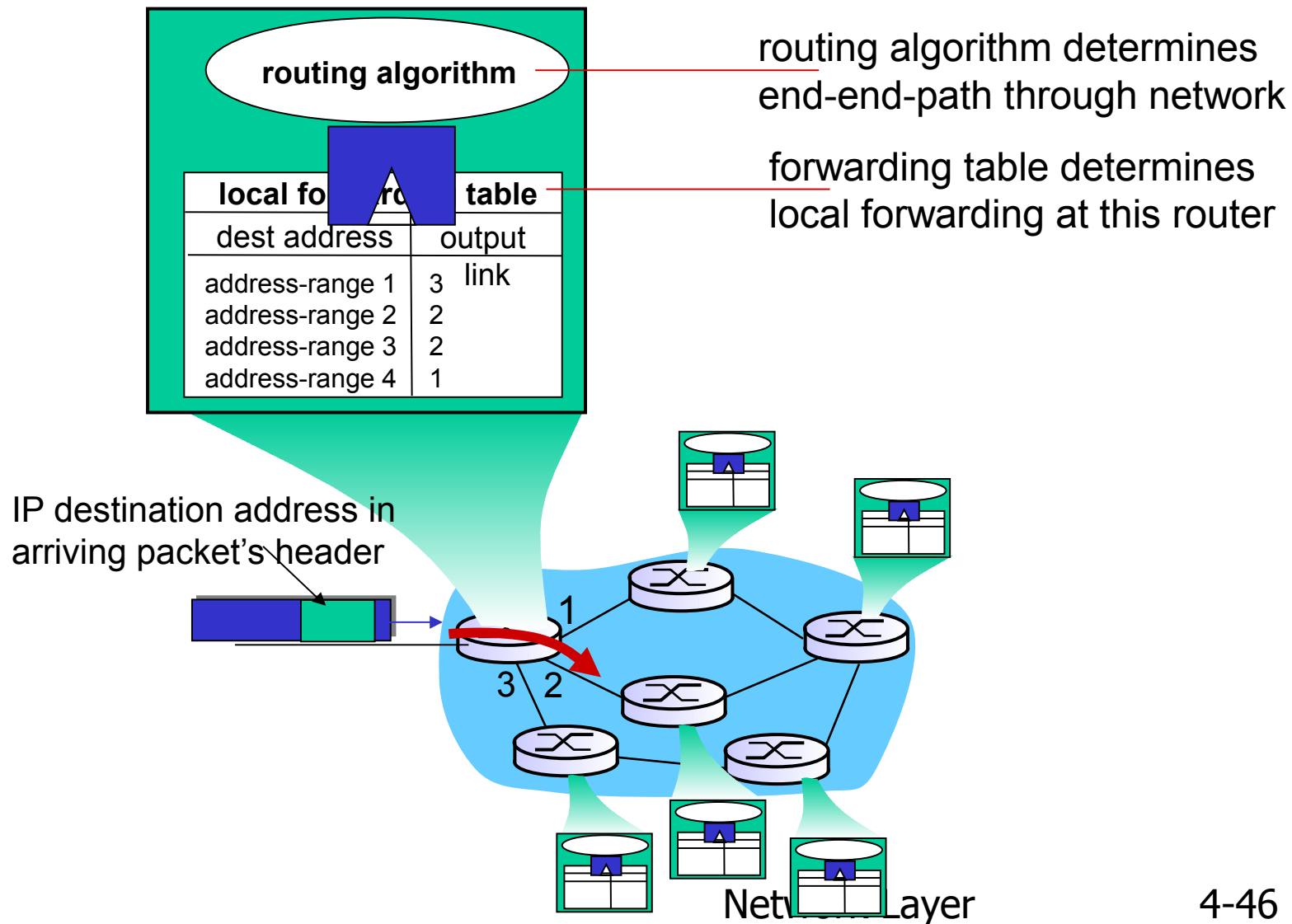- when ICMP messages arrives, source records RTTs

*stopping criteria:*
- ❖ UDP segment eventually arrives at destination host
- ❖ destination returns ICMP "port unreachable" message (type 3, code 3)
- ❖ source stops

3 probes    3 probes

3 probes

# Agenda

❐ Link Layer
  ○ Link-Layer Addressing and ARP
  ○ Switches
  ○ VLANs
❐ Network Layer
  ○ IPv4 Addressing & ICMP
  ○ Routing algorithms
  ○ Routing in the Internet
❐ A Day in the Life of a Web Request

# Interplay between routing, forwarding



routing algorithm
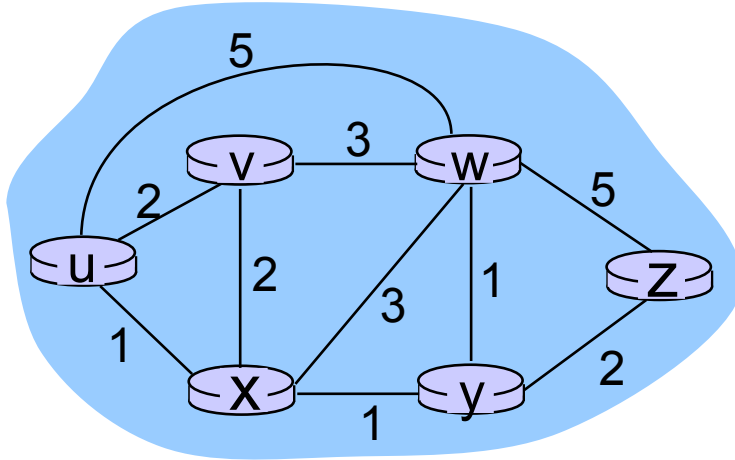
**local forwarding table**

| dest address | output link |
|---|---|
| address-range 1 | 3 |
| address-range 2 | 2 |
| address-range 3 | 2 |
| address-range 4 | 1 |

routing algorithm determines end-end-path through network

forwarding table determines local forwarding at this router

IP destination address in arriving packet's header

1
3  2

# Graph abstraction



graph: G = (N,E)

N = set of routers = { u, v, w, x, y, z }

E = set of links ={ (u,v), (u,x), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z) }

# Graph abstraction: costs



c(x,x') = cost of link (x,x')
     e.g., c(w,z) = 5

cost could always be 1, or inversely related to bandwidth, or inversely related to congestion

cost of path (x1, x2, x3,…, xp) = c(x1,x2) + c(x2,x3) + … + c(xp-1,xp)

*key question:* what is the least-cost path between u and z ?
*routing algorithm:* algorithm that finds that least cost path

# Hierarchical routing

our routing study thus far - idealization
- ❖ all routers identical
- ❖ network "flat"
- … *not* true in practice

*scale:* with 600 million destinations:
- ❖ can't store all dest's in routing tables!
- ❖ routing table exchange would swamp links!

*administrative autonomy*
- ❖ internet = network of networks
- ❖ each network admin may want to control routing in its own network

# Hierarchical routing

- aggregate routers into regions, "autonomous systems" (AS)
- routers in same AS run same routing protocol
  - "intra-AS" routing protocol
  - routers in different AS can run different intra-AS routing protocol

*gateway router:*
- at "edge" of its own AS
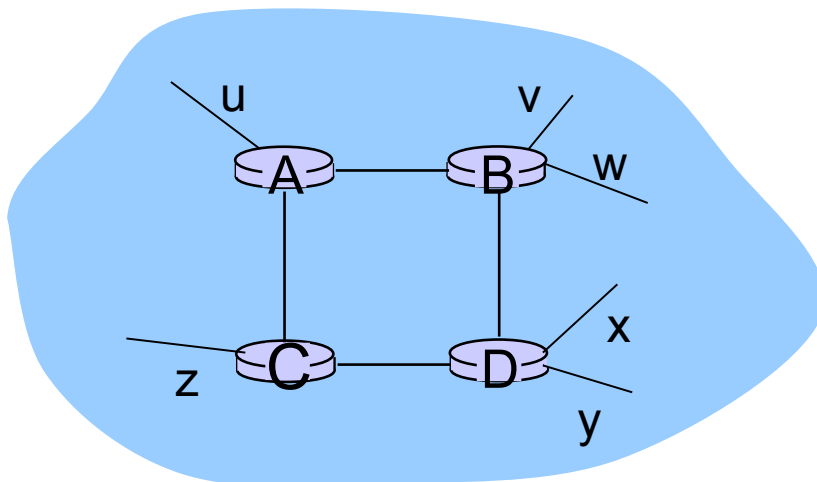- has link to router in another AS

# Interconnected ASes



- forwarding table configured by both intra- and inter-AS routing algorithm
  - intra-AS sets entries for internal dests
  - inter-AS & intra-AS sets entries for external dests

# Intra-AS Routing

☐ also known as *interior gateway protocols (IGP)*

☐ most common intra-AS routing protocols:

○ RIP: Routing Information Protocol

○ OSPF: Open Shortest Path First

○ IGRP: Interior Gateway Routing Protocol (Cisco proprietary)
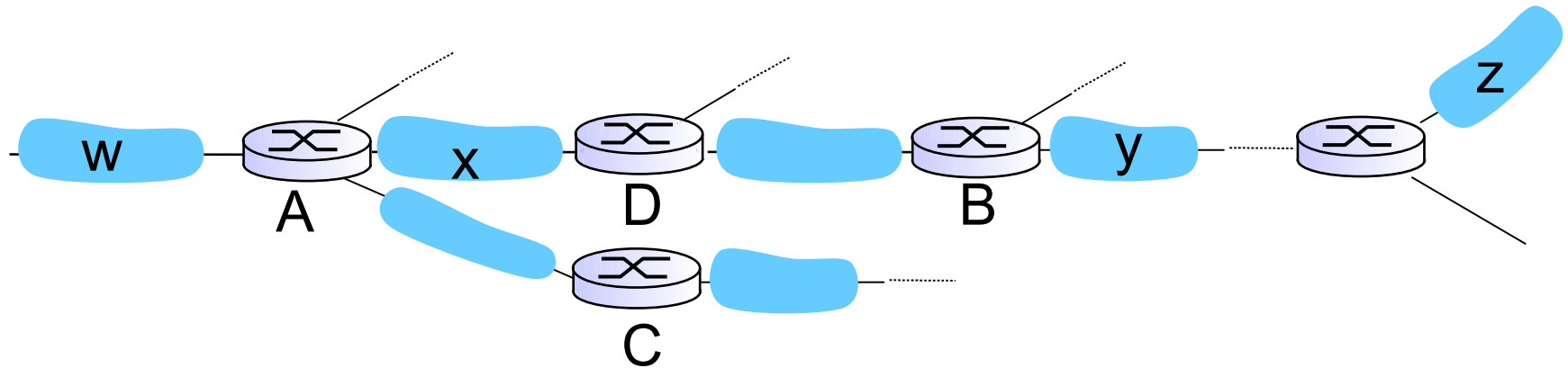
# RIP ( Routing Information Protocol)

❖ included in BSD-UNIX distribution in 1982

❖ distance vector algorithm
  ▪ distance metric: # hops (max = 15 hops), each link has cost 1
  ▪ DVs exchanged with neighbors every 30 sec in response message (aka advertisement)
  ▪ each advertisement: list of up to 25 destination *subnets (in IP addressing sense)*

from router A to destination *subnets:*

| subnet | hops |
|--------|------|
| u | 1 |
| v | 2 |
| w | 2 |
| x | 3 |
| y | 3 |
| z | 2 |

# RIP: example



routing table in router D

| destination subnet | next router | # hops to dest |
|---|---|---|
| w | A | 2 |
| y | B | 2 |
| | | |
| x | -- | 1 |
| …. | …. | .... |

# RIP: example

A-to-D advertisement

| dest | next | hops |
|------|------|------|
| w | - | 1 |
| x | - | 1 |
| z | C | 4 |
| .... | ... | ... |



routing table in router D

| destination subnet | next router | # hops to dest |
|--------------------|-------------|----------------|
| w | A | 2 |
| y | B | 2 |
| z | B | 7 |
| x | -- | 1 |
| .... | .... | .... |

A          5

# RIP: link failure, recovery

if no advertisement heard after 180 sec --> neighbor/link declared dead

- routes via neighbor invalidated
- new advertisements sent to neighbors
- neighbors in turn send out new advertisements (if tables changed)
- link failure info quickly (?) propagates to entire net
- *poison reverse* used to prevent ping-pong loops (infinite distance = 16 hops)

# Internet inter-AS routing: BGP

❖ **BGP (Border Gateway Protocol):** *the* de facto inter-domain routing protocol
  ▪ "glue that holds the Internet together"
❖ BGP provides each AS a means to:
  ▪ **eBGP:** obtain subnet reachability information from neighboring ASs.
  ▪ **iBGP:** propagate reachability information to all AS-internal routers.
  ▪ determine "good" routes to other networks based on reachability information and policy.
❖ allows subnet to advertise its existence to rest of Internet: *"I am here"*

# Agenda

❖ Link Layer
- Link-Layer Addressing and ARP
- Switches
- VLANs

❖ Network Layer
- IPv4 Addressing
- Routing algorithms
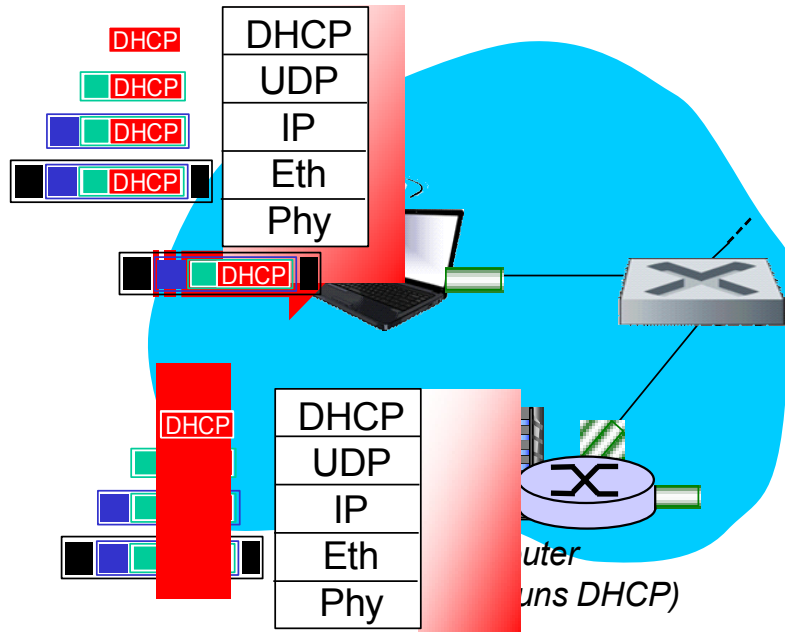- Routing in the Internet

❖ 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
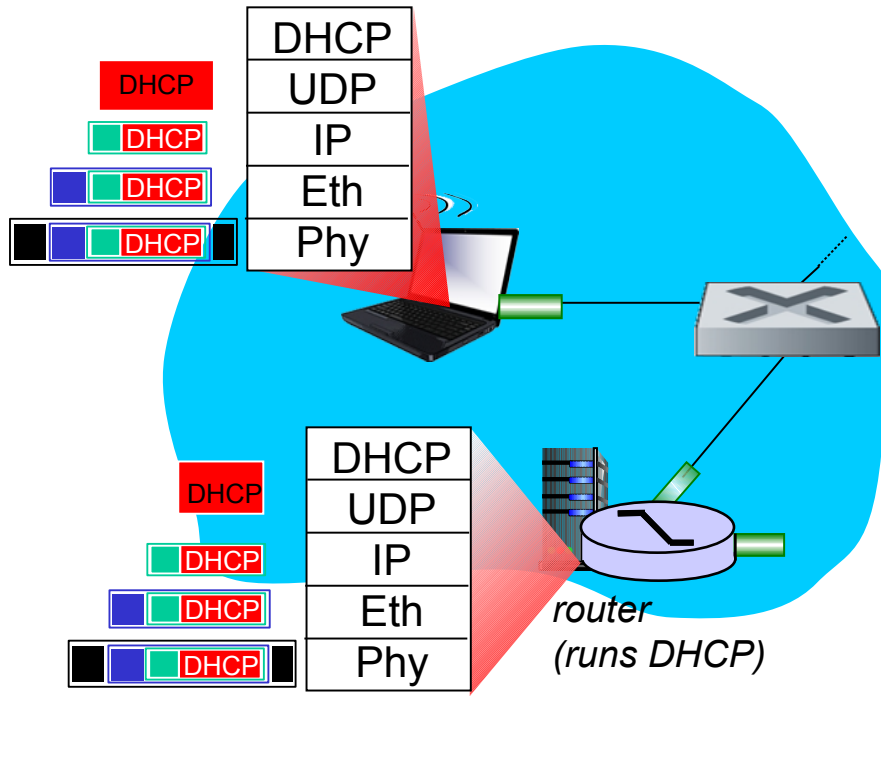
# A day in the life: scenario



browser

DNS server

Comcast network
68.80.0.0/13

school network
68.80.2.0/24

web page

Google

| | |
|---|---|
| Google Search | I'm Feeling Lucky |

Advanced Search
Preferences
Language Tools

Advertising Programs - Business Solutions - About Google

©2009 - Privacy

web server
64.233.169.105

Google's network
64.233.160.0/19

Link Layer

# A day in the life… connecting to the Internet



r 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
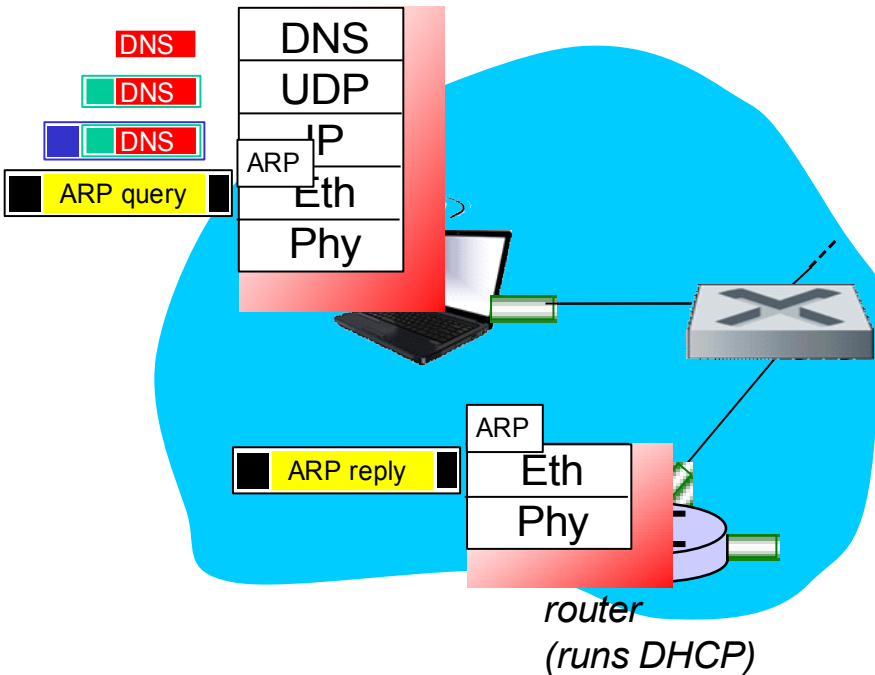
❖ Ethernet demuxed to IP demuxed, UDP demuxed to DHCP

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

DHCP
UDP
IP
Eth
Phy

DHCP
DHCP
DHCP
DHCP

DHCP
UDP
IP
Eth
Phy

DHCP
DHCP
DHCP
DHCP

*router*
*(runs DHCP)*

❐ 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
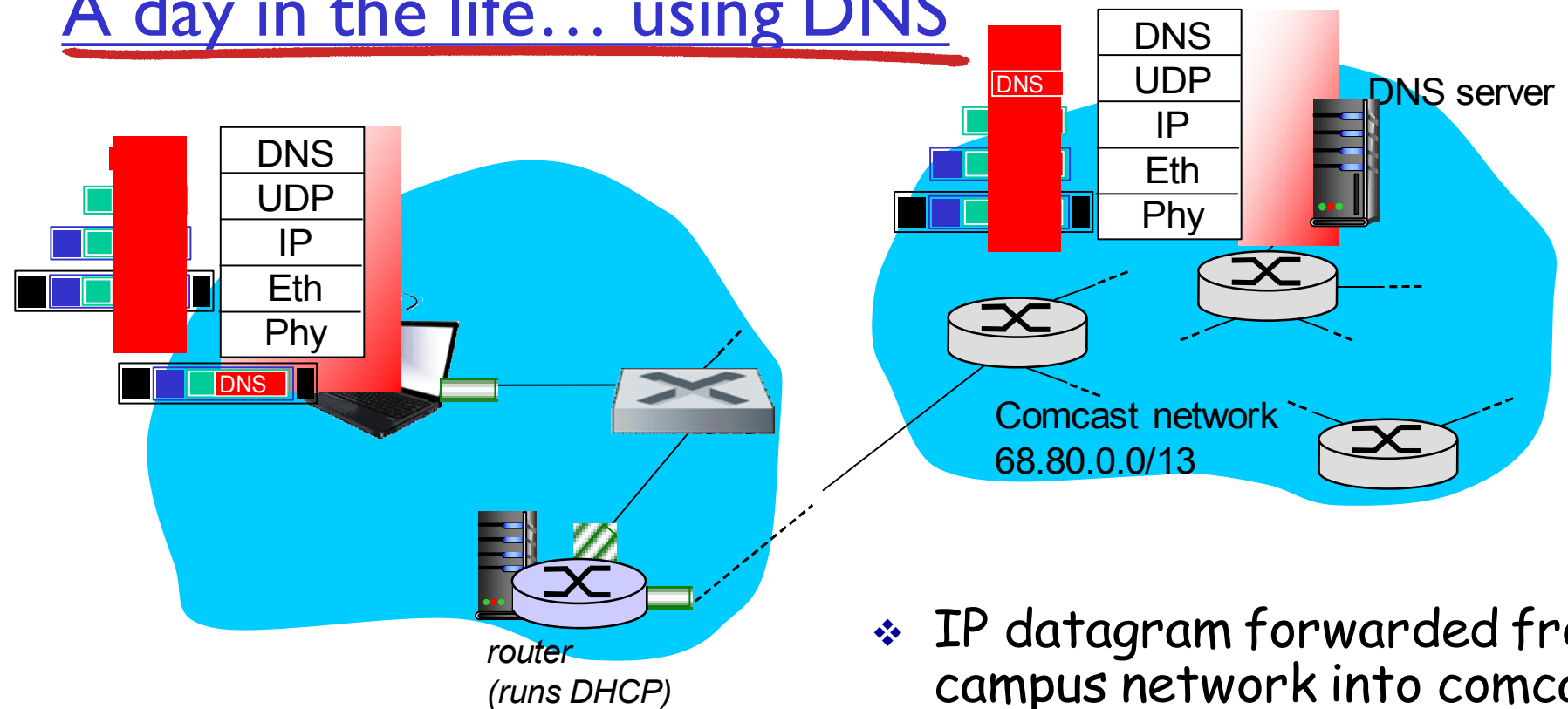
❖ DHCP client receives DHCP ACK reply

*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)



*router*
*(runs DHCP)*

r 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
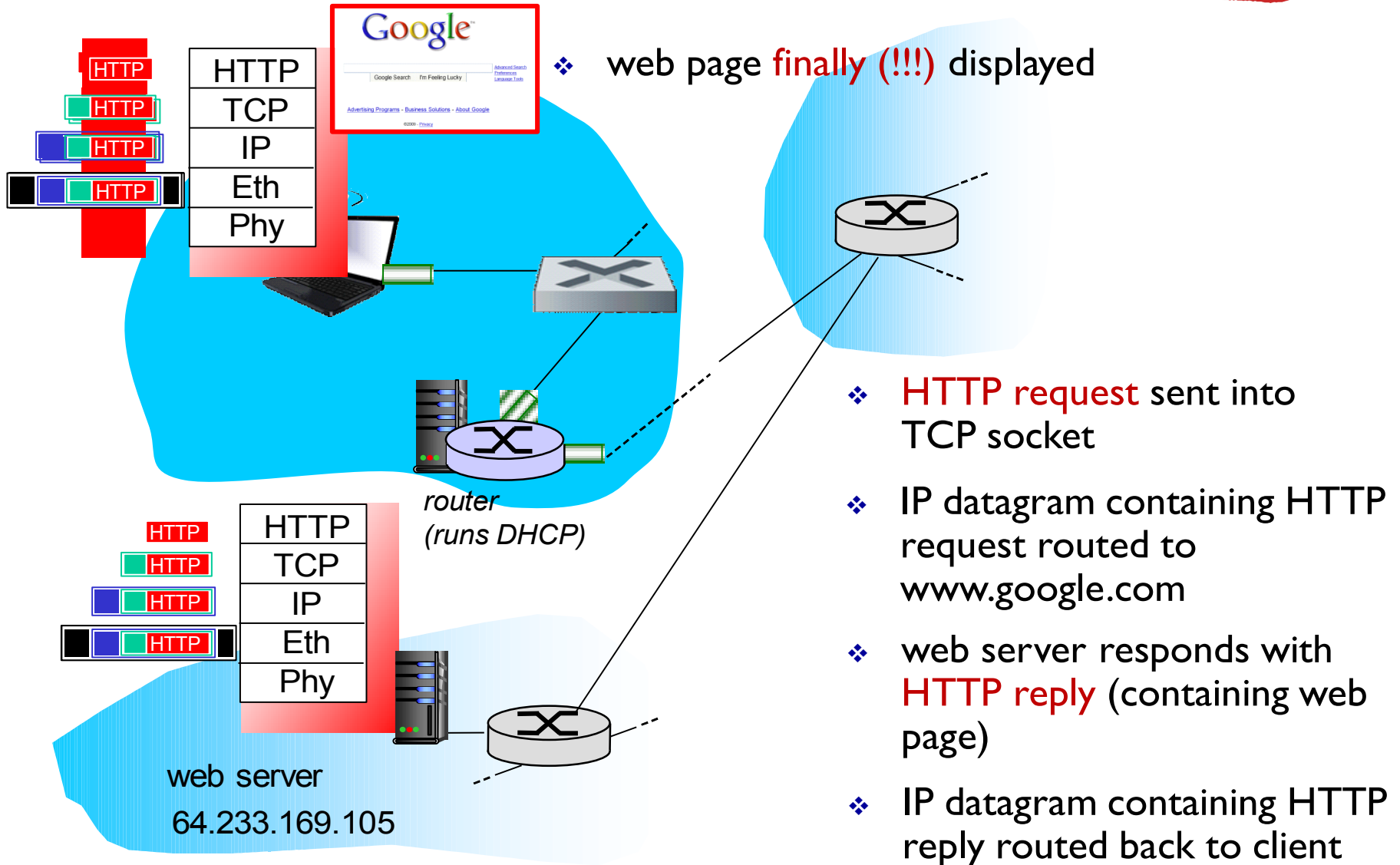
# A day in the life… using DNS



router
(runs DHCP)

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

Comcast network
68.80.0.0/13

❖ 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

| HTTP |
|------|
| TCP  |
| IP   |
| Eth  |
| Phy  |

SYNACK
SYNACK
SYNACK

router
(runs DHCP)

| TCP |
|-----|
| IP  |
| Eth |
| Phy |

SYNACK
SYNACK
SYNACK

web server
64.233.169.105

- ❖ to send HTTP request, client first opens TCP socket to web server

- ❖ TCP SYN segment (step 1 in 3-way handshake) inter-domain routed to web server

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

- ❖ TCP connection established!

# A day in the life… HTTP request/reply



❖ 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

# Τέλος Ενότητας

# Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στα πλαίσια του εκπαιδευτικού έργου του διδάσκοντα.

- Το έργο «**Ανοικτά Ακαδημαϊκά Μαθήματα στο Πανεπιστήμιο Κρήτης**» έχει χρηματοδοτήσει μόνο τη αναδιαμόρφωση του εκπαιδευτικού υλικού.

- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «Εκπαίδευση και Δια Βίου Μάθηση» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.



1

# Σημειώματα

# Σημείωμα αδειοδότησης

# Σημείωμα Αναφοράς

# Διατήρηση Σημειωμάτων

Οποιαδήποτε αναπαραγωγή ή διασκευή του υλικού θα πρέπει να συμπεριλαμβάνει:

- το Σημείωμα Αναφοράς

- το Σημείωμα Αδειοδότησης

- τη δήλωση Διατήρησης Σημειωμάτων

- το Σημείωμα Χρήσης Έργων Τρίτων (εφόσον υπάρχει)

μαζί με τους συνοδευόμενους υπερσυνδέσμους.