# Chapter 4
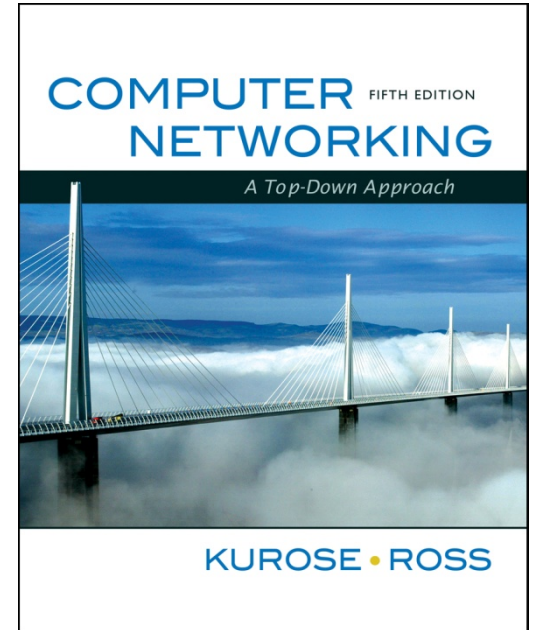# Network Layer

*Computer Networking: A Top Down Approach*
5th edition.
Jim Kurose, Keith Ross
Addison-Wesley, April 2009.
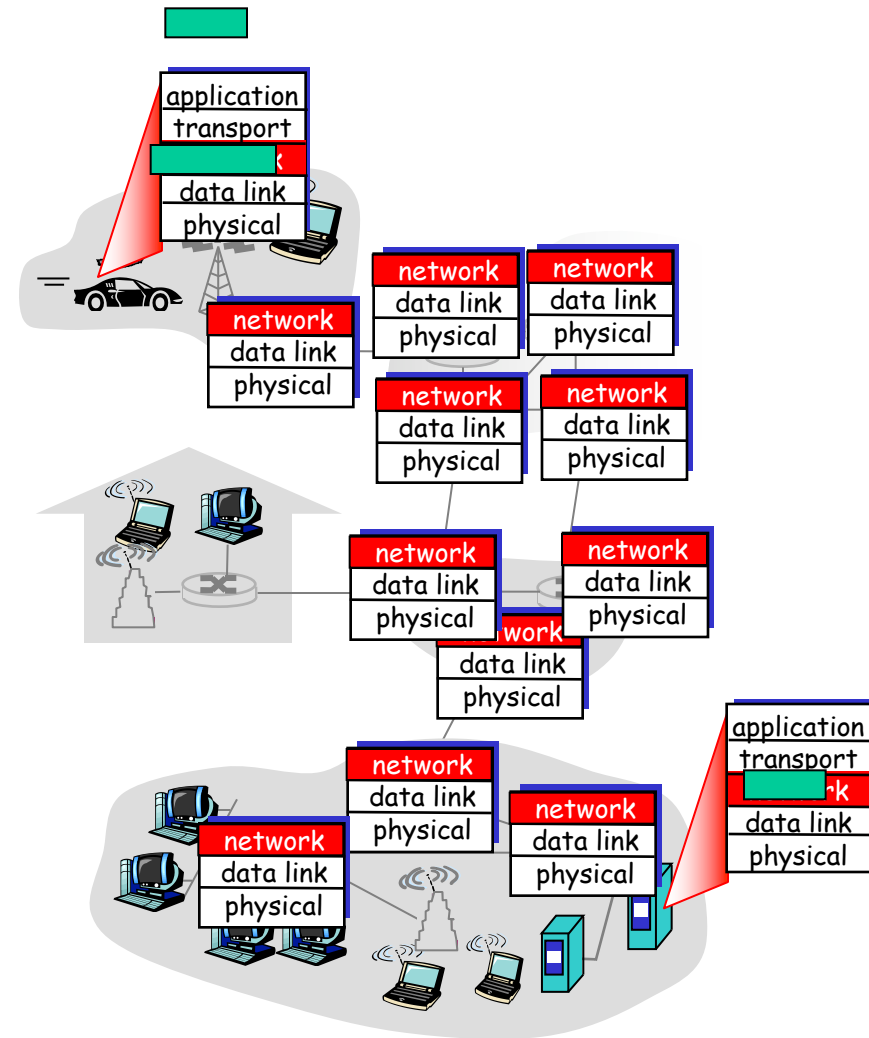
# Chapter 4: Network Layer

## Chapter goals:

❑ understand principles behind network layer services:

- ❖ network layer service models
- ❖ forwarding versus routing
- ❖ how a router works
- ❖ routing (path selection)
- ❖ dealing with scale
- ❖ advanced topics: IPv6, mobility

❑ instantiation, implementation in the Internet

# Chapter 4: Network Layer

# Network layer

- transport segment from sending to receiving host

- on sending side encapsulates segments into datagrams

- on rcving side, delivers segments to transport layer

- network layer protocols in *every* host, router

- router examines header fields in all IP datagrams passing through it

# Two Key Network-Layer Functions

❑ *forwarding:* move packets from router's input to appropriate router output

❑ *routing:* determine route taken by packets from source to dest.
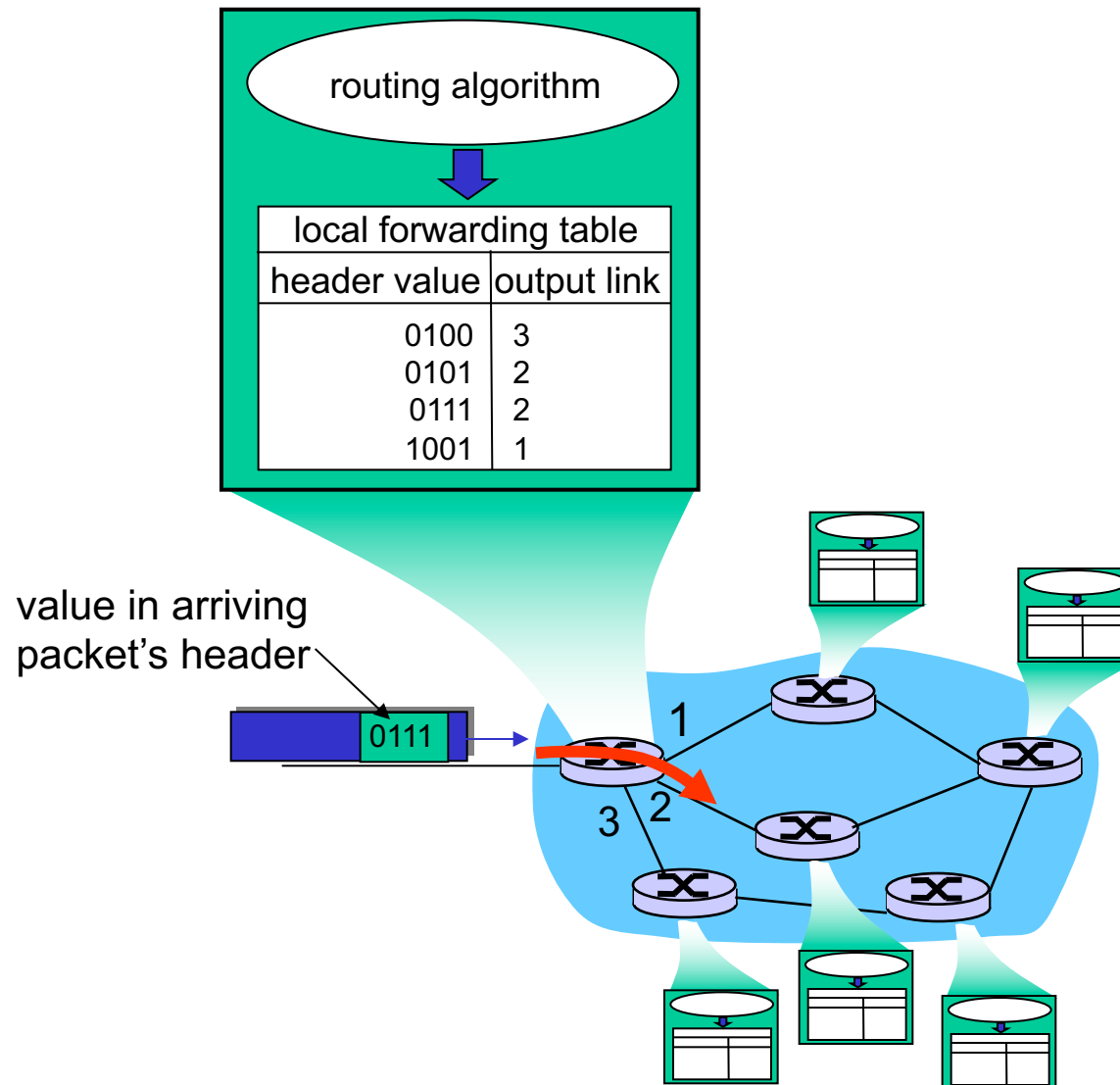
❖ *routing algorithms*

❑ routing: process of planning trip from source to dest

❑ forwarding: process of getting through single interchange

# Interplay between routing and forwarding

routing algorithm

| local forwarding table | |
|---|---|
| header value | output link |
| 0100 | 3 |
| 0101 | 2 |
| 0111 | 2 |
| 1001 | 1 |

value in arriving
packet's header

0111

1

3   2

# Third important function: connection setup

❑ 3<sup>rd</sup> important function in *some* network architectures:
  ❖ ATM, frame relay, X.25
❑ before datagrams flow, two end hosts *and* intervening routers establish virtual connection
  ❖ routers get involved
❑ network vs transport layer connection service:
  ❖ network: between two hosts (may also involve intervening routers in case of VCs)
  ❖ transport: between two processes

# Chapter 4: Network Layer

# Network layer connection and connection-less service

❑ **datagram network** provides network-layer *connectionless* service

❑ **VC network** provides network-layer *connection* service

❑ analogous to the transport-layer services, but:
  - ❖ **service:** host-to-host
  - ❖ **no choice:** network provides one or the other
  - ❖ **implementation:** in network core

# Virtual circuits

"source-to-dest path behaves much like telephone circuit"

  ❖ performance-wise
  ❖ network actions along source-to-dest path

❑ call setup, teardown for each call *before* data can flow
❑ each packet carries VC identifier (not destination host address)
❑ *every* router on source-dest path maintains "state" for each passing connection
❑ link, router resources (bandwidth, buffers) may be *allocated* to VC (dedicated resources = predictable service)

# VC implementation

a VC consists of:

1. path from source to destination
2. VC numbers, one number for each link along path
3. entries in forwarding tables in routers along path

❑ packet belonging to VC carries VC number (rather than dest address)

❑ VC number can be changed on each link.

❖ New VC number comes from forwarding table

# Forwarding table



VC number

Router 1

interface number

## Forwarding table in Router 1:

| Incoming interface | Incoming VC # | Outgoing interface | Outgoing VC # |
|---|---|---|---|
| 1 | 12 | 3 | 22 |
| 2 | 63 | 1 | 18 |
| 3 | 7 | 2 | 17 |
| 1 | 97 | 3 | 87 |
| … | … | … | … |

Routers maintain connection state information!

# Virtual circuits: signaling protocols

❑ used to setup, maintain  teardown VC

❑ used in ATM, frame-relay, X.25

❑ not used in today's Internet

*What's the difference between VC setup and TCP three-way handshake?*

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

5. Data flow begins

4. Call connected

1. Initiate call

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

6. Receive data

3. Accept call

2. incoming call

# Datagram networks

❑ no call setup at network layer

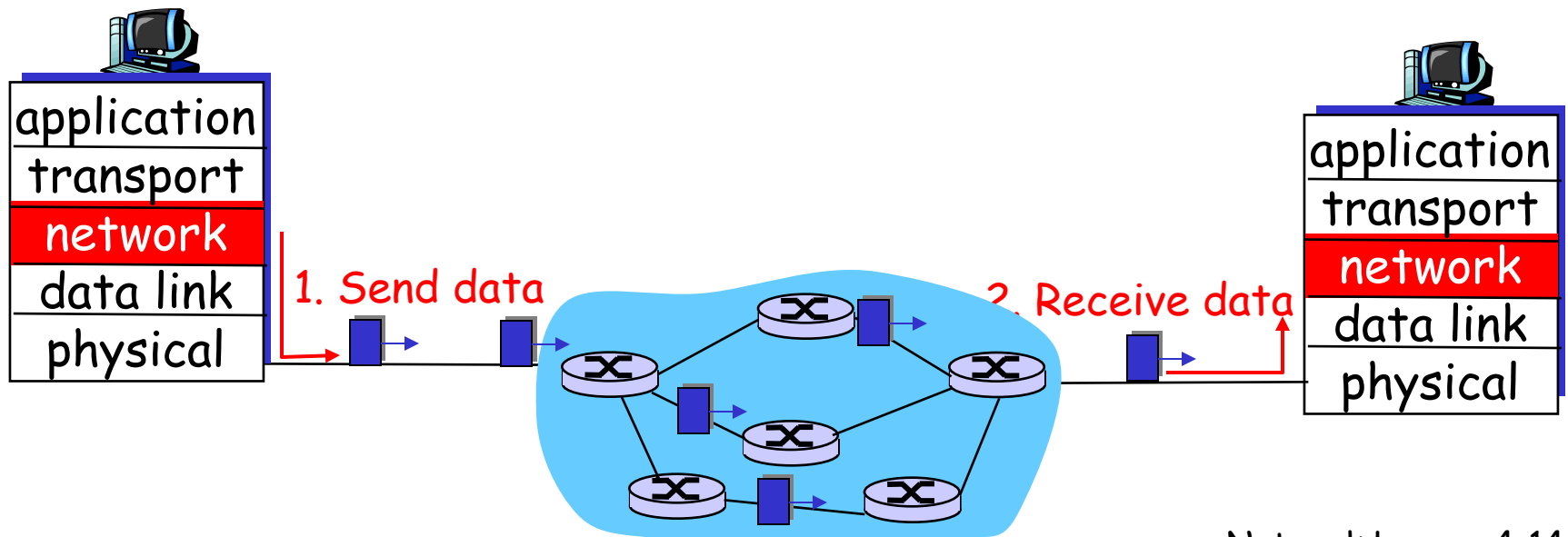❑ routers: no state about end-to-end connections
  ❖ no network-level concept of "connection"

❑ packets forwarded using destination host address
  ❖ packets between same source-dest pair may take different paths

| application |
|-------------|
| transport |
| **network** |
| data link |
| physical |

1. Send data

2. Receive data

| application |
|-------------|
| transport |
| **network** |
| data link |
| physical |

# Forwarding table

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

# Forwarding table

| Destination Address Range | Link Interface |
|---|---|
| **11001000 00010111 00010**000 00000000<br>through<br>**11001000 00010111 00010**111 11111111 | 0 |
| **11001000 00010111 00011000** 00000000<br>through<br>**11001000 00010111 00011000** 11111111 | 1 |
| **11001000 00010111 00011**001 00000000<br>through<br>**11001000 00010111 00011**111 11111111 | 2 |
| otherwise | 3 |

# Longest prefix matching

|  Prefix Match | Link Interface |
|---|---|
| 11001000 00010111 00010 | 0 |
| 11001000 00010111 00011000 | 1 |
| 11001000 00010111 00011 | 2 |
| otherwise | 3 |

Examples

DA: 11001000  00010111  00010110  10100001        Which interface?
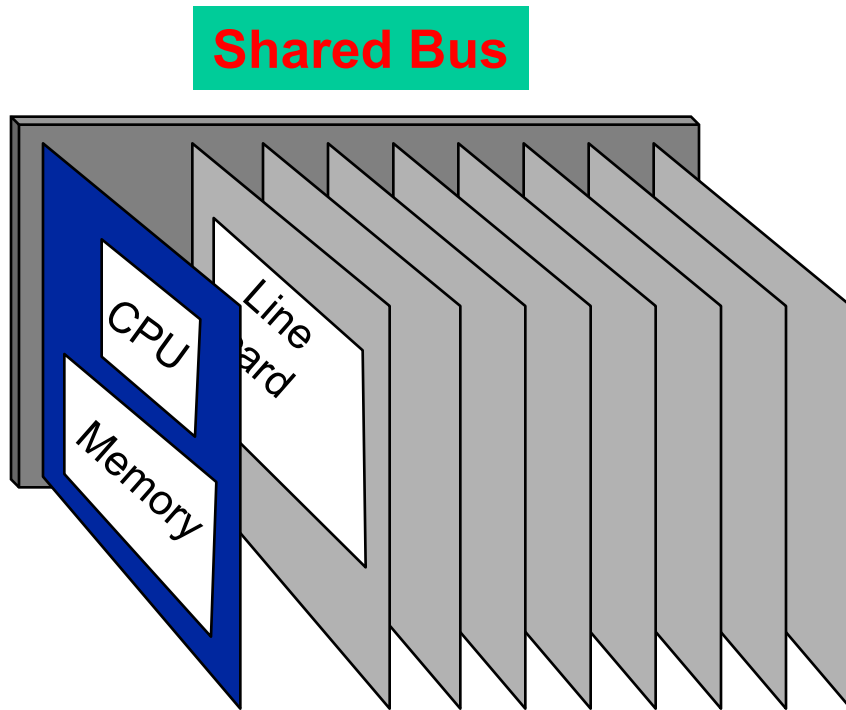
DA: 11001000  00010111  00011000  10101010        Which interface?

DA: 11001000  00010111  00111000  10101010

# Chapter 4: Network Layer

# First Generation Routers

**Shared Bus**

**Off-chip Buffer**

CPU

Memory

Line Card

CPU  Route Table  Buffer Memory

Line cards  Line cards  ......  Line cards

MAC  MAC  MAC

Line cards: enables router to attach to different networks that employ different Data link technologies

# First Generation Routers

- Traditional computers with switching under direct control of CPU
- Packet copied to system's shared memory
- Speed limited by memory bandwidth
- Every packet needs two transfers between line cards and memory
- Does not scale to too many line cards
- Suffices for Low speed routers

# Second Generation Routers



CPU    Route Table    **Buffer Memory**

**Line Card**    **Line Card**    **Line Card**

Buffer Memory

Fwding Cache

MAC

Cache update

Fast path

Slow path

# Second Generation IP routers

- ❑ Each line card has a route cache
- ❑ On a hit, forward directly
  - ❖ Fast path
- ❑ On a miss, via CPU bus, memory
  - ❖ Slow path
- ❑ Bus contention:  switching speed limited by bus bandwidth

# Third Generation Routers

**"Crossbar": Switched Backplane**

Line
CPU
Memory

Line Card

CPU Card

Line Card

Local Buffer Memory

Routing Table

Local Buffer Memory

Fwding Table

Fwding Table

MAC

MAC

23

# Crossbar

❑ Every input port has a connection to every output port
  ❖ N inputs, N outputs →actually, inputs are also outputs

❑ During each timeslot, each input connected to zero or one outputs

❑ **Advantage:** Exploits parallelism
❑ **Disadvantage:** Need scheduling algorithm
❑ Scheduling:
  ❖ Scheduling is pipelined (4 state): get bids, allocate schedule, notify pairings, ship bits
  ❖ Maximize switch throughput rather than bounded latency



crossbar

# Where does the queuing occur?

□ Q: do we need queue if:
- ❖ All line cards have the identical speeds (input/output)
- ❖ $n$ input line cards and $n$ output line cards
- ❖ Switching fabric is at least $n$ times the line speed

# Output port queuing



Output Port Contention at Time t

One Packet Time Later

- ❑ buffering when arrival rate via switch exceeds output line speed
- ❑ *queuing (delay) and loss due to output port buffer overflow!*

# Input Port Queuing

☐ Fabric slower that input ports combined -> queueing may occur at input queues

☐ Head-of-the-Line (HOL) blocking: queued datagram at front of queue prevents others in queue from moving forward

☐ *queueing delay and loss due to input buffer overflow!*

output port contention
at time t – only one red
packet can be transferred

green packet
experiences HOL blocking

# Solution: Virtual Output Queues

❑ Maintain N virtual queues at each input
  ❖ one per output

Input 1

Input 2

Input 3

Output 1

Output 2

Output 3

# Chapter 4: Network Layer

- ❑ 4. 1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - ❖ Datagram format
  - ❖ IPv4 addressing
  - ❖ ICMP
  - ❖ IPv6

- ❑ 4.5 Routing algorithms
  - ❖ Link state
  - ❖ Distance Vector
  - ❖ Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - ❖ RIP
  - ❖ OSPF
  - ❖ BGP
- ❑ 4.7 Broadcast and multicast routing

# The Internet Network layer

Host, router network layer functions:

Network layer

| |
| Transport layer: TCP, UDP |

**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

# Chapter 4: Network Layer

- ❑ 4. 1 Introduction
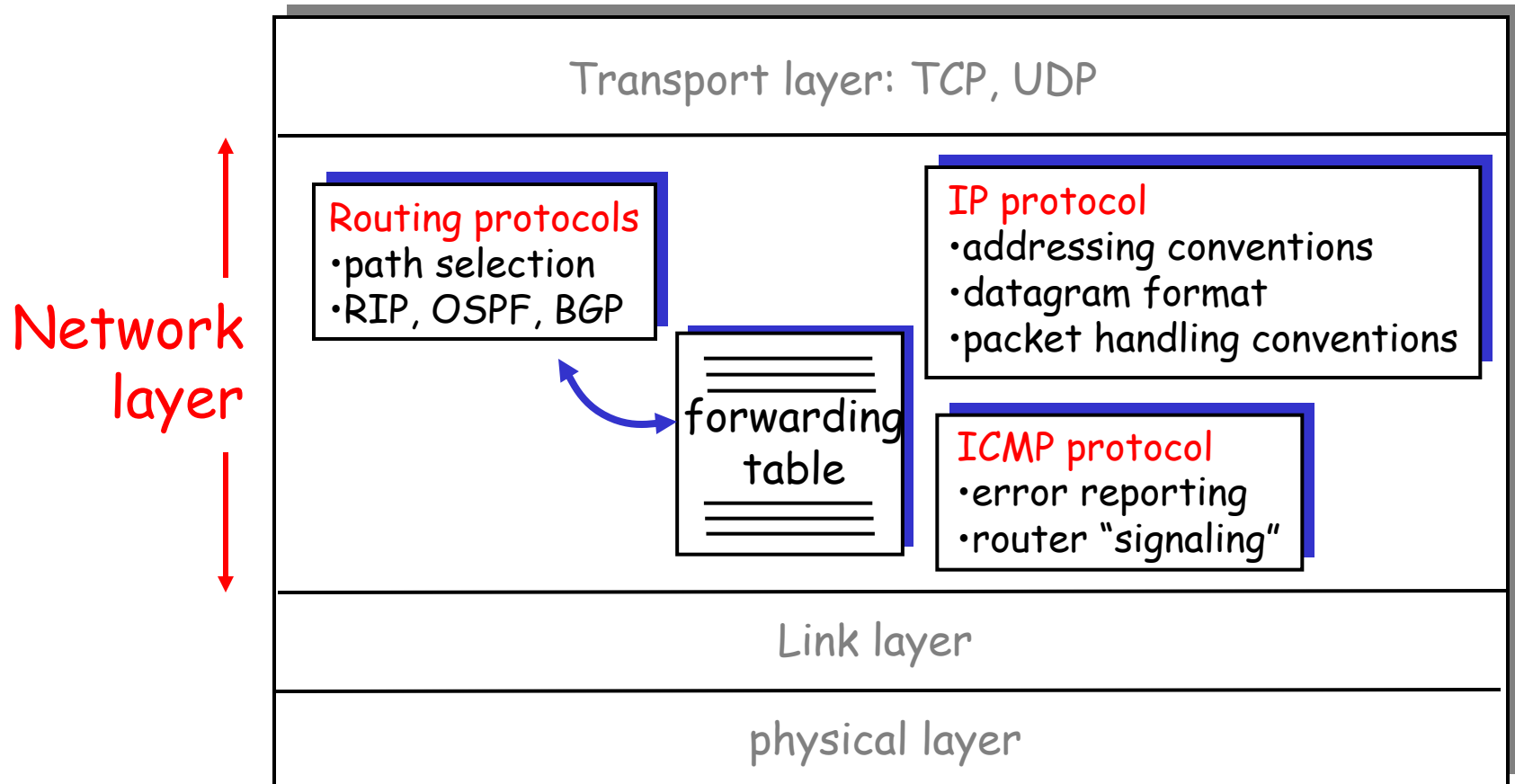- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - ❖ Datagram format
  - ❖ IPv4 addressing
  - ❖ ICMP
  - ❖ IPv6

- ❑ 4.5 Routing algorithms
  - ❖ Link state
  - ❖ Distance Vector
  - ❖ Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - ❖ RIP
  - ❖ OSPF
  - ❖ BGP
- ❑ 4.7 Broadcast and multicast routing

# IP datagram format

IP protocol version number

header length (bytes)

"type" of data

max number remaining hops (decremented at each router)

upper layer protocol to deliver payload to

total datagram length (bytes)

for fragmentation/ reassembly

E.g. timestamp, record route taken, specify list of routers to visit.

32 bits

| ver | head. len | type of service | length |
| 16-bit identifier | | flgs | fragment offset |
| time to live | upper layer | header checksum |
| 32 bit source IP address | | | |
| 32 bit destination IP address | | | |
| Options (if any) | | | |
| data (variable length, typically a TCP or UDP segment) | | | |

how much overhead with TCP?

❑ 20 bytes of TCP

❑ 20 bytes of IP

❑ = 40 bytes + app layer overhead

# IP Fragmentation & Reassembly

❑ network links have MTU (max.transfer size) - largest possible link-level frame.

  ❖ different link types, different MTUs

❑ large IP datagram divided ("fragmented") within net

  ❖ one datagram becomes several datagrams

  ❖ "reassembled" only at final destination

  ❖ IP header bits used to identify, order related fragments

fragmentation:
in: one large datagram
out: 3 smaller datagrams

reassembly

# IP Fragmentation and Reassembly

| | length =4000 | ID =x | fragflag =0 | offset =0 | |
|---|---|---|---|---|---|

**Example**

- 4000 byte datagram
- MTU = 1500 bytes

One large datagram becomes several smaller datagrams

1480 bytes in data field

offset = 1480/8

| | length =1500 | ID =x | fragflag =1 | offset =0 | |
|---|---|---|---|---|---|

| | length =1500 | ID =x | fragflag =1 | offset =185 | |
|---|---|---|---|---|---|

| | length =1040 | ID =x | fragflag =0 | offset =370 | |
|---|---|---|---|---|---|

# Chapter 4: Network Layer

- 4. 1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
  - ❖ Datagram format
  - ❖ IPv4 addressing
  - ❖ ICMP
  - ❖ IPv6

- 4.5 Routing algorithms
  - ❖ Link state
  - ❖ Distance Vector
  - ❖ Hierarchical routing
- 4.6 Routing in the Internet
  - ❖ RIP
  - ❖ OSPF
  - ❖ BGP
- 4.7 Broadcast and multicast routing

# 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 interface
  ❖ IP addresses associated with each interface
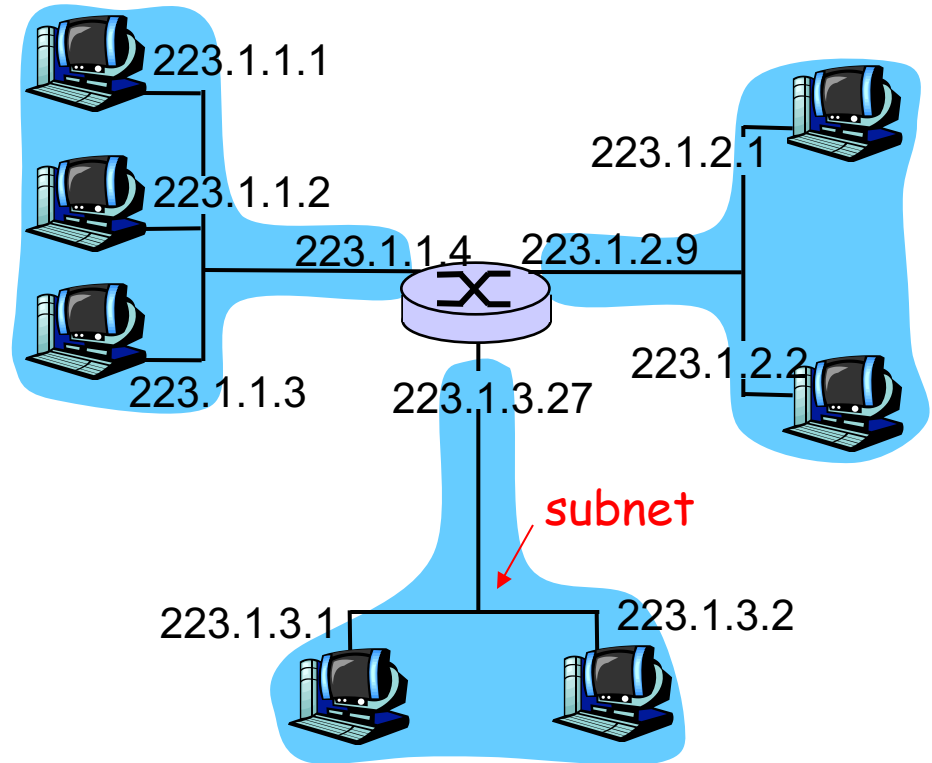
223.1.1.1

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

223.1.3.1                     223.1.3.2

223.1.1.1 = 11011111 00000001 00000001 00000001

     223         1         1         1

# 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

223.1.1.3

223.1.2.1

223.1.2.9

223.1.2.2

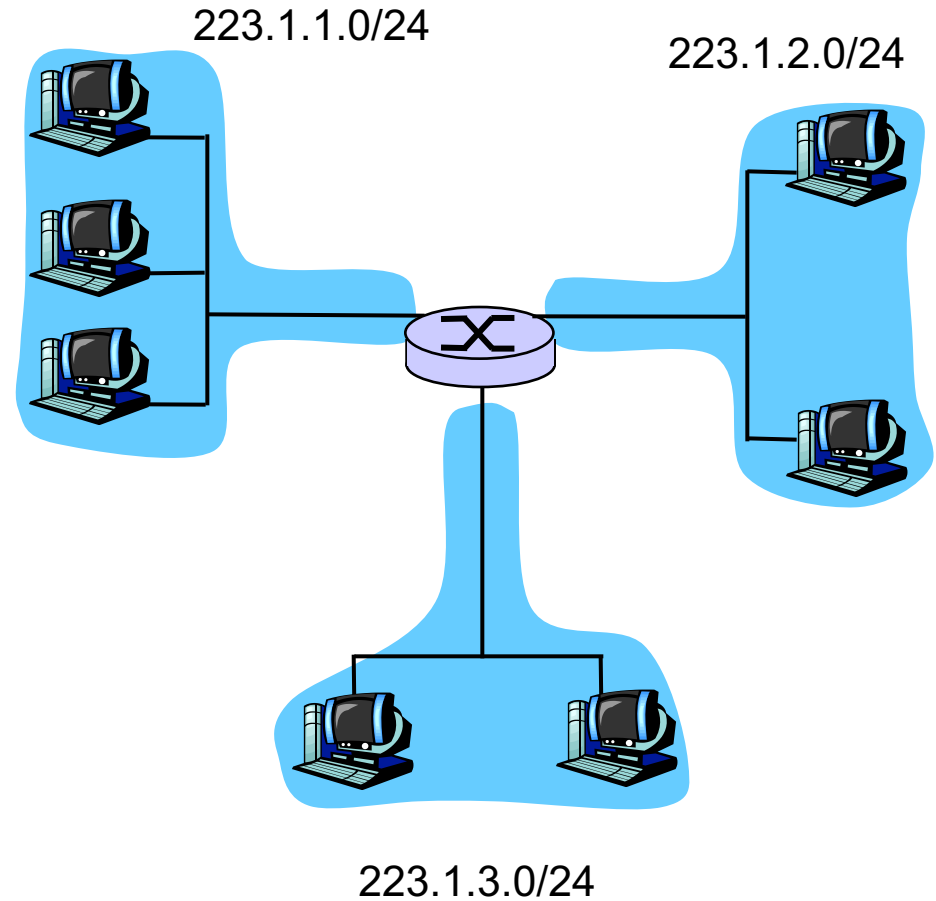223.1.3.27

subnet

223.1.3.1

223.1.3.2

network consisting of 3 subnets

# Subnets

## Recipe

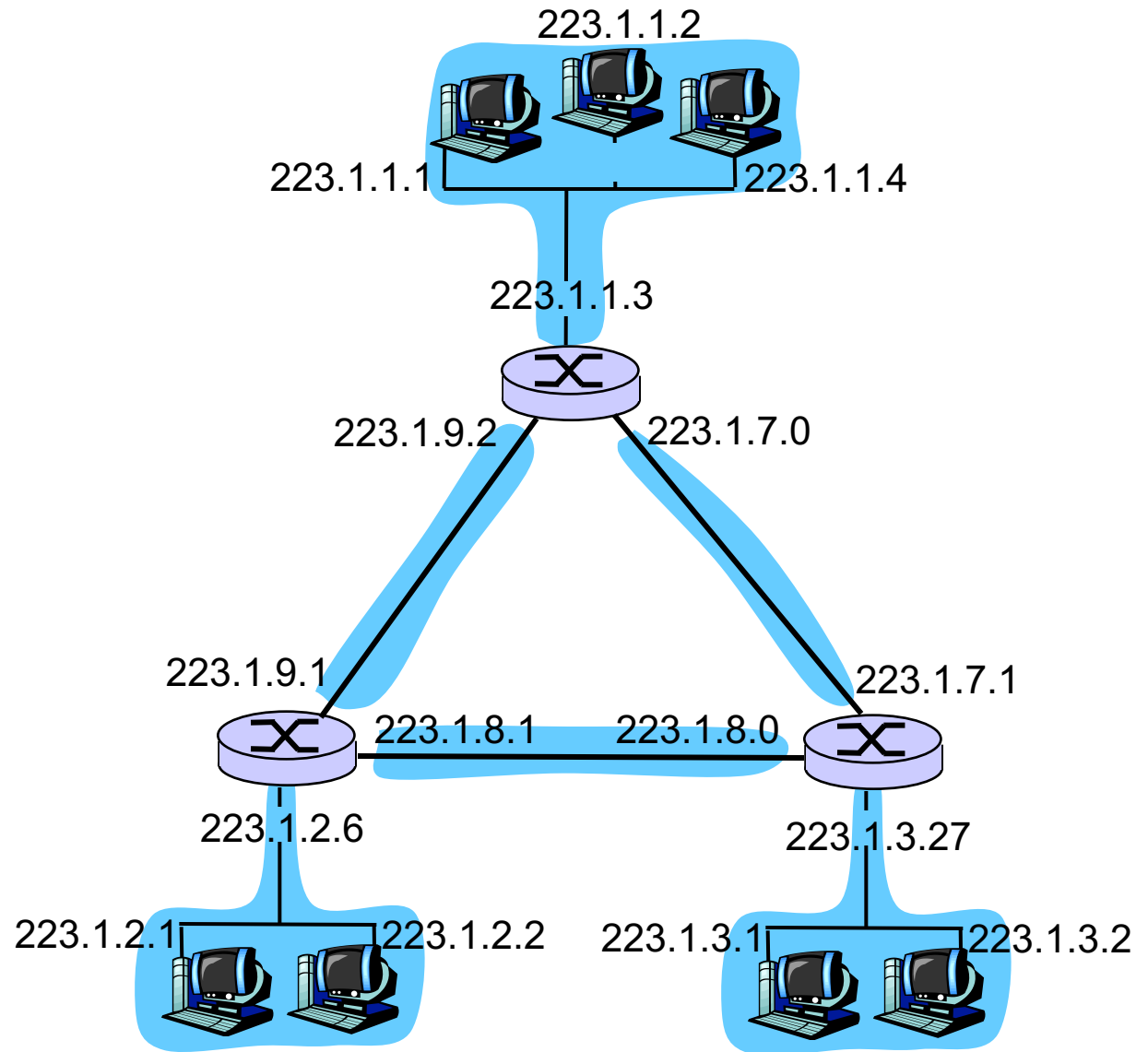❑ To determine the subnets, detach each interface from its host or router, creating islands of isolated networks. Each isolated network is called a subnet.

223.1.1.0/24

223.1.2.0/24

223.1.3.0/24

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
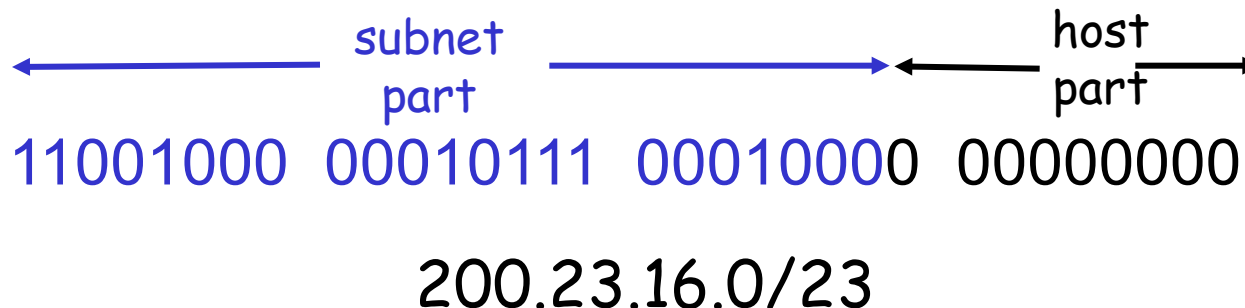
223.1.2.6

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
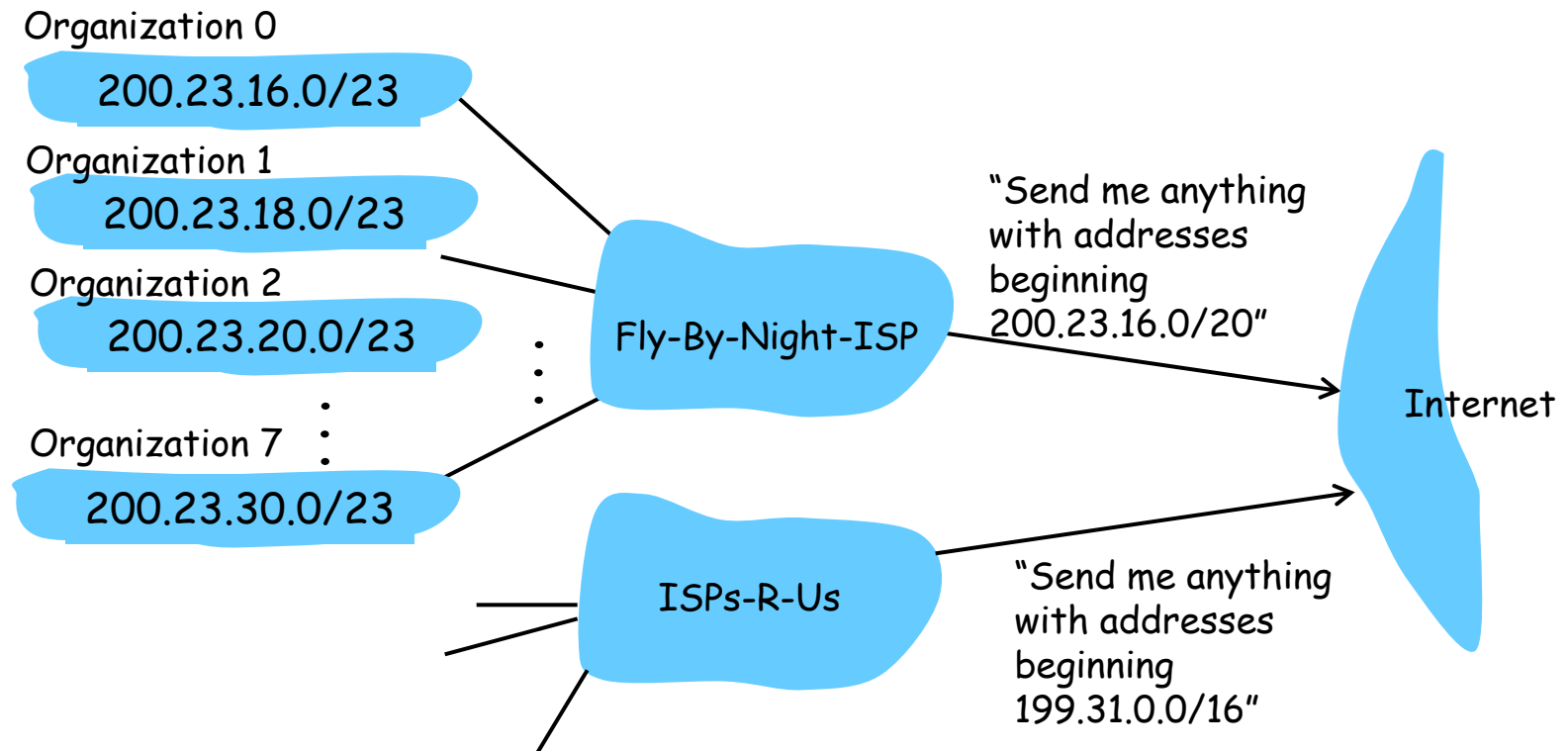❖ address format: a.b.c.d/x, where x is # bits in subnet portion of address

subnet part ← → | host part ← →

11001000 00010111 00010000 00000000

200.23.16.0/23

# IP addresses: how to get one?

**Q:** How does *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 |

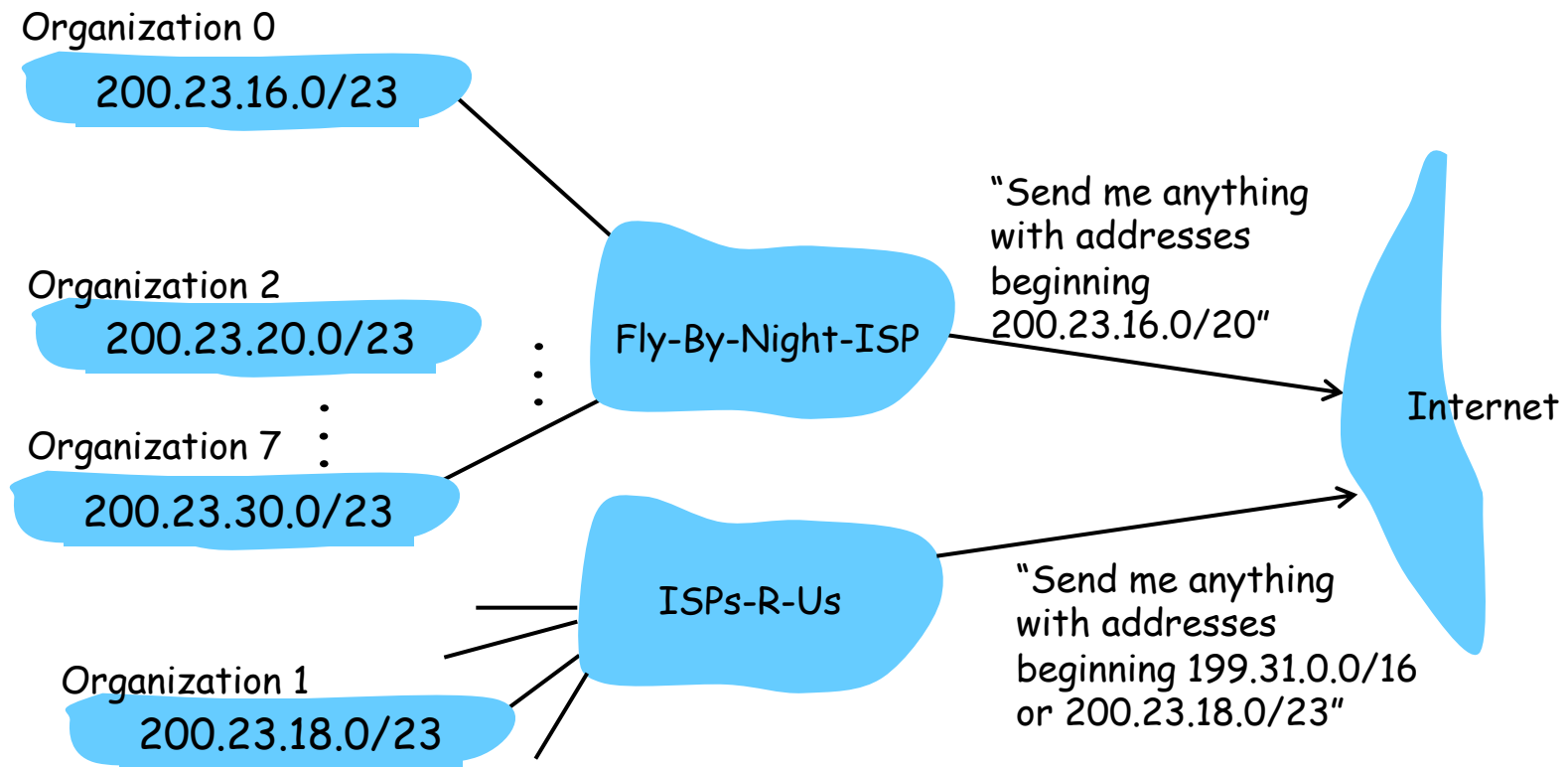# Hierarchical addressing: route aggregation

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

Organization 1
200.23.18.0/23

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 or 200.23.18.0/23"

Internet

# IP addressing: the last word...

Q: How does an ISP get block of addresses?

A: ICANN: Internet Corporation for Assigned Names and Numbers

- ❖ allocates addresses
- ❖ manages DNS
- ❖ assigns domain names, resolves disputes

# 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
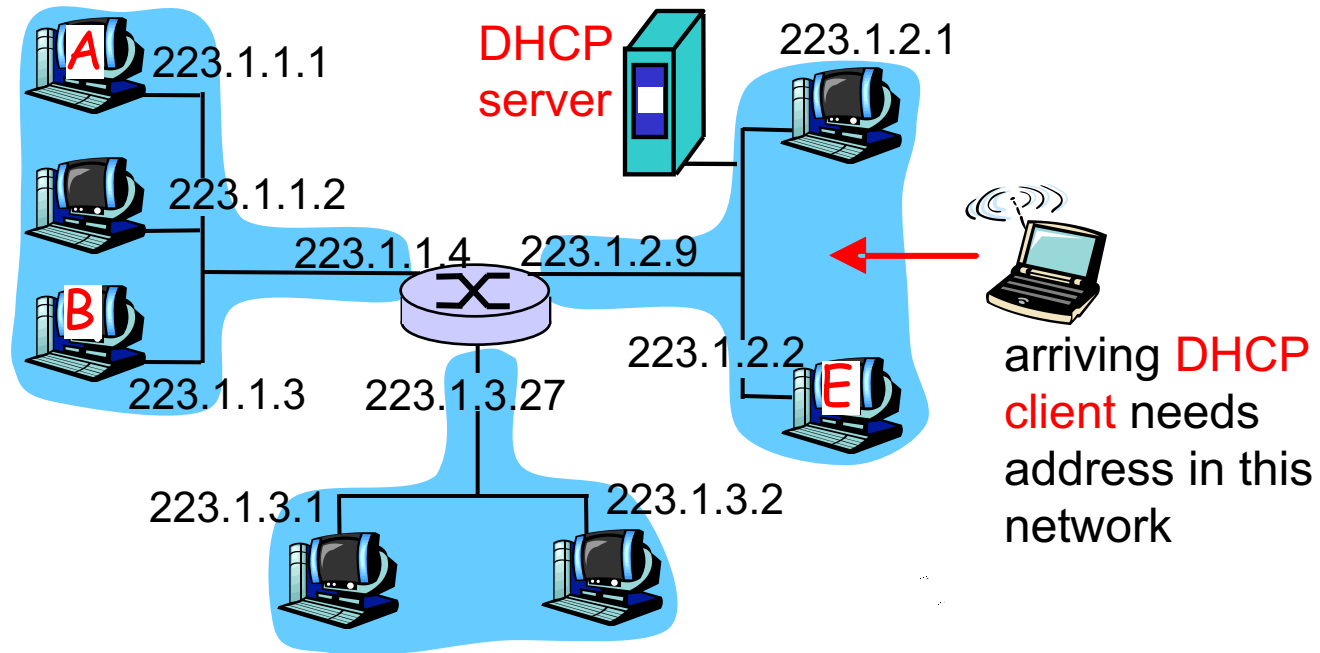
Can renew its lease on address in use

Allows reuse of addresses (only hold address while connected an "on")

Support for mobile users who want to join network (more shortly)

DHCP overview:

❖ host broadcasts "DHCP discover" msg

❖ DHCP server responds with "DHCP offer" msg

❖ host requests IP address: "DHCP request" msg

❖ DHCP server sends address: "DHCP ack" msg

# DHCP client-server scenario



A 223.1.1.1

223.1.1.2

B 223.1.1.3

223.1.1.4    223.1.2.9

223.1.3.27

223.1.3.1    223.1.3.2

DHCP server

223.1.2.1

223.1.2.2

E

arriving DHCP client needs address in this network

# 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
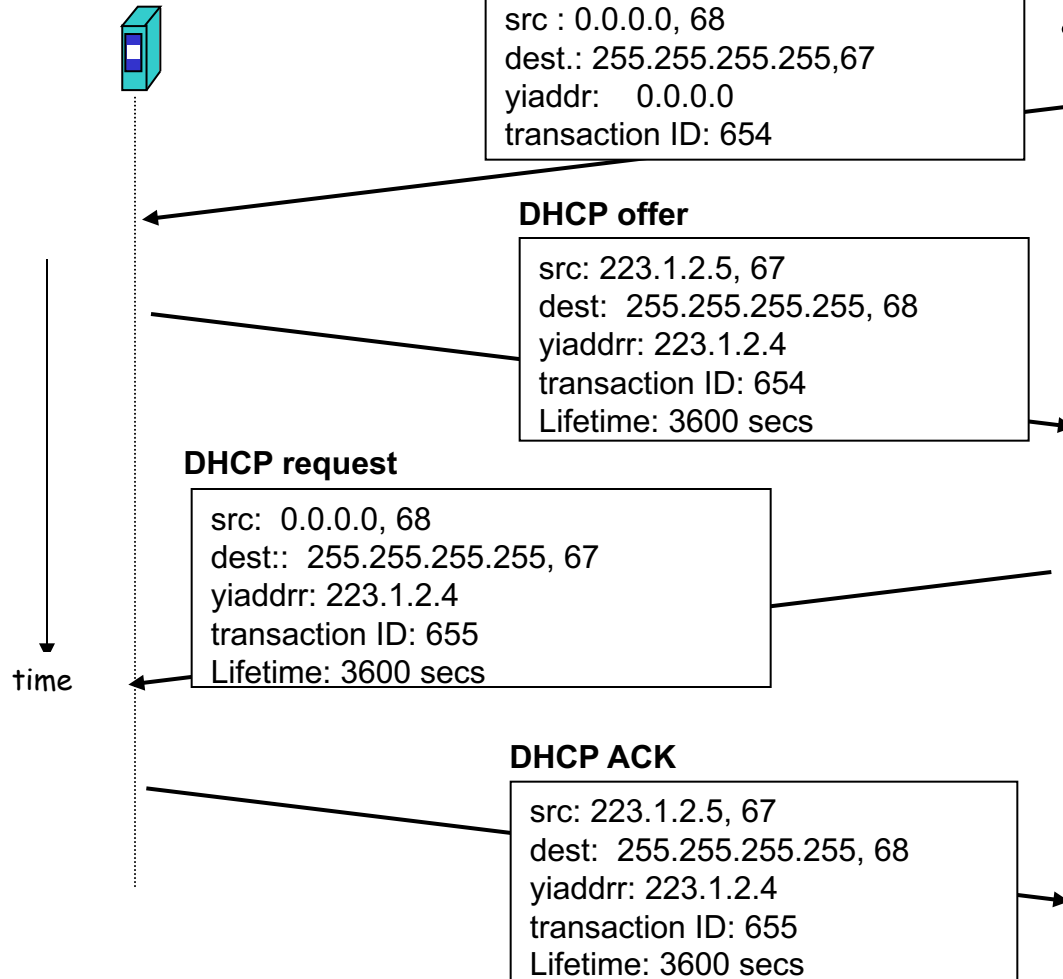yiaddr:   0.0.0.0
transaction ID: 654

**DHCP offer**

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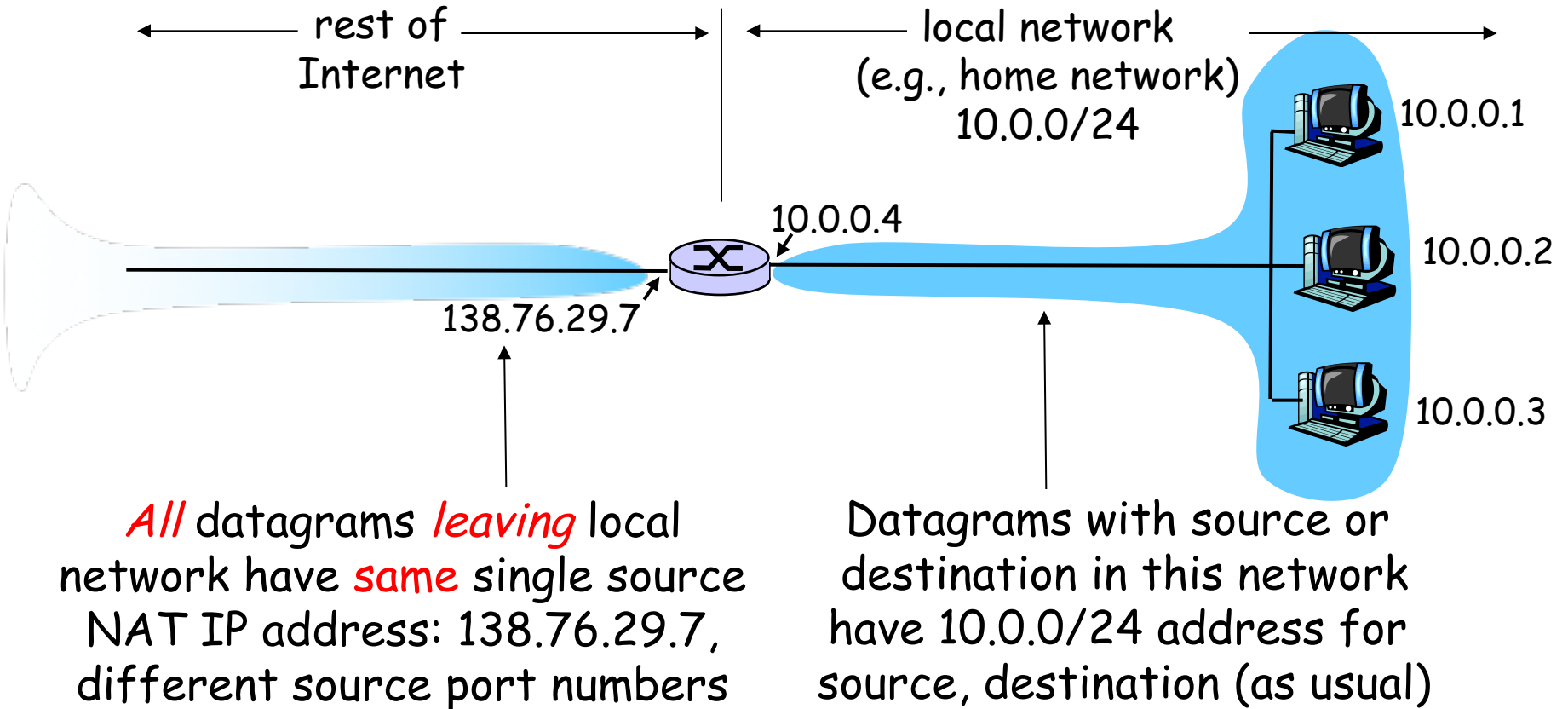transaction ID: 654
Lifetime: 3600 secs

**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

time

# 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)
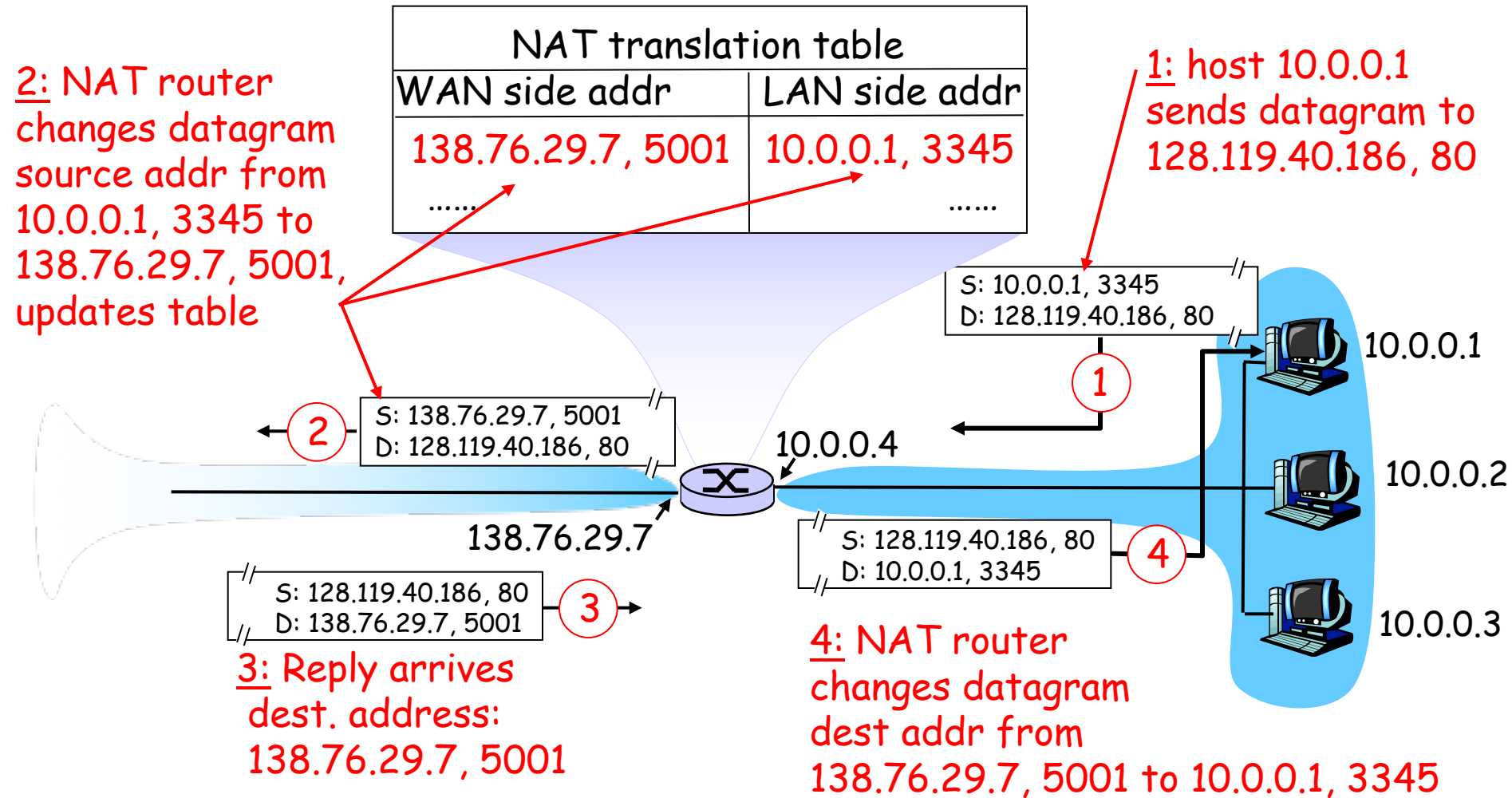
# NAT: Network Address Translation

❑ 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

10.0.0.1

2

S: 138.76.29.7, 5001
D: 128.119.40.186, 80

10.0.0.4

10.0.0.2

138.76.29.7

S: 128.119.40.186, 80
D: 10.0.0.1, 3345

4

3

S: 128.119.40.186, 80
D: 138.76.29.7, 5001

3: Reply arrives dest. address: 138.76.29.7, 5001

10.0.0.3

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

# 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, eg, P2P applications

   ❖ address shortage should instead be solved by IPv6

# Chapter 4: Network Layer

# 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 |

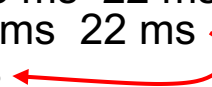| IP header | ICMP Message |
|---|---|

20 bytes

# Traceroute

traceroute: gaia.cs.umass.edu to www.eurecom.fr

Three delay measurements from
gaia.cs.umass.edu to cs-gw.cs.umass.edu

```
1  cs-gw (128.119.240.254)  1 ms  1 ms  2 ms
2  border1-rt-fa5-1-0.gw.umass.edu (128.119.3.145)  1 ms  1 ms  2 ms
3  cht-vbns.gw.umass.edu (128.119.3.130)  6 ms 5 ms 5 ms
4  jn1-at1-0-0-19.wor.vbns.net (204.147.132.129)  16 ms 11 ms 13 ms
5  jn1-so7-0-0-0.wae.vbns.net (204.147.136.136)  21 ms 18 ms 18 ms
6  abilene-vbns.abilene.ucaid.edu (198.32.11.9)  22 ms  18 ms  22 ms
7  nycm-wash.abilene.ucaid.edu (198.32.8.46)  22 ms  22 ms  22 ms
8  62.40.103.253 (62.40.103.253)  104 ms 109 ms 106 ms
9  de2-1.de1.de.geant.net (62.40.96.129)  109 ms 102 ms 104 ms
10  de.fr1.fr.geant.net (62.40.96.50)  113 ms 121 ms 114 ms
11  renater-gw.fr1.fr.geant.net (62.40.103.54)  112 ms  114 ms  112 ms
12  nio-n2.cssi.renater.fr (193.51.206.13)  111 ms  114 ms  116 ms
13  nice.cssi.renater.fr (195.220.98.102)  123 ms  125 ms  124 ms
14  r3t2-nice.cssi.renater.fr (195.220.98.110)  126 ms  126 ms  124 ms
15  eurecom-valbonne.r3t2.ft.net (193.48.50.54)  135 ms  128 ms  133 ms
16  194.214.211.25 (194.214.211.25)  126 ms  128 ms  126 ms
17  * * *
18  * * *
19  fantasia.eurecom.fr (193.55.113.142)  132 ms  128 ms  136 ms
```

trans-oceanic link

* means no response (probe lost, router not replying)

# Traceroute and ICMP

- Source sends series of UDP segments to dest
  - First has TTL =1
  - Second has TTL=2, etc.
  - Unlikely port number
- When nth datagram arrives to nth router:
  - Router discards datagram
  - And sends to source an ICMP message (type 11, code 0)
  - Message includes name of router& IP address

- When ICMP message arrives, source calculates RTT
- Traceroute does this 3 times

Stopping criterion

- UDP segment eventually arrives at destination host
- Destination returns ICMP "port unreachable" packet (type 3, code 3)
- When source gets this ICMP, stops.

# Chapter 4: Network Layer

# IPv6

❑ Initial motivation: 32-bit address space soon to be completely allocated.

❑ Additional motivation:

  ❖ header format helps speed processing/forwarding

  ❖ header changes to facilitate QoS

  IPv6 datagram format:

  ❖ fixed-length 40 byte header
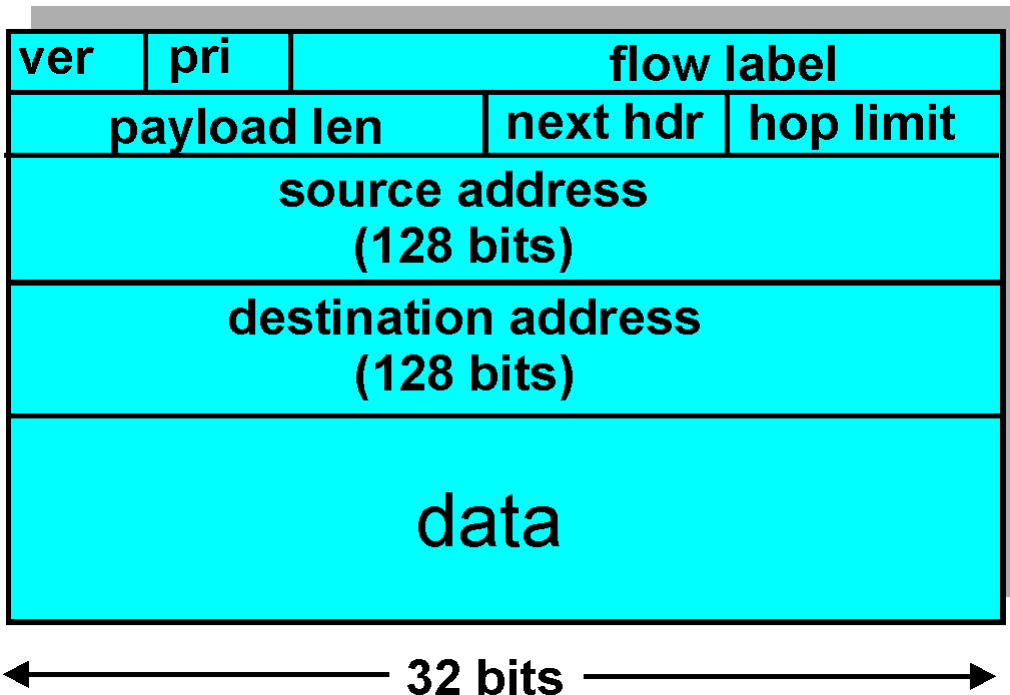
  ❖ no fragmentation allowed

# IP4 datagram format

IP protocol version number

header length (bytes)

"type" of data

max number remaining hops (decremented at each router)

upper layer protocol to deliver payload to

← 32 bits →

total datagram length (bytes)

for fragmentation/ reassembly

E.g. timestamp, record route taken, specify list of routers to visit.

| ver | head. len | type of service | length | |
|---|---|---|---|---|
| 16-bit identifier | | | flgs | fragment offset |
| time to live | upper layer | | header checksum | |
| 32 bit source IP address | | | | |
| 32 bit destination IP address | | | | |
| Options (if any) | | | | |
| data (variable length, typically a TCP or UDP segment) | | | | |

# IPv6 Header (Cont)

*Priority:* identify priority among datagrams in flow
*Flow Label:* identify datagrams in same "flow."
                    (concept of"flow" not well defined).
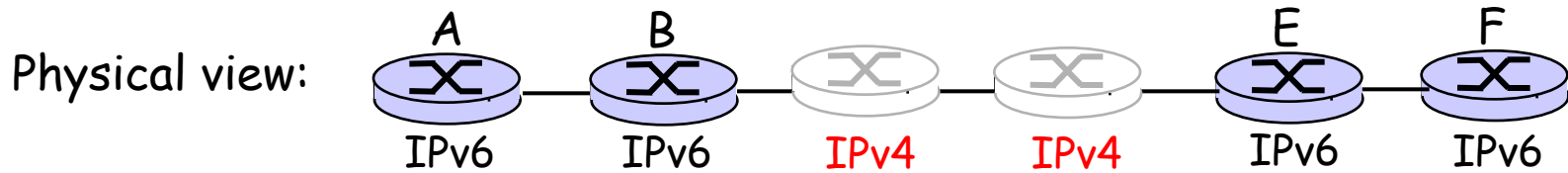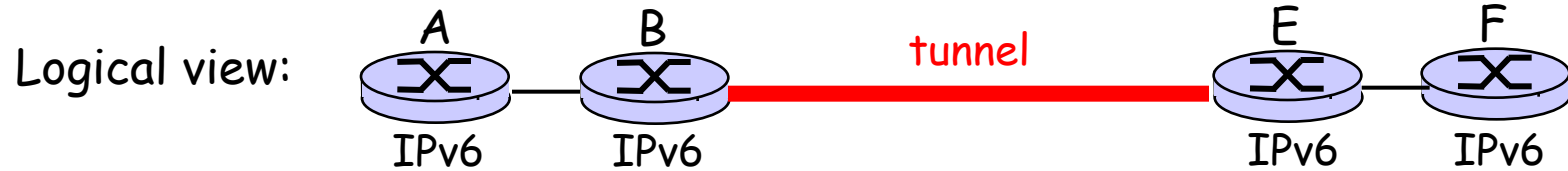*Next header:* identify upper layer protocol for data

| ver | pri | flow label | | |
|-----|-----|------------|---|---|
| payload len | | | next hdr | hop limit |
| source address (128 bits) | | | | |
| destination address (128 bits) | | | | |
| data | | | | |

← **32 bits** →

# Other Changes from IPv4

❑ *Checksum*: removed entirely to reduce processing time at each hop

❑ *Options:* allowed, but outside of header, indicated by "Next Header" field

❑ *ICMPv6:* new version of ICMP

  ❖ additional message types, e.g. "Packet Too Big"
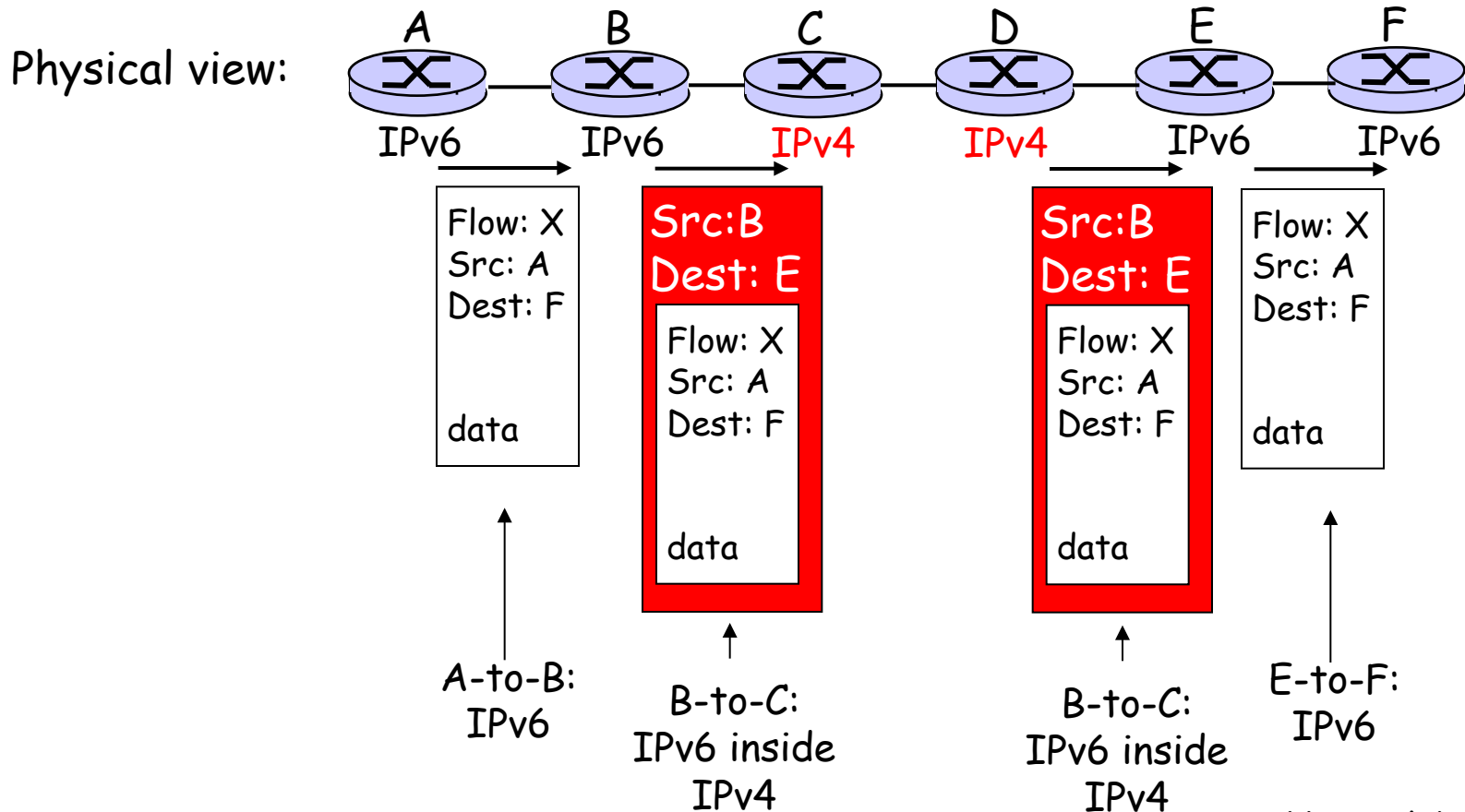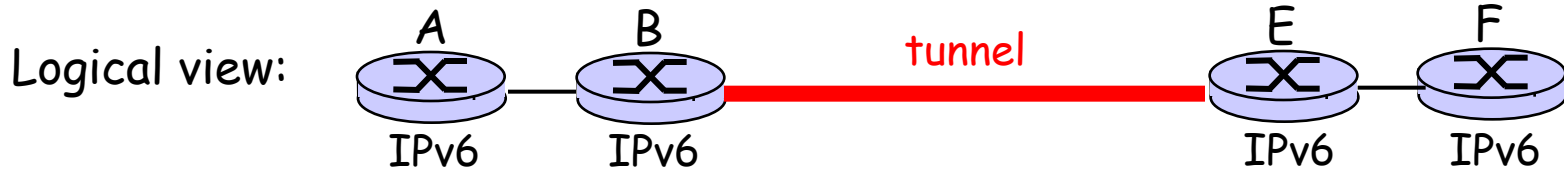
  ❖ multicast group management functions

# Transition From IPv4 To IPv6

❑ Not all routers can be upgraded simultaneous
- ❖ no "flag days"
- ❖ How will the network operate with mixed IPv4 and IPv6 routers?

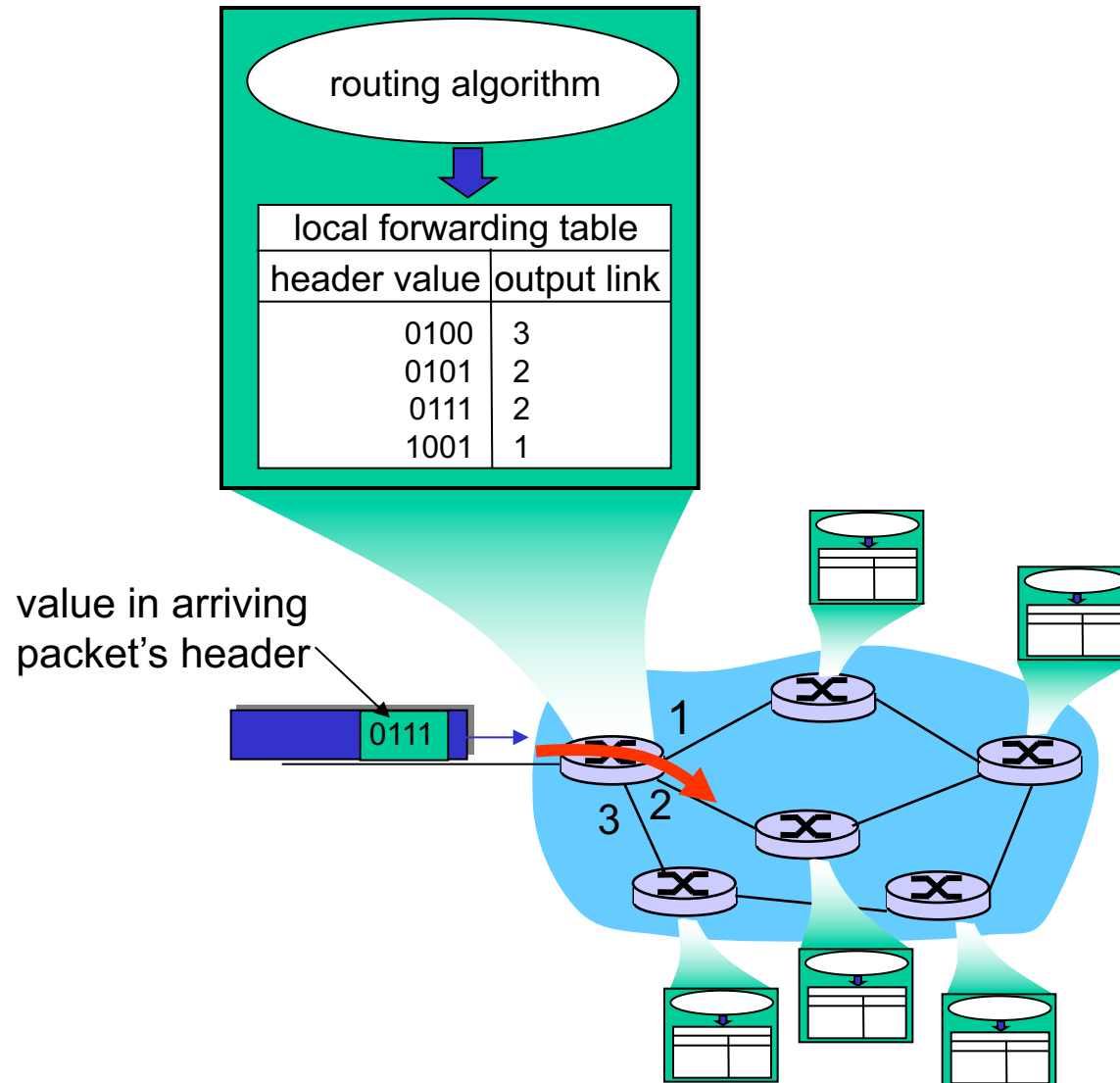❑ *Tunneling:* IPv6 carried as payload in IPv4 datagram among IPv4 routers

# Tunneling

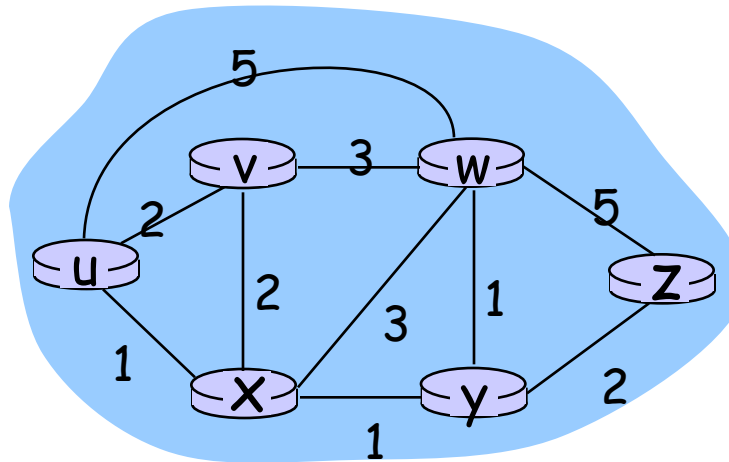Logical view:

A      B          tunnel        E     F

IPv6    IPv6                IPv6    IPv6

Physical view:

A      B                  E     F

IPv6    IPv6    IPv4    IPv4    IPv6    IPv6

# Tunneling

Logical view:

A — B === tunnel === E — F

IPv6    IPv6                    IPv6    IPv6

Physical view:

A — B — C — D — E — F

IPv6   IPv6   IPv4   IPv4   IPv6   IPv6

| Flow: X<br>Src: A<br>Dest: F<br><br>data | Src:B<br>Dest: E<br>Flow: X<br>Src: A<br>Dest: F<br><br>data | Src:B<br>Dest: E<br>Flow: X<br>Src: A<br>Dest: F<br><br>data | Flow: X<br>Src: A<br>Dest: F<br><br>data |

A-to-B:
IPv6

B-to-C:
IPv6 inside
IPv4

B-to-C:
IPv6 inside
IPv4

E-to-F:
IPv6

# Chapter 4: Network Layer

# Interplay between routing, forwarding



routing algorithm

| local forwarding table | |
|---|---|
| header value | output link |
| 0100 | 3 |
| 0101 | 2 |
| 0111 | 2 |
| 1001 | 1 |

value in arriving packet's header
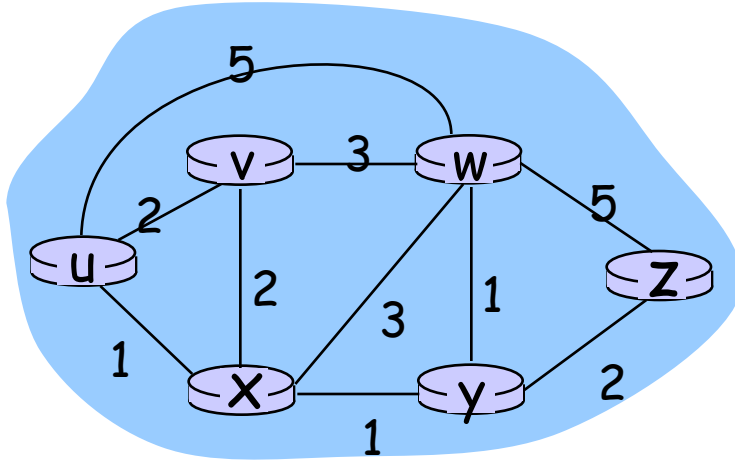
0111

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) }

Remark: Graph abstraction is useful in other network contexts

Example: P2P, where N is set of peers and E is set of TCP connections

# 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 $(x_1, x_2, x_3,..., x_p) = c(x_1,x_2) + c(x_2,x_3) + ... + c(x_{p-1},x_p)$

Question: What's the least-cost path between u and z ?

Routing algorithm: algorithm that finds least-cost path

# Routing Algorithm classification

**Global or decentralized information?**

Global:

- all routers have complete topology, link cost info
- "link state" algorithms

Decentralized:

- router knows physically-connected neighbors, link costs to neighbors
- iterative process of computation, exchange of info with neighbors
- "distance vector" algorithms

**Static or dynamic?**

Static:

- routes change slowly over time

Dynamic:

- routes change more quickly
  - ❖ periodic update
  - ❖ in response to link cost changes

# Chapter 4: Network Layer

- ❑ 4. 1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - ❖ Datagram format
  - ❖ IPv4 addressing
  - ❖ ICMP
  - ❖ IPv6

- ❑ 4.5 Routing algorithms
  - ❖ Link state
  - ❖ Distance Vector
  - ❖ Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - ❖ RIP
  - ❖ OSPF
  - ❖ BGP
- ❑ 4.7 Broadcast and multicast routing

# A Link-State Routing Algorithm

## Dijkstra's algorithm

❑ net topology, link costs known to all nodes
  ❖ accomplished via "link state broadcast"
  ❖ all nodes have same info
❑ computes least cost paths from one node ('source") to all other nodes
  ❖ gives forwarding table for that node
❑ iterative: after k iterations, know least cost path to k dest.'s

## Notation:

❑ c(x,y): link cost from node x to y;  = ∞ if not direct neighbors
❑ D(v): current value of cost of path from source to dest. v
❑ p(v): predecessor node along path from source to v
❑ N': set of nodes whose least cost path definitively known

# Dijkstra's Algorithm

```
1  Initialization:
2     N' = {u}
3   for all nodes v
4      if v adjacent to u
5          then D(v) = c(u,v)
6      else D(v) = ∞
7
8   Loop
9     find w not in N' such that D(w) is a minimum
10     add w to N'
11     update D(v) for all v adjacent to w and not in N' :
12        D(v) = min( D(v), D(w) + c(w,v) )
13     /* new cost to v is either old cost to v or known
14       shortest path cost to w plus cost from w to v */
15  until all nodes in N'
```
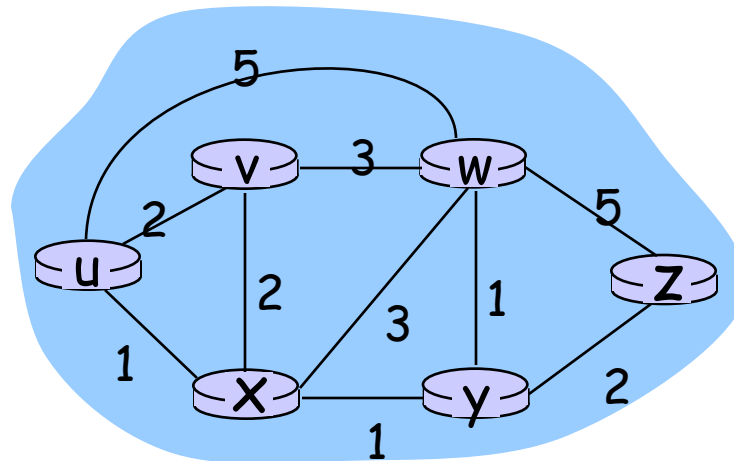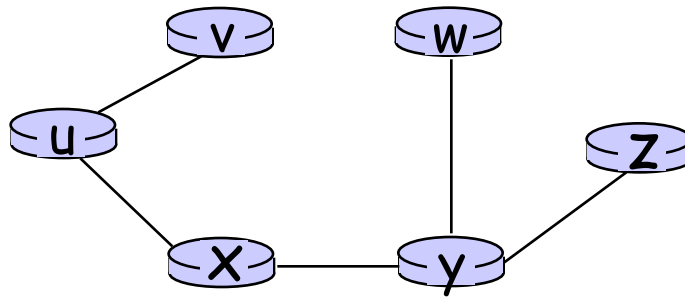
# Dijkstra's algorithm: example

| Step | N' | D(v),p(v) | D(w),p(w) | D(x),p(x) | D(y),p(y) | D(z),p(z) |
|------|-------|-----------|-----------|-----------|-----------|-----------|
| 0 | u | 2,u | 5,u | 1,u | ∞ | ∞ |
| 1 | ux | 2,u | 4,x | | 2,x | ∞ |
| 2 | uxy | 2,u | 3,y | | | 4,y |
| 3 | uxyv | | 3,y | | | 4,y |
| 4 | uxyvw | | | | | 4,y |
| 5 | uxyvwz | | | | | |

# Dijkstra's algorithm: example (2)

Resulting shortest-path tree from u:



Resulting forwarding table in u:

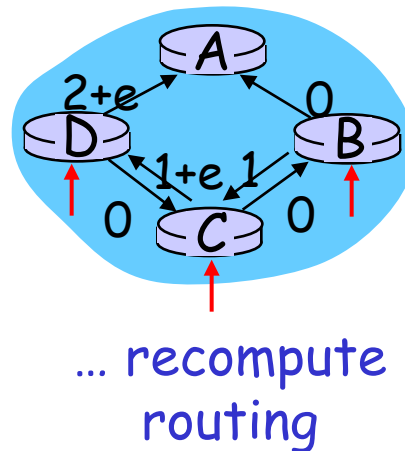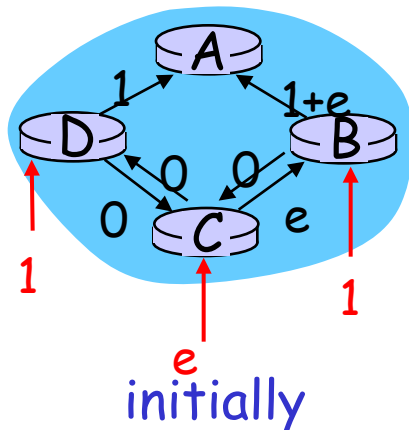| destination | link |
|---|---|
| v | (u,v) |
| x | (u,x) |
| y | (u,x) |
| w | (u,x) |
| z | (u,x) |

# Dijkstra's algorithm, discussion

Algorithm complexity: n nodes

- ❑ each iteration: need to check all nodes, w, not in N
- ❑ $n(n+1)/2$ comparisons: $O(n^2)$
- ❑ more efficient implementations possible: $O(n \log n)$

Oscillations possible:

- ❑ e.g., link cost = amount of carried traffic



... recompute routing    ... recompute    ... recompute

initially

# Chapter 4: Network Layer

❑ 4. 1 Introduction

❑ 4.2 Virtual circuit and datagram networks

❑ 4.3 What's inside a router

❑ 4.4 IP: Internet Protocol
  ❖ Datagram format
  ❖ IPv4 addressing
  ❖ ICMP
  ❖ IPv6

❑ 4.5 Routing algorithms
  ❖ Link state
  ❖ Distance Vector
  ❖ Hierarchical routing

❑ 4.6 Routing in the Internet
  ❖ RIP
  ❖ OSPF
  ❖ BGP

❑ 4.7 Broadcast and multicast routing

# Distance Vector Algorithm

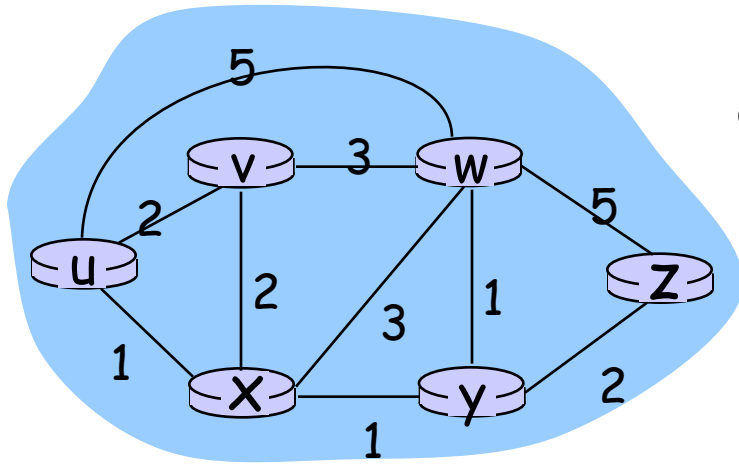Bellman-Ford Equation (dynamic programming)

Define

$d_u(y) :=$ cost of least-cost path from u to y

Then

$$d_u(y) = \min_v \{c(u,v) + d_v(y)\}$$

where min is taken over all neighbors v of u

# Bellman-Ford example



Clearly, $d_v(z) = 5$, $d_x(z) = 3$, $d_w(z) = 3$

B-F equation says:

$$d_u(z) = \min \{ c(u,v) + d_v(z),$$
$$c(u,x) + d_x(z),$$
$$c(u,w) + d_w(z) \}$$
$$= \min \{2 + 5,$$
$$1 + 3,$$
$$5 + 3\} = 4$$

Node that achieves minimum is next
hop in shortest path ➜ forwarding table

# Distance Vector Algorithm

- ❑ $D_x(y)$ = estimate of least cost from x to y
- ❑ Node x knows cost to each neighbor v: $c(x,v)$
- ❑ Node x maintains its own distance vector $\mathbf{D_x} = [D_x(y): y \in N\ ]$
- ❑ Node x also maintains its neighbors' distance vectors
  - ❖ For each neighbor v, x maintains $\mathbf{D_v} = [D_v(y): y \in N\ ]$

# Distance vector algorithm (4)

**Basic idea:**

❑ From time-to-time, each node sends its own distance vector estimate to neighbors

❑ Asynchronous

❑ When a node x receives new DV estimate from neighbor, it updates its own DV using B-F equation:

$$D_x(y) \leftarrow min_v\{c(x,v) + D_v(y)\} \quad \text{for each node } y \in N$$

❑ Under minor, natural conditions, the estimate $D_x(y)$ converge to the actual least cost $d_x(y)$

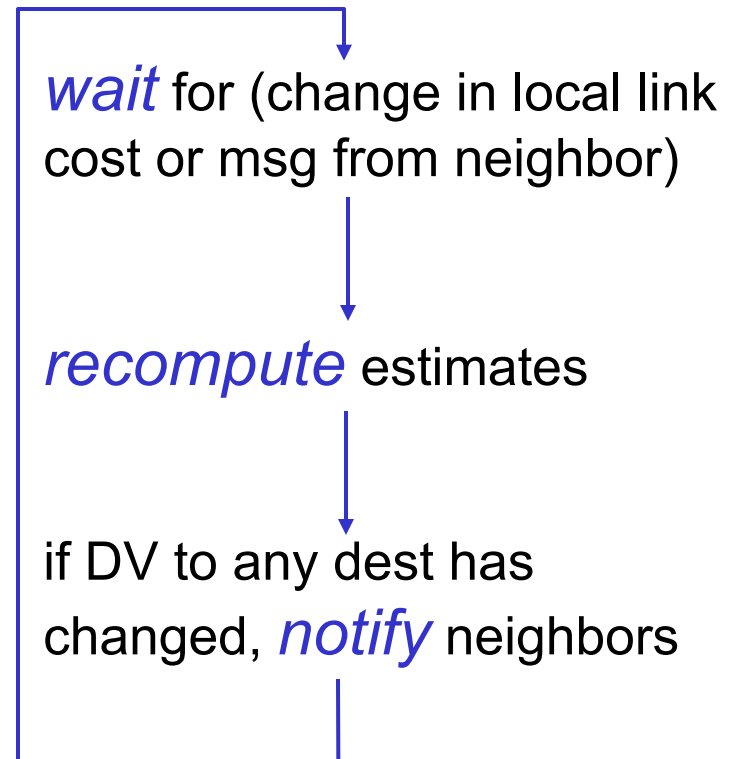# Distance Vector Algorithm (5)

## Iterative, asynchronous:
each local iteration caused by:

❑ local link cost change
❑ DV update message from neighbor

## Distributed:
❑ each node notifies neighbors *only* when its DV changes
   ❖ neighbors then notify their neighbors if necessary

## Each node:

*wait* for (change in local link cost or msg from neighbor)

*recompute* estimates

if DV to any dest has changed, *notify* neighbors

$$D_x(y) = \min\{c(x,y) + D_y(y), c(x,z) + D_z(y)\}$$
$$= \min\{2+0, 7+1\} = 2$$

$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\}$$
$$= \min\{2+1, 7+0\} = 3$$

**node x table**

cost to

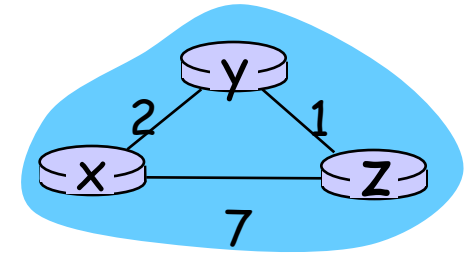|   | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 7 |
| y | ∞ | ∞ | ∞ |
| z | ∞ | ∞ | ∞ |

from

cost to

|   | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 3 |
| y | 2 | 0 | 1 |
| z | 7 | 1 | 0 |

from

**node y table**

cost to

|   | x | y | z |
|---|---|---|---|
| x | ∞ | ∞ | ∞ |
| y | 2 | 0 | 1 |
| z | ∞ | ∞ | ∞ |

from

**node z table**

cost to

|   | x | y | z |
|---|---|---|---|
| x | ∞ | ∞ | ∞ |
| y | ∞ | ∞ | ∞ |
| z | 7 | 1 | 0 |

from

time

$$D_x(y) = \min\{c(x,y) + D_y(y), c(x,z) + D_z(y)\}$$
$$= \min\{2+0, 7+1\} = 2$$

$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\}$$
$$= \min\{2+1, 7+0\} = 3$$

**node x table**

cost to

|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 7 |
| y | ∞ | ∞ | ∞ |
| z | ∞ | ∞ | ∞ |

from

cost to

|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 3 |
| y | 2 | 0 | 1 |
| z | 7 | 1 | 0 |

from

cost to

|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 3 |
| y | 2 | 0 | 1 |
| z | 3 | 1 | 0 |

from

**node y table**

cost to

|  | x | y | z |
|---|---|---|---|
| x | ∞ | ∞ | ∞ |
| y | 2 | 0 | 1 |
| z | ∞ | ∞ | ∞ |

from

cost to

|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 7 |
| y | 2 | 0 | 1 |
| z | 7 | 1 | 0 |

from

cost to

|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 3 |
| y | 2 | 0 | 1 |
| z | 3 | 1 | 0 |

from

**node z table**

cost to

|  | x | y | z |
|---|---|---|---|
| x | ∞ | ∞ | ∞ |
| y | ∞ | ∞ | ∞ |
| z | 7 | 1 | 0 |

from

cost to

|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 7 |
| y | 2 | 0 | 1 |
| z | 3 | 1 | 0 |

from

cost to

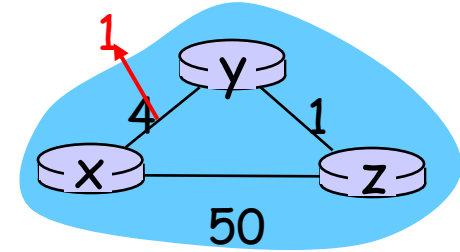|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 3 |
| y | 2 | 0 | 1 |
| z | 3 | 1 | 0 |

from

time

# Distance Vector: link cost changes

**Link cost changes:**

- ❑ node detects local link cost change
- ❑ updates routing info, recalculates distance vector
- ❑ if DV changes, notify neighbors



**"good news travels fast"**

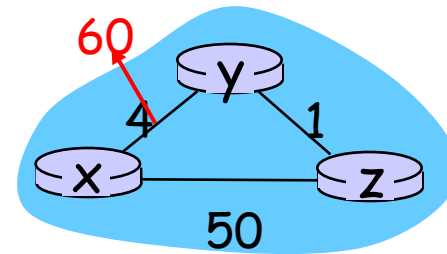At time $t_0$, $y$ detects the link-cost change, updates its DV, and informs its neighbors.

At time $t_1$, $z$ receives the update from $y$ and updates its table. It computes a new least cost to $x$ and sends its neighbors its DV.

At time $t_2$, $y$ receives $z$'s update and updates its distance table. $y$'s least costs do not change and hence $y$ does *not* send any message to $z$.

# Distance Vector: link cost changes

Link cost changes:

❑ good news travels fast

❑ bad news travels slow - "count to infinity" problem!

❑ 44 iterations before algorithm stabilizes: see text

Poisoned reverse:

❑ If Z routes through Y to get to X :

  ❑ Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)

❑ will this completely solve count to infinity problem?

# Comparison of LS and DV algorithms

## Message complexity

- **LS:** with n nodes, E links, O(nE) msgs sent
- **DV:** exchange between neighbors only
  - ❖ convergence time varies

## Speed of Convergence

- **LS:** O(n²) algorithm requires O(nE) msgs
  - ❖ may have oscillations
- **DV**: convergence time varies
  - ❖ may be routing loops
  - ❖ count-to-infinity problem

## Robustness: what happens if router malfunctions?

## LS:

- ❖ node can advertise incorrect *link* cost
- ❖ each node computes only its *own* table

## DV:

- ❖ DV node can advertise incorrect *path* cost
- ❖ each node's table used by others
  - • error propagate thru network

# Chapter 4: Network Layer

# Hierarchical Routing

Our routing study thus far - idealization
- ❑ all routers identical
- ❑ network "flat"

... *not* true in practice

**scale:** with 200 million destinations:
- ❑ can't store all dest's in routing tables!
- ❑ routing table exchange would swamp links!

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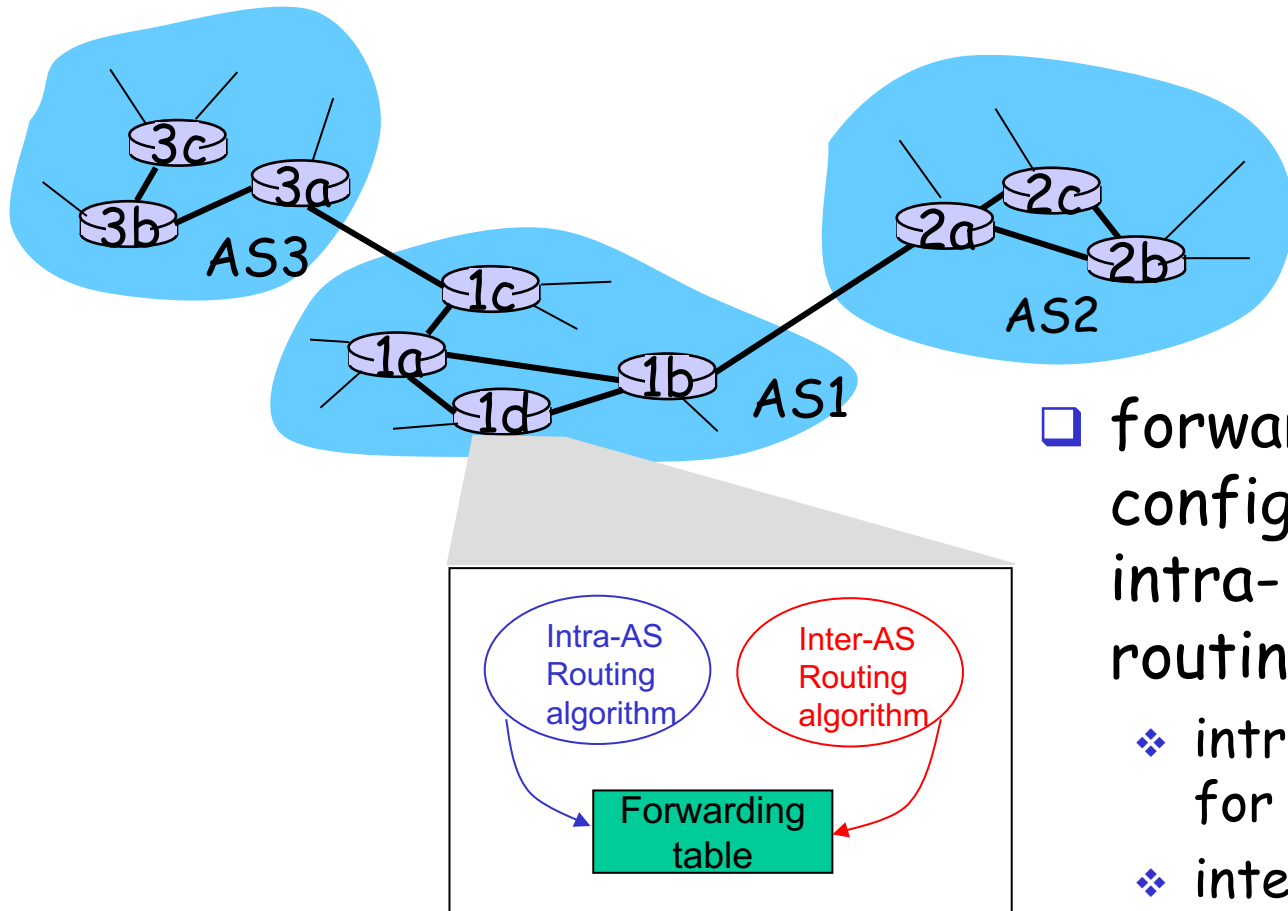- ❑ each network admin may want to control routing in its own network

# 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

□ Direct 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
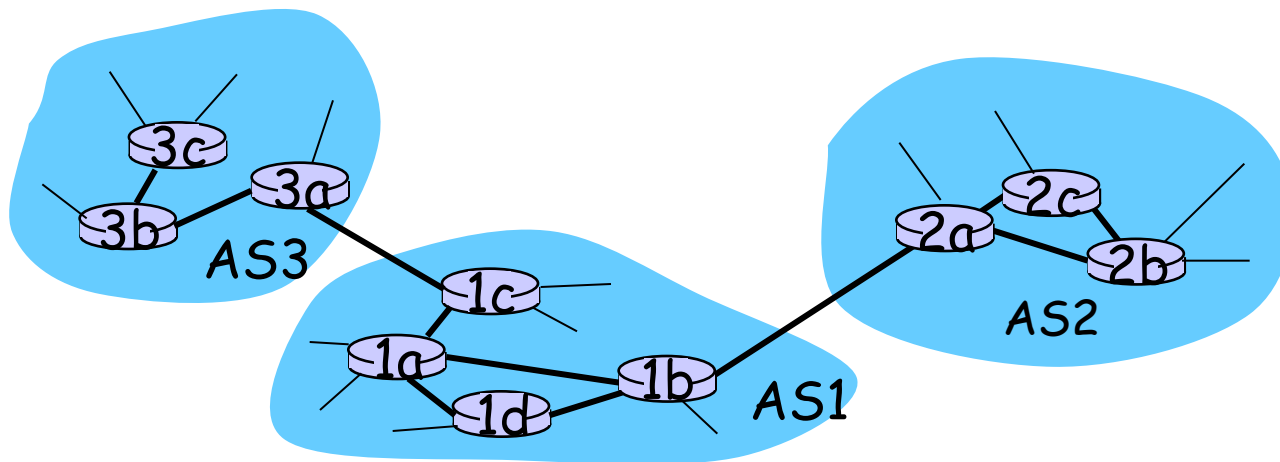- ❖ inter-AS & intra-As sets entries for external dests

# Inter-AS tasks

- suppose router in AS1 receives datagram destined outside of AS1:
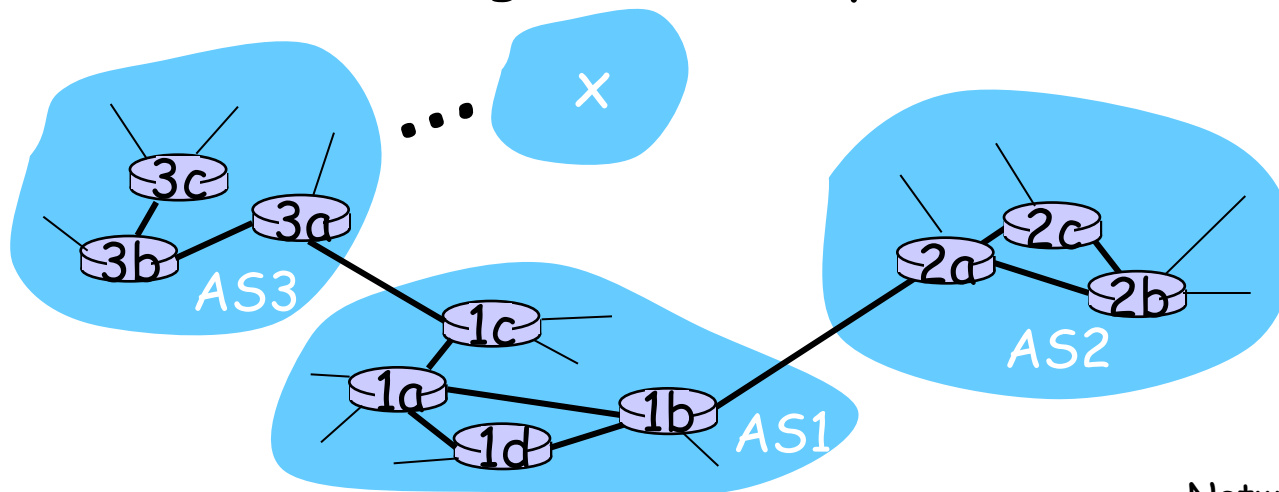  - router should forward packet to gateway router, but which one?

AS1 must:

1. learn which dests are reachable through AS2, which through AS3

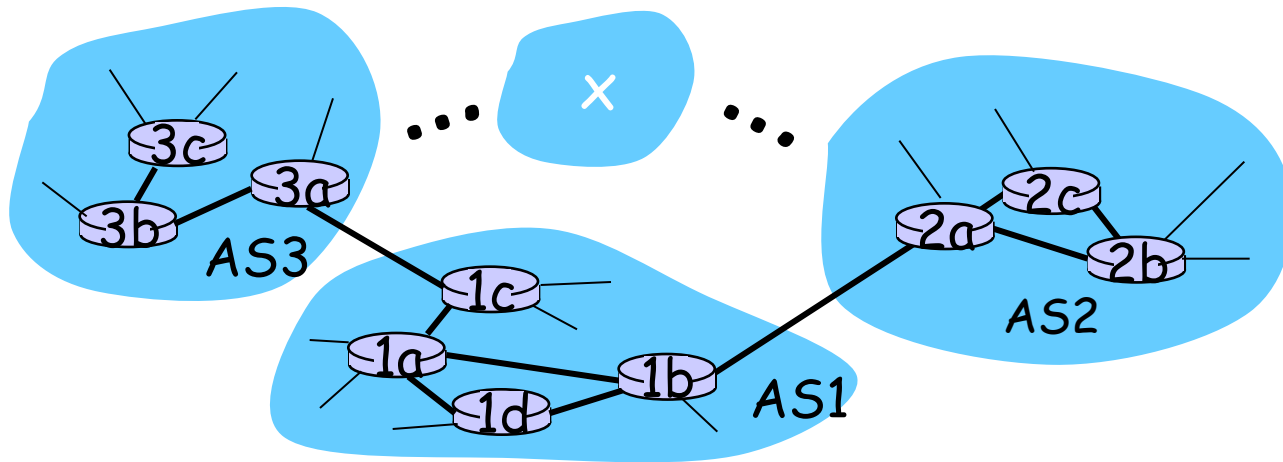2. propagate this reachability info to all routers in AS1

Job of inter-AS routing!

# Example: Setting forwarding table in router 1d

❑ suppose AS1 learns (via inter-AS protocol) that subnet *x* reachable via AS3 (gateway 1c) but not via AS2.

❑ inter-AS protocol propagates reachability info to all internal routers.

❑ router 1d determines from intra-AS routing info that its interface *I* is on the least cost path to 1c.
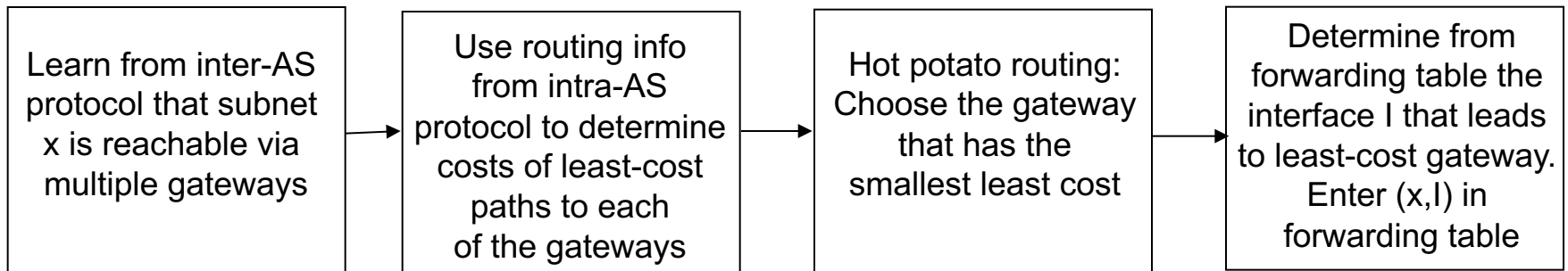
  ❖ installs forwarding table entry *(x,I)*

# Example: Choosing among multiple ASes

❑ now suppose AS1 learns from inter-AS protocol that subnet *x* is reachable from AS3 *and* from AS2.

❑ to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest *x*.

  ❖ this is also job of inter-AS routing protocol!

# Example: Choosing among multiple ASes

❑ now suppose AS1 learns from inter-AS protocol that subnet $x$ is reachable from AS3 *and* from AS2.

❑ to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest $x$.
  ❖ this is also job of inter-AS routing protocol!

❑ hot potato routing: send packet towards closest of two routers.

| Learn from inter-AS protocol that subnet x is reachable via multiple gateways | Use routing info from intra-AS protocol to determine costs of least-cost paths to each of the gateways | Hot potato routing: Choose the gateway that has the smallest least cost | Determine from forwarding table the interface I that leads to least-cost gateway. Enter (x,I) in forwarding table |
|---|---|---|---|

# Chapter 4: Network Layer

- 4. 1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6

- 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- 4.7 Broadcast and multicast routing

# 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)
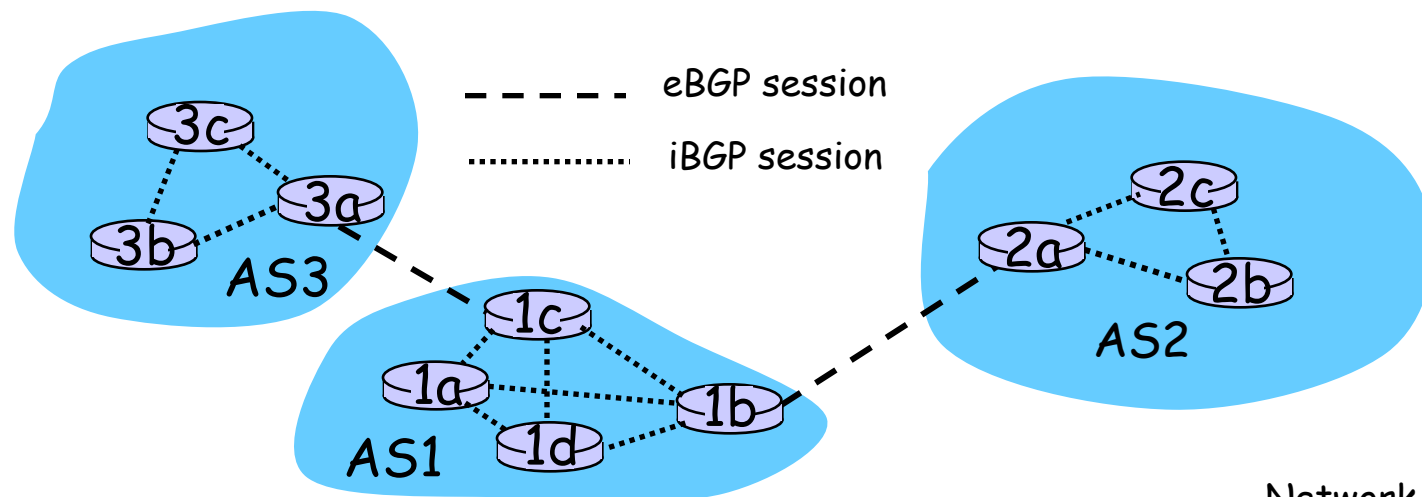
# Chapter 4: Network Layer

# Internet inter-AS routing: BGP

❑ BGP (Border Gateway Protocol): *the* de facto standard

❑ BGP provides each AS a means to:
   1. Obtain subnet reachability information from neighboring ASs.
   2. Propagate reachability information to all AS-internal routers.
   3. Determine "good" routes to subnets based on reachability information and policy.

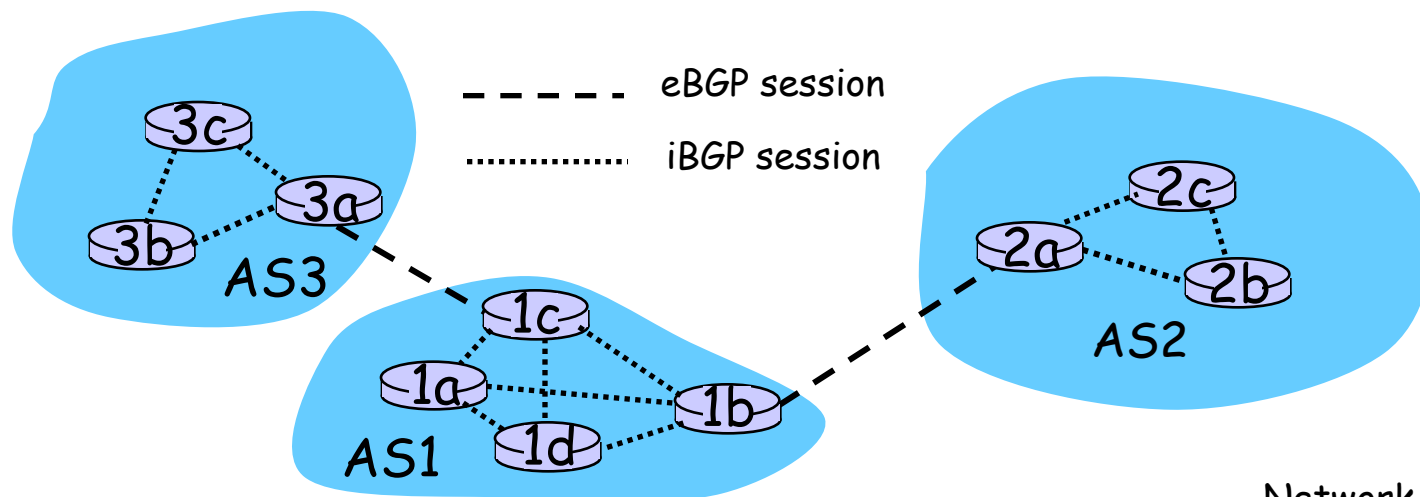❑ allows subnet to advertise its existence to rest of Internet: *"I am here"*

# BGP basics

❑ pairs of routers (BGP peers) exchange routing info over semi-permanent TCP connections: <span style="color:red">BGP sessions</span>
  ❖ BGP sessions need not correspond to physical links.
❑ when AS2 advertises a prefix to AS1:
  ❖ AS2 <span style="color:red">*promises*</span> it will forward datagrams towards that prefix.
  ❖ AS2 can aggregate prefixes in its advertisement

- - - - - eBGP session

............... iBGP session

3c
3a
3b
AS3

2c
2a
2b
AS2

1c
1a
1b
1d
AS1

# Distributing reachability info

❑ using eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.

  ❖ 1c can then use iBGP do distribute new prefix info to all routers in AS1

  ❖ 1b can then re-advertise new reachability info to AS2 over 1b-to-2a eBGP session

❑ when router learns of new prefix, it creates entry for prefix in its forwarding table.



eBGP session

iBGP session

3c

3a

3b

AS3

1c

1a

1b

1d

AS1

2c

2a

2b

AS2

# Path attributes & BGP routes

❑ advertised prefix includes BGP attributes.
   ❖ prefix + attributes = "route"
❑ two important attributes:
   ❖ AS-PATH: contains ASs through which prefix advertisement has passed: e.g, AS 67, AS 17
   ❖ NEXT-HOP: indicates specific internal-AS router to next-hop AS. (may be multiple links from current AS to next-hop-AS)
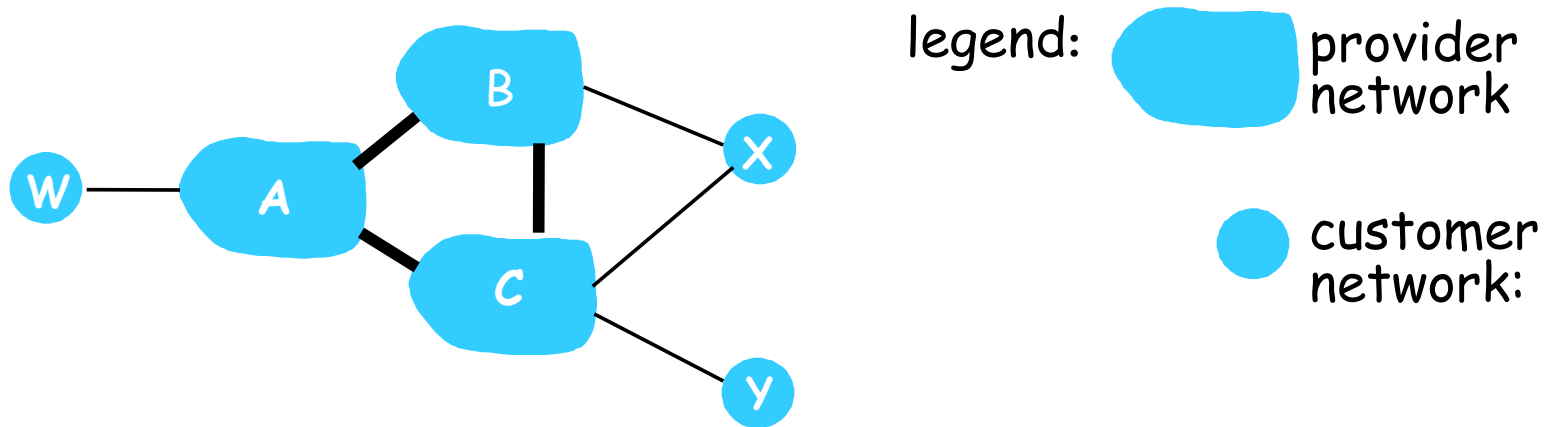❑ when gateway router receives route advertisement, uses import policy to accept/decline.

# BGP route selection

❑ router may learn about more than 1 route to some prefix. Router must select route.

❑ elimination rules:

1. local preference value attribute: policy decision
2. shortest AS-PATH
3. closest NEXT-HOP router: hot potato routing
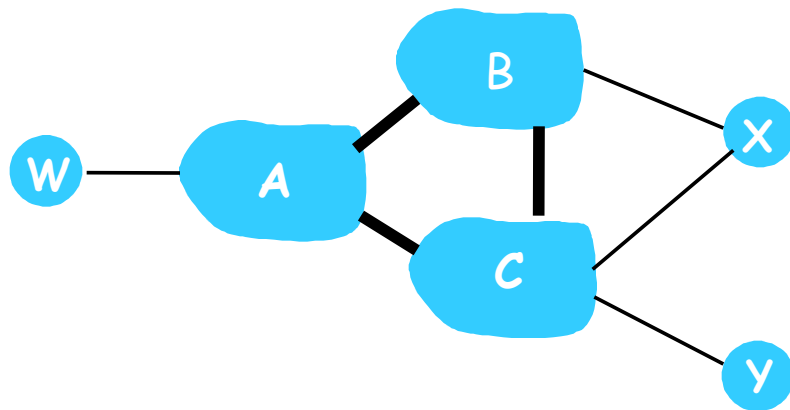4. additional criteria

# BGP messages

❑ BGP messages exchanged using TCP.

❑ BGP messages:
  ❖ OPEN: opens TCP connection to peer and authenticates sender
  ❖ UPDATE: advertises new path (or withdraws old)
  ❖ KEEPALIVE keeps connection alive in absence of UPDATES; also ACKs OPEN request
  ❖ NOTIFICATION: reports errors in previous msg; also used to close connection

# BGP routing policy



legend:

provider network

customer network:

❑ A,B,C are provider networks

❑ X,W,Y are customer (of provider networks)

❑ X is dual-homed: attached to two networks

  ❑ X does not want to route from B via X to C

  ❑ .. so X will not advertise to B a route to C

# BGP routing policy (2)



legend:

provider
network

customer
network:

❑ A advertises path AW  to B
❑ B advertises path BAW to X
❑ Should B advertise path BAW to C?
- ❑ No way! B gets no "revenue" for routing CBAW since neither W nor C are B's customers
- ❑ B wants to force C to route to w via A
- ❑ B wants to route *only* to/from its customers!

# Why different Intra- and Inter-AS routing ?

## Policy:

- ❑ Inter-AS: admin wants control over how its traffic routed, who routes through its net.
- ❑ Intra-AS: single admin, so no policy decisions needed

## Scale:

- ❑ hierarchical routing saves table size, reduced update traffic

## Performance:

- ❑ Intra-AS: can focus on performance
- ❑ Inter-AS: policy may dominate over performance

# Chapter 4: summary

- 4. 1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- 4.7 Broadcast and multicast routing

# Exercise questions

❑ What is the difference between forwarding and routing?

❑ Do the routers in both datagram networks and virtual-circuit networks use forwarding tables?

❑ What is the 32-bit binary equivalent of the IP address 128.64.0.1?

❑ Define the following terms: subnet, prefix.

# Exercise questions

❑ How does BGP use the AS-PATH attribute?

❑ Compare and contrast link-state and distance-vector routing algorithms.