# CSB051 –Computer Networks
## 電腦網路

# Chapter 4
# Network Layer
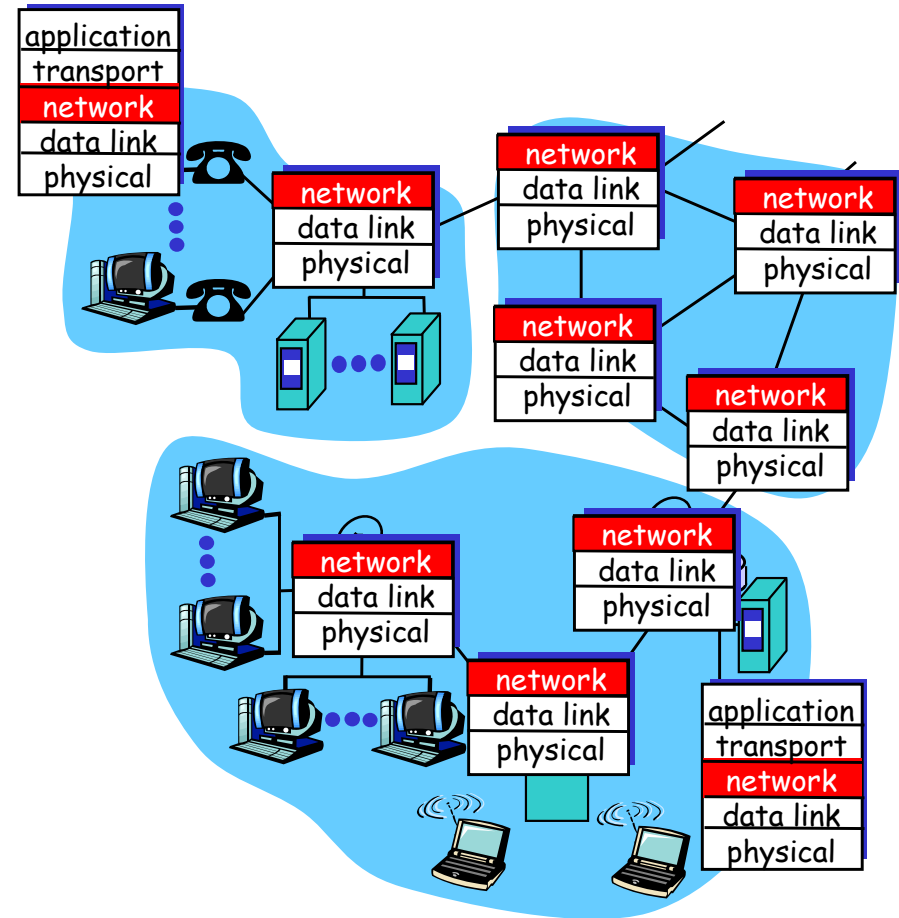
吳俊興

國立高雄大學 資訊工程學系

# 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

# 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*

analogy:

- routing: process of planning trip from source to dest

- forwarding: process of getting through single interchange

# Chapter 4: Network Layer

# Network layer connection and connection-less service
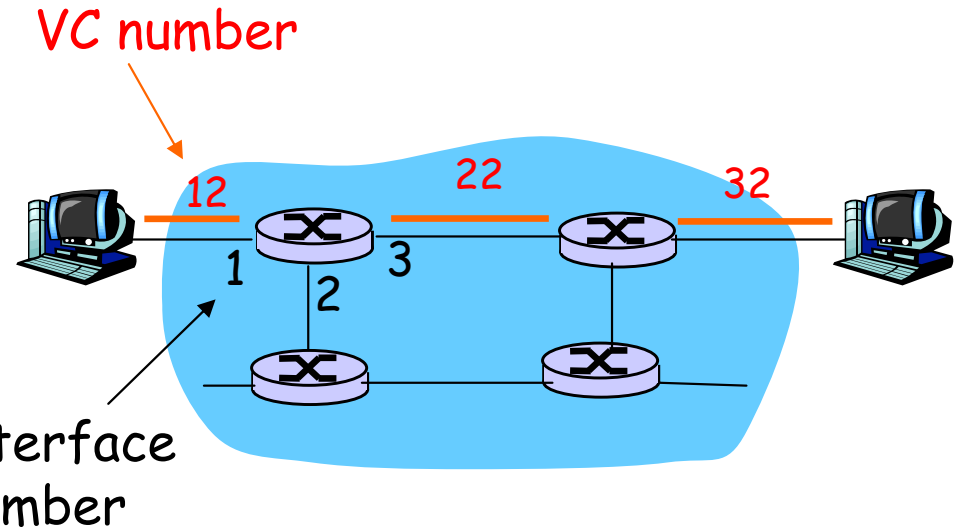
☐ Datagram network provides network-layer connectionless service

☐ 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 the 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

# Forwarding table

VC number



interface number

## Forwarding table in northwest router:

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

Routers maintain connection state information!

# 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<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 |

# 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?

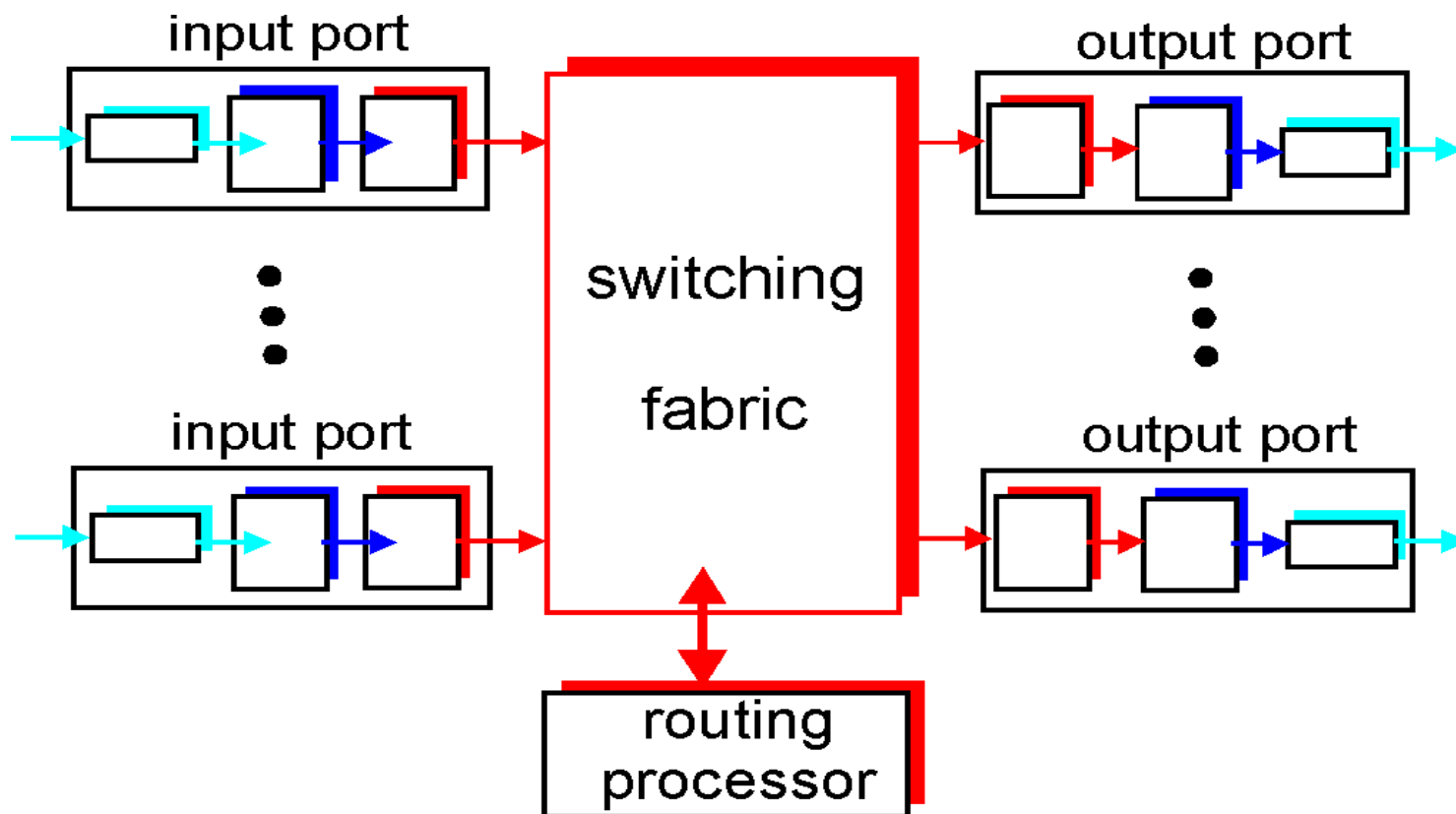DA: 11001000  00010111  00011000  10101010          Which interface?

# 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
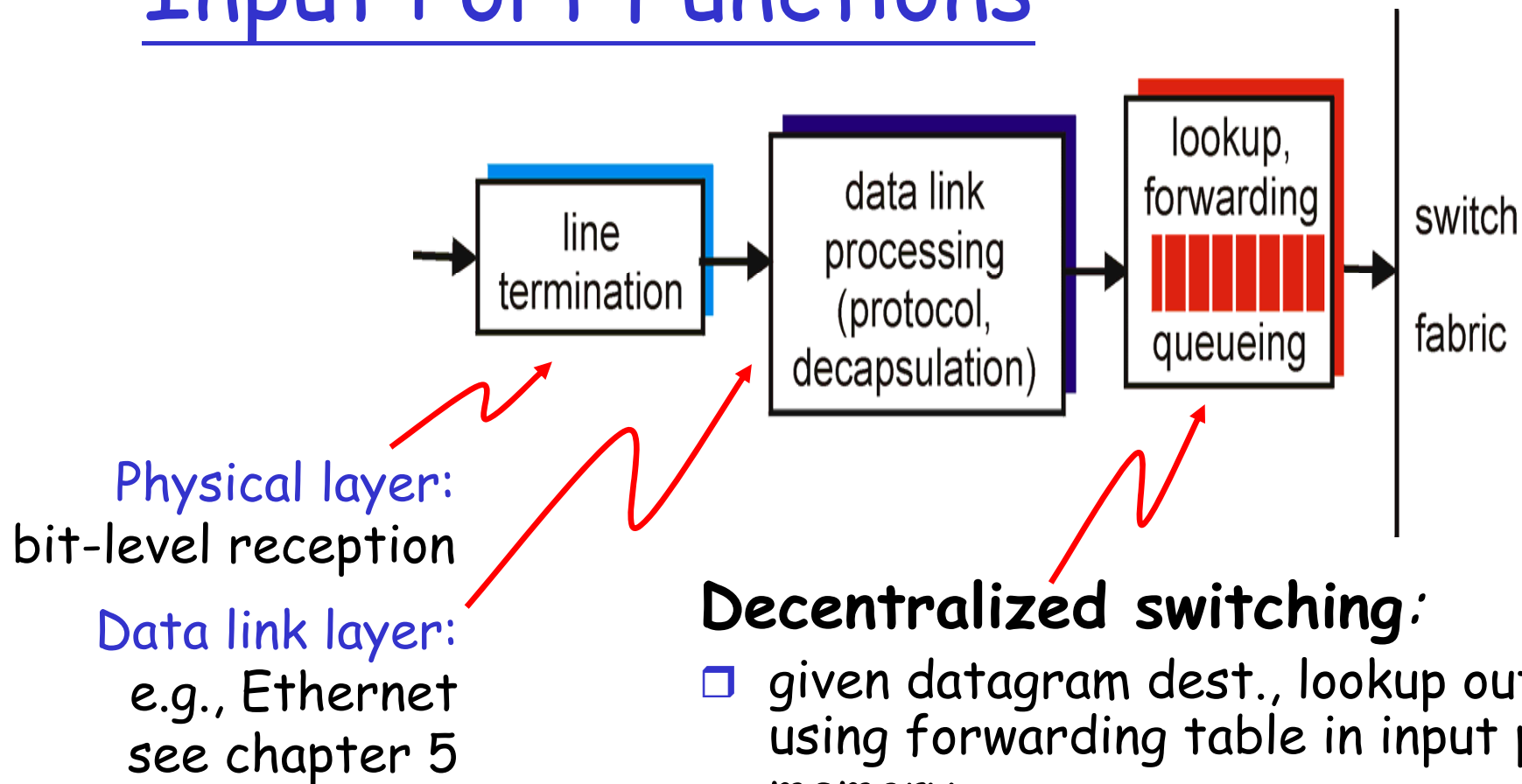  - OSPF
  - BGP
- 4.7 Broadcast and multicast routing

# Router Architecture Overview

Two key router functions:

- run routing algorithms/protocol (RIP, OSPF, BGP)
- *forwarding* datagrams from incoming to outgoing link

# Input Port Functions



Physical layer:
bit-level reception

Data link layer:
e.g., Ethernet
see chapter 5

**Decentralized switching:**

- given datagram dest., lookup output port using forwarding table in input port memory
- goal: complete input port processing at 'line speed'
- queuing: if datagrams arrive faster than forwarding rate into switch fabric

# Three types of switching fabrics

# Output Ports



□ *Buffering* required when datagrams arrive from fabric faster than the transmission rate

□ *Scheduling discipline* chooses among queued datagrams for transmission

# Chapter 4: Network Layer

# 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

# 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 (1:ICMP, 6:TCP, 17:UDP)

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 | | Internet 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<br>=4000 | ID<br>=x | fragflag<br>=0 | offset<br>=0 | |
|---|---|---|---|---|---|

**Example**

- ☐ 4000 byte datagram
- ☐ MTU = 1500 bytes

One large datagram becomes several smaller datagrams

| | length<br>=1500 | ID<br>=x | fragflag<br>=1 | offset<br>=0 | |
|---|---|---|---|---|---|

| | length<br>=1500 | ID<br>=x | fragflag<br>=1 | offset<br>=185 | |
|---|---|---|---|---|---|

| | length<br>=1040 | ID<br>=x | fragflag<br>=0 | offset<br>=370 | |
|---|---|---|---|---|---|

1480 bytes in data field

offset = 1480/8

# 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 may have multiple interfaces

  ○ IP addresses associated with each interface

223.1.1.1

223.1.2.1

223.1.1.2

223.1.1.4   223.1.2.9

223.1.2.2

223.1.1.3   223.1.3.27

223.1.3.1       223.1.3.2

223.1.1.1 = 11011111 00000001 00000001 00000001

      223         1          1          1
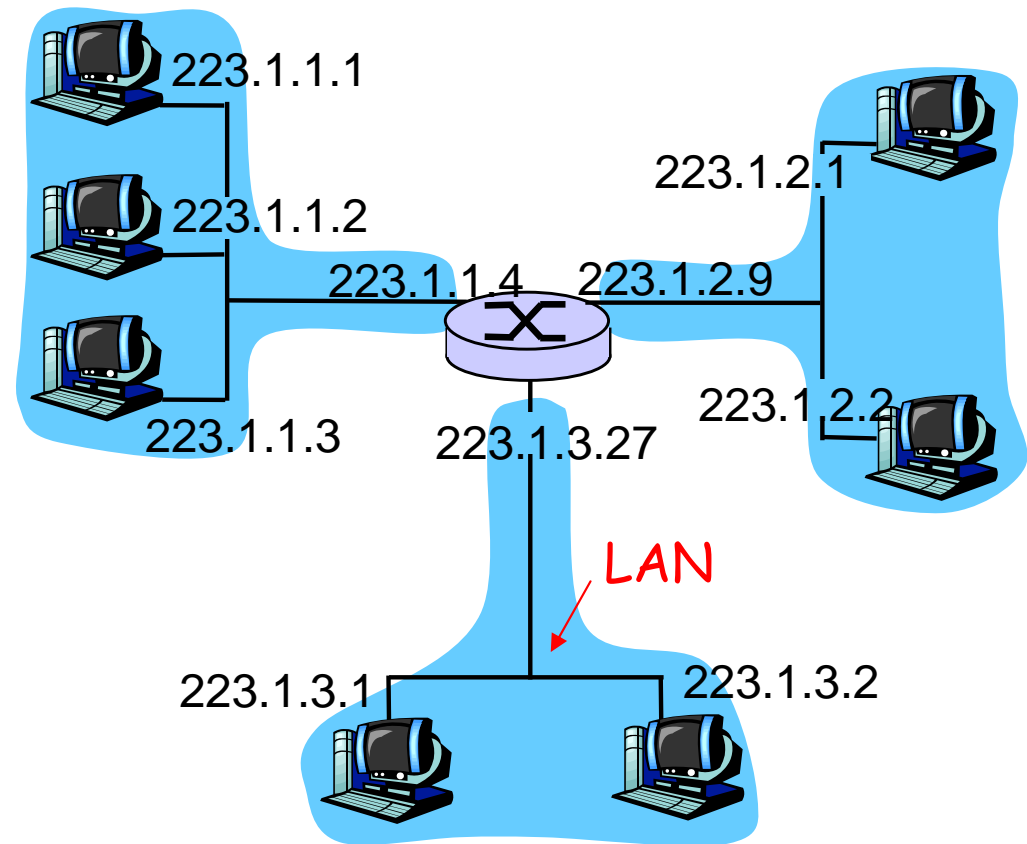
# Subnets

- 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

223.1.2.1

223.1.1.4    223.1.2.9

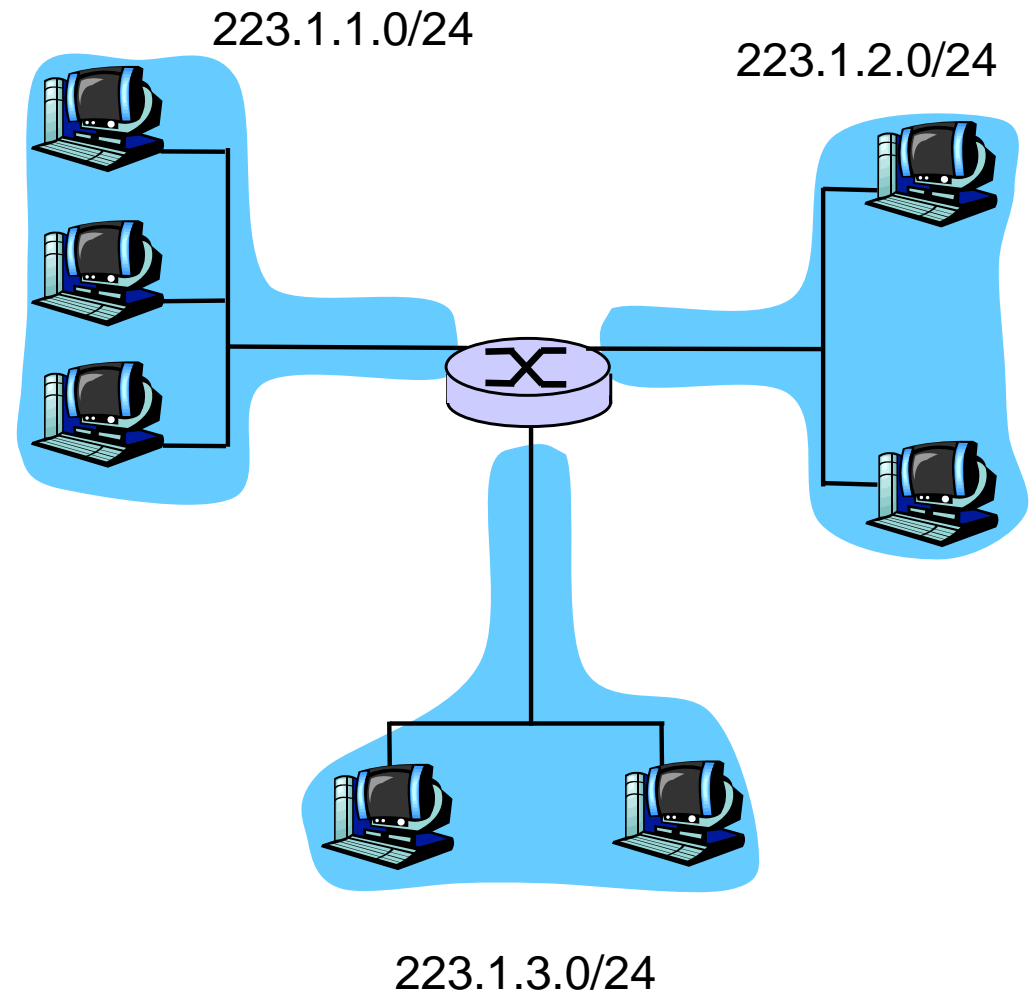223.1.2.2

223.1.1.3    223.1.3.27

LAN

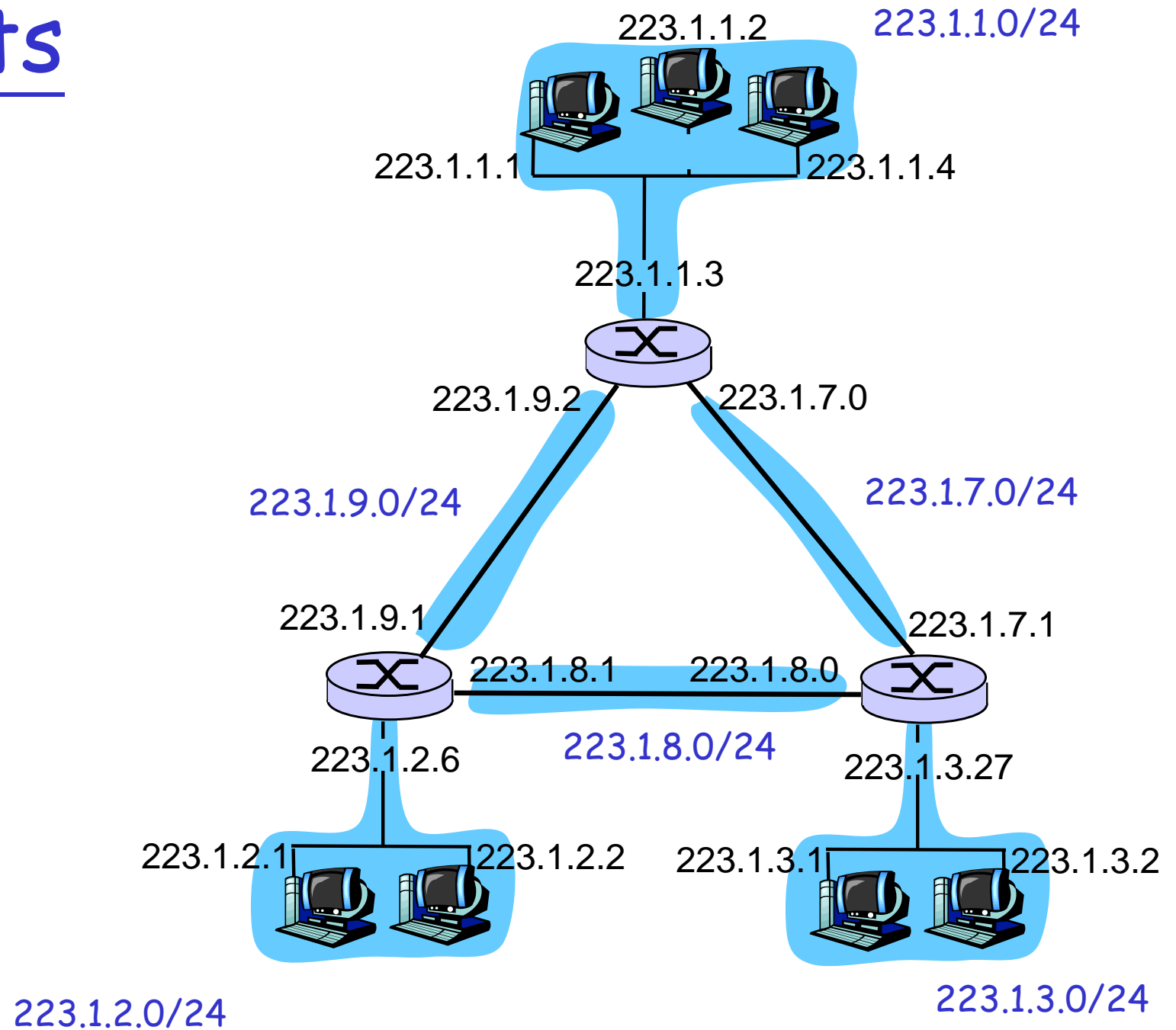223.1.3.1    223.1.3.2

network consisting of 3 subnets

# Subnets

□ 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?

6 subnets

223.1.1.2                     223.1.1.0/24

223.1.1.1                     223.1.1.4

223.1.1.3

223.1.9.2        223.1.7.0

223.1.9.0/24                  223.1.7.0/24

223.1.9.1                     223.1.7.1

223.1.8.1        223.1.8.0

223.1.2.6        223.1.8.0/24        223.1.3.27

223.1.2.1    223.1.2.2    223.1.3.1    223.1.3.2

223.1.2.0/24                  223.1.3.0/24

# IP addressing: CIDR

**CIDR:** **C**lassless **I**nter**D**omain **R**outing

- 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

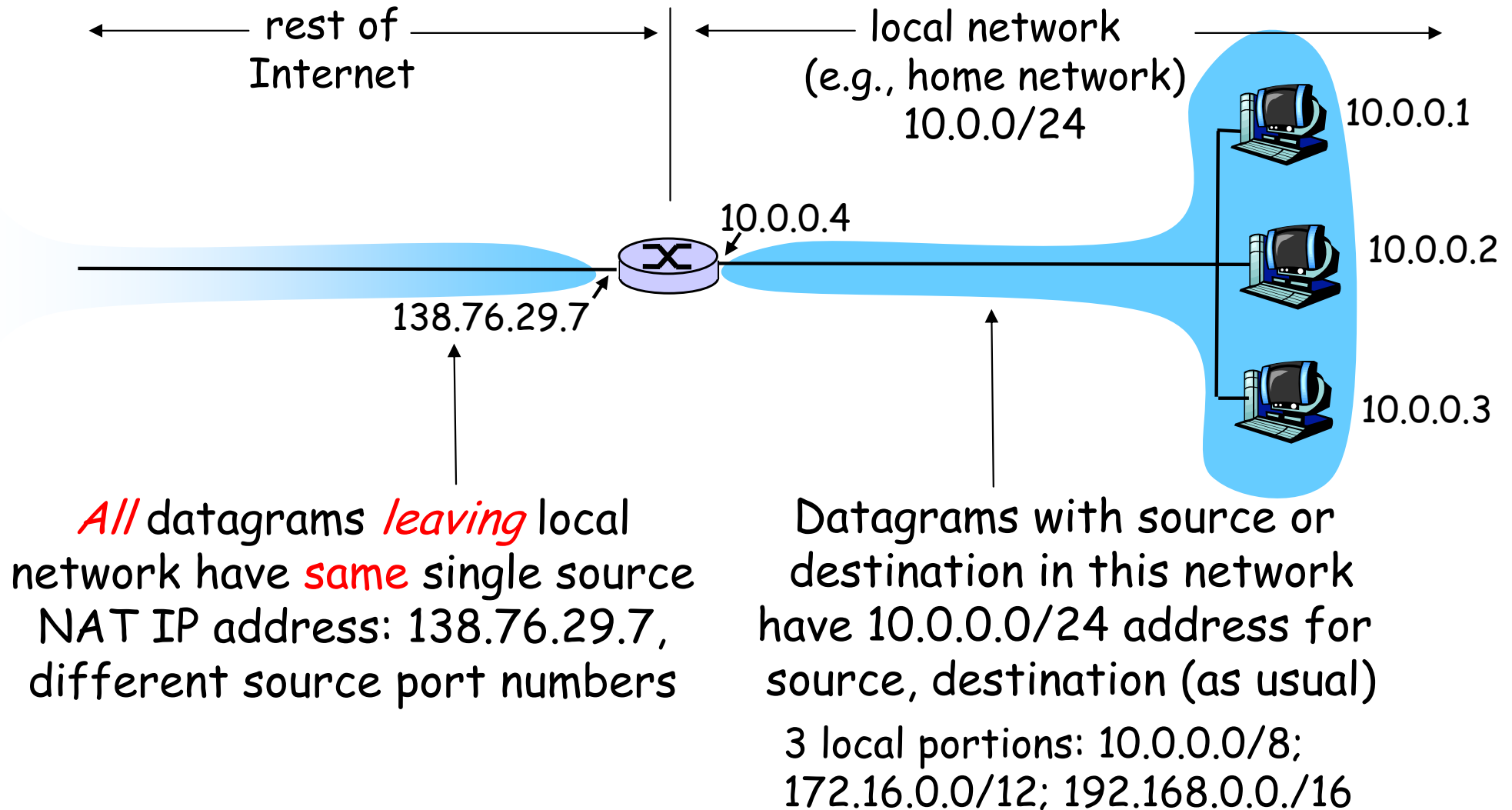11001000 00010111 00010000 00000000

200.23.16.0/23

# IP addresses: how to get one?

Q: How does *host* get IP address?

□ hard-coded by system admin in a file
- Wintel: 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"

(more in next chapter)

# NAT: Network Address Translation

rest of Internet

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

10.0.0.4

138.76.29.7

10.0.0.1

10.0.0.2

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 10.0.0.0/24 address for source, destination (as usual)

3 local portions: 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 word is concerned:

- ○ no need to be allocated range of addresses from ISP:
  - just one IP address is used 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
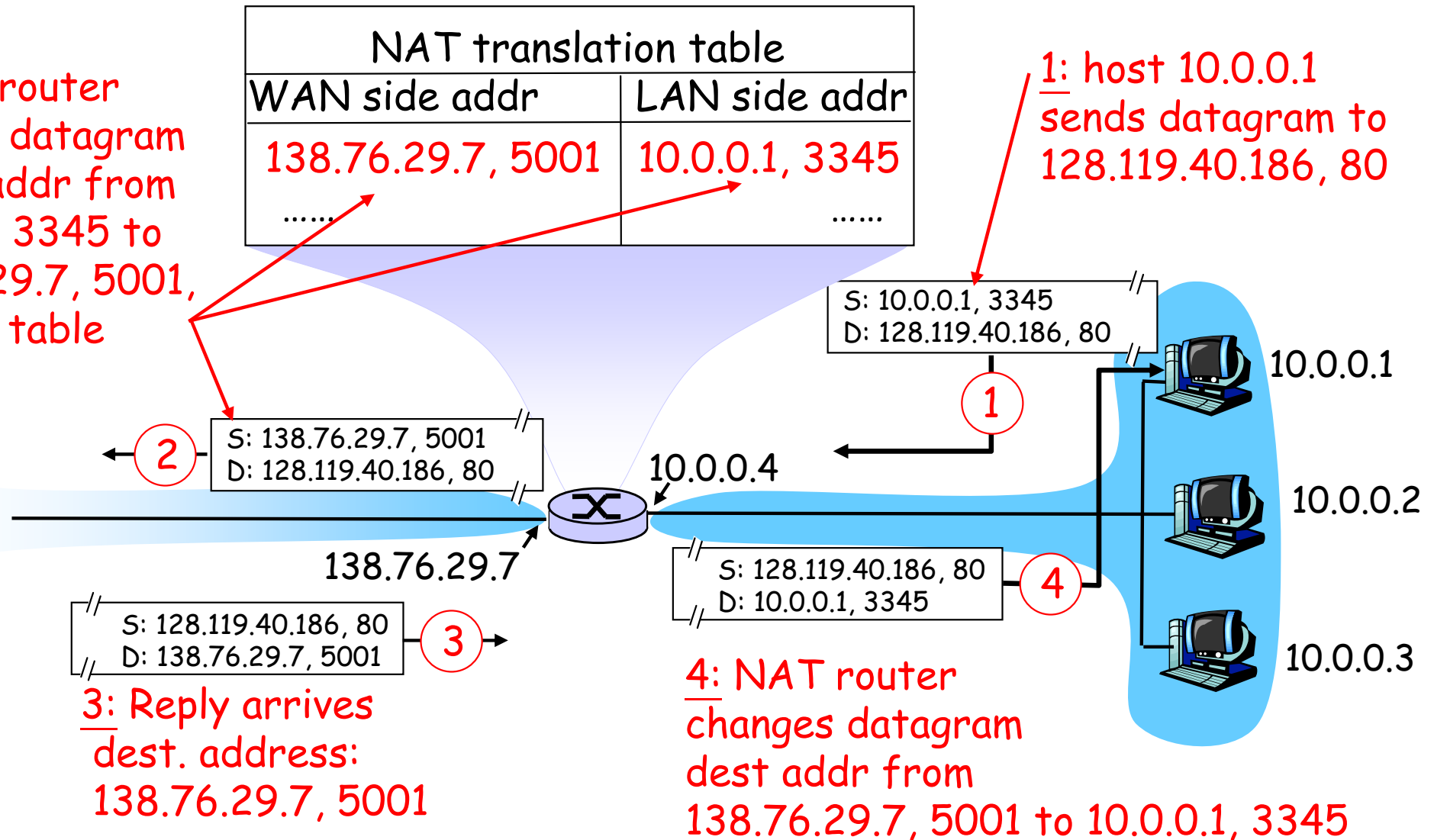- ○ devices inside local net not explicitly addressable, visible by outside world (a security plus).

# 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

**NAT translation table**

| WAN side addr | LAN side addr |
|---|---|
| 138.76.29.7, 5001 | 10.0.0.1, 3345 |
| ...... | ...... |

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

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)

S: 138.76.29.7, 5001
D: 128.119.40.186, 80

(2)

10.0.0.4

10.0.0.1

10.0.0.2

10.0.0.3

138.76.29.7

S: 128.119.40.186, 80
D: 10.0.0.1, 3345

(4)

S: 128.119.40.186, 80
D: 138.76.29.7, 5001

(3)

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

# 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

# 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 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 "host unreachable" packet (type 3, code 3)
□ When source gets this ICMP, stops.

4-35

# 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

# IPv6 Header (Cont)

*Traffic Class (Priority):* identify priority among datagrams in flow

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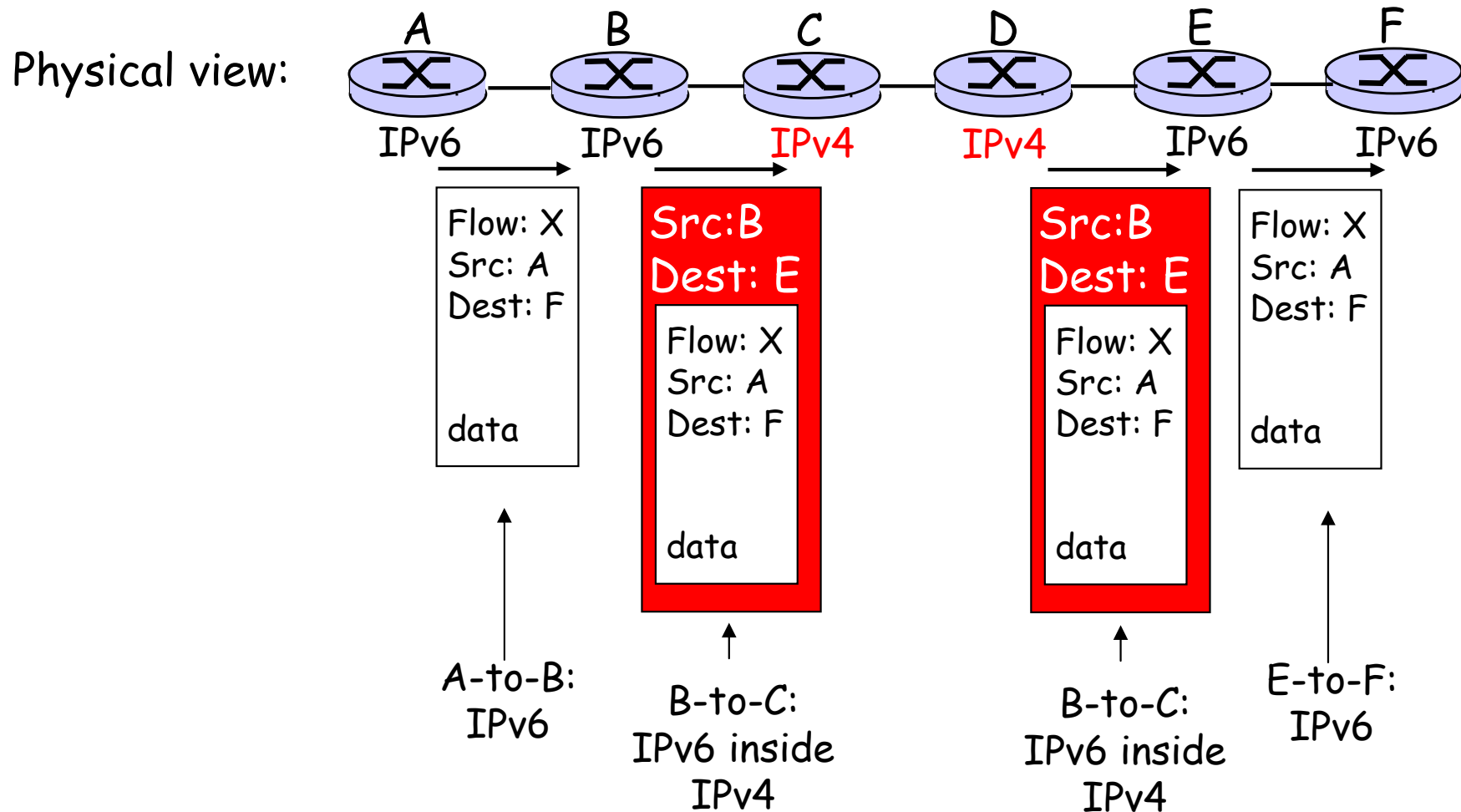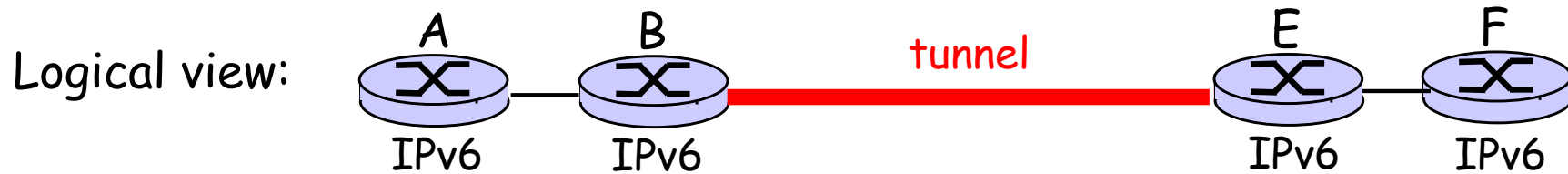
*Next header:* identify upper layer protocol for data

32 bits

| Version | Traffic class | Flow label | |
|---|---|---|---|
| Payload length | | Next hdr | Hop limit |
| Source address (128 bits) | | | |
| Destination address (128 bits) | | | |
| Data | | | |

32 bits

| Version | Header length | Type of service | Datagram length (bytes) | |
|---|---|---|---|---|
| 16-bit Identifier | | | Flags | 13-bit Fragmentation offset |
| Time-to-live | Upper-layer protocol | | Header checksum | |
| 32-bit Source IP address | | | | |
| 32-bit Destination IP address | | | | |
| Options (if any) | | | | |
| Data | | | | |

# 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:



Physical view:

# Chapter 4: Network Layer

# 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

# 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

# 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; $= \infty$ 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

# Dijsktra's Algorithm

```
1  Initialization: source u
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'
```
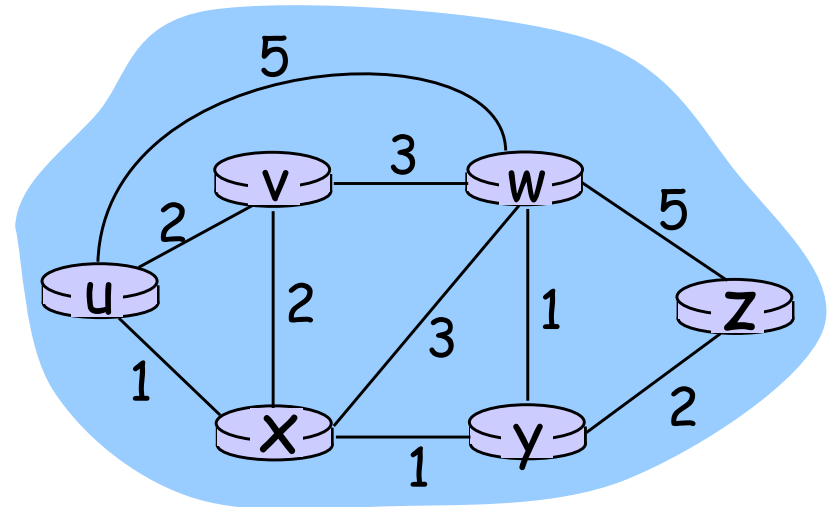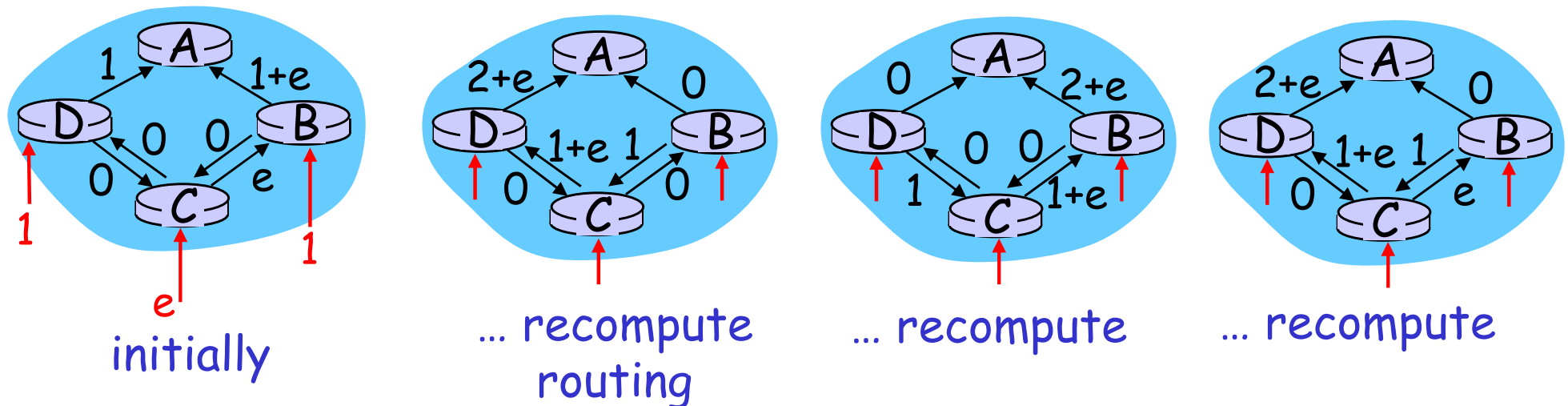
# Dijkstra's algorithm: example

| 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 | | | | | |

1  **Initialization:** source u
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, 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



initially ... recompute routing ... recompute ... recompute

Self-synchronize: even if start at different times but with the same period
Solution: randomize the time it sends out a link advertisement

# Distance Vector Algorithm (1)

Bellman-Ford Equation (dynamic programming)
Define
$d_x(y)$ := cost of least-cost path from x to y

Then

$$d_x(y) = \min_v \{c(x,v) + d_v(y) \}$$

where min is taken over all neighbors of x

Intuitive: the least cost from x to y is the minimum of $c(x,v)+d_v(y)$
taken over all neighbors v

# Bellman-Ford example (2)



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 (3)

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

# Distance vector algorithm (4)

Basic idea:

□ Each node periodically sends its own distance vector estimate to neighbors

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

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

□ Under minor, natural conditions, the estimate $D_x(y)$ *converge the actual least cost* $d_x(y)$
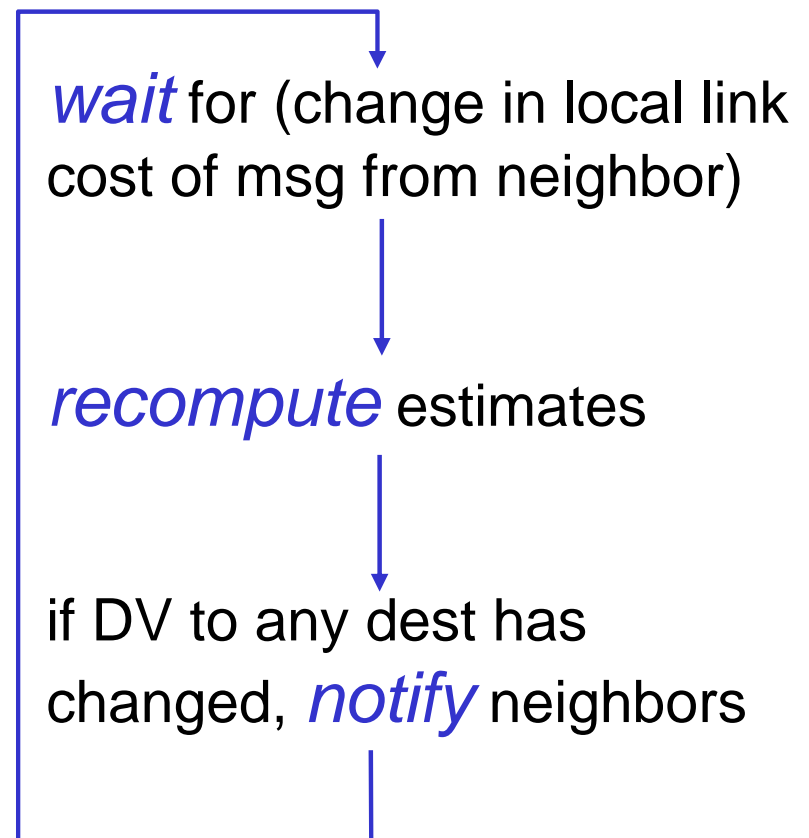
# Distance Vector Algorithm (5)

**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 of 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

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

cost to

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

cost to

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

**node y table**

cost to

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

cost to

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

cost to

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

**node z table**

cost to

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

cost to

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

cost to

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

time

2  y  1
x  7  z

# 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 (msg xchg between y and z) before algorithm stabilizes: see text

## Count-to-infinity problem:

- Initially, $D_y(x)$=4, $D_y(z)$=1, $D_z(y)$=1, $D_z(x)$=5
- Assume y detects the link-cost change
  - new $D_y(x)$=min{$c(y,x)$ +$D_x(x)$, $c(y,z)+D_z(x)$} = {60+0, 1+5} = 6

  wrong value ! (correct route via z)

- Routing loop: Old $D_z(x)$ = 5 = $c(z,y)+D_y(x)$
  - new $D_z(x)$=min{$c(z,y)+D_y(x)$,$c(z,x)+D_x(x)$} ={1+6, 50+0} = 7

  wrong value and wrong route via y (route thru the node with changing cost)

- Loop again: y recomputes the DV
  - new $D_y(x)$=min{$c(y,x)$ +$D_x(x)$, $c(y,z)+D_z(x)$} = {60+0, 1+7} = 8

60

4    y    1

x         z

50

## Poissoned 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? NO!

# 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

# 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



- 3c
- 3a
- 3b
- AS3
- 1c
- 1a
- 1d
- 1b
- AS1
- 2a
- 2c
- 2b
- AS2

Intra-AS Routing algorithm

Inter-AS Routing algorithm

Forwarding table

□ Forwarding table is 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 for which dest is outside of AS1
  - Router should forward packet towards on of the gateway routers, but which one?

AS1 needs:

1. to learn which dests are reachable through AS2 and which through AS3
2. to propagate this reachability info to all routers in AS1

Job of inter-AS routing!

# Chapter 4: Network Layer

# 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)

□ Distance vector algorithm

□ Included in BSD-UNIX Distribution in 1982

□ Distance metric: # of hops (max = 15 hops)



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

# RIP advertisements

- Distance vectors: exchanged among neighbors every 30 sec via Response Message (also called advertisement)
- Each advertisement: list of up to 25 destination nets within AS

# RIP: Example



| Destination Network | Next Router | Num. of hops to dest. |
|---|---|---|
| w | A | 2 |
| y | B | 2 |
| z | B | 7 |
| x | -- | 1 |
| …. | …. | …. |

Routing table in D

# RIP: Example

| Dest | Next | hops |
|------|------|------|
| w | - | - |
| x | - | - |
| z | C | 4 |
| …. | … | … |

Advertisement from A to D

W   A   x   D   B   y   z

C

| Destination Network | Next Router | Num. of hops to dest. |
|---------------------|-------------|-----------------------|
| w | A | 2 |
| y | B | 2 |
| z | ~~B~~ A | ~~7~~ 5 |
| x | -- | 1 |
| …. | …. | …. |

Routing table in D

# RIP: Link Failure and 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)

# RIP Table processing

□ RIP routing tables managed by **application-level** process called route-d (daemon)

□ advertisements sent in UDP packets, periodically repeated

# OSPF (Open Shortest Path First)

□ "open": publicly available

□ Uses Link State algorithm
   ○ LS packet dissemination
   ○ Topology map at each node
   ○ Route computation using Dijkstra's algorithm

□ OSPF advertisement carries one entry per neighbor router

□ Advertisements disseminated to entire AS (via flooding)
   ○ Carried in OSPF messages directly over IP (rather than TCP or UDP

# OSPF "advanced" features (not in RIP)

- Security: all OSPF messages authenticated (to prevent malicious intrusion)
- Multiple same-cost paths allowed (only one path in RIP)
- For each link, multiple cost metrics for different TOS (e.g., satellite link cost set "low" for best effort; high for real time)
- Integrated uni- and multicast support:
  - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- Hierarchical OSPF in large domains.

# 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 the reachability information to all routers internal to the AS.
   3. Determine "good" routes to subnets based on reachability information and policy.

□ Allows a subnet to advertise its existence to rest of the Internet: *"I am here"*

# BGP basics

□ Pairs of routers (BGP peers) exchange routing info over semi-permanent TCP conctns: <span style="color:red">BGP sessions</span>

□ Note that BGP sessions do not correspond to physical links.

□ When AS2 advertises a prefix to AS1, AS2 is <span style="color:red">*promising*</span> it will forward any datagrams destined to that prefix towards the prefix.

   ○ AS2 can aggregate prefixes in its advertisement



— — — — —  eBGP session

·············  iBGP session

# Distributing reachability info

- With eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
- 1c can then use iBGP do distribute this new prefix reach info to all routers in AS1
- 1b can then re-advertise the new reach info to AS2 over the 1b-to-2a eBGP session
- When router learns about a new prefix, it creates an entry for the prefix in its forwarding table.



eBGP session

iBGP session

# Path attributes & BGP routes

- When advertising a prefix, advert includes BGP attributes.
  - prefix + attributes = "route"
- Two important attributes:
  - AS-PATH: contains the ASs through which the advert for the prefix passed: AS 67 AS 17
  - NEXT-HOP: Indicates the specific internal-AS router to next-hop AS. (There may be multiple links from current AS to next-hop-AS.)
- When gateway router receives route advert, 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

# 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: Network Layer

# Broadcast Routing Algorithms

□ N-way-unicast
- ○ source duplication
- ○ in-network duplication

□ Controlled flooding
- ○ Uncontrolled flooding: broadcast storm
- ○ sequence-number-controlled flooding
- ○ reverse path forwarding

□ Spanning-Tree broadcast
- ○ Center-based approach

# N-way-unicast to broadcasting



**Figure 4.40** Source-duplication versus in-network duplication.
(a) source duplication, (b) in-network duplication

# Controlled Flooding

□ Sequence-number-controlloed:
  ○ Put node address and a broadcast sequence number into a broadcast packet for checking duplicate

□ Reverse path forwarding
  ○ Know the neighbor on its unicast shortest path to the sender
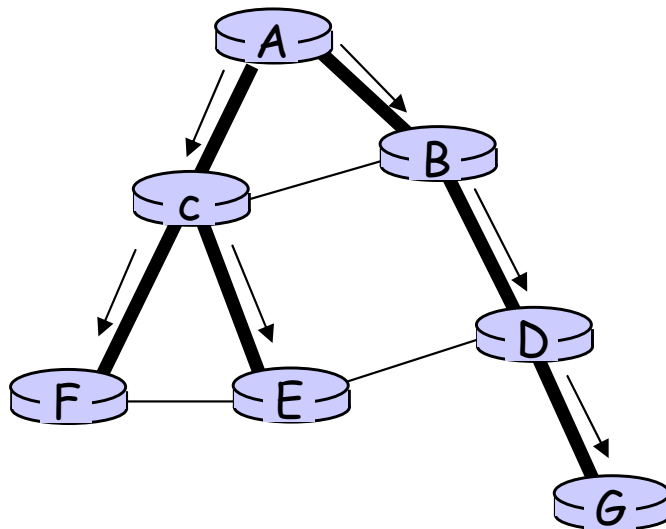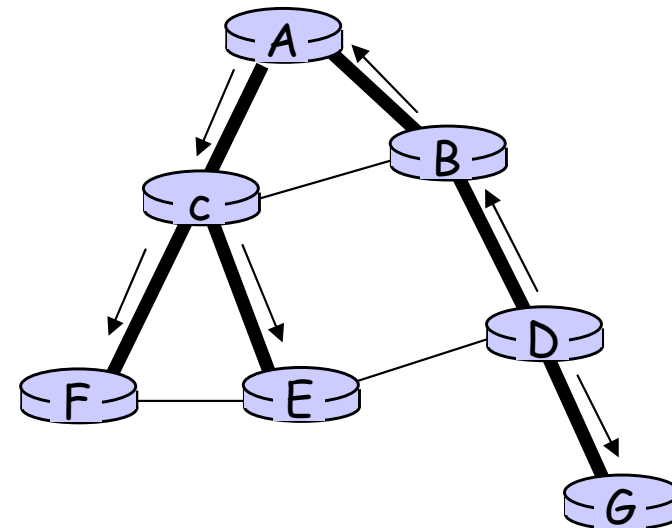  ○ Only the broadcast packet from the above neighbor is transmitted on all of its outgoing links

**Figure 4.41**: Reverse path forwarding

# Spanning-Tree Broadcast

☐ Construct a (minimum) spanning tree to broadcast packets
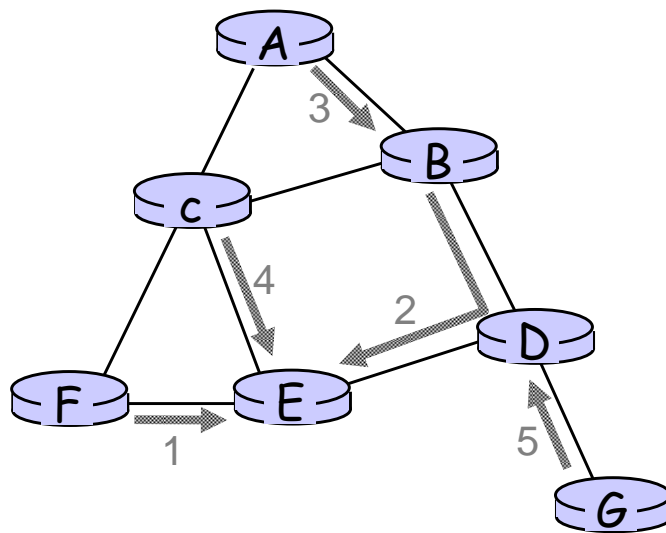


(a) Broadcast initiated at A
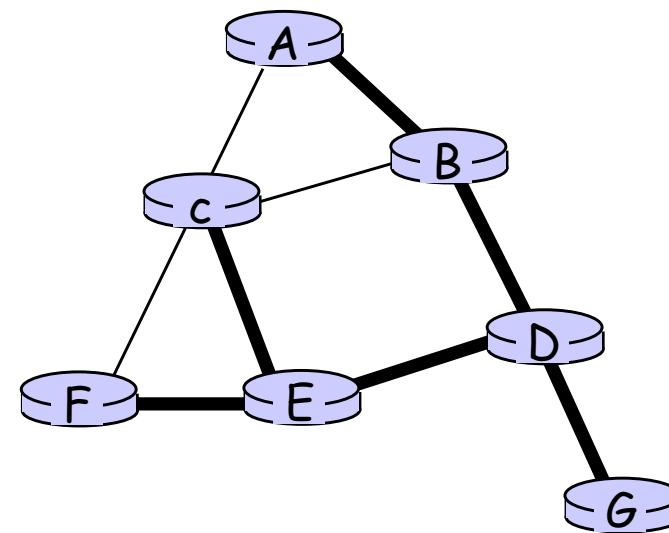
(b) Broadcast initiated at D

**Figure 4.42**: Broadcast along a spanning tree

# Center-based construction of a spanning tree

☐ Nodes unicast tree-join messages addressed to the center node

○ Assume E is the center node
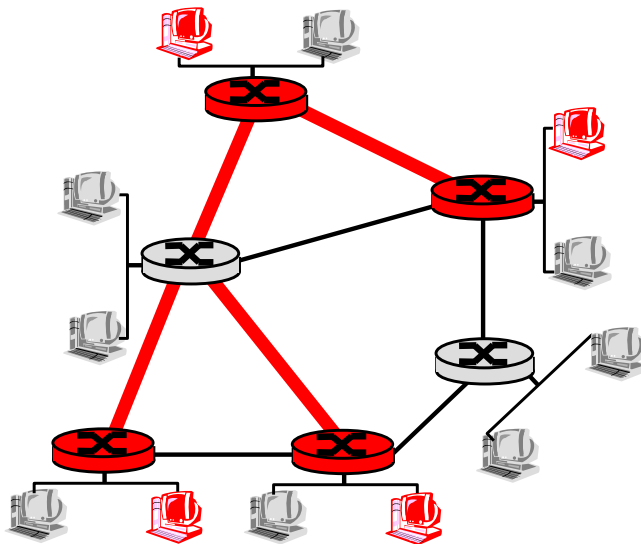


**(a) Stepwise construction of spanning tree**

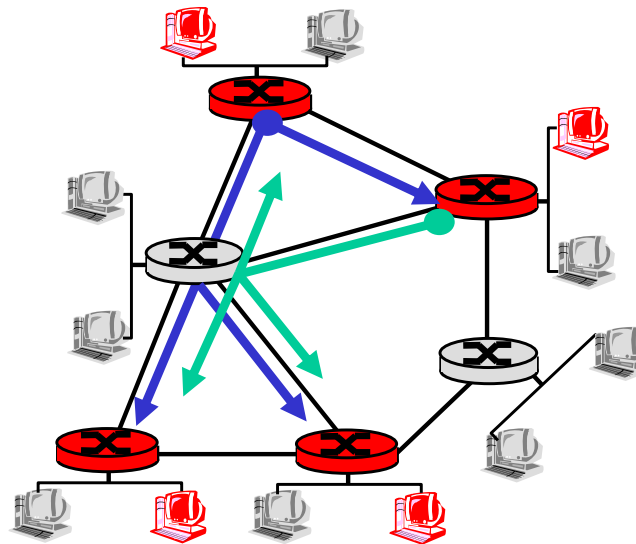**(b) Constructed spanning tree**

**Figure 4.43**: Center-based construction of a spanning tree

# Multicast Routing: Problem Statement

□ **_Goal:_** find a tree (or trees) connecting routers having local mcast group members

  ○ _tree:_ not all paths between routers used
  ○ _source-based:_ different tree from each sender to rcvrs
  ○ _shared-tree:_ same tree used by all group members

Shared tree                    Source-based trees
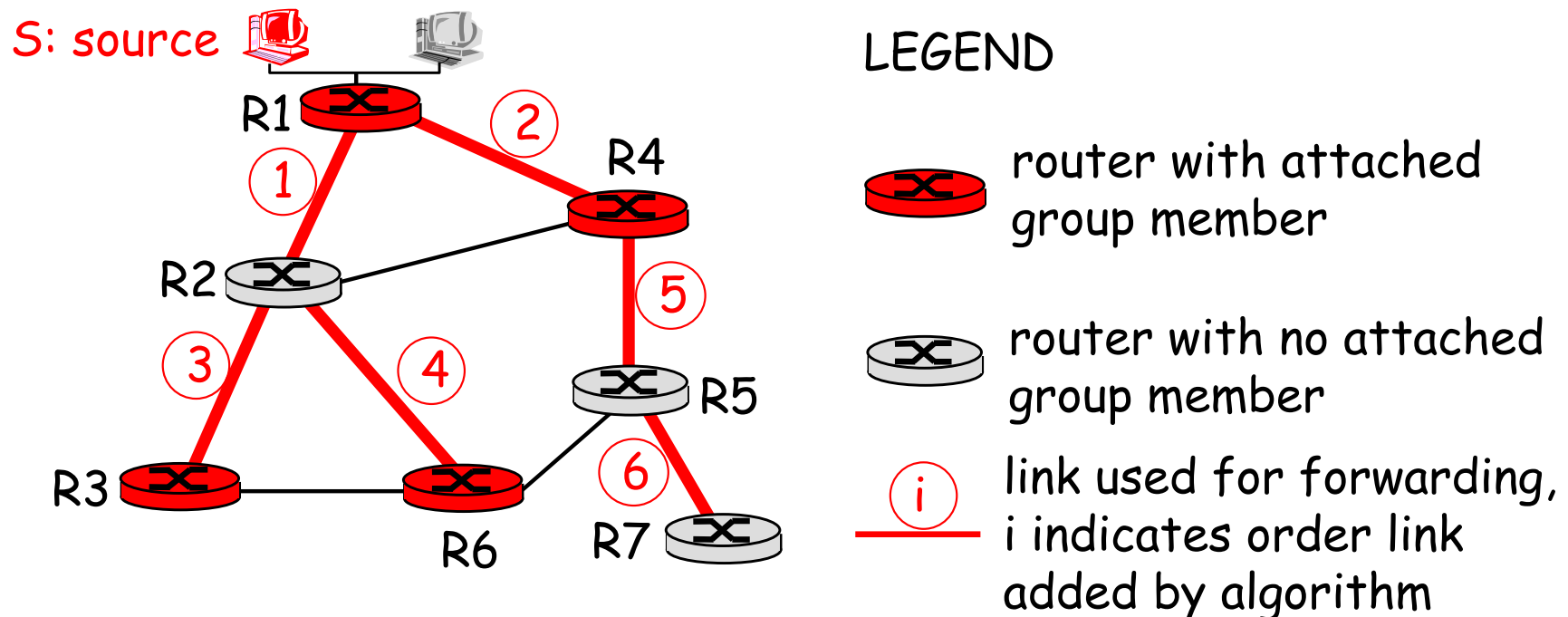
# Approaches for building mcast trees

Approaches:

□ source-based tree: one tree per source
  ○ shortest path trees
  ○ reverse path forwarding

□ group-shared tree: group uses one tree
  ○ minimal spanning (Steiner)
  ○ center-based trees

...we first look at basic approaches, then specific protocols adopting these approaches

# Shortest Path Tree

□ mcast forwarding tree: tree of shortest path routes from source to all receivers
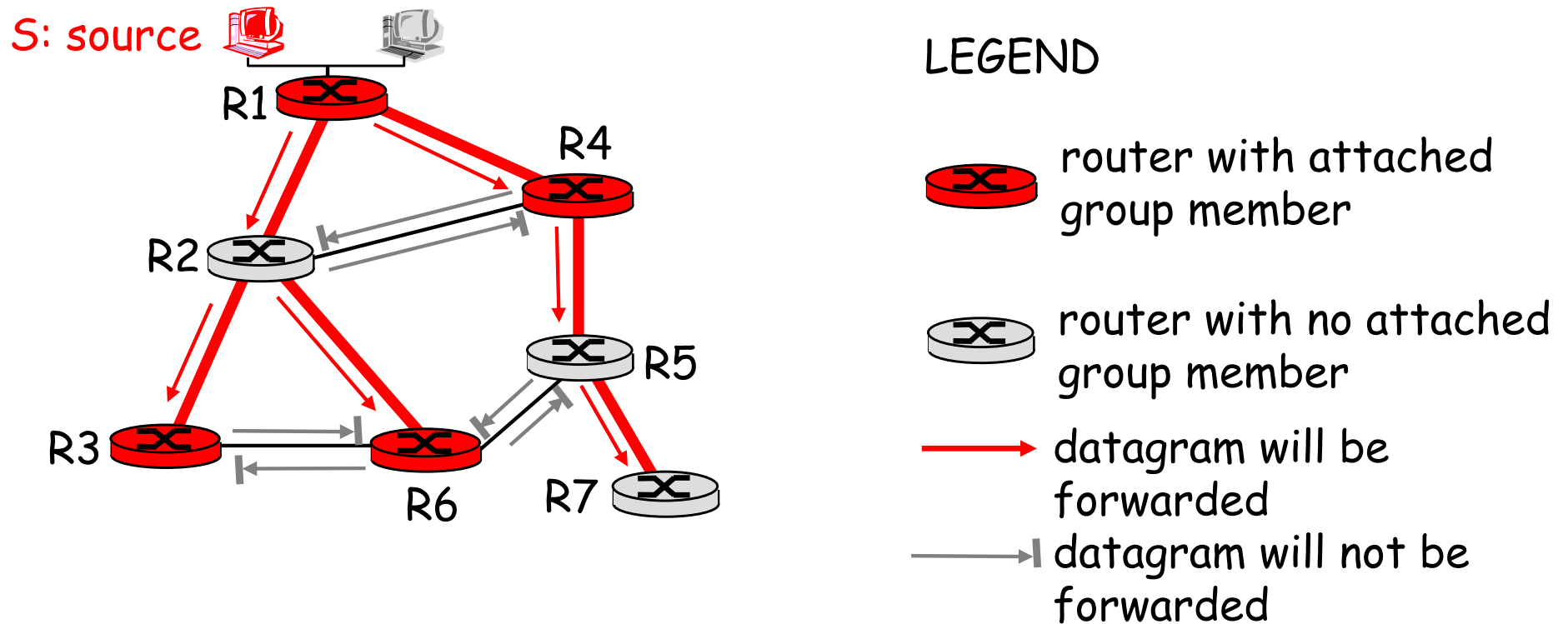  ○ Dijkstra's algorithm

# Reverse Path Forwarding

- ❑ rely on router's knowledge of unicast shortest path from it to sender
- ❑ each router has simple forwarding behavior:

*if* (mcast datagram received on incoming link
   on shortest path back to center)
   *then* flood datagram onto all outgoing links
   *else* ignore datagram

# Reverse Path Forwarding: example

S: source

LEGEND

router with attached group member

router with no attached group member

→ datagram will be forwarded

→ datagram will not be forwarded

R1
R2
R3
R4
R5
R6
R7

- result is a source-specific *reverse* SPT
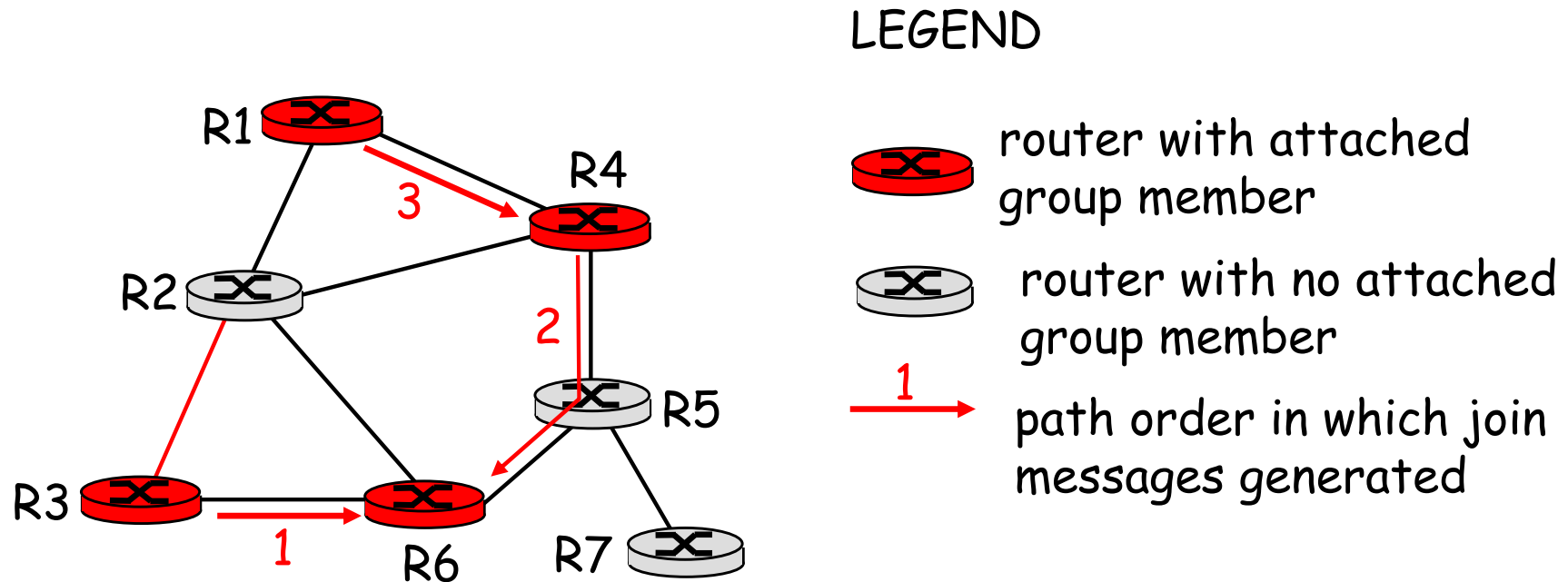  - may be a bad choice with asymmetric links

# Shared-Tree: Steiner Tree

❑ **Steiner Tree:** minimum cost tree connecting all routers with attached group members

❑ problem is NP-complete

❑ excellent heuristics exists

❑ not used in practice:

  ❍ computational complexity

  ❍ information about entire network needed

  ❍ monolithic: rerun whenever a router needs to join/leave

# Center-based trees

□ single delivery tree shared by all

□ one router identified as *"center"* of tree

□ to join:

  ○ edge router sends unicast *join-msg* addressed to center router

  ○ *join-msg* "processed" by intermediate routers and forwarded towards center

  ○ *join-msg* either hits existing tree branch for this center, or arrives at center

  ○ path taken by *join-msg* becomes new branch of tree for this router

# Center-based trees: an example

Suppose R6 chosen as center:



LEGEND

router with attached group member

router with no attached group member

1 → path order in which join messages generated
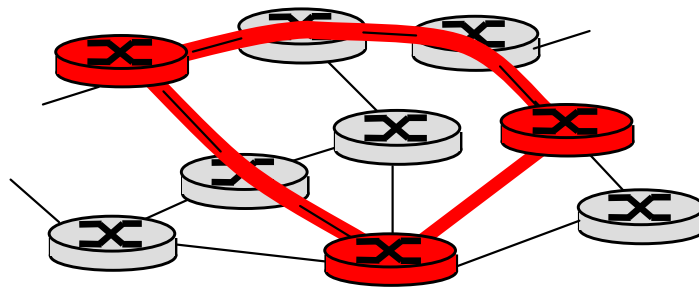
# Internet Multicasting Routing: DVMRP

□ **DVMRP:** distance vector multicast routing protocol, RFC1075

□ *flood and prune:* reverse path forwarding, source-based tree

  ○ RPF tree based on DVMRP's own routing tables constructed by communicating DVMRP routers

  ○ no assumptions about underlying unicast

  ○ initial datagram to mcast group flooded everywhere via RPF

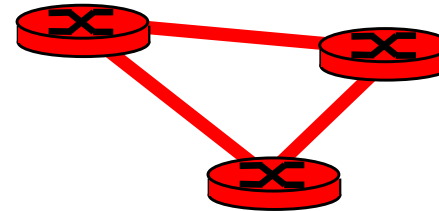  ○ routers not wanting group: send upstream prune msgs

# DVMRP: continued...

□ *soft state:* DVMRP router periodically (1 min.) "forgets" branches are pruned:

- mcast data again flows down unpruned branch
- downstream router: reprune or else continue to receive data

□ routers can quickly regraft to tree

- following IGMP join at leaf

□ odds and ends

- commonly implemented in commercial routers
- Mbone routing done using DVMRP

# Tunneling

**Q:** How to connect "islands" of multicast routers in a "sea" of unicast routers?



physical topology    logical topology

❑ mcast datagram encapsulated inside "normal" (non-multicast-addressed) datagram

❑ normal IP datagram sent thru "tunnel" via regular IP unicast to receiving mcast router

❑ receiving mcast router unencapsulates to get mcast datagram

# PIM: Protocol Independent Multicast

□ not dependent on any specific underlying unicast routing algorithm (works with all)

□ two different multicast distribution scenarios :

*Dense*:

❑ group members densely packed, in "close" proximity.

❑ bandwidth more plentiful

*Sparse:*

❑ # networks with group members small wrt # interconnected networks

❑ group members "widely dispersed"

❑ bandwidth not plentiful

# Consequences of Sparse-Dense Dichotomy:

## Dense
- group membership by routers *assumed* until routers explicitly prune
- *data-driven* construction on mcast tree (e.g., RPF)
- bandwidth and non-group-router processing *profligate*

## Sparse:
- no membership until routers explicitly join
- *receiver- driven* construction of mcast tree (e.g., center-based)
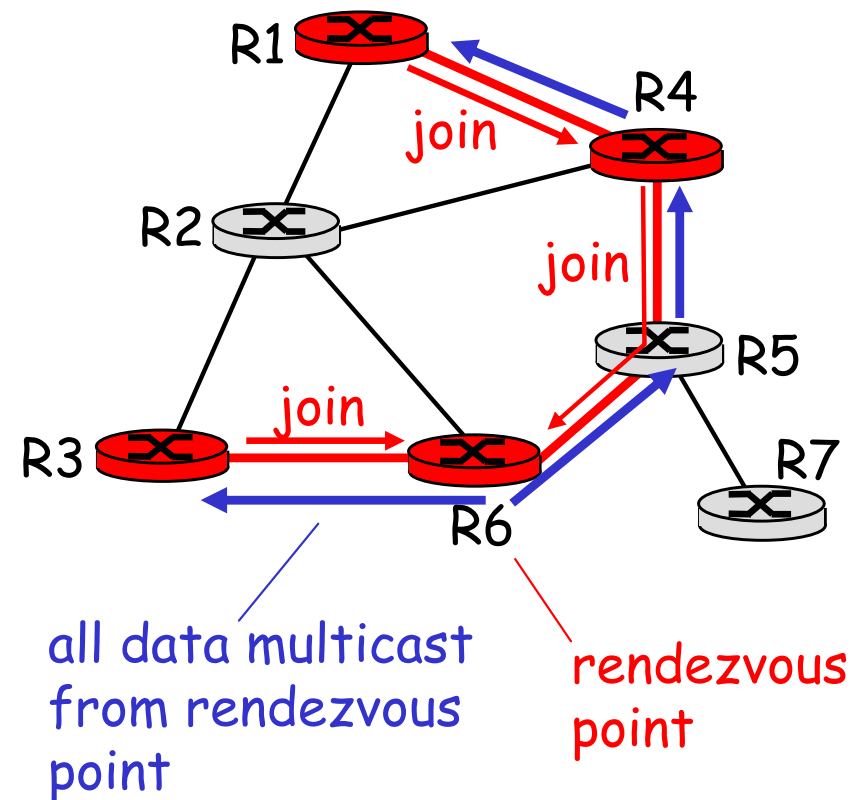- bandwidth and non-group-router processing *conservative*

# PIM- Dense Mode

flood-and-prune RPF, similar to DVMRP but

- ❑ underlying unicast protocol provides RPF info for incoming datagram

- ❑ less complicated (less efficient) downstream flood than DVMRP reduces reliance on underlying routing algorithm

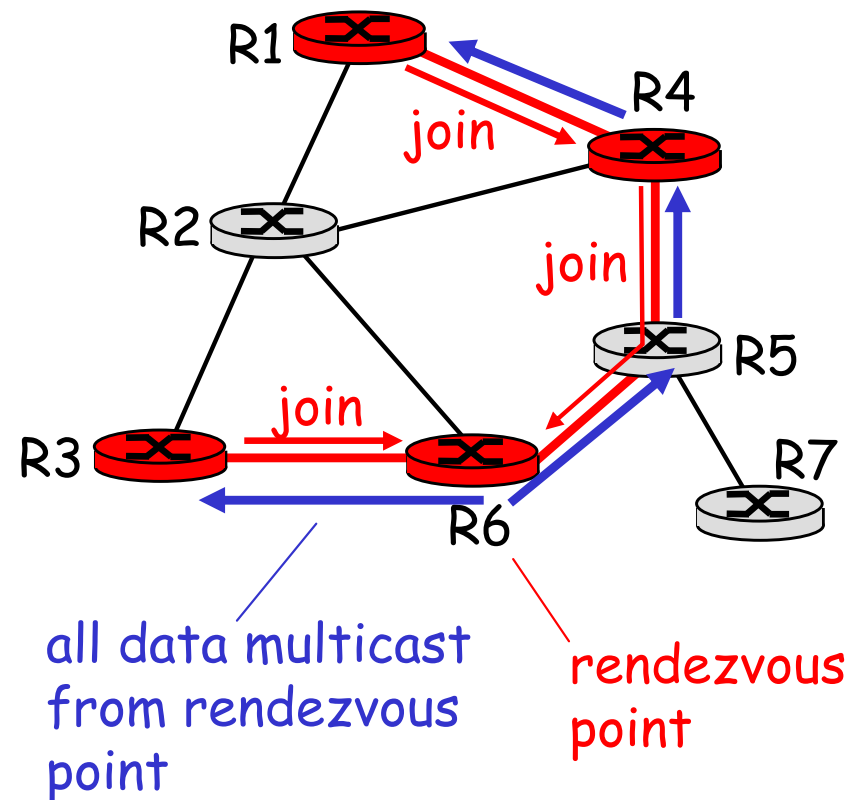- ❑ has protocol mechanism for router to detect it is a leaf-node router

# PIM - Sparse Mode

□ center-based approach

□ router sends *join* msg to rendezvous point (RP)

  ▫ intermediate routers update state and forward *join*

□ after joining via RP, router can switch to source-specific tree

  ▫ increased performance: less concentration, shorter paths

R1

R4

join

R2

join

R5

R3    join    R7

R6

all data multicast from rendezvous point

rendezvous point

# PIM - Sparse Mode

## sender(s):

□ unicast data to RP, which distributes down RP-rooted tree

□ RP can extend mcast tree upstream to source

□ RP can send *stop* msg if no attached receivers
   ○ "no one is listening!"

R1

R4

join

R2

join

R5

R3

join

R6

R7

all data multicast from rendezvous point

rendezvous point

# Network Layer: summary

<span style="color:red">What we've covered:</span>

- network layer services
- routing principles: link state and distance vector
- hierarchical routing
- IP
- Internet routing protocols RIP, OSPF, BGP
- what's inside a router?
- IPv6