# Chapter 4
# Network Layer

Computer Networking
*A Top-Down Approach*

SIXTH EDITION

James F. Kurose • Keith W. Ross

INTERNATIONAL EDITION

ALWAYS LEARNING

PEARSON

# Chapter 4: Network Layer

# Network layer

r   network layer protocols in *every* host, router

r   router examines header fields in all IP datagrams passing through it
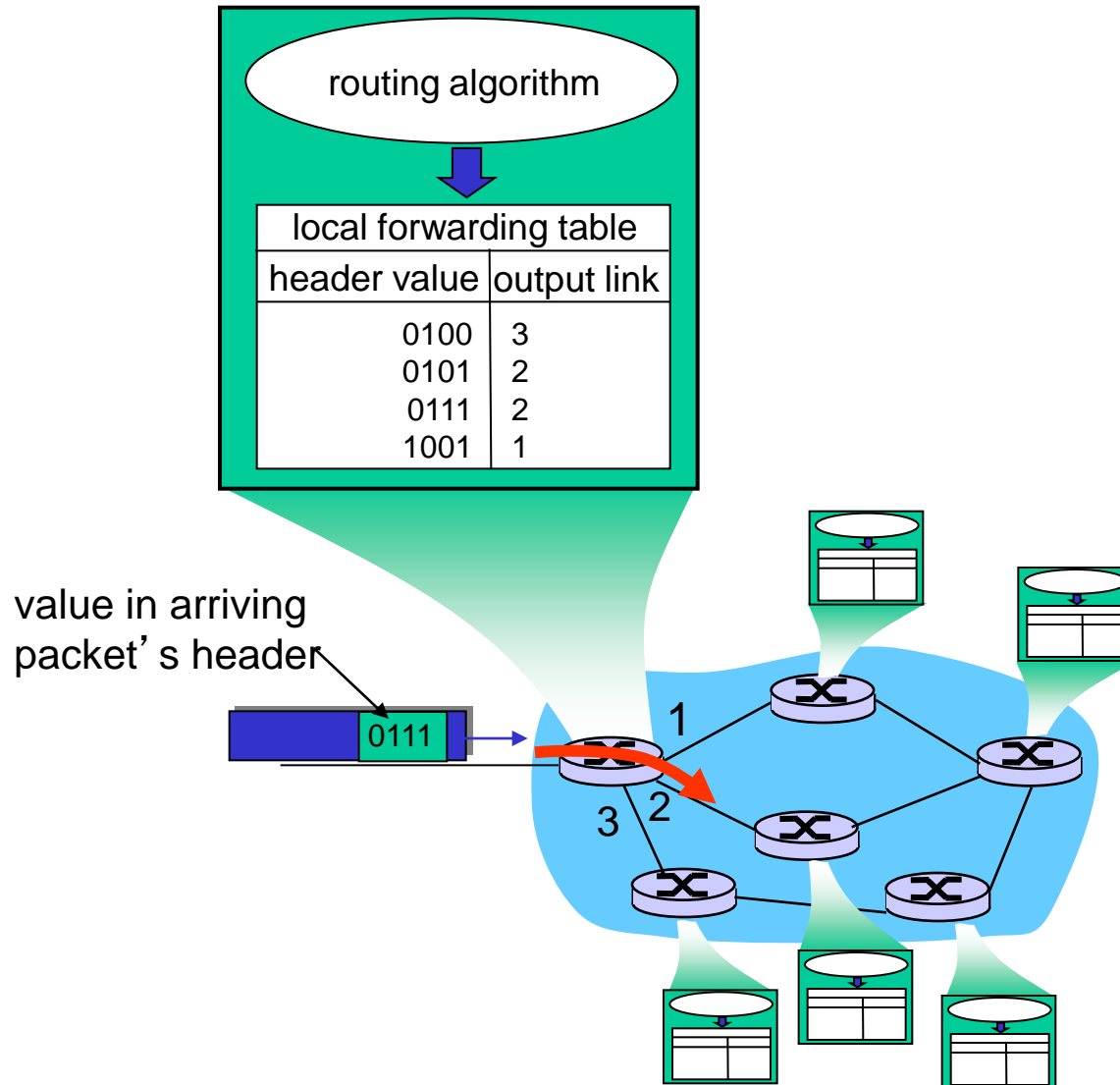
# Two Key Network-Layer Functions

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

r *routing:* determine route taken by packets from source to dest.

m *routing algorithms*

r routing: process of planning trip from source to dest

r 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

# Chapter 4: Network Layer

# Network layer connection and connection-less service

r datagram network provides network-layer connectionless service

r VC network provides network-layer connection service

# Virtual circuits

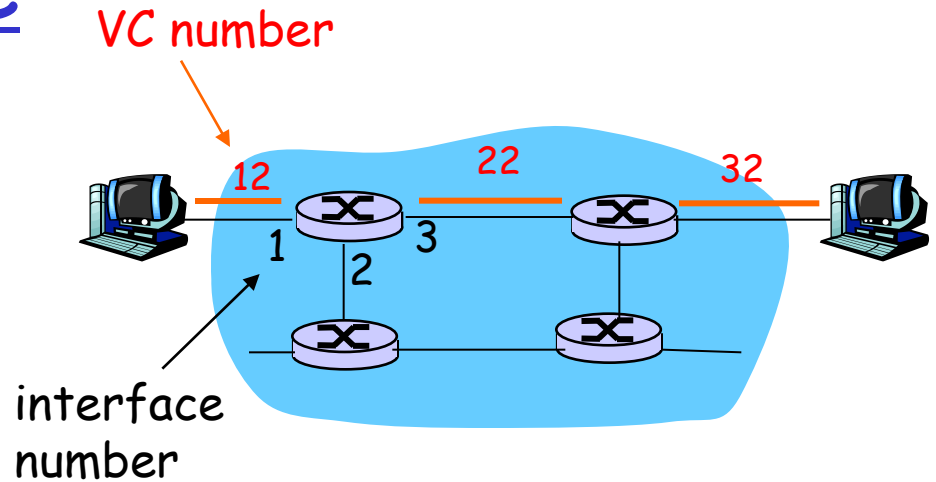"source-to-dest path behaves much like telephone circuit"

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

r each packet carries VC identifier (not destination host address)

r *every* router on source-dest path maintains "state" for each passing connection

r link, router resources (bandwidth, buffers) may be *allocated* to VC (dedicated resources = predictable service)

# Forwarding table



VC number

interface number

## Forwarding table in northwest router:

| 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

r   used in ATM, frame-relay, X.25
r   not used in today's Internet



5. Data flow begins
4. Call connected
1. Initiate call

6. Receive data
3. Accept call
2. incoming call

application
transport
network
data link
physical

application
transport
network
data link
physical

# Datagram networks

r no call setup at network layer

r routers: no state about end-to-end connections

    m no network-level concept of "connection"

r packets forwarded using destination host address

    m packets between same source-dest pair may take different paths

# 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?
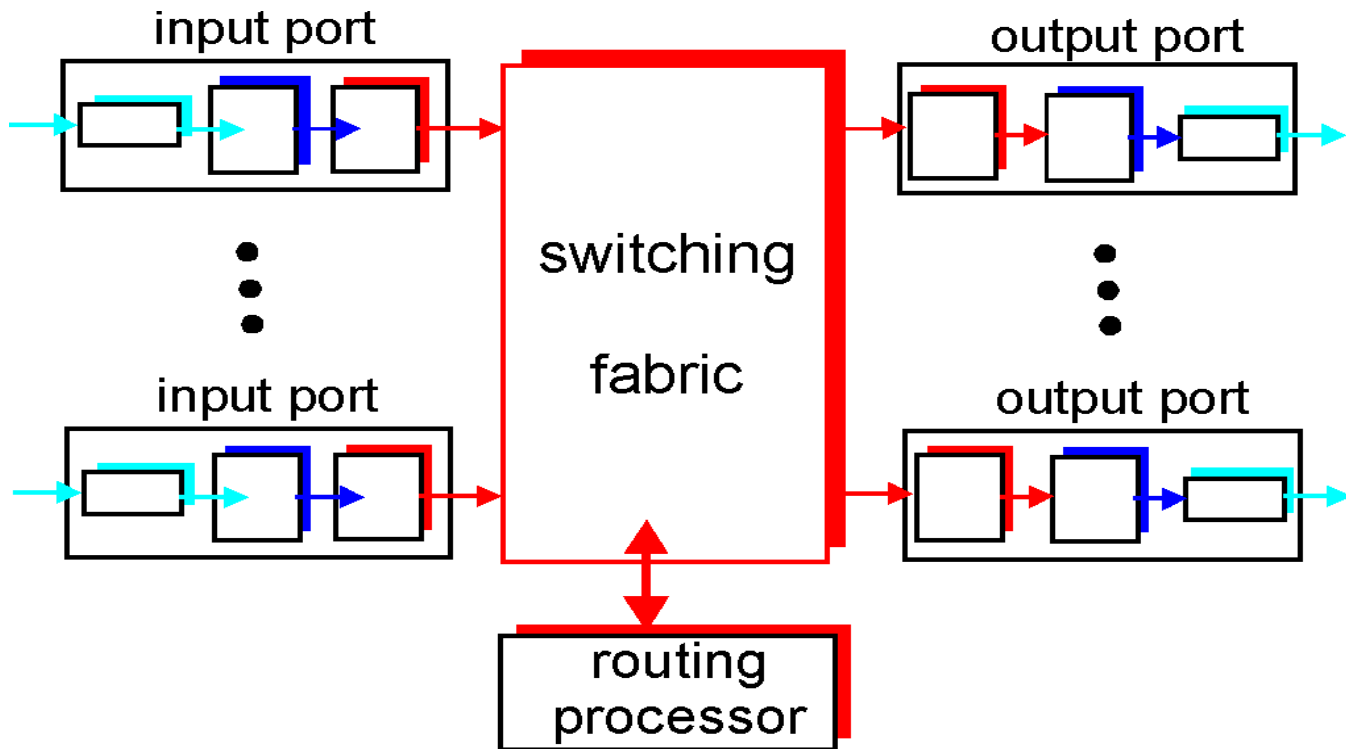
DA: 11001000  00010111  00011000  10101010        Which interface?
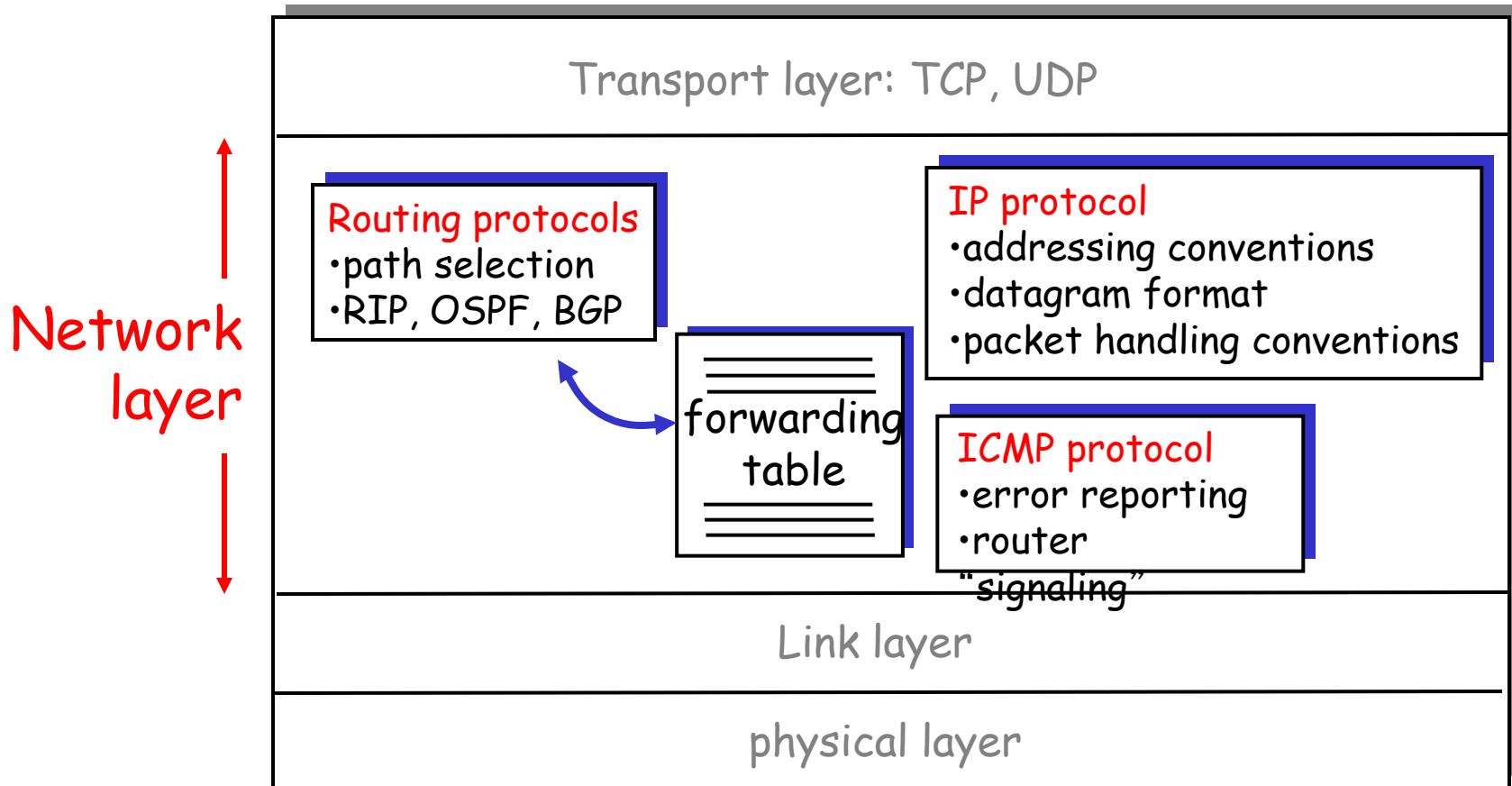
# Router Architecture Overview

Two key router functions:

r   run routing algorithms/protocol (RIP, OSPF, BGP)

r   *forwarding* datagrams from incoming to outgoing link

# 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

total datagram length (bytes)

for fragmentation/ reassembly

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) |

# IP Fragmentation & Reassembly

r network links have MTU (max.transfer size)

- m largest possible link-level frame.

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

- m one datagram becomes several datagrams
- m "reassembled" only at final destination
- m IP header bits used to identify, order related fragments

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

reassembly

# IP Fragmentation and Reassembly

**Example**

r  4000 byte datagram

r  MTU = 1500 bytes

1480 bytes in data field

offset = 1480/8

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

One large datagram becomes several smaller datagrams

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

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

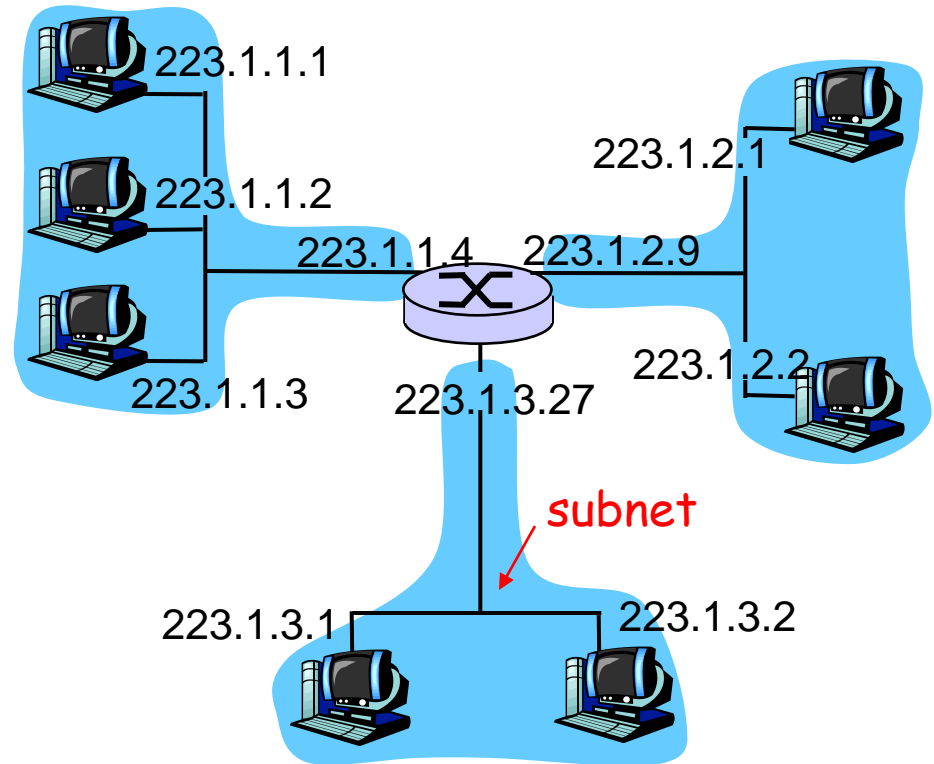| | length =1040 | ID =x | fragflag =0 | offset =370 | |
|---|---|---|---|---|---|

# Chapter 4: Network Layer

# IP Addressing: introduction

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

r *interface:* connection between host/router and physical link

 m router's typically have multiple interfaces

 m host typically has one interface

 m 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

223.1.2.2

223.1.1.3  223.1.3.27

223.1.3.1   223.1.3.2

223.1.1.1 = 11011111 00000001 00000001 00000001
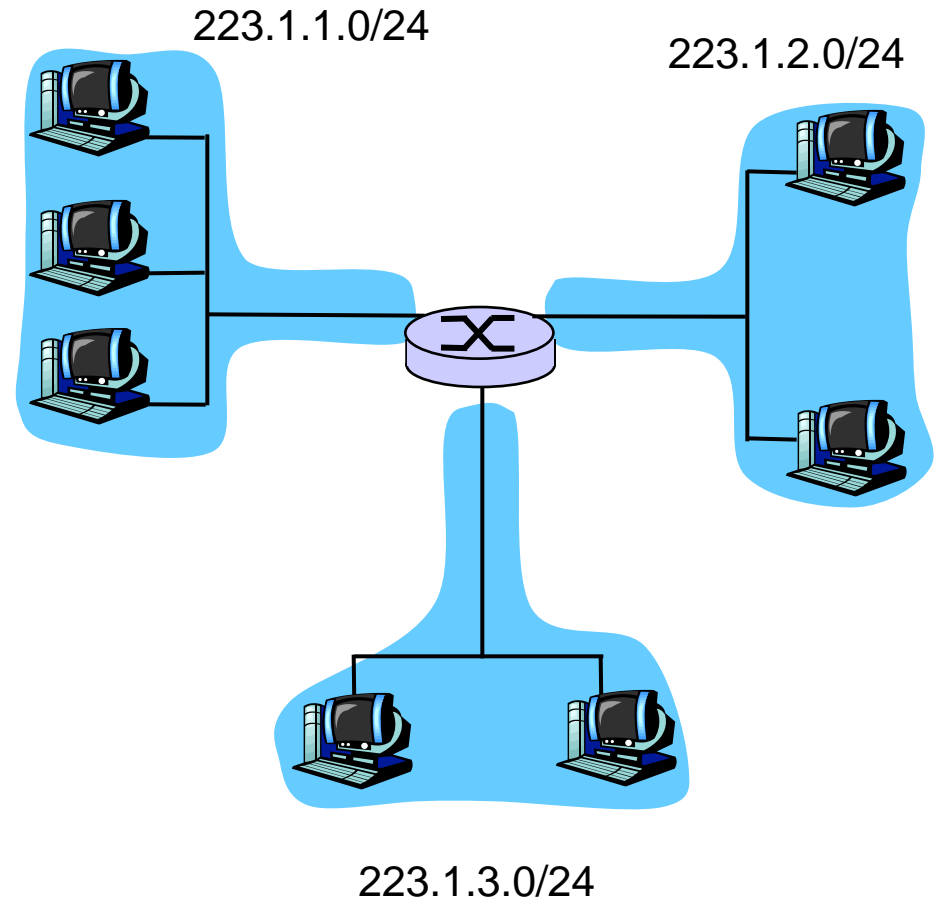
      223    1     1    1

# Subnets

r **IP address:**

  m subnet part (high order bits)

  m host part (low order bits)

r *What's a subnet ?*

  m device interfaces with same subnet part of IP address

  m can physically reach each other without intervening router

223.1.1.1

223.1.1.2

223.1.1.4

223.1.2.1

223.1.2.9

223.1.1.3

223.1.3.27

223.1.2.2

subnet

223.1.3.1
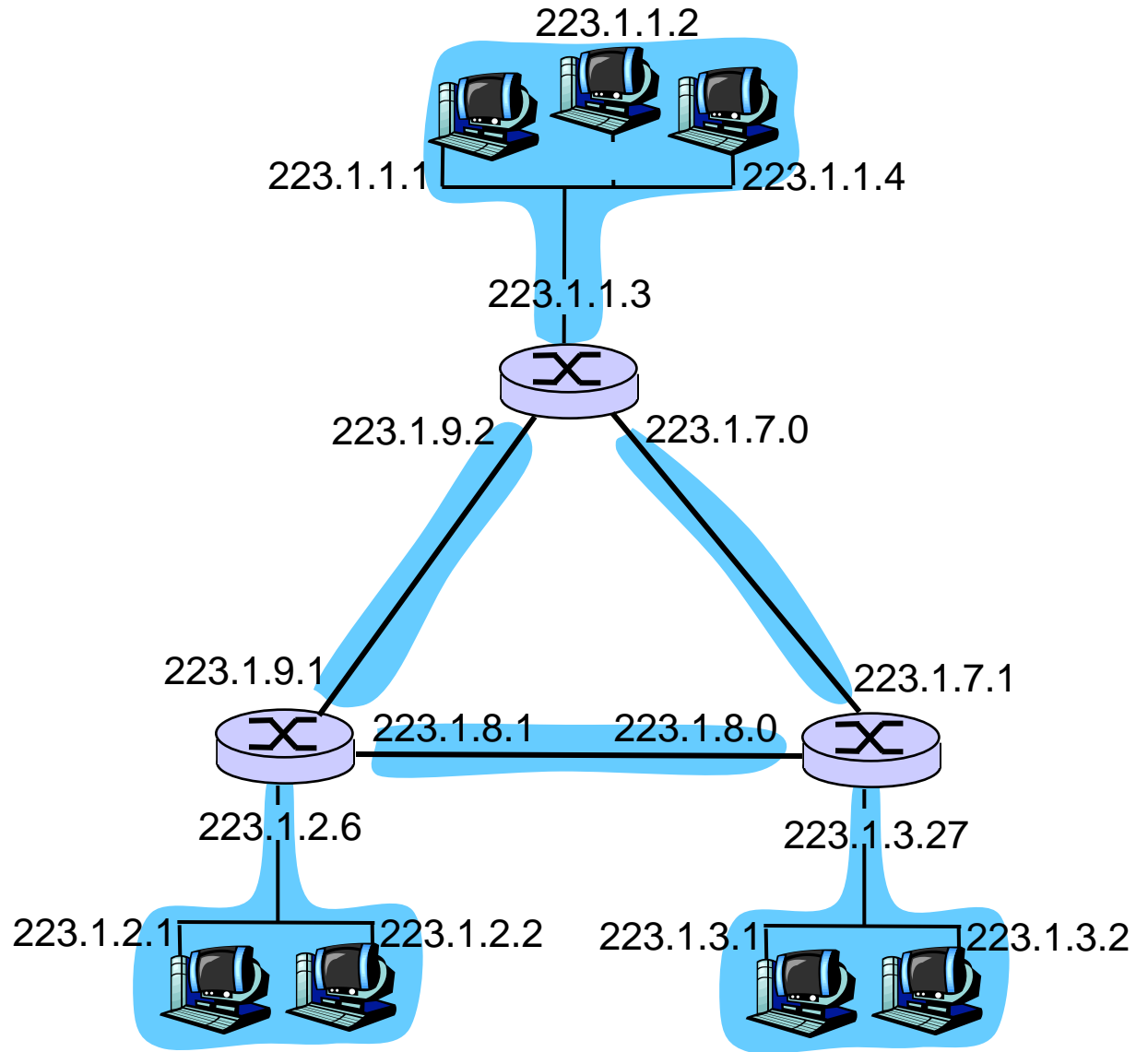
223.1.3.2

network consisting of 3 subnets

# Subnets

r  To determine the subnets, detach each interface from its host or router, creating islands of isolated networks. Each isolated network is called a <span style="color:red">subnet</span>.

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
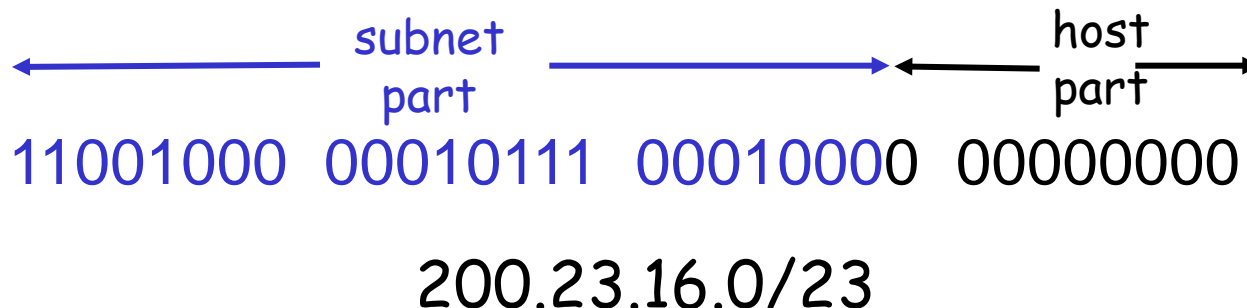
# IP addressing: CIDR

## 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 →|

subnet part ────────────────────→ ←── host part ──→

11001000 00010111 00010000 00000000
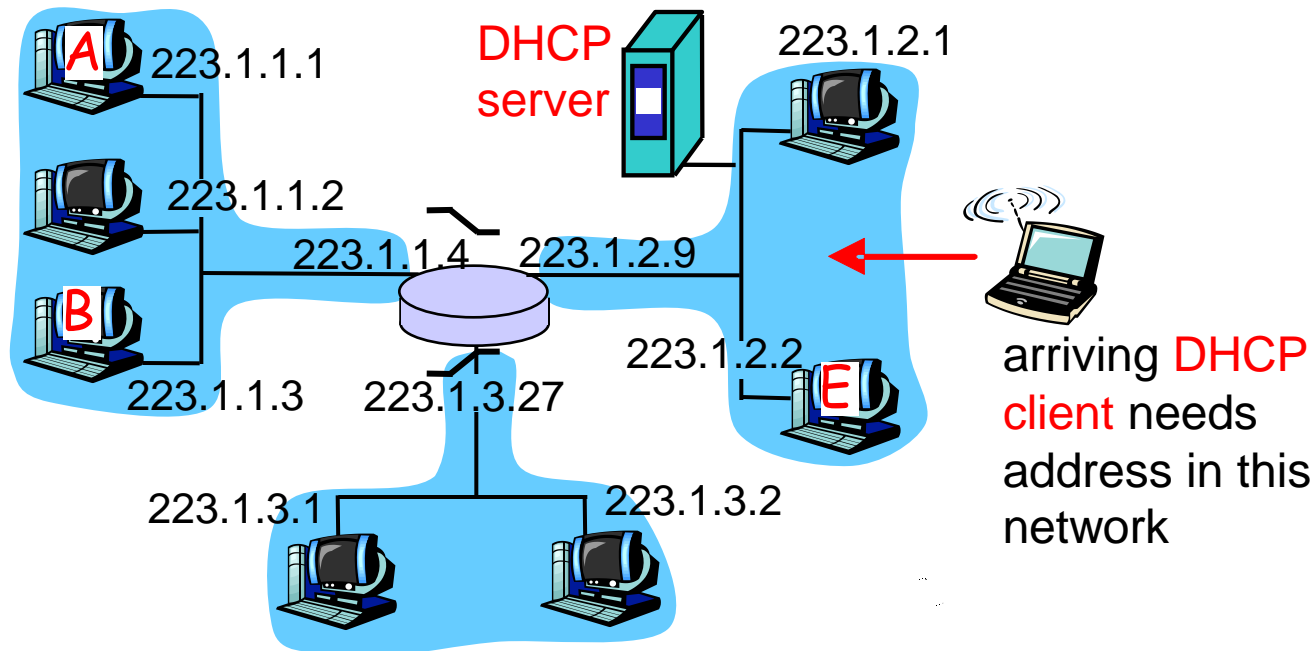
200.23.16.0/23

# IP addresses: how to get one?

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

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

# DHCP: Dynamic Host Configuration Protocol

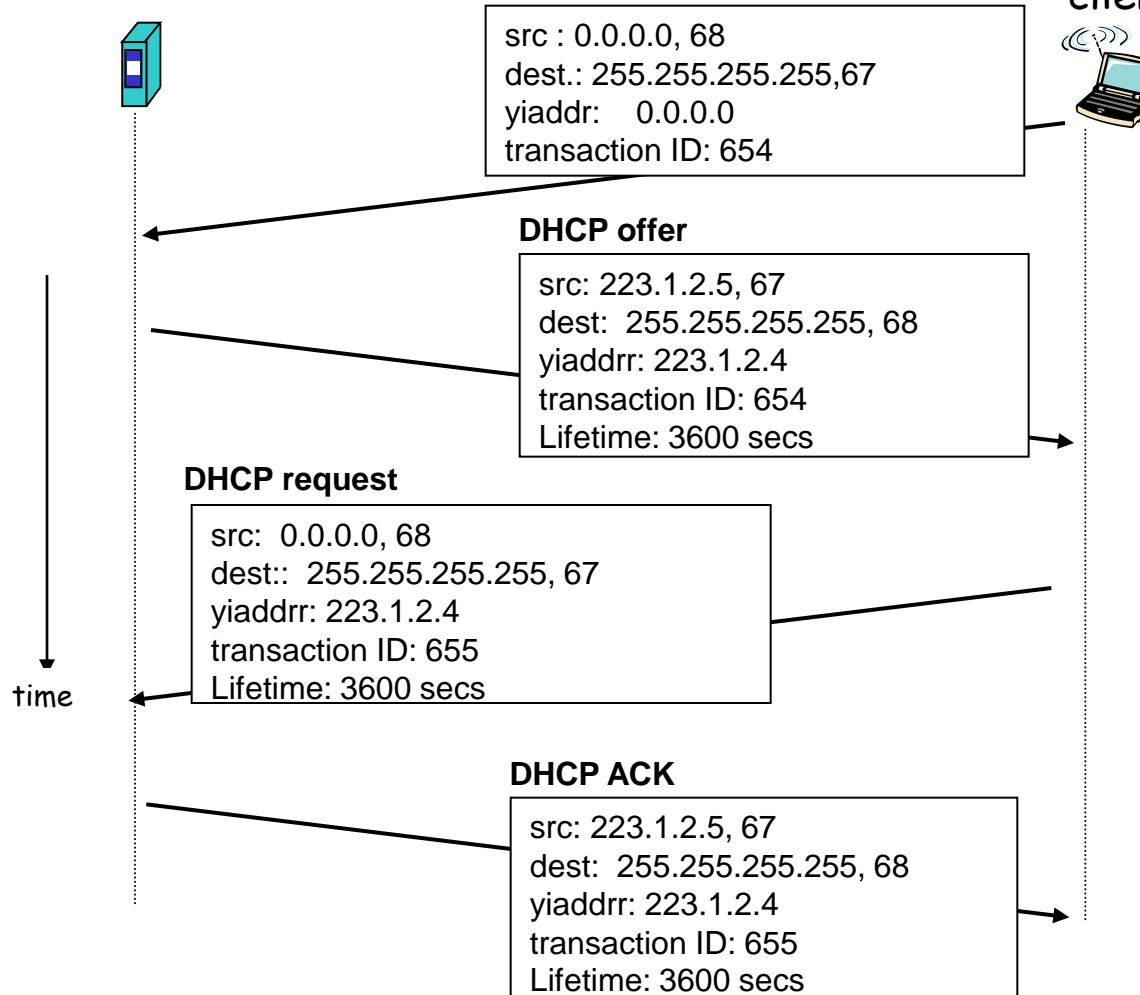Goal: allow host to *dynamically* obtain its IP address from network server when it joins network

- m  Allows reuse of addresses

# 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

time

**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

# 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
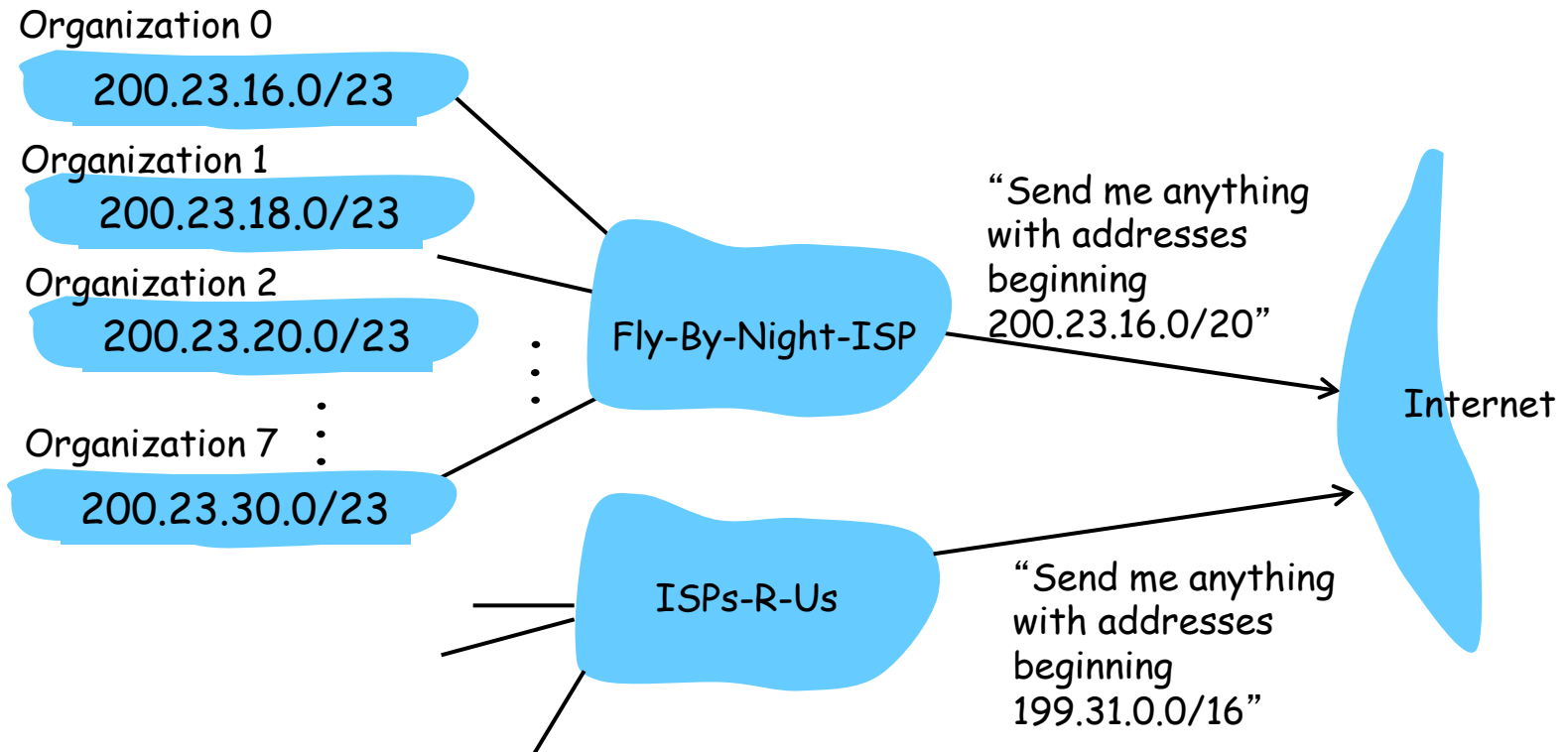
# IP addresses: how to get one?

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

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

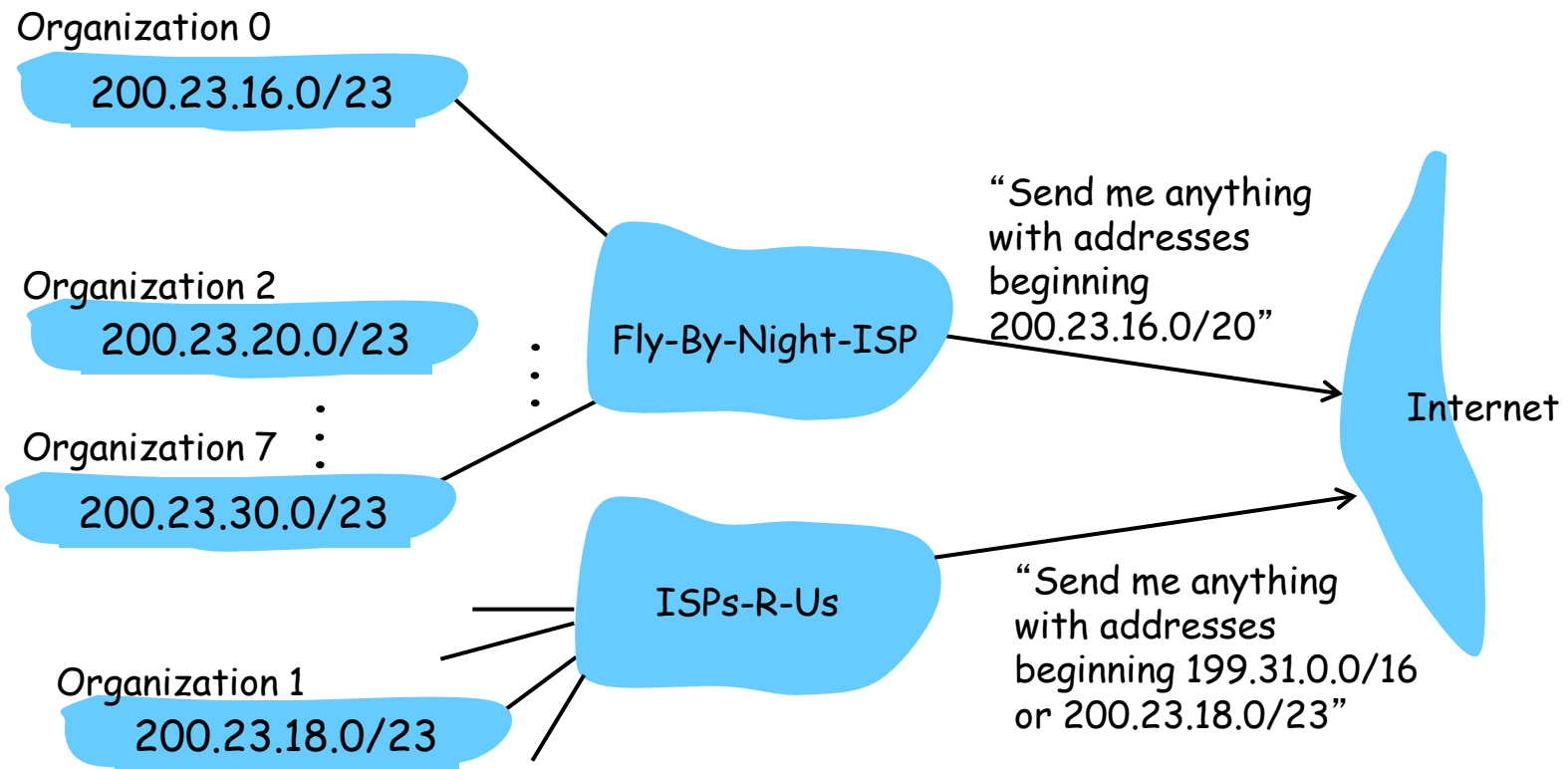| | | | |
|---|---|---|---|
| ISP's block | 11001000  00010111  0001<u>0000</u> | 00000000 | 200.23.16.0/20 |
| | | | |
| Organization 0 | 11001000  00010111  0001000<u>0</u> | 00000000 | 200.23.16.0/23 |
| Organization 1 | 11001000  00010111  0001001<u>0</u> | 00000000 | 200.23.18.0/23 |
| Organization 2 | 11001000  00010111  0001010<u>0</u> | 00000000 | 200.23.20.0/23 |
| ... | ….. | …. | …. |
| Organization 7 | 11001000  00010111  0001111<u>0</u> | 00000000 | 200.23.30.0/23 |

# Hierarchical addressing: route aggregation

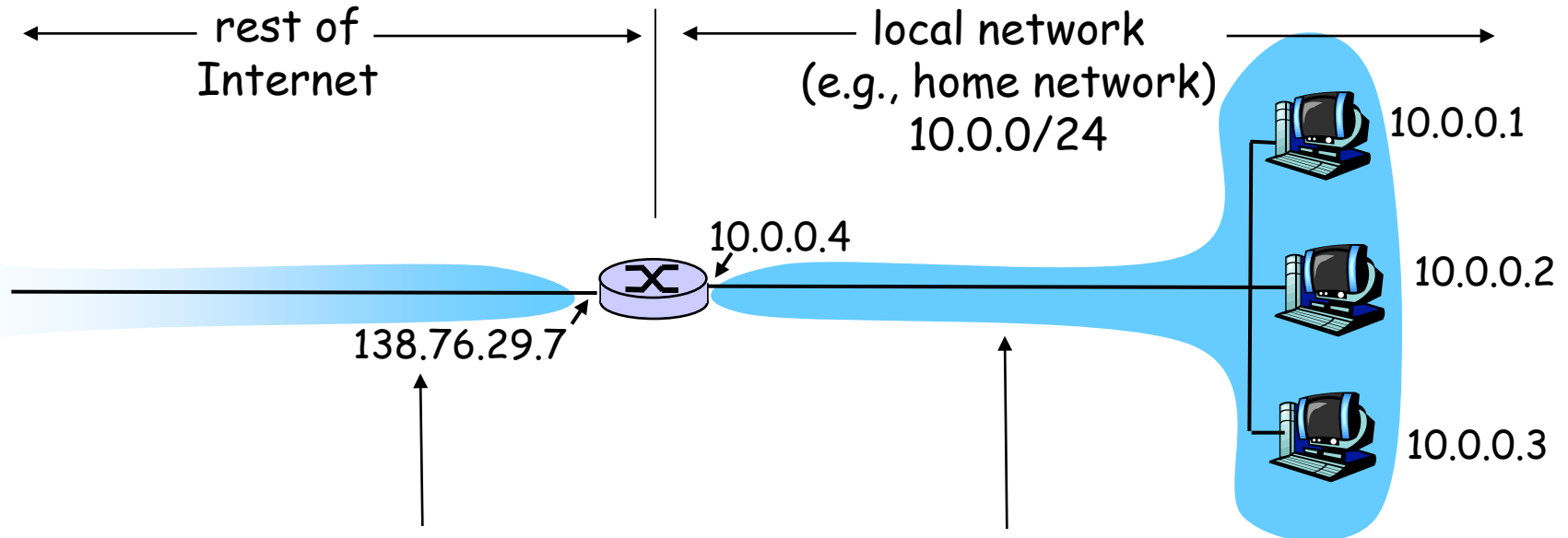Hierarchical addressing allows efficient advertisement of routing information:

Organization 0
200.23.16.0/23

Organization 1
200.23.18.0/23

Organization 2
200.23.20.0/23

Organization 7
200.23.30.0/23

Fly-By-Night-ISP

ISPs-R-Us

"Send me anything with addresses beginning 200.23.16.0/20"

"Send me anything with addresses beginning 199.31.0.0/16"

Internet

# Hierarchical addressing: more specific routes

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



Organization 0
200.23.16.0/23

Organization 2
200.23.20.0/23

Organization 7
200.23.30.0/23

Fly-By-Night-ISP

"Send me anything with addresses beginning 200.23.16.0/20"

Internet

ISPs-R-Us

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

Organization 1
200.23.18.0/23

# 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

r Motivation: local network uses just one IP address as far as outside world is concerned:

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

# NAT: Network Address Translation

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

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

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

S: 10.0.0.1, 3345
D: 128.119.40.186, 80

1

10.0.0.1

S: 138.76.29.7, 5001
D: 128.119.40.186, 80

2

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

S: 128.119.40.186, 80
D: 138.76.29.7, 5001

3

**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**

# Chapter 4: Network Layer

# ICMP: Internet Control Message Protocol

r used by hosts & routers to communicate network-level information

- m error reporting: unreachable host, network, port, protocol

- m echo request/reply (used by ping)

| 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 |

# Chapter 4: Network Layer

# IPv6

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

r Additional motivation:

  m header format helps speed processing/forwarding

  m header changes to facilitate QoS

  IPv6 datagram format:

  m fixed-length 40 byte header

  m no fragmentation allowed

# 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 →

# Transition From IPv4 To IPv6

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

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

# 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:

r   all routers have complete topology, link cost info

r   "link state" algorithms

Decentralized:

r   router knows physically-connected neighbors, link costs to neighbors

r   iterative process of computation, exchange of info with neighbors

r   "distance vector" algorithms

## Static or dynamic?

Static:

r   routes change slowly over time

Dynamic:

r   routes change more quickly

m   periodic update

m   in response to link cost changes

# Chapter 4: Network Layer

# A Link-State Routing Algorithm

## Dijkstra's algorithm

r    net topology, link costs known to all nodes

    m    accomplished via "link state broadcast"

    m    all nodes have same info

r    computes least cost paths from one node ('source") to all other nodes

    m    gives forwarding table for that node

r    iterative: after k iterations, know least cost path to k dest.'s

## Notation:

r    $c(x,y)$: link cost from node x to y;  = ∞ if not direct neighbors

r    $D(v)$: current value of cost of path from source to dest. v

r    $p(v)$: predecessor node along path from source to v

r    $N'$: set of nodes whose least cost path definitively known

# Dijsktra'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'
```
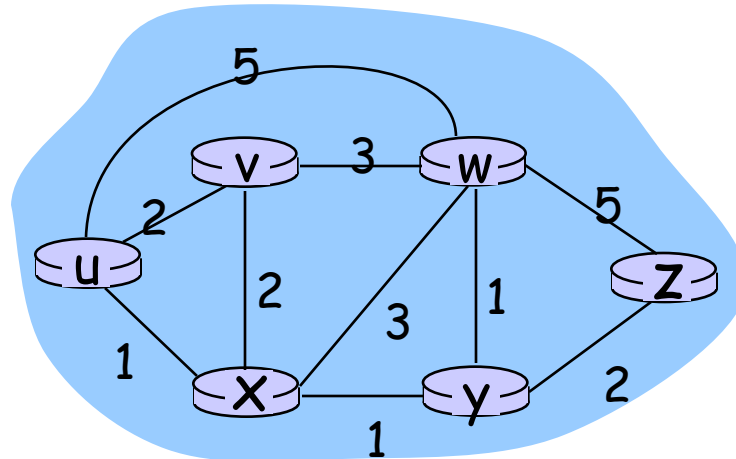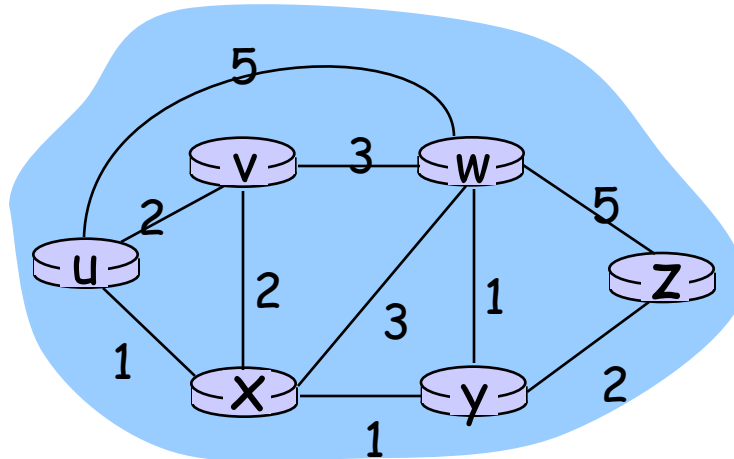
# Dijkstra's algorithm: example

| 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 | | | | | | |
| 3 | | | | | | |
| 4 | | | | | | |
| 5 | | | | | | |

# 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) |

# Chapter 4: Network Layer

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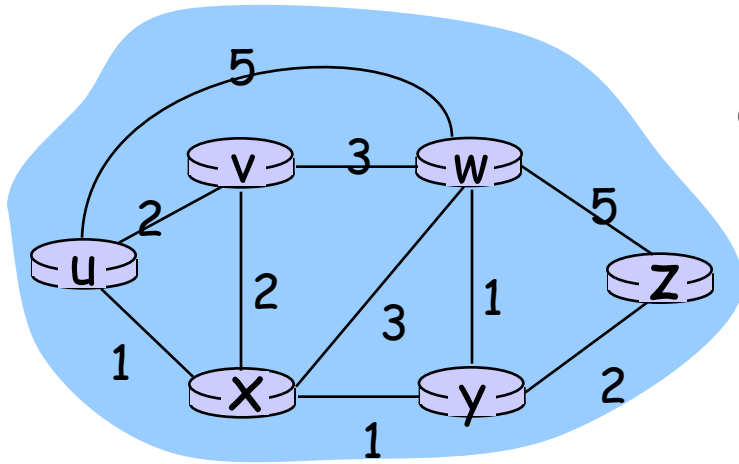
Bellman-Ford Equation (dynamic programming)
Define
$d_x(y)$ := cost of least-cost path from x to y

Then

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

where min is taken over all neighbors v of x

# 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

r  $D_x(y)$ = estimate of least cost from x to y

r  Node x knows cost to each neighbor v: $c(x,v)$

r  Node x maintains  distance vector $D_x$ = $[D_x(y): y \in N ]$

r  Node x also maintains its neighbors' distance vectors

  m  For each neighbor v, x maintains $D_v = [D_v(y): y \in N ]$

# Distance vector algorithm (4)

**Basic idea:**

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

r  Asynchronous

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

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

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

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

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

**node x table**

cost to

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

cost to

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

**node y table**

cost to

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

**node z table**

cost to

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

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

| 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

# Chapter 4: Network Layer

# Hierarchical Routing

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

... *not* true in practice

scale: with 200 million destinations:

- r can't store all dest's in routing tables!
- r routing table exchange would swamp links!

administrative autonomy

- r internet = network of networks
- r each network admin may want to control routing in its own network

# Hierarchical Routing

r  aggregate routers into regions, "autonomous systems" (AS)

r  routers in same AS run same routing protocol

  m  "intra-AS" routing protocol

  m  routers in different AS can run different intra-AS routing protocol

Gateway router

r  Direct link to router in another AS

# Inter-AS tasks

r  suppose router in AS1 receives datagram destined outside of AS1:

  m  router should forward packet to gateway router, but which one?

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!

# Chapter 4: Network Layer

# Intra-AS Routing

r   also known as Interior Gateway Protocols (IGP)

r   most common Intra-AS routing protocols:

   m RIP: Routing Information Protocol

   m OSPF: Open Shortest Path First

   m IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

# RIP ( Routing Information Protocol)

r   distance vector algorithm

r   included in BSD-UNIX Distribution in 1982
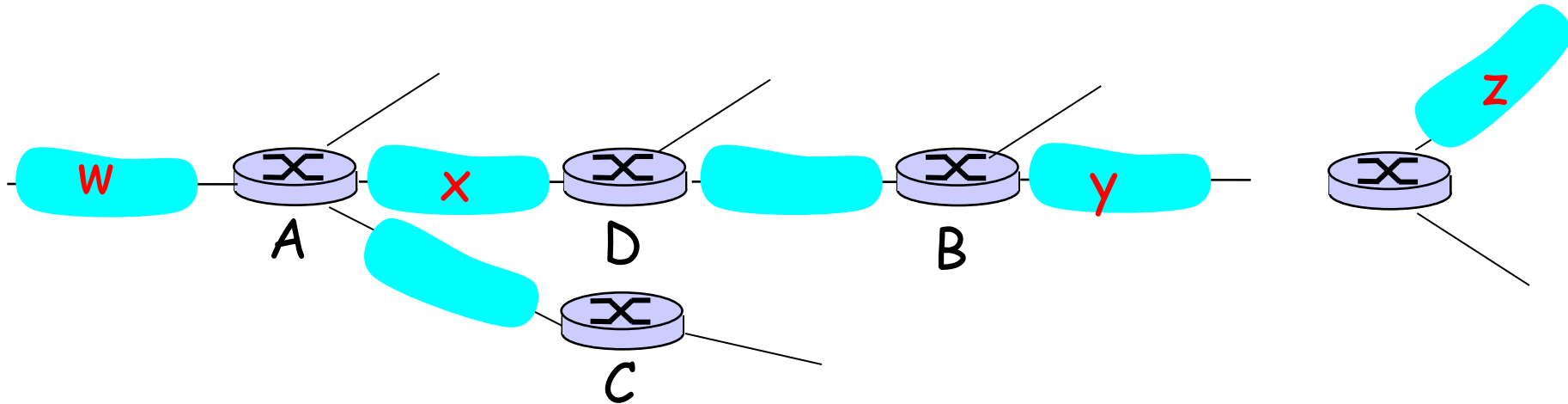
r   distance metric: # of hops (max = 15 hops)



From router A to subnets:

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

# RIP advertisements

r *distance vectors:* exchanged among neighbors every 30 sec via Response Message (also called advertisement)
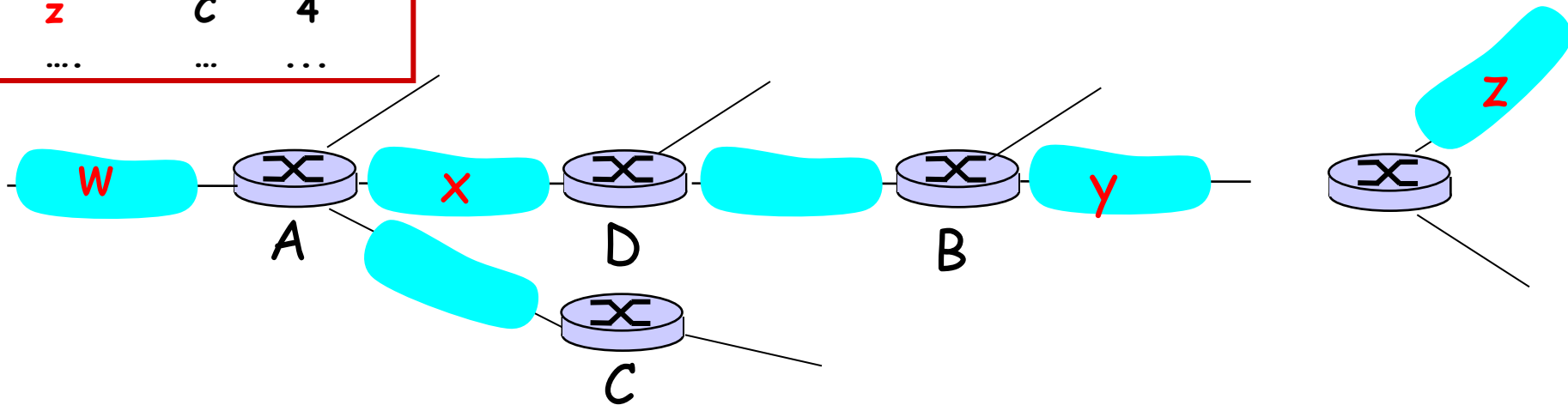
r each advertisement: list of up to 25 destination subnets within AS

# RIP: Example



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

Routing/Forwarding table in D

# RIP: Example

| Dest | Next | hops |
|------|------|------|
| w | – | 1 |
| x | – | 1 |
| z | C | 4 |
| …. | … | … |

Advertisement from A to D



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

Routing/Forwarding table in D

# RIP: Link Failure and Recovery

If no advertisement heard after 180 sec --> neighbor/link declared dead
- m routes via neighbor invalidated
- m new advertisements sent to neighbors
- m neighbors in turn send out new advertisements (if tables changed)
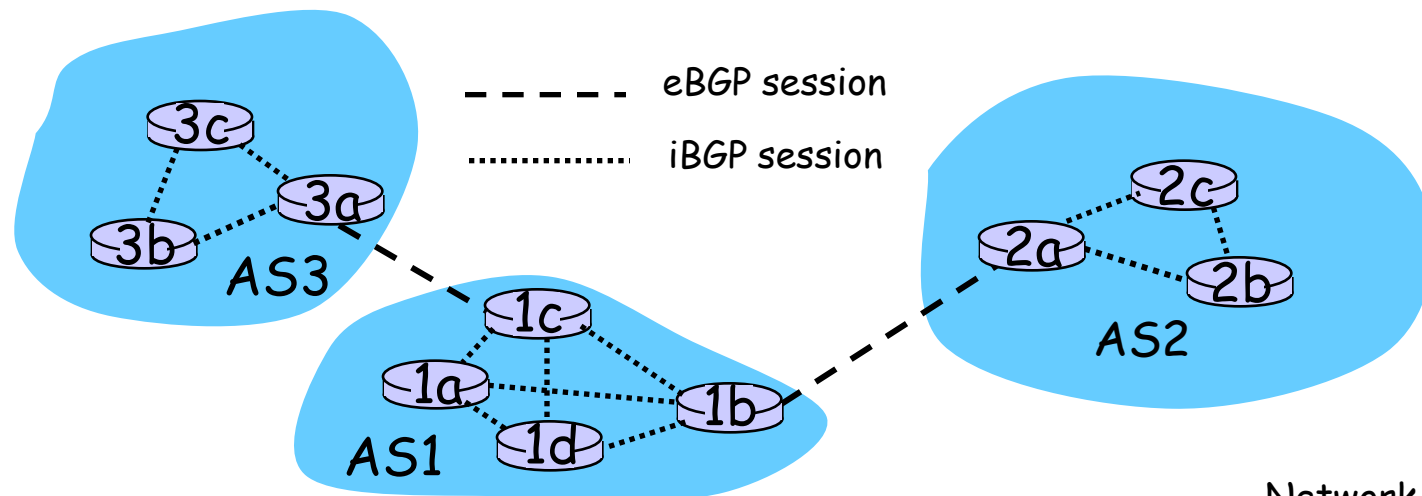- m link failure info propagates quickly to entire net

# Chapter 4: Network Layer

r 4. 1 Introduction

r 4.2 Virtual circuit and datagram networks

r 4.3 What's inside a router

r 4.4 IP: Internet Protocol

- m Datagram format
- m IPv4 addressing
- m ICMP
- m IPv6

r 4.5 Routing algorithms

- m Link state
- m Distance Vector
- m Hierarchical routing

r 4.6 Routing in the Internet

- m RIP
- m OSPF
- m BGP

r 4.7 Broadcast and multicast routing

# OSPF "advanced" features (not in RIP)

r   security: all OSPF messages authenticated (to prevent malicious intrusion)

r   multiple same-cost paths allowed (only one path in RIP)

r   integrated uni- and multicast support:

   m  Multicast OSPF (MOSPF) uses same topology data base as OSPF

r   hierarchical OSPF in large domains.

# Hierarchical OSPF

# Hierarchical OSPF

r  two-level hierarchy: local area, backbone.

   m Link-state advertisements only in area

   m each nodes has detailed area topology; only know direction (shortest path) to nets in other areas.

r  *area border routers:* "summarize" distances to nets in own area, advertise to other Area Border routers.

r  *backbone routers:* run OSPF routing limited to backbone.

r  *boundary routers:* connect to other AS's.

# Chapter 4: Network Layer

# Internet inter-AS routing: BGP

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

r 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.

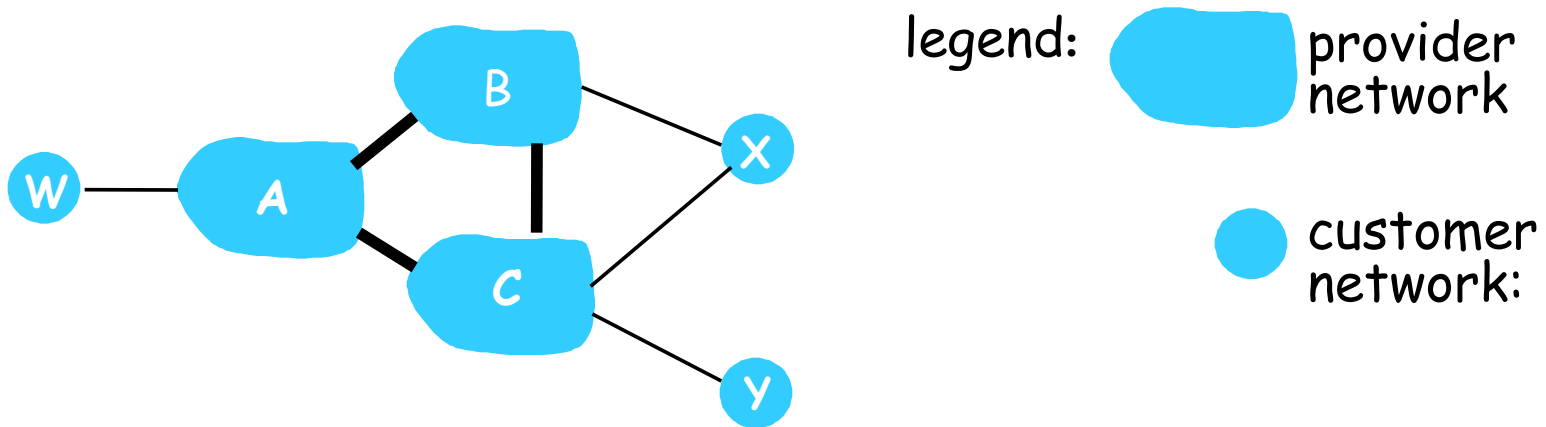r allows subnet to advertise its existence to rest of Internet: *"I am here"*

# BGP basics

r   pairs of routers (BGP peers) exchange routing info over semi-permanent TCP connections: BGP sessions

   m   BGP sessions need not correspond to physical links.

r   when AS2 advertises a prefix to AS1:

   m   AS2 *promises* it will forward datagrams towards that prefix.

   m   AS2 can aggregate prefixes in its advertisement



- - - - -   eBGP session

· · · · · · · · ·   iBGP session

# BGP route selection

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

r    elimination rules:

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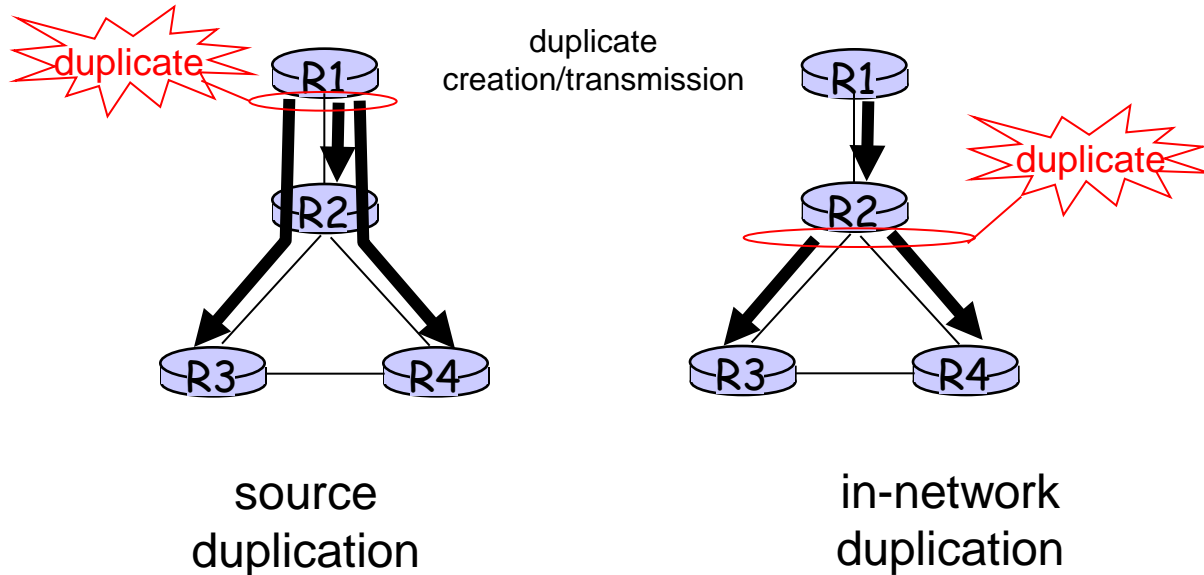4. additional criteria

# BGP routing policy



legend:

provider network

customer network:

r   A,B,C are provider networks

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

r   X is dual-homed: attached to two networks

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

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

# Chapter 4: Network Layer

# Broadcast Routing

r deliver packets from source to all other nodes

r source duplication is inefficient:
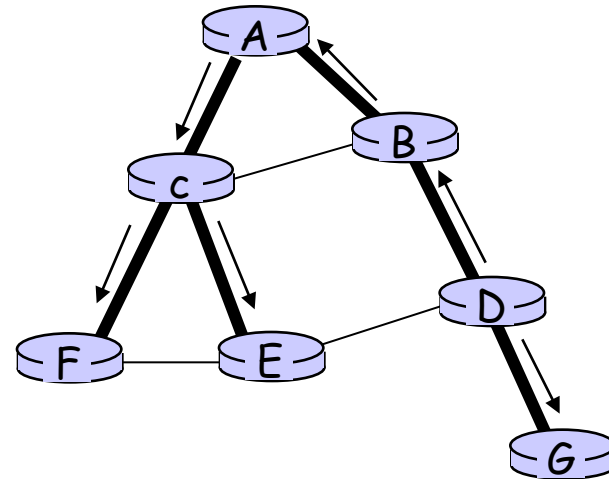


source
duplication

in-network
duplication

r source duplication: how does source determine recipient addresses?

# Spanning Tree

r  First construct a spanning tree
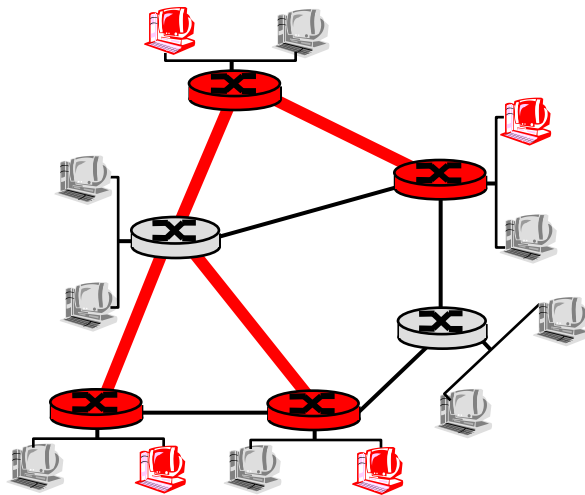r  Nodes forward copies only along spanning tree
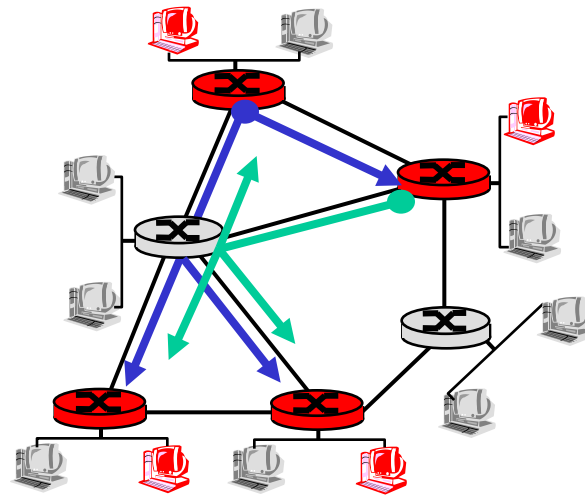


(a) Broadcast initiated at A          (b) Broadcast initiated at D

# Multicast Routing: Problem Statement

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

  m  _tree:_ not all paths between routers used

  m  _source-based:_ different tree from each sender to rcvrs

  m  _shared-tree:_ same tree used by all group members



Shared tree                    Source-based trees

# Chapter 4: summary

r 4. 1 Introduction

r 4.2 Virtual circuit and datagram networks

r 4.3 What's inside a router

r 4.4 IP: Internet Protocol
  m Datagram format
  m IPv4 addressing
  m ICMP
  m IPv6

r 4.5 Routing algorithms
  m Link state
  m Distance Vector
  m Hierarchical routing

r 4.6 Routing in the Internet
  m RIP
  m OSPF
  m BGP

r 4.7 Broadcast and multicast routing