# Chapter 4
# Network Layer

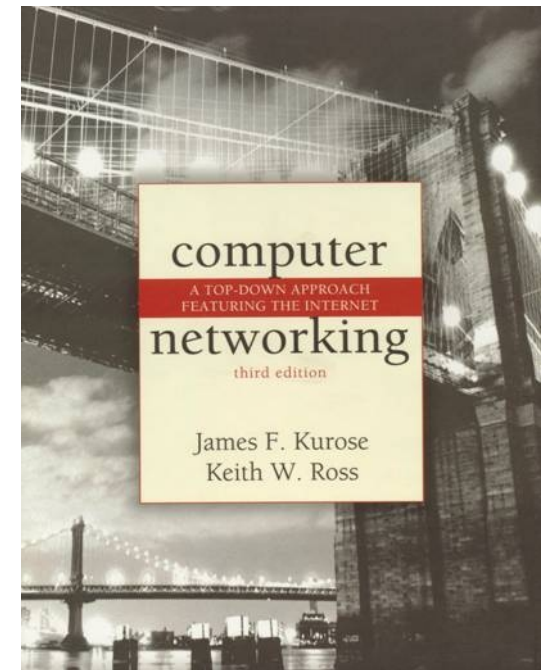## A note on the use of these ppt slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

❑ If you use these slides (e.g., in a class) in substantially unaltered form, that you mention their source (after all, we'd like people to use our book!)
❑ If you post any slides in substantially unaltered form on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy!  JFK/KWR

*Computer Networking: A Top Down Approach Featuring the Internet,*

3rd edition.
Jim Kurose, Keith Ross
Addison-Wesley, July 2004.

# Chapter 4: Network Layer
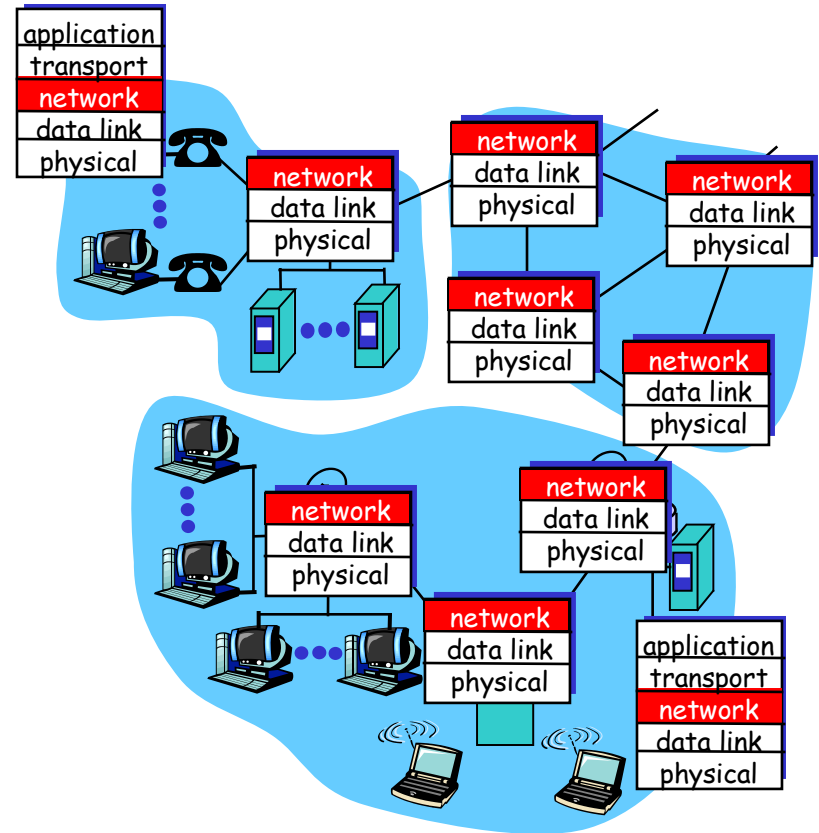
## Chapter goals:

- understand principles behind network layer services:
  - routing (path selection)
  - dealing with scale
  - how a router works
  - advanced topics: IPv6, mobility
- instantiation and 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 receiving 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 destination
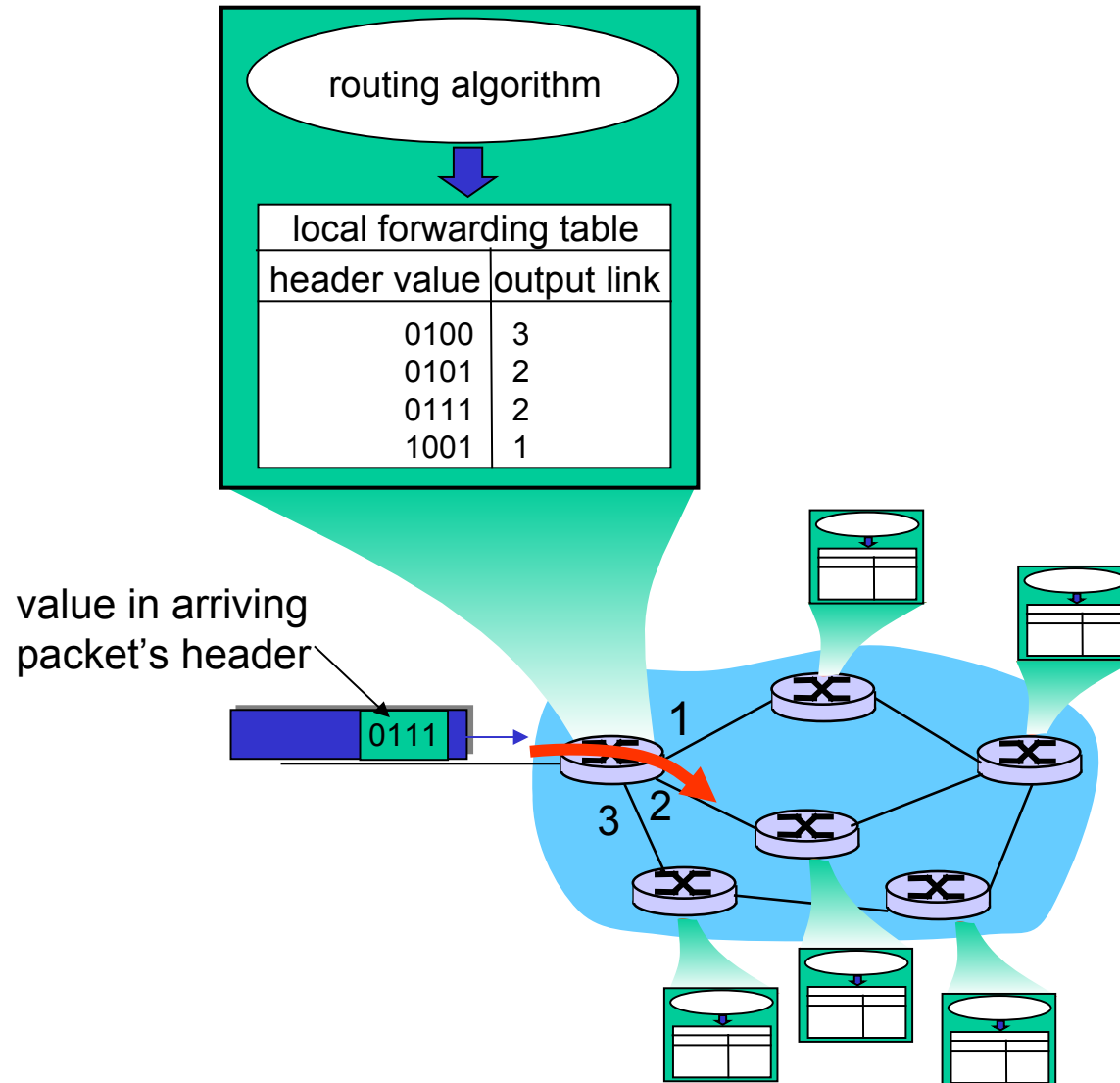
  - *Routing algorithms*

analogy:

- routing: process of planning trip from source to destination

- 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

# Connection setup

- 3rd important function in *some* network architectures:
  - ATM, frame relay, X.25
- Before datagrams flow, two hosts and intervening routers establish virtual connection
  - Routers get involved
- Network and transport layer connection service:
  - Network: between two hosts
  - Transport: between two processes

# Network service model

Q: What *service model* for "channel" transporting datagrams from sender to receiver?

## Example services for individual datagrams:

- guaranteed delivery
- Guaranteed delivery with less than 40 msec delay

## Example services for a flow of datagrams:

- In-order datagram delivery
- Guaranteed minimum bandwidth to flow
- Restrictions on changes in inter-packet spacing

# Network layer service models:

| Network Architecture | Service Model | Guarantees ? | | | | Congestion feedback |
|---|---|---|---|---|---|---|
| | | Bandwidth | Loss | Order | Timing | |
| Internet | best effort | none | no | no | no | no (inferred via loss) |
| ATM | CBR (constant) | constant rate | yes | yes | yes | no congestion |
| ATM | VBR (variable) | guaranteed rate | yes | yes | yes | no congestion |
| ATM | ABR (available) | guaranteed minimum | no | yes | no | yes |
| ATM | UBR (unspecific) | none | no | yes | no | no |

# 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

# 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

<div style="border: 2px solid red; padding: 10px;">

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

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

</div>

- call setup 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
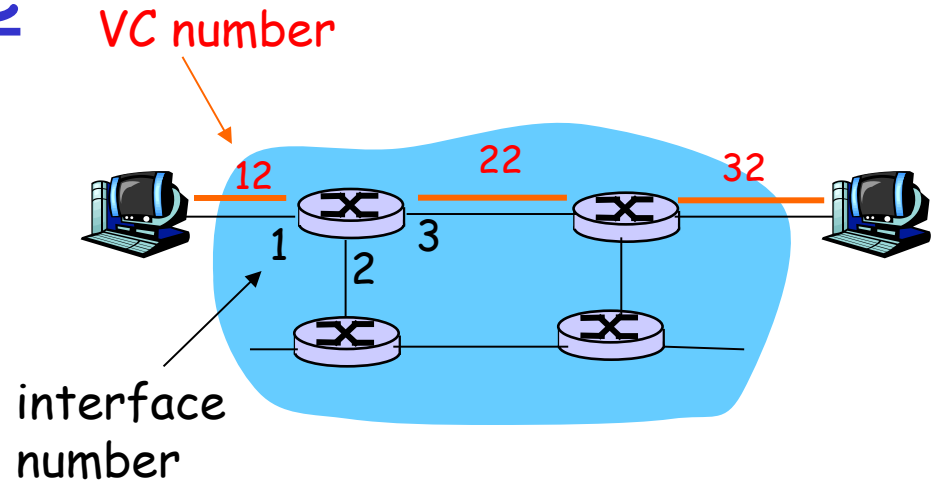- link, router resources (bandwidth, buffers) may be *allocated* to VC
- VC teardown

# 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 a VC number.

- VC number must be changed on each link.
  - New VC number comes from forwarding table
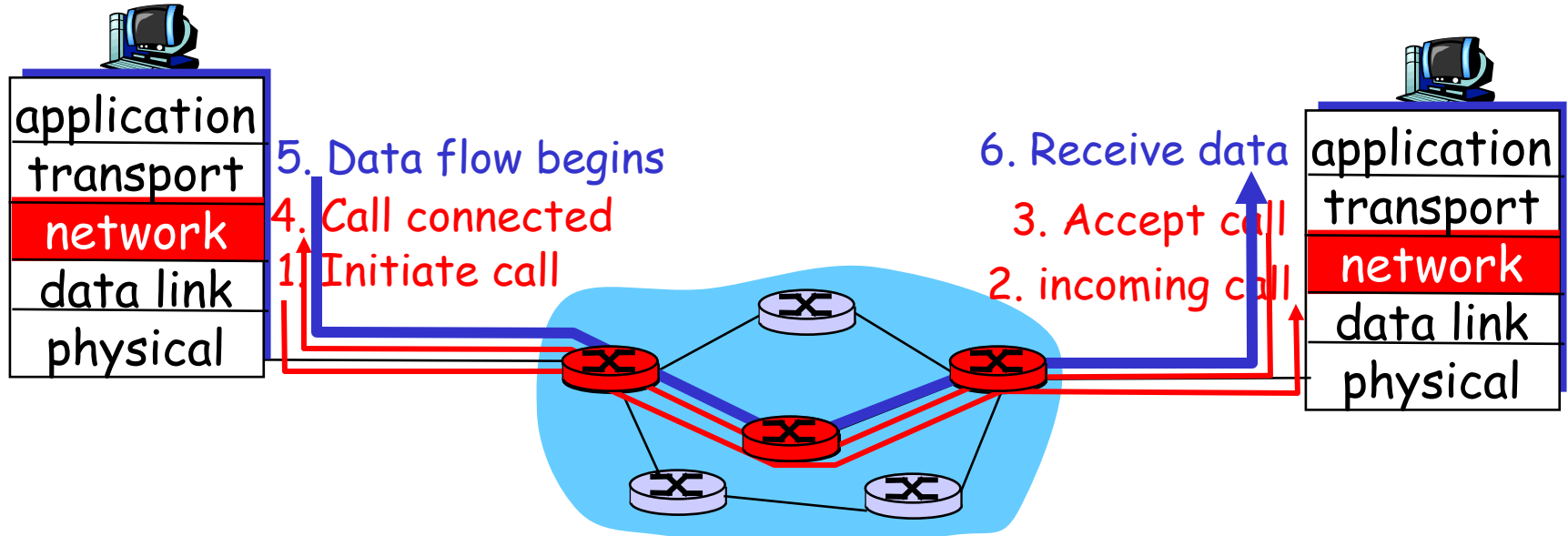
# Forwarding table

VC number



interface number

## Forwarding table in northwest router:

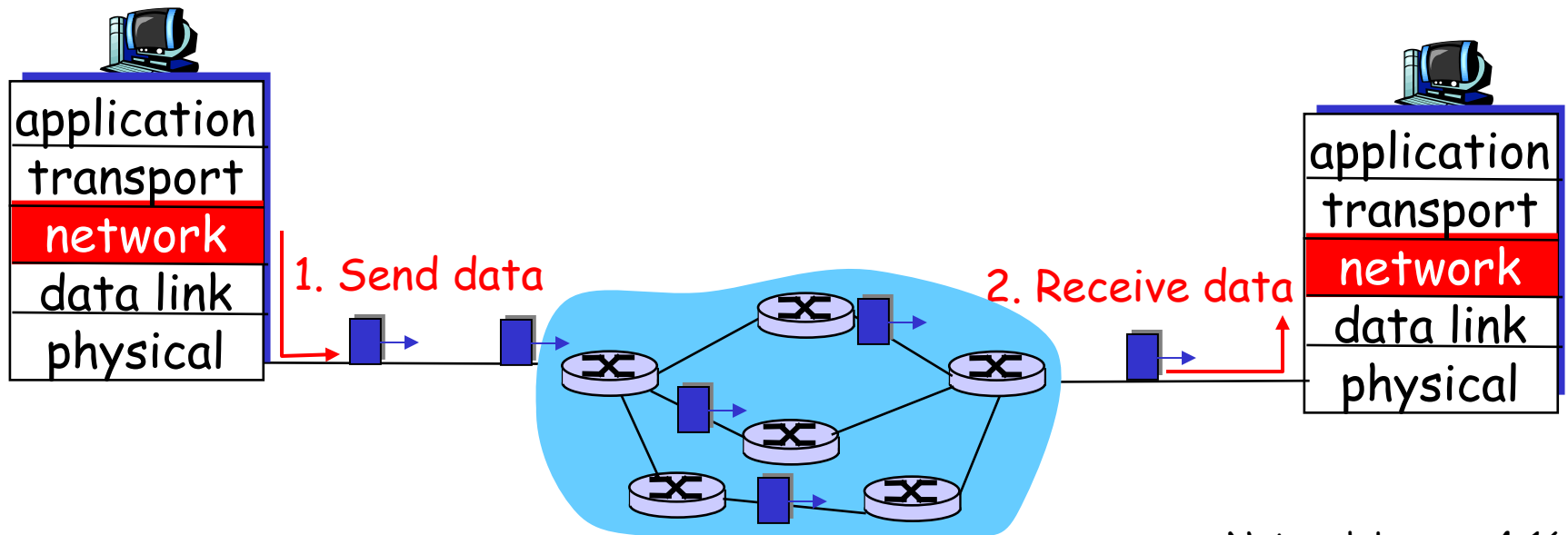| Incoming interface | Incoming VC # | Outgoing interface | Outgoing VC # |
|---|---|---|---|
| 1 | 12 | 3 | 22 |
| 2 | 63 | 1 | 18 |
| 3 | 7 | 2 | 17 |
| 1 | 97 | 3 | 87 |
| … | … | … | … |

Routers maintain connection state information!

# Virtual circuits: signaling protocols

- used to setup, maintain  teardown VC
- used in ATM, frame-relay, X.25
- 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

# 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



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?

# Datagram or VC network: why?

## Internet

- data exchange among computers
  - "elastic" service, no strict timing req.
- "smart" end systems (computers)
  - can adapt, perform control, error recovery
  - simple inside network, complexity at "edge"
- Easier to interconnect networks that use different link types (satellite, Ethernet, fiber, radio)
  - different characteristics
  - uniform service difficult

## ATM

- evolved from telephony
- human conversation:
  - strict timing, reliability requirements
  - need for guaranteed service
- "dumb" end systems
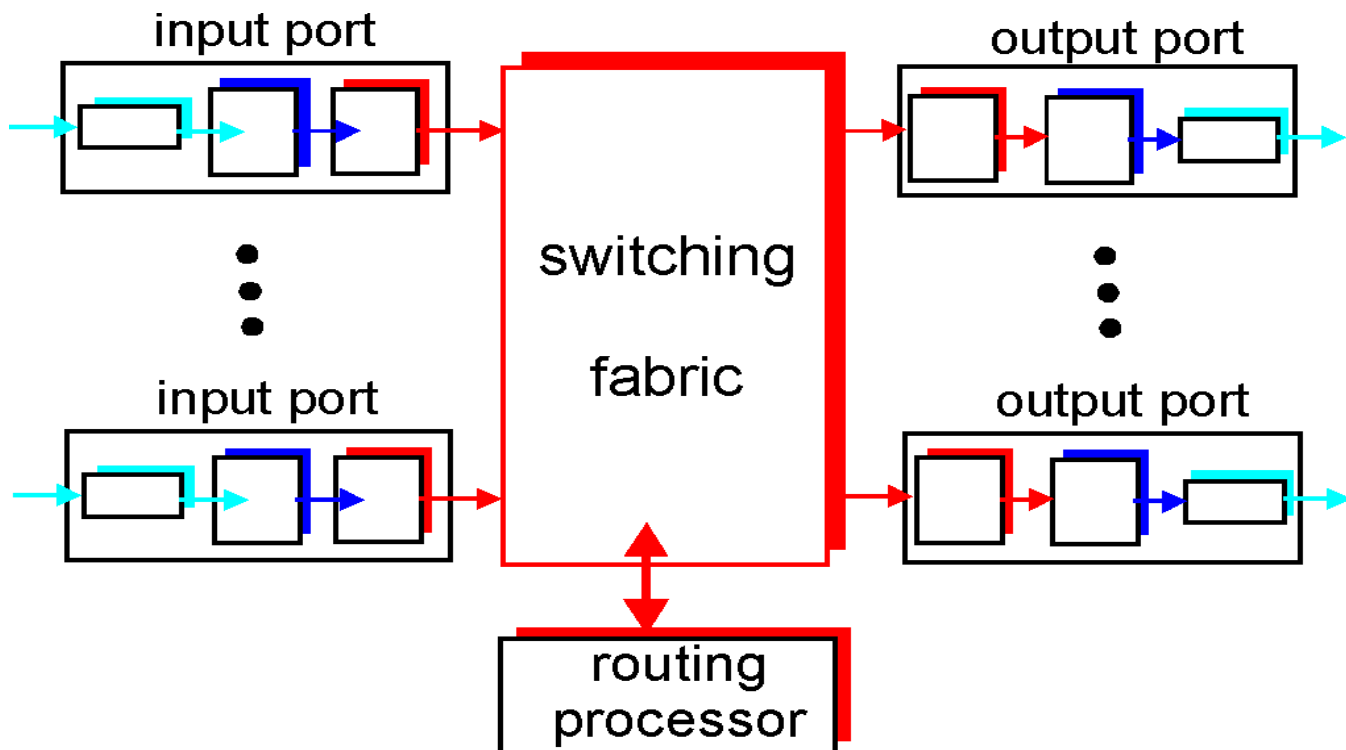  - telephones
  - complexity inside network

# 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
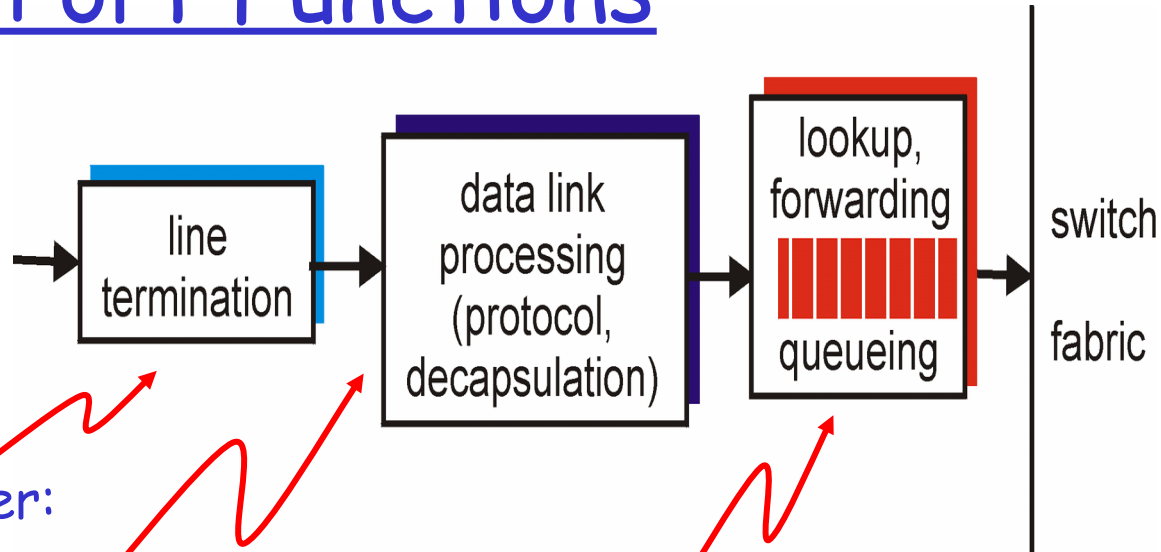- 4.7 Broadcast and multicast routing

# Router Architecture Overview

Two key router functions:

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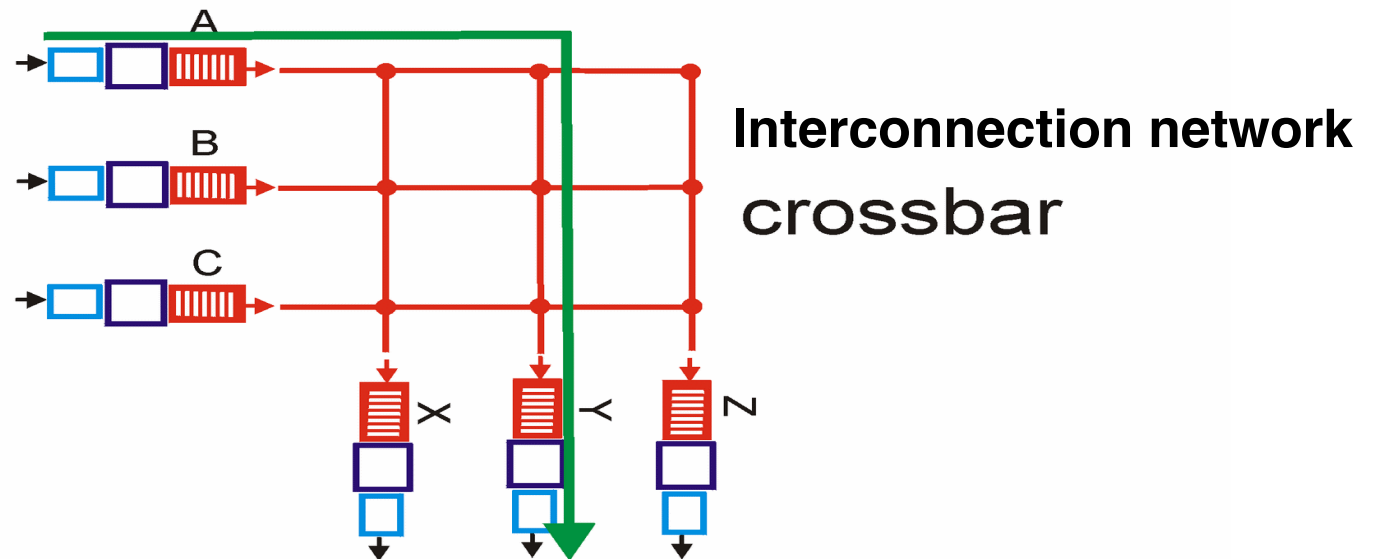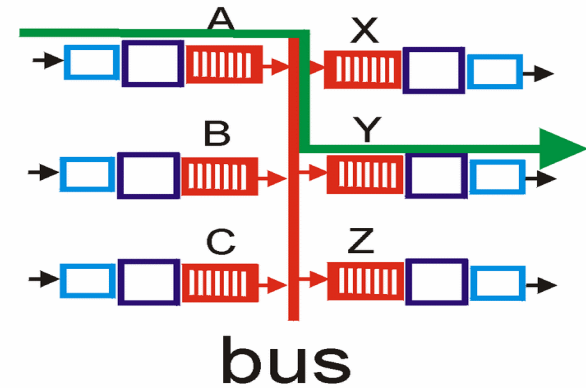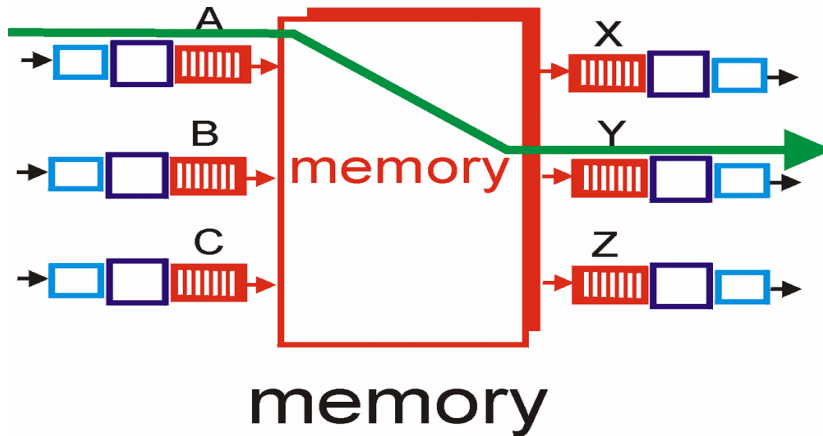
# Input Port Functions



**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'
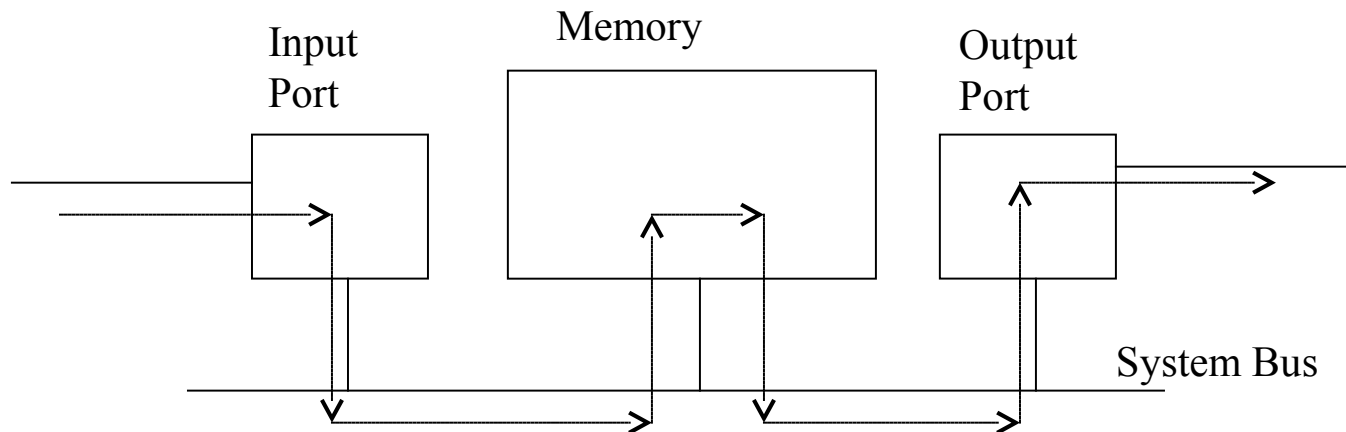- queuing: if datagrams arrive faster than forwarding rate into switch fabric

# Three types of switching fabrics

memory

bus

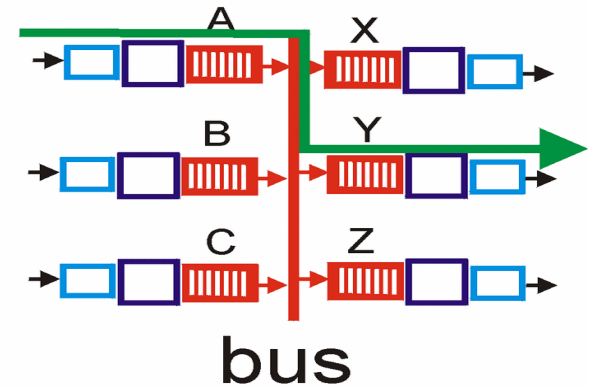**Interconnection network**

crossbar

# Switching Via Memory

First generation routers:

☐ traditional computers with switching under direct control of CPU

☐ packet copied to system's memory

☐ speed limited by memory bandwidth (2 bus crossings per datagram)

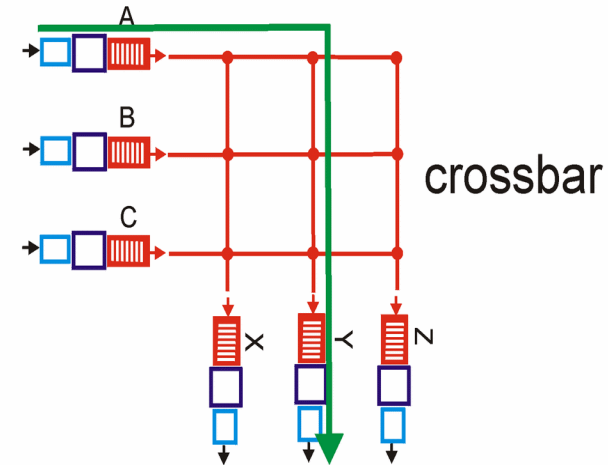Input Port          Memory          Output Port

System Bus

# Switching Via a Bus



bus

- datagram from input port memory to output port memory via a shared bus

- bus contention: switching speed limited by bus bandwidth

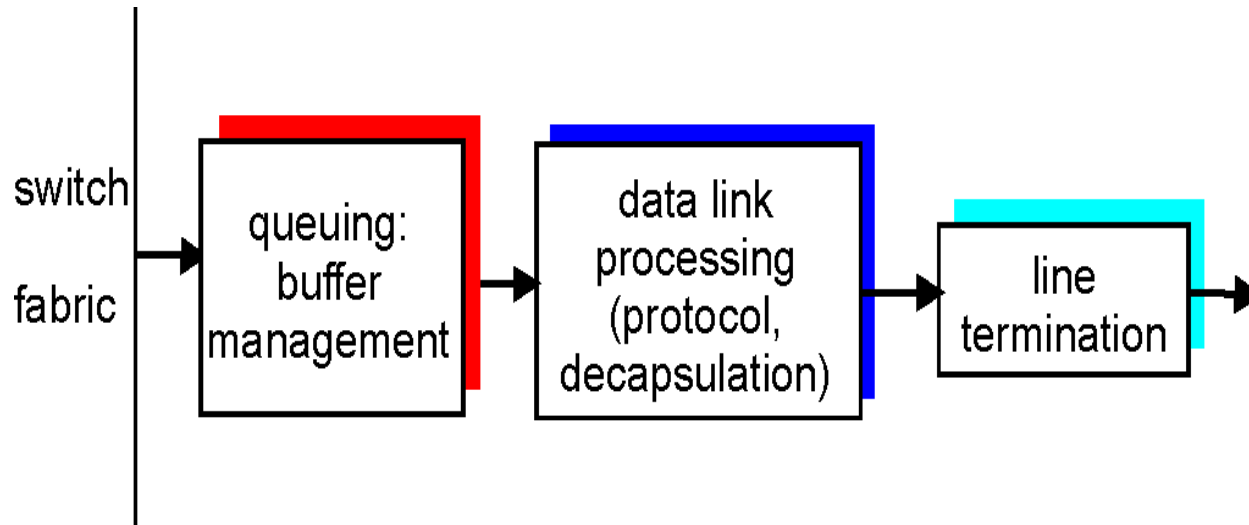- 1 Gbps bus, Cisco 1900: sufficient speed for access and enterprise routers (not regional or backbone)

# Switching Via An Interconnection Network

Interconnection network



crossbar

- overcome  bus bandwidth limitations
- Banyan networks, other interconnection networks initially developed to connect processors in multiprocessor
- Advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- Cisco 12000: switches Gbps through the interconnection network

# Output Ports



- *Buffering* required when datagrams arrive from fabric faster than the transmission rate
- *Scheduling discipline* chooses among queued datagrams for transmission

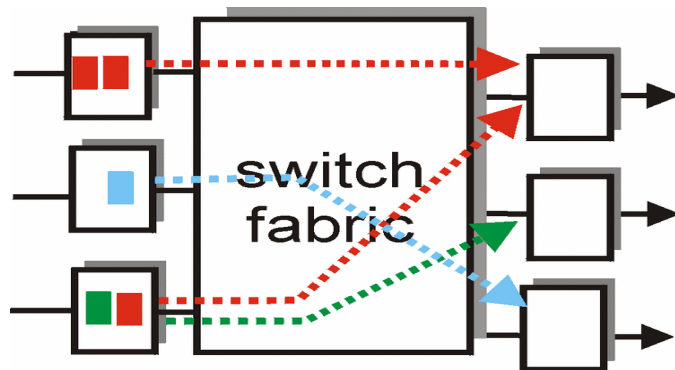# Output port queueing



Output Port Contention at Time $t$

One Packet Time Later

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

# Input Port Queuing

- Fabric slower than 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

# 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:

Three major components: IP protocol, Routing protocols, ICMP protocol
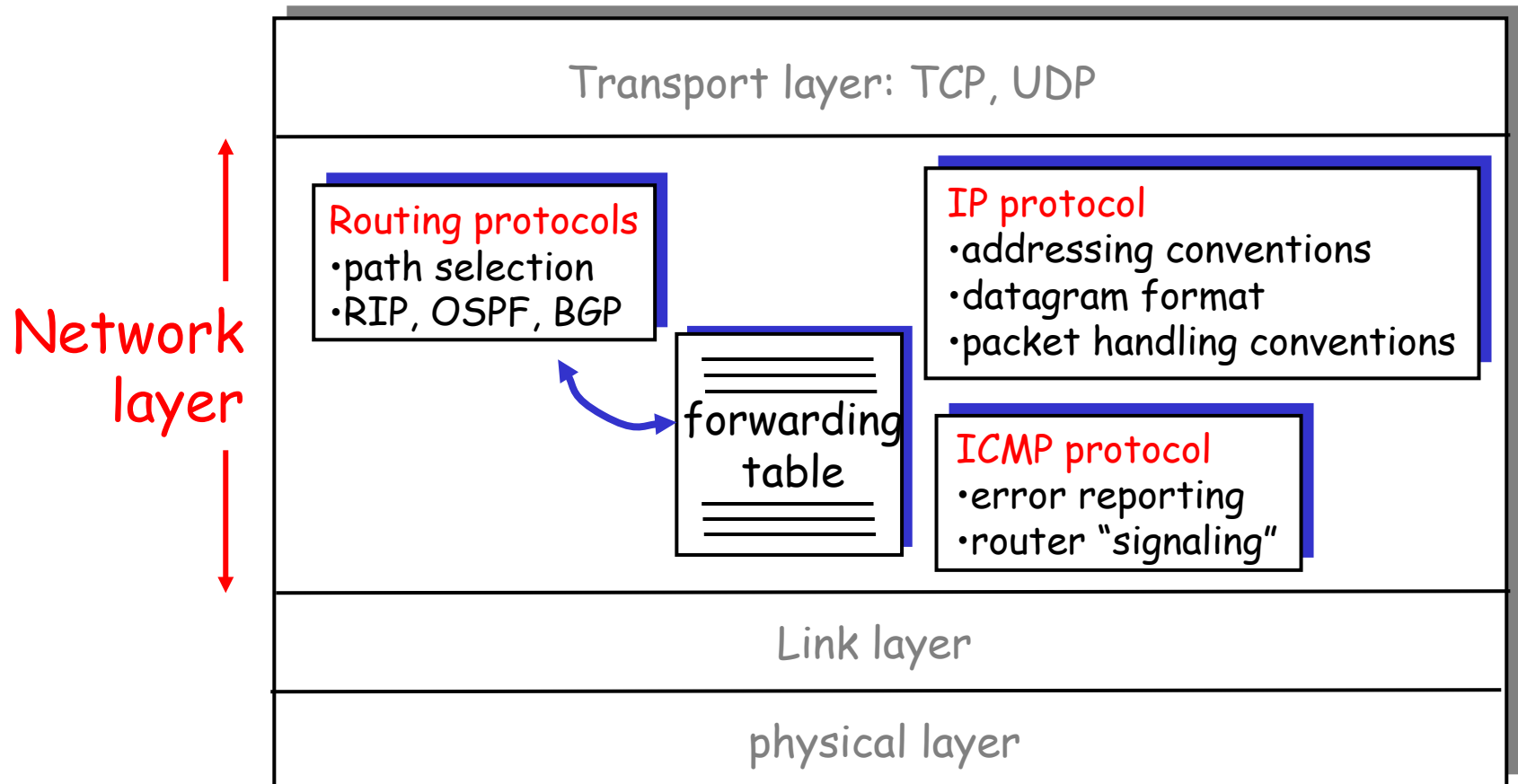
# Chapter 4: Network Layer

- 4. 1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- <span style="color:red">4.4 IP: Internet Protocol</span>
  - <span style="color:red">Datagram format</span>
  - 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
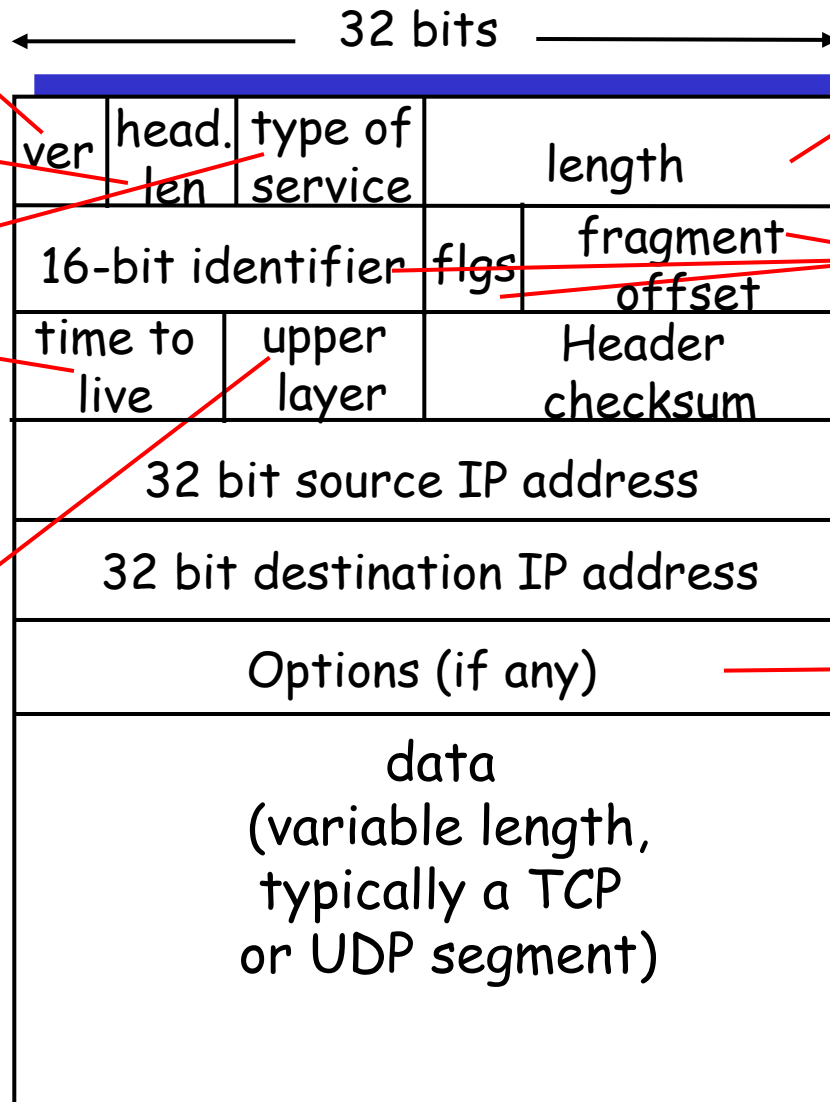
# IP datagram format

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

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) | | | | |

## 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.
    - The links on a route can use different link types, different link-layer protocols, and different MTUs
- large IP datagram are divided ("fragmented") within network
    - 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

Three fields for fragmentation and reassembly: identifier, flag, fragmentation offset

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

One large datagram becomes several smaller datagrams

Identifier: created by sender, all fragments have the same identification number as the original datagram

Flag = 1 there is more fragment

Flag = 0 this is the last fragment

Offset: byte number of the 1st byte of the fragment (specified in units of 8-byte chunks)

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

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

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

## Example

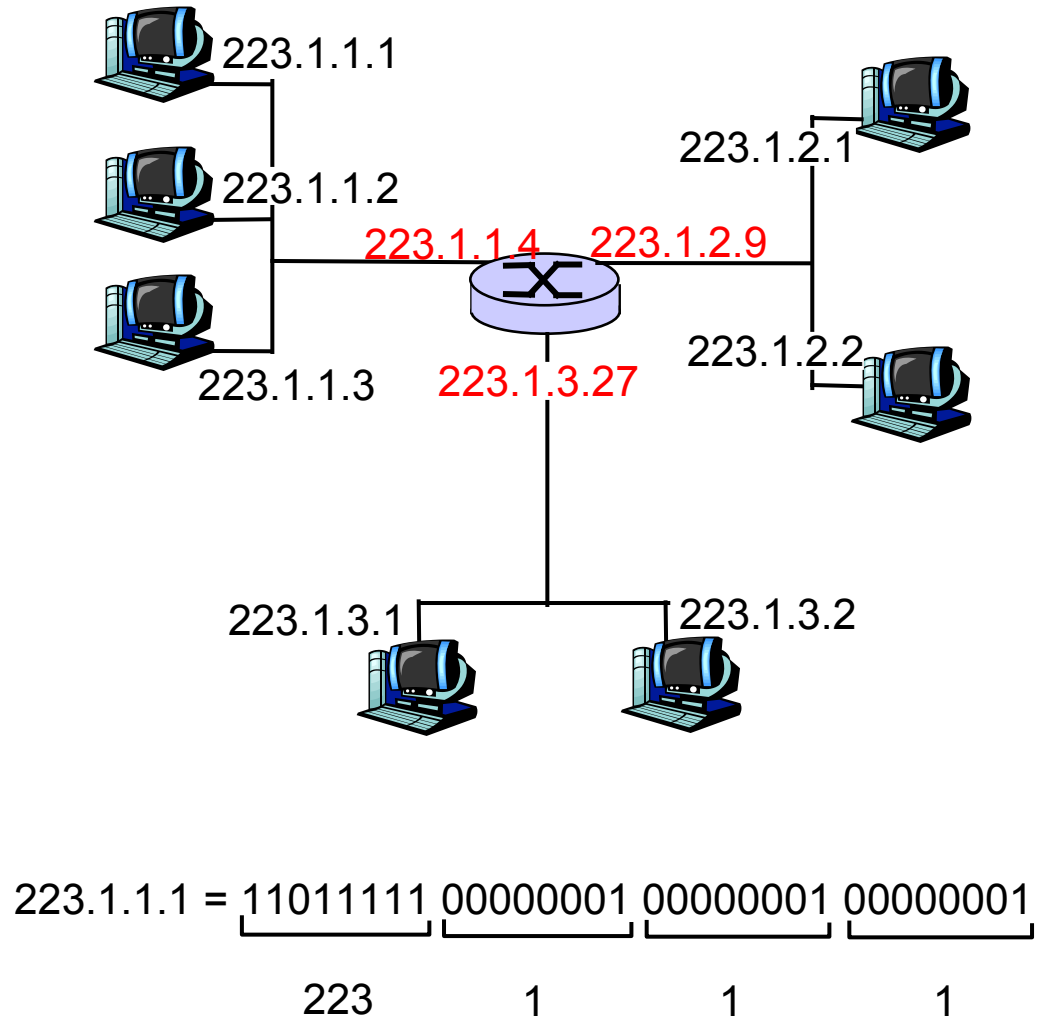- 4000 byte datagram
- MTU = 1500 bytes

1480 bytes in data field

offset = 1480/8

# 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
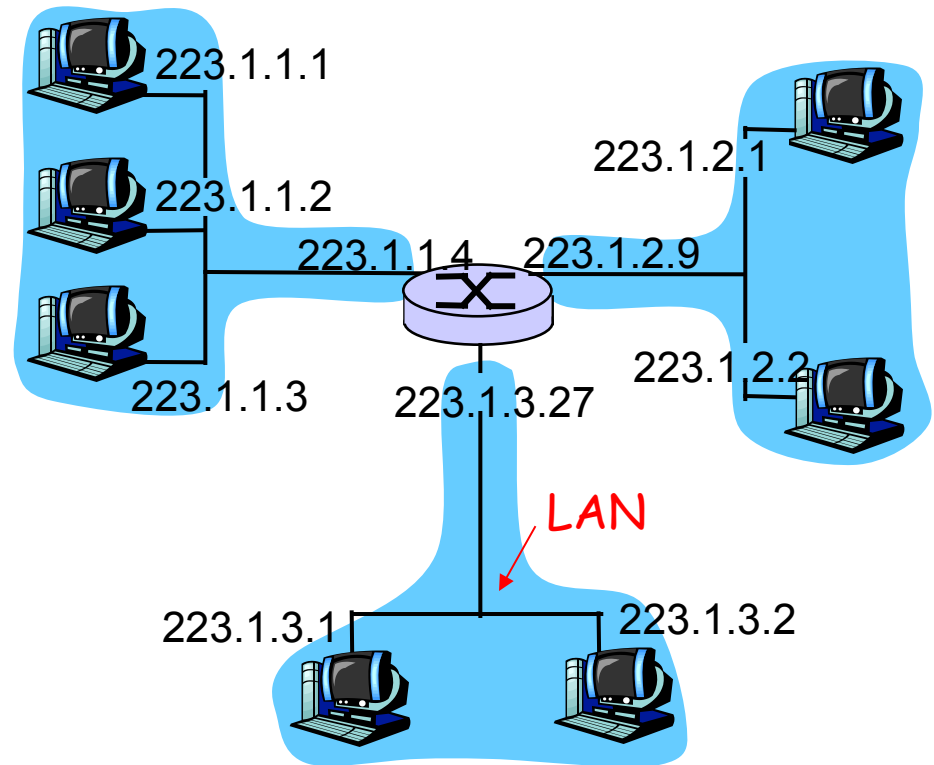- 4.7 Broadcast and multicast routing

# 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 may have multiple interfaces
  - IP addresses associated with each interface

223.1.1.1

223.1.1.2

223.1.1.4    223.1.2.9

223.1.1.3    223.1.3.27

223.1.2.1

223.1.2.2

223.1.3.1    223.1.3.2

223.1.1.1 = 11011111 00000001 00000001 00000001

    223        1        1        1

# Subnets

- IP address:
  - subnet part (high order bits)
  - host part (low order bits)
- *What's a subnet ?*
  - device interfaces with same subnet part of IP address
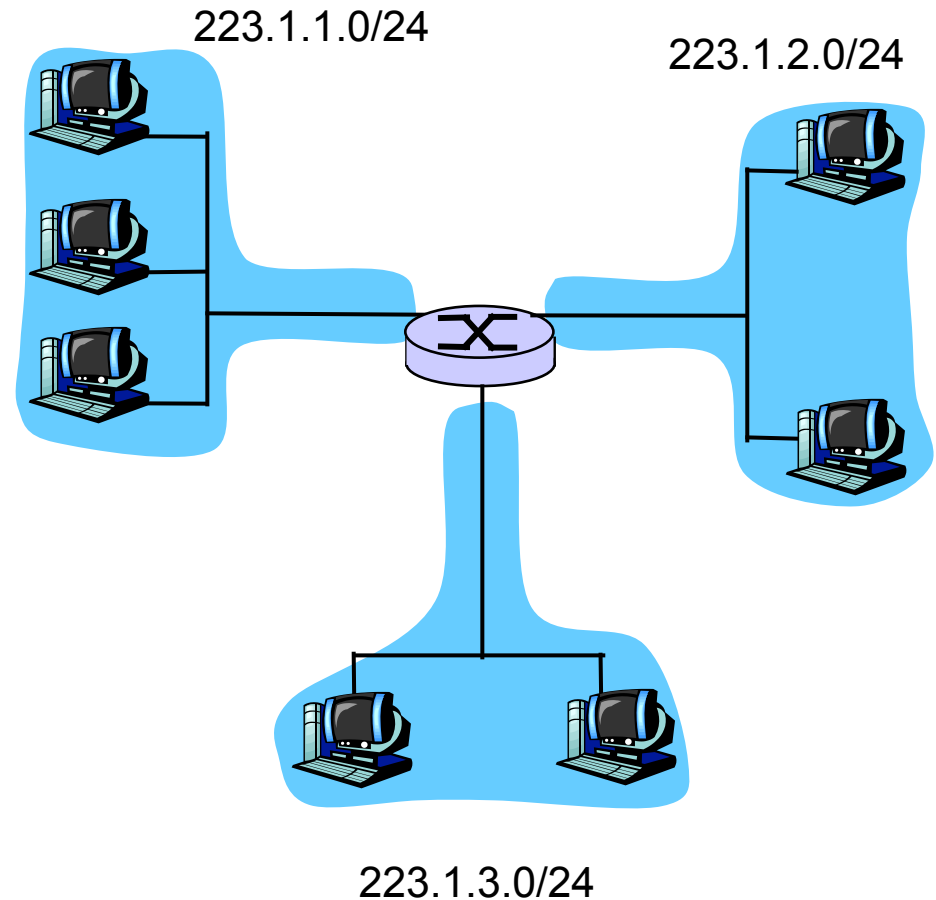  - can physically reach each other without intervening router

223.1.1.1

223.1.1.2

223.1.1.4    223.1.2.9

223.1.1.3    223.1.3.27

223.1.2.1

223.1.2.2

LAN

223.1.3.1    223.1.3.2

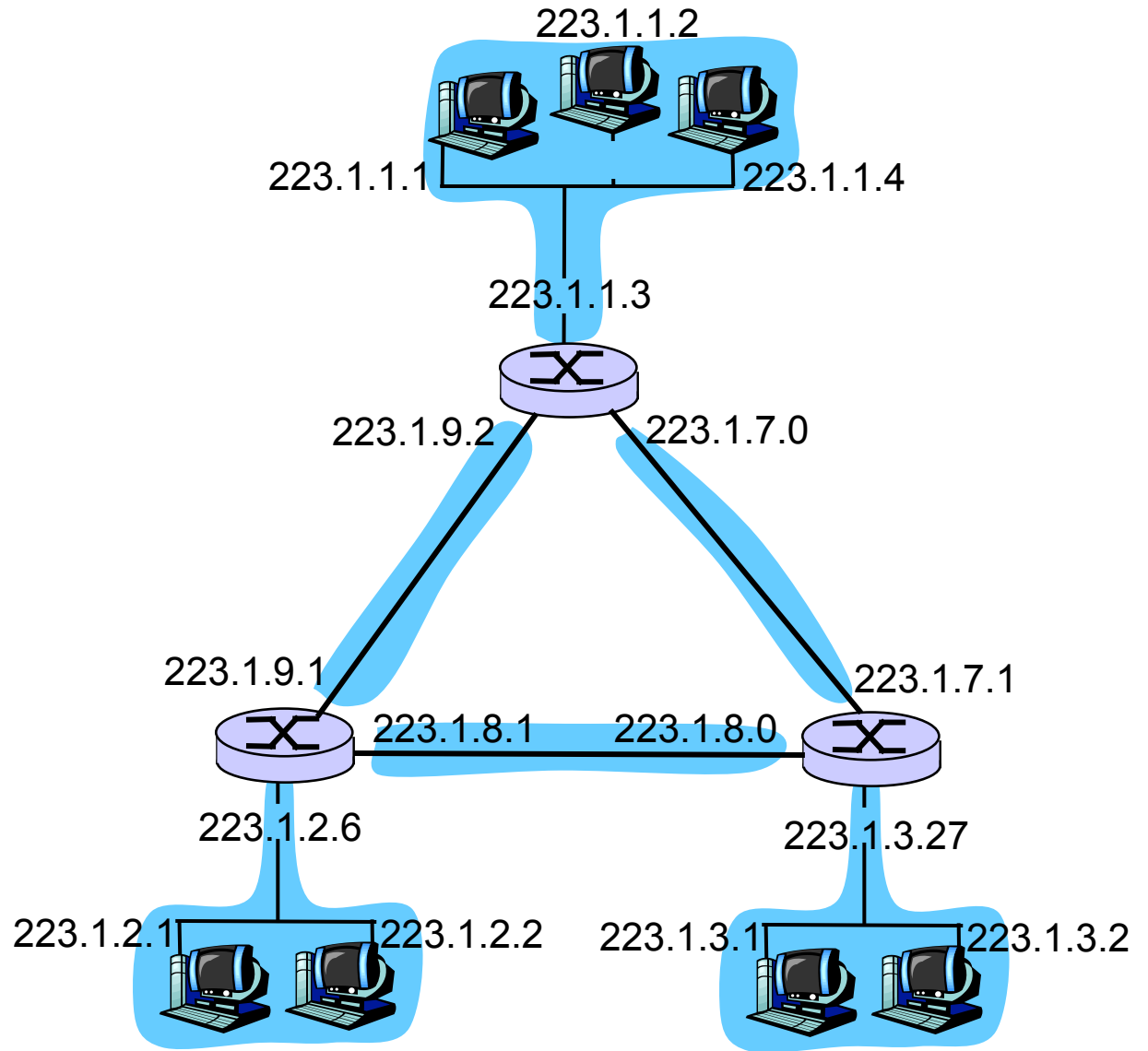network consisting of 3 subnets

# Subnets

223.1.2.0/24

## 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.

223.1.3.0/24

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                      host
         part                       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 a server
  - "plug-and-play"
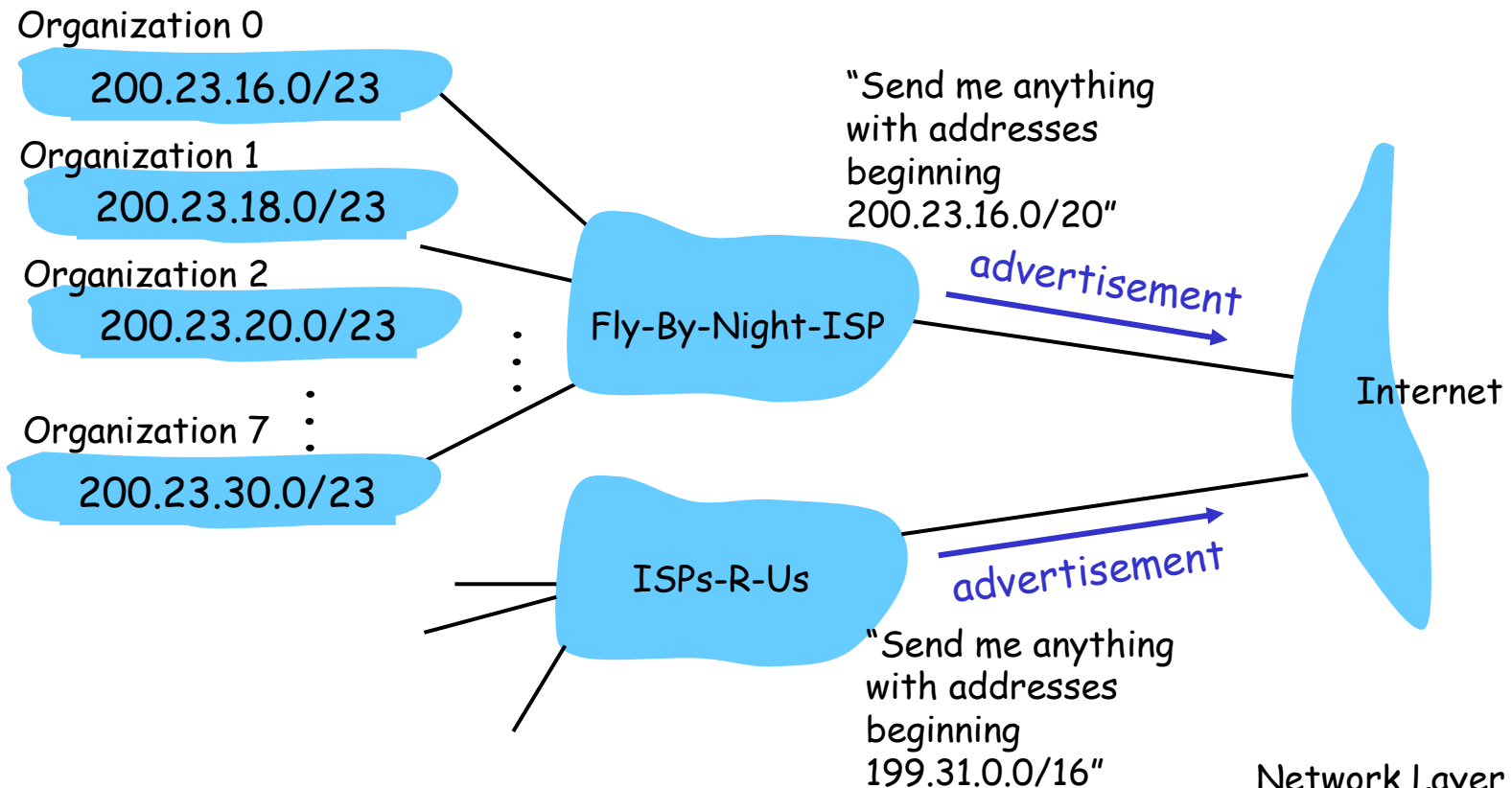  (more in next chapter)

# IP addresses: how to get one?

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

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

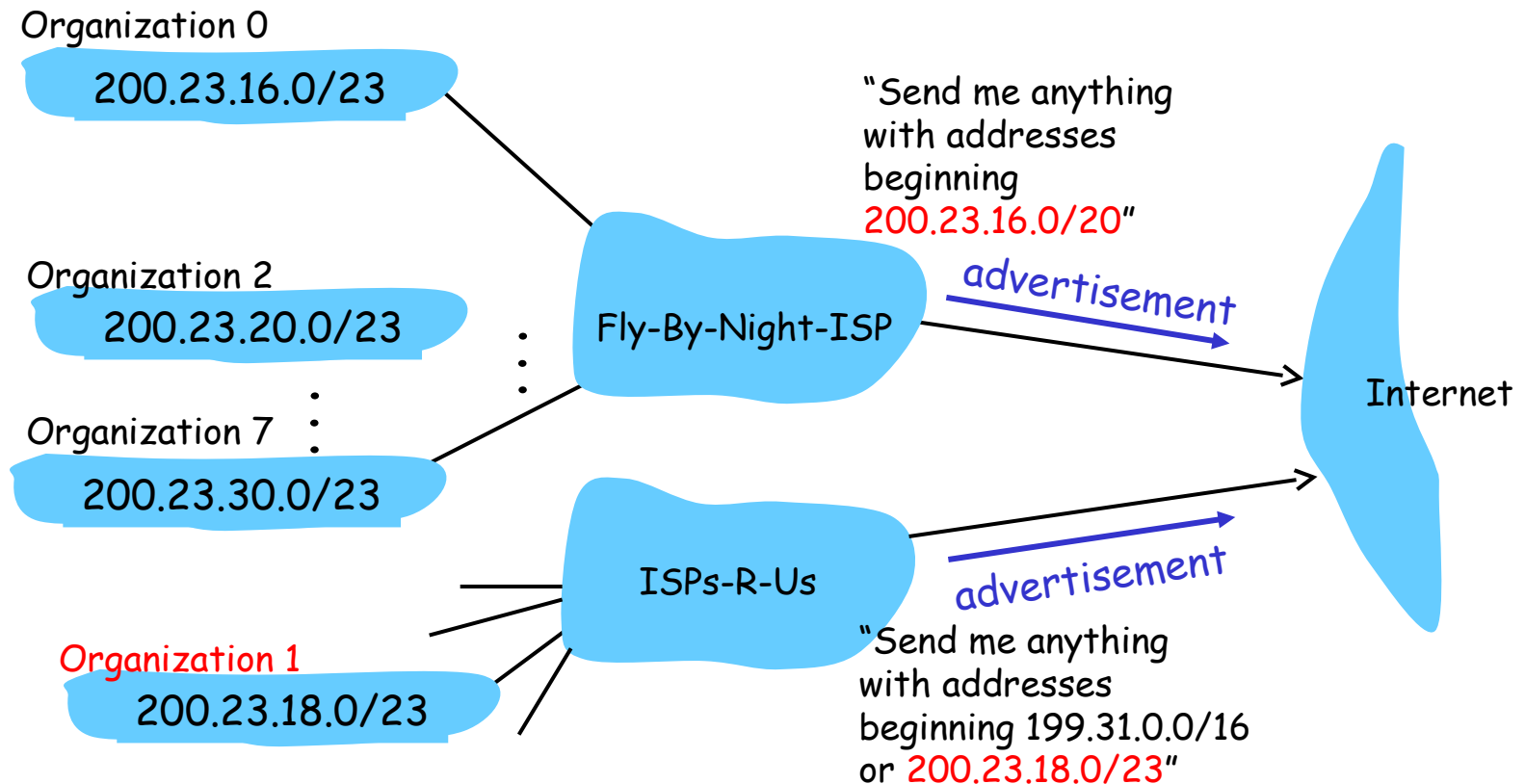| | | | |
|---|---|---|---|
| ISP's block | 11001000 00010111 0001|0000 00000000 | 200.23.16.0/20 |
| | | | |
| Organization 0 | 11001000 00010111 00010000 | 00000000 | 200.23.16.0/23 |
| Organization 1 | 11001000 00010111 00010010 | 00000000 | 200.23.18.0/23 |
| Organization 2 | 11001000 00010111 00010100 | 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

❑ Route aggregation: use a single network prefix to advertise multiple networks

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

"Send me anything with addresses beginning 200.23.16.0/20"

advertisement

ISPs-R-Us

"Send me anything with addresses beginning 199.31.0.0/16"

advertisement

Internet

# Hierarchical addressing: more specific routes

Organization 1 disconnects from Fly-By-Night-ISP and connects to ISPs-R-Us

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

"Send me anything with addresses beginning 200.23.16.0/20"

Fly-By-Night-ISP

advertisement

ISPs-R-Us

advertisement

Internet

"Send me anything with addresses beginning 199.31.0.0/16 or 200.23.18.0/23"
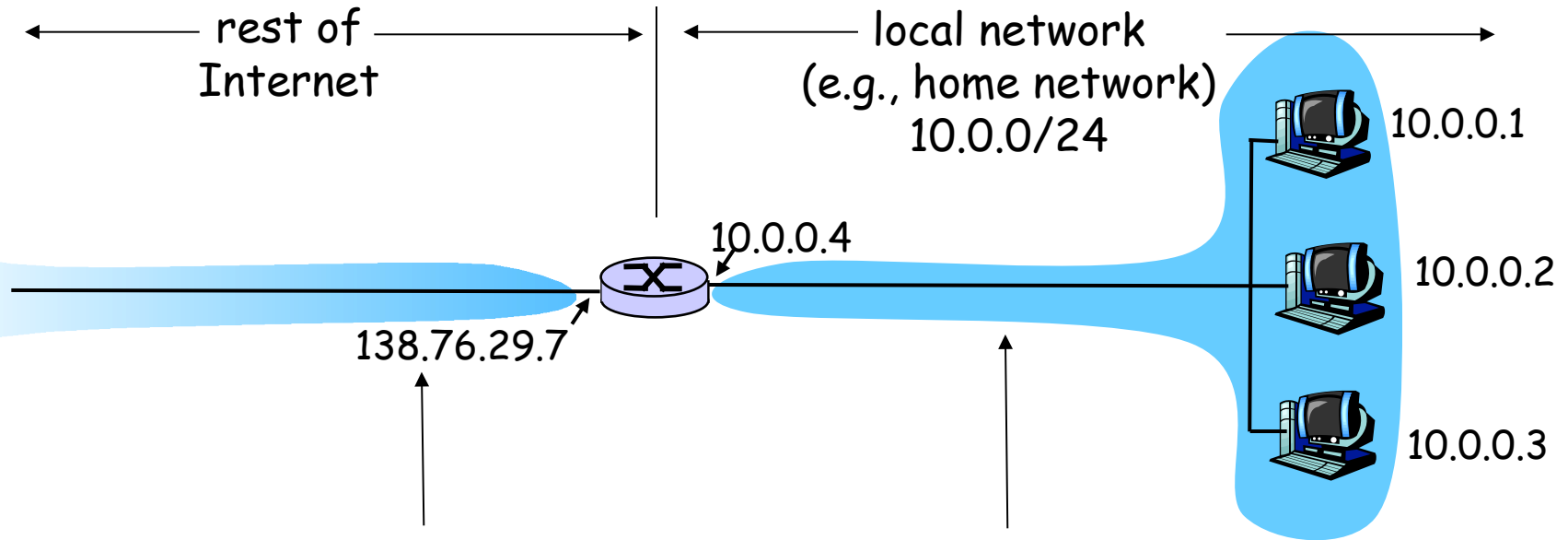
200.23.18.0/23 ⊂ 200.23.16.0/20 ⇒ use longest prefix matching

# 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

# 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)

Private IP addresses: 10.0.0.0 - 10.255.255.255
172.16.0.0 - 172.31.255.255
192.168.0.0 - 192.168.255.255

# NAT: Network Address Translation

◊ Motivation: local network uses just one IP address as far as outside world 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

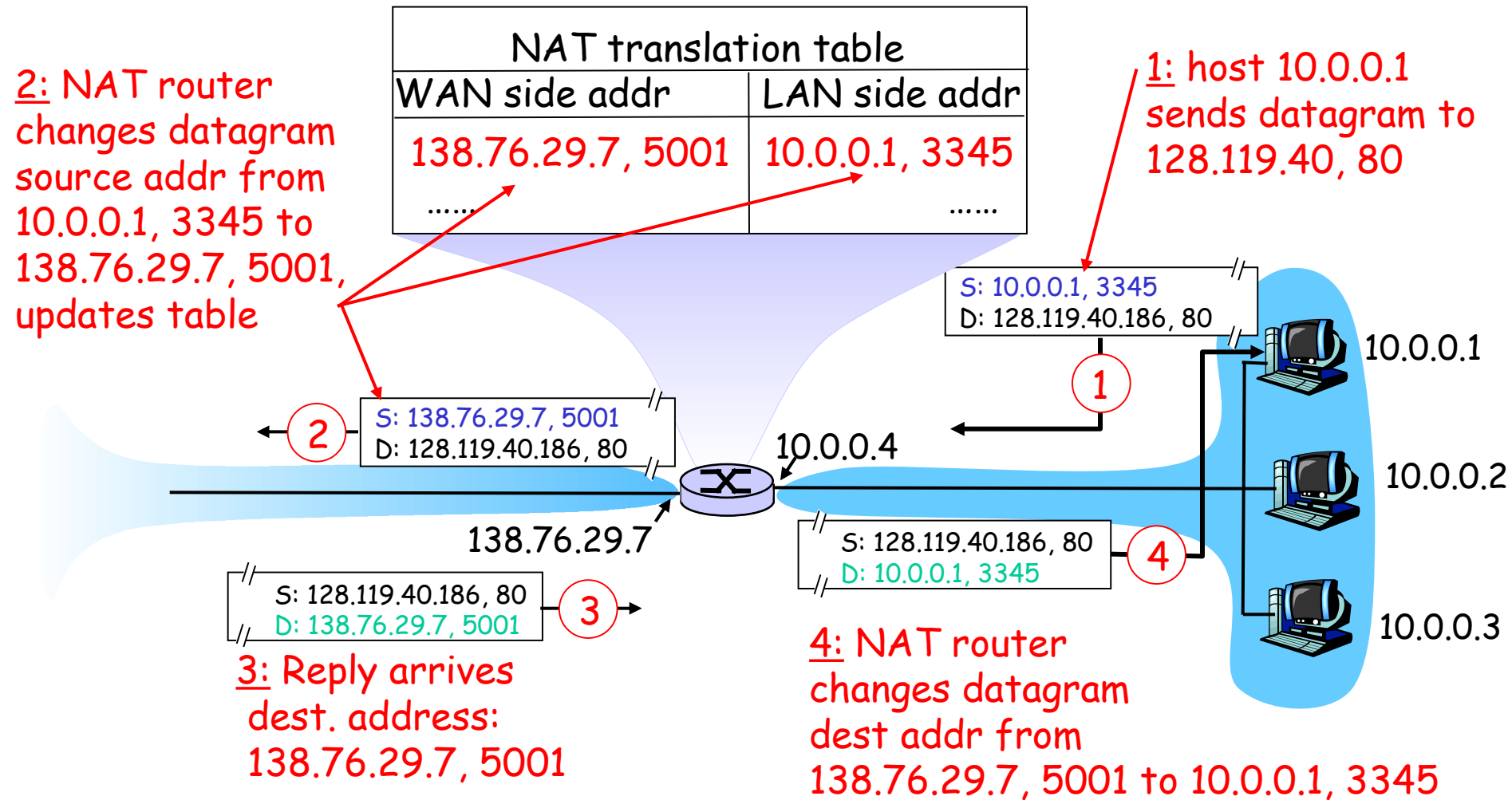  ◊ devices inside local net not explicitly addressable, visible by outside world (a security plus).

# 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

**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, 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

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**

# NAT: Network Address Translation

- 16-bit port-number field:
  - More than 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

- 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

# 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 messages 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 | router 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 TTL expired 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)
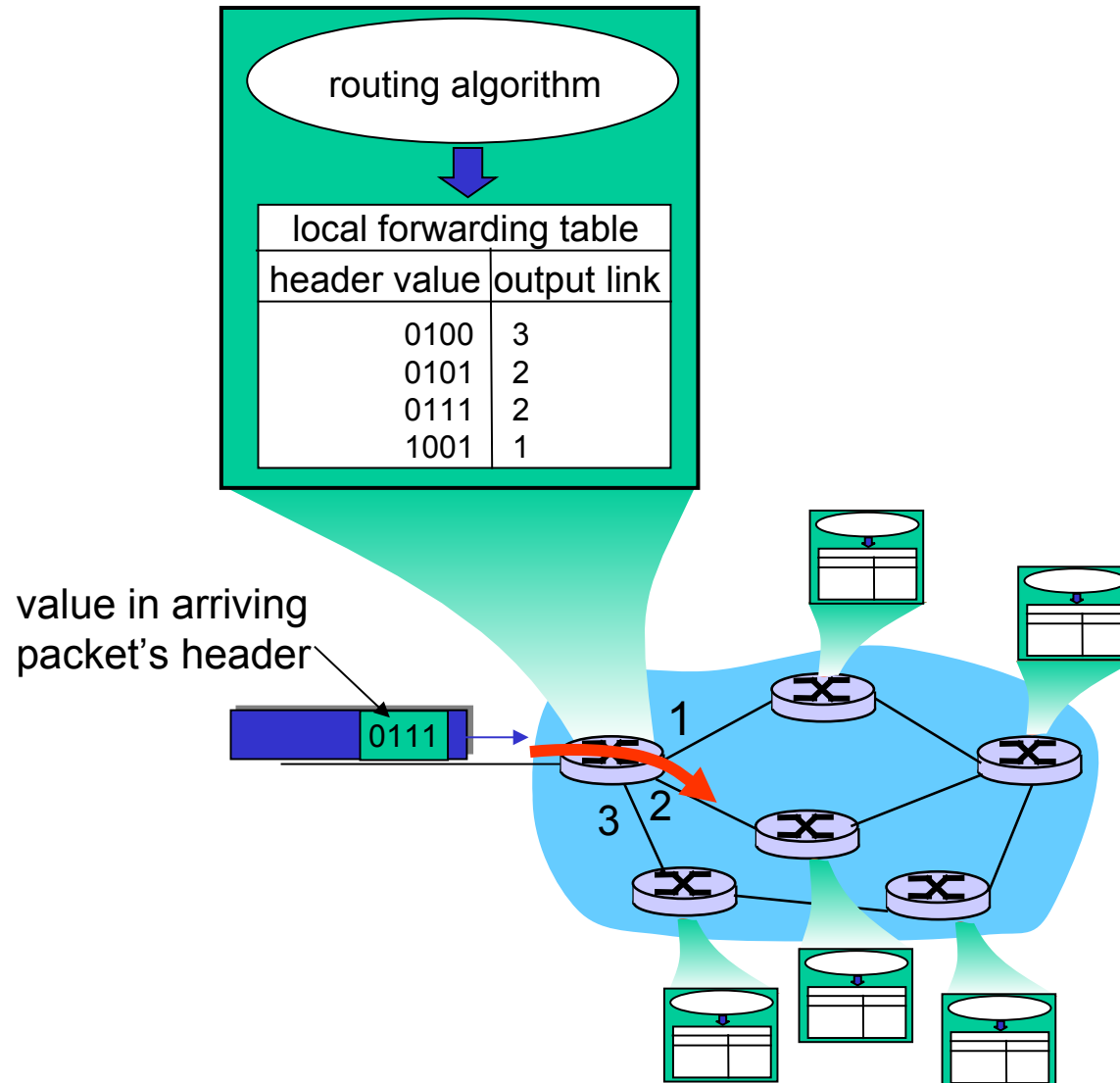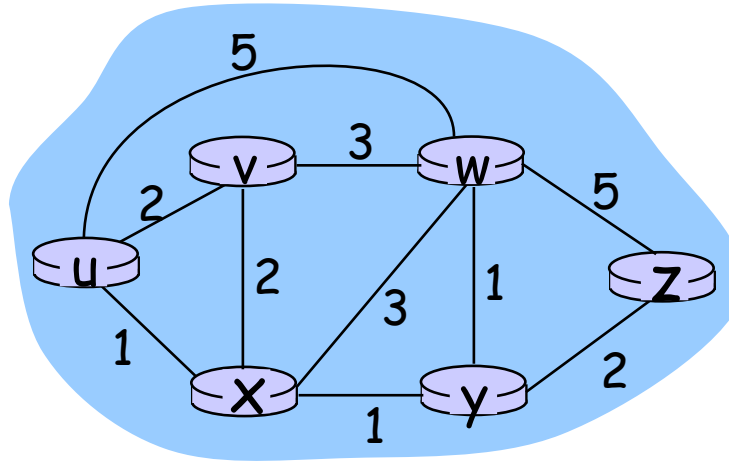* When source gets this ICMP, stops.

# Chapter 4: Network Layer

- 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

# IPv6

- **Initial motivation:** 32-bit address space soon to be completely allocated.
  - IPv6: 128-bits IP address
- Additional motivation:
  - header format helps speed processing/forwarding
  - header changes to facilitate QoS

  IPv6 datagram format:
  - fixed-length 40 byte header
    - Faster processing of IP datagram
  - 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 →

# 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 — 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><br>Flow: X<br>Src: A<br>Dest: F<br><br>data | Src:B<br>Dest: E<br><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:<br>IPv6 | B-to-C:<br>IPv6 inside<br>IPv4 | D-to-E:<br>IPv6 inside<br>IPv4 | E-to-F:<br>IPv6 |

# 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 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 information 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 destinations

**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 destination 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:
2    N' = {u}                        u: source node
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 | | | | | |

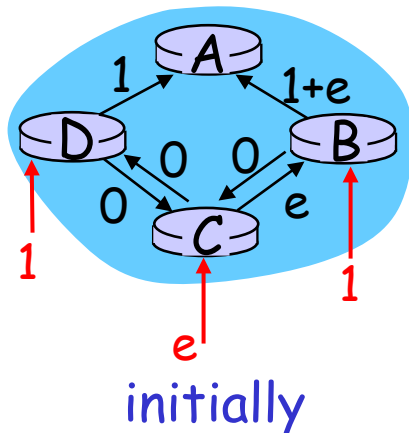# Dijkstra's algorithm: example

# Dijkstra's algorithm: example

# 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

# Chapter 4: Network Layer

# Distance Vector Algorithm (1)

 Iterative, asynchronous, distributed
 Bellman-Ford Equation (dynamic programming)

Define
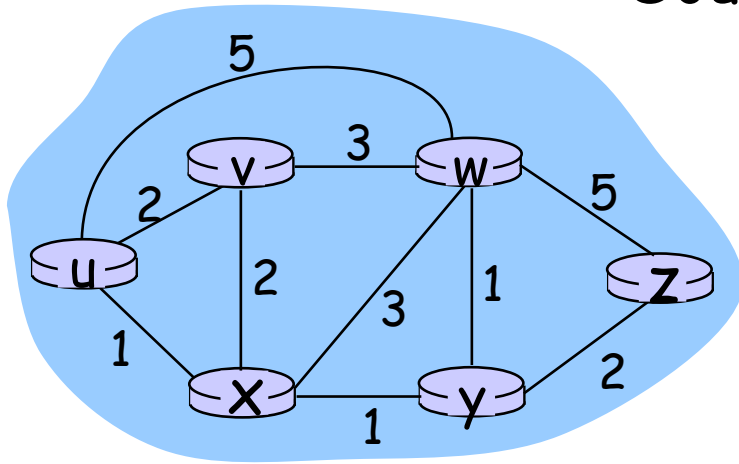
$d_x(y)$ := cost of least-cost path from x to y

Then

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

where min is taken over all neighbors of x

# Bellman-Ford example (2)

Source: U    Destination: Z



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 the next
hop in shortest path in 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 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 $D_x(y)$ *converges to the actual least cost* $d_x(y)$

# Distance Vector Algorithm (5)
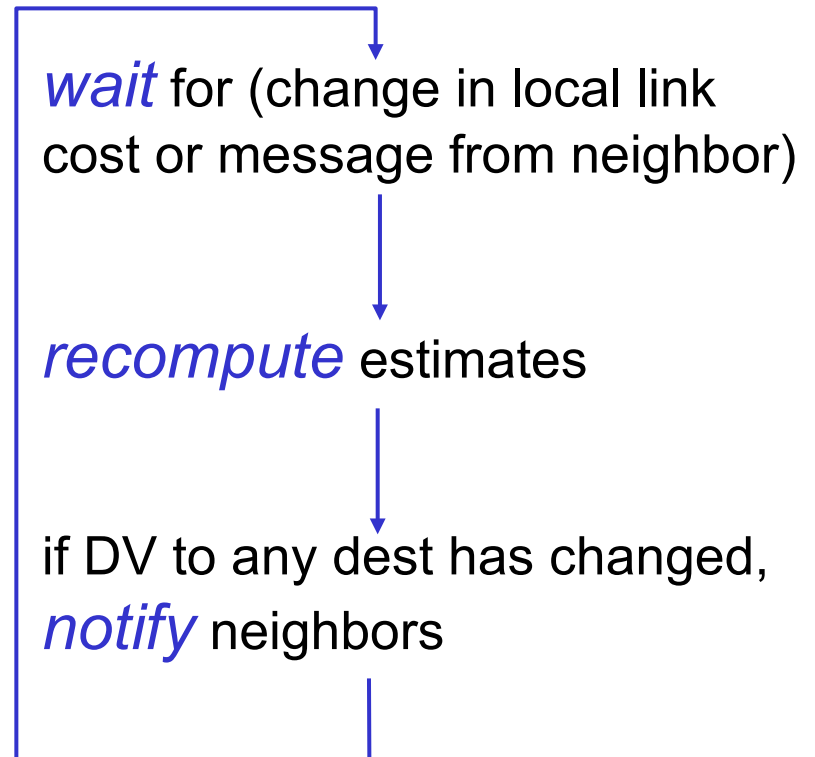
**Iterative, asynchronous:**
each local iteration caused by:

- 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 message 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

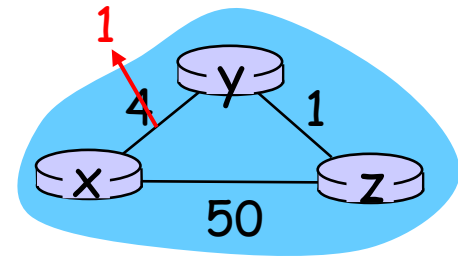| from | x | y | z |
|------|---|---|---|
| x | 0 | 2 | 3 |
| y | 2 | 0 | 1 |
| z | 3 | 1 | 0 |

time

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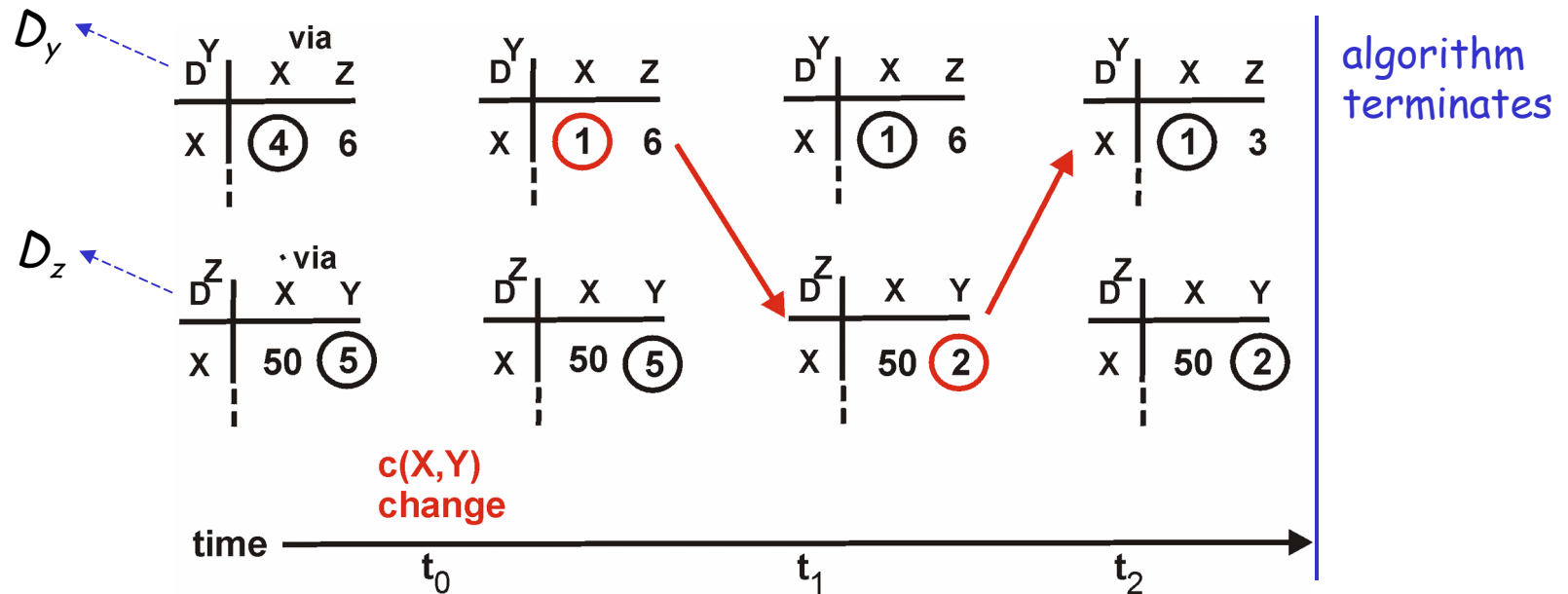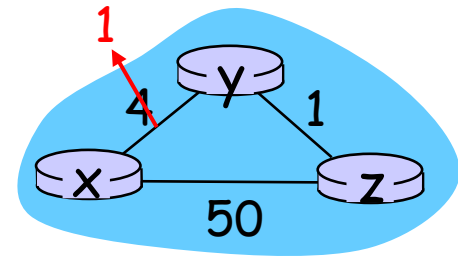
At time $t_2$, $y$ receives $z$'s update and updates its distance table. $y$'s least costs do not change and hence $y$ does *not* send any message to $z$.
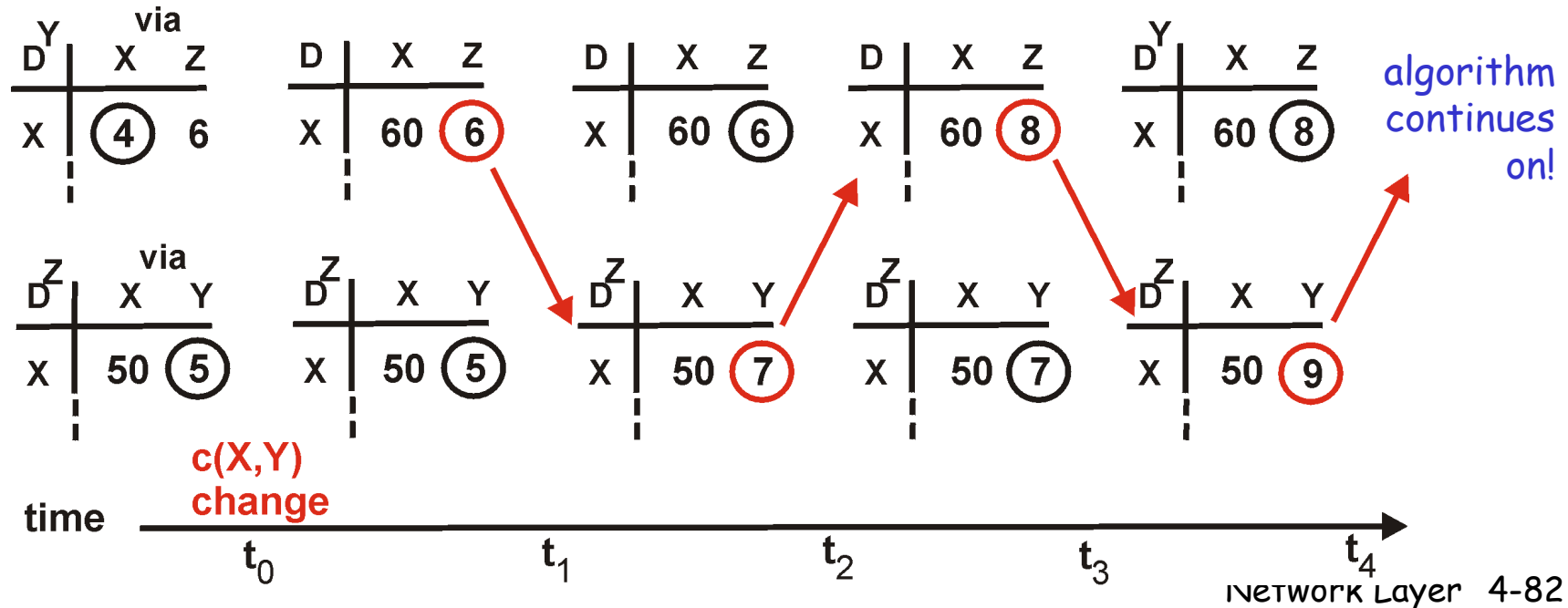
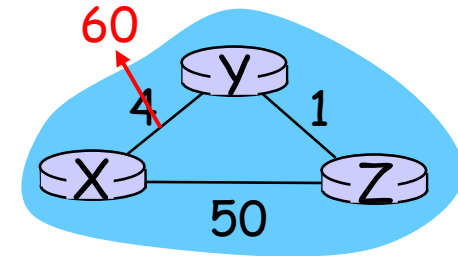# Distance Vector: link cost changes

"good news Travels fast"

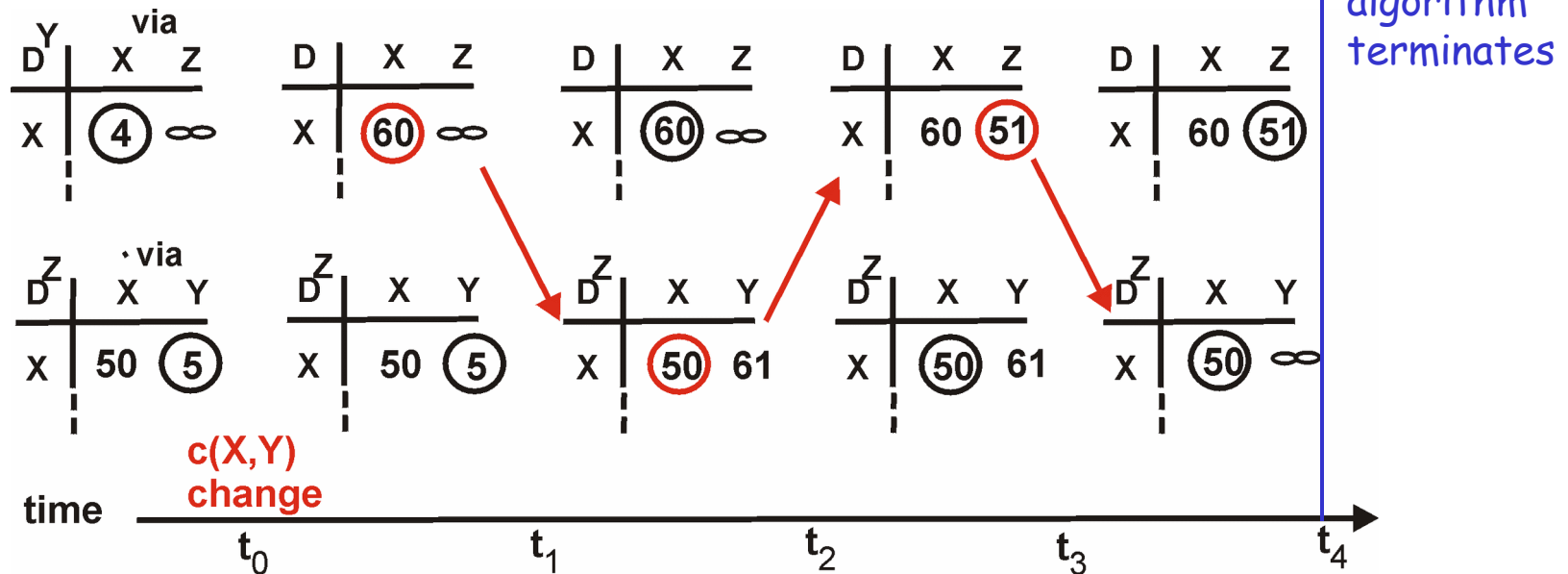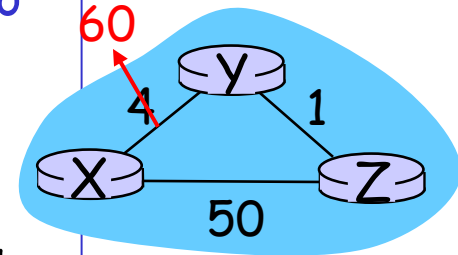# Distance Vector: link cost changes

**Link cost changes:**

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

* 44 iterations before algorithm stabilizes

* z (y) does not know that the least distance from y (z) to x that y (z) tells z (y) is the distance of the path y-z-y-x (z-y-x)

60

4    Y    1

X         Z

50

| $D^Y$ | via X | Z |
|---|---|---|
| X | ④ | 6 |

| D | X | Z |
|---|---|---|
| X | 60 | ⑥ |

| D | X | Z |
|---|---|---|
| X | 60 | ⑥ |

| D | X | Z |
|---|---|---|
| X | 60 | ⑧ |

| $D^Y$ | X | Z |
|---|---|---|
| X | 60 | ⑧ |

algorithm continues on!

| $D^Z$ | via X | Y |
|---|---|---|
| X | 50 | ⑤ |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ⑤ |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ⑦ |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ⑦ |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ⑨ |

c(X,Y) change

time

$t_0$          $t_1$          $t_2$          $t_3$          $t_4$

# Distance Vector: 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?
- ☐ Loops involving three or more nodes cannot be solved using the technique

60

Y

4    1

X         Z

50

algorithm terminates

| $D^Y$ | via X | Z |
|---|---|---|
| X | ④ | ∞ |

| D | X | Z |
|---|---|---|
| X | ⑥⓪ | ∞ |

| D | X | Z |
|---|---|---|
| X | ⑥⓪ | ∞ |

| D | X | Z |
|---|---|---|
| X | 60 | ㊟51 |

| D | X | Z |
|---|---|---|
| X | 60 | ㊟51 |

| $D^Z$ | via X | Y |
|---|---|---|
| X | 50 | ⑤ |

| D | X | Y |
|---|---|---|
| X | 50 | ⑤ |

| D | X | Y |
|---|---|---|
| X | ㊟50 | 61 |

| D | X | Y |
|---|---|---|
| X | ㊟50 | 61 |

| D | X | Y |
|---|---|---|
| X | ㊟50 | ∞ |

**c(X,Y) change**

**time**

$t_0$      $t_1$      $t_2$      $t_3$      $t_4$

# Comparison of LS and DV algorithms

## Message complexity

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

## Speed of Convergence

- **LS:** O(n²) algorithm requires O(nE) messages
  - 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

 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

# 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 destinations in routing tables!
- routing table exchange would swamp links!

**administrative autonomy**
- internet = network of networks
- each network administrator 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 is configured by both intra- and inter-AS routing algorithm
  - Intra-AS sets entries for internal destinations
  - Inter-AS & Intra-As sets entries for external destinations

# Inter-AS tasks

○ Suppose router in AS1 receives datagram for which destination is outside of AS1

  ○ Router should forward packet towards one of the gateway routers, but which one?

<span style="color:red">AS1 needs:</span>

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

<span style="color:red">Job of inter-AS routing!</span>

# Example: Setting forwarding table in router 1d

- Suppose AS1 learns from the inter-AS protocol that subnet *x* is reachable from AS3 (gateway 1c) but not from AS2.

- Inter-AS protocol propagates reachability information to all internal routers.

- Router 1d determines from intra-AS routing information that its interface $I$ is on the least cost path to 1c.

- Puts in forwarding table entry *(x,I)*.

# Example: Choosing among multiple ASes

- Now suppose AS1 learns from the 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 destination *x*.
- This is also the job on inter-AS routing protocol!
- Hot potato routing: send packet towards the closest of the 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

# Intra-AS Routing

- Also known as <span style="color:red">Interior Gateway Protocols (IGP)</span>
- 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

- 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

# RIP ( Routing Information Protocol)

- ☐ Distance vector algorithm
- ☐ Included in BSD-UNIX Distribution in 1982
- ☐ Distance metric: # of hops (max = 15 hops)

Source node: A

| 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: a 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 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 or 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

| routed | | | | | routed | |
|---|---|---|---|---|---|---|

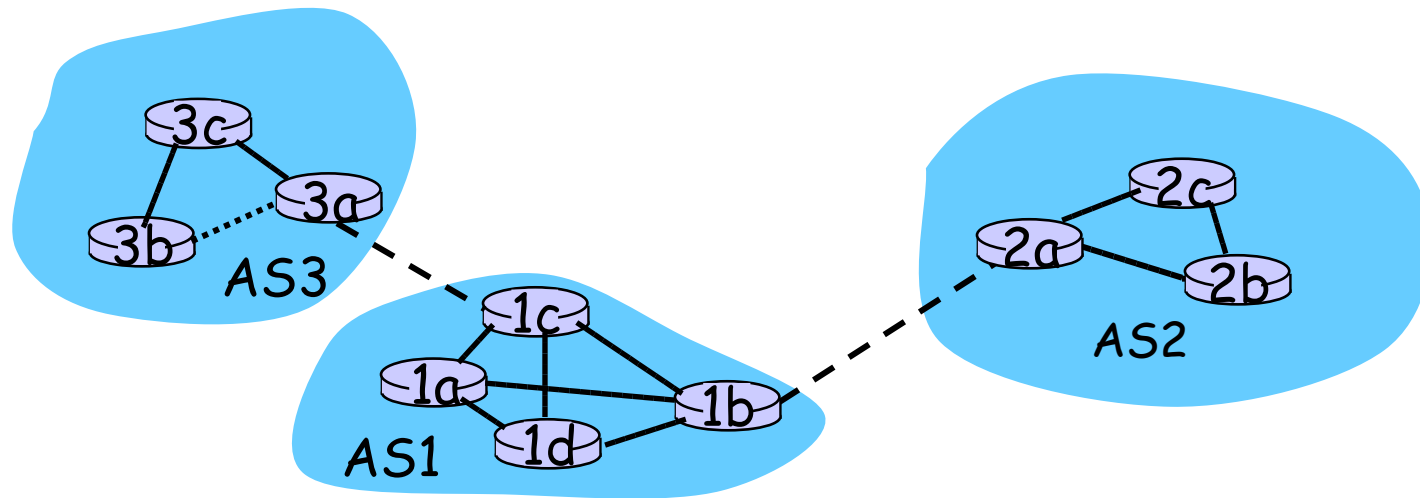| Transprt (UDP) | | | | | | Transprt (UDP) |
|---|---|---|---|---|---|---|
| network (IP) | forwarding table | | | forwarding table | | network (IP) |
| link | | | | | | link |
| physical | | | | | | physical |

# 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

# OSPF (Open Shortest Path First)

- "open": publicly available – defined in RFC 2328
- Uses Link State algorithm
  - Link-State 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.

# Hierarchical OSPF

- An OSPF autonomous system (AS) can be configured into areas
- Exactly one OSPF area in the AS is configured to be the backbone area
- Each area runs its own OSPF link-state routing algorithm
- Two-level hierarchy: local area, backbone.
  - Link-state advertisements only in area
  - each nodes has detailed area topology; only know direction (shortest path) to nets in other areas.

# Hierarchical OSPF



boundary router

backbone router

Backbone

area border routers

internal routers

Area 1

Area 2

Area 3

# Hierarchical OSPF

## Four types of routers

- Internal routers: perform only intra AS routing
- Area border routers: belong to both an area and the backbone
- Backbone routers: run OSPF routing limited to backbone.
- Boundary routers: connect to other AS's.

# 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

# 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 information over semi-permanent TCP connections: BGP sessions
- Note that BGP sessions do not correspond to physical links.
- When AS2 advertises a prefix to AS1, AS2 is *promising* it will forward any datagrams destined to that prefix towards the prefix.
  - AS2 can aggregate prefixes in its advertisement



3c
3a
3b
AS3

2c
2a
2b
AS2

1c
1a
1b
1d
AS1

– – – – – External BGP (eBGP) session

————— Internal BGP (iBGP) session

# Aggregation of prefixes

138.16.64/24

138.16.65/24

138.16.66/24    = >    138.16.64/22

138.16.67/24

# Distributing reachability info

- With eBGP session between 3a and 1c, AS3 sends prefix reachability information to AS1.
- 1c can then use iBGP to distribute this new prefix reachability information to all routers in AS1
- 1b can then re-advertise the new reachability information 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.



AS3

AS1

AS2

– – – – – eBGP session

——— iBGP session

# Path attributes & BGP routes

- When advertising a prefix, advertisement includes BGP attributes.
  - prefix + attributes = "route"
- Two important attributes:
  - AS-PATH: contains the ASs through which the advertisement for the prefix passed: AS 67 AS 17
    - used to detect and prevent looping advertisement
    - also use in choosing among multiple path to the same prefix
  - 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 advertisement, uses import policy to accept/decline.

# BGP route selection

- Router may learn about more than 1 route to any one prefix. Router must select route.
- Elimination rules invoked sequentially until one route remains:
  1. Local preference value attribute: policy decision – AS's network administrator
  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 message; 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 to B the path AW
- B advertises to X the path BAW
- Should B advertise to C the path BAW?
  - 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: administrator 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 and multicast routing

 Broadcast routing –- deliver a packet from a source node to all other nodes

 Multicast routing – deliver a packet from a source node to a subset of other nodes

# Source-duplication versus in-network duplication



duplicate creation/transmission

(a) source duplication, (b) in-network duplication

# Broadcast routing algorithms

 Uncontrolled flooding

 Controlled flooding

 Spanning-tree broadcast

# Uncontrolled flooding

* The source node sends a copy of the packet to all of its neighbors

* When a node receives a broadcast packet, it duplicates the packet and forwards it to all of its neighbors (except the neighbor from which it receive the packet)

Problems:

* I f the graph has cycles, then one or more copies of each broadcast packet will cycle indefinitely

* Broadcast storm

# Controlled flooding

- Sequence-number-controlled flooding
  - Source node puts its address and a broadcast sequence number into a broadcast packet
  - Each node maintains a list of the source address and sequence number of each packet it has received
  - When a node receives a broadcast packet
    - If the packet is in the list, the packet is dropped
    - Otherwise, the packet is duplicated and forwarded

# Controlled flooding

⬜ Reverse path forwarding

⬜ When a router receives a broadcast packet, it duplicates and forwards the packet only if the packet arrives on the link that is on its own shortest unicast path back to the source



→ Packet will be forwarded

→| Packet not forwarded beyond receiving router

# Controlled flooding

Drawback

* Some of the nodes receive redundant packets



Redundant packets

Ideally, every node should receive only one copy of the broadcast packet.

# Spanning-tree broadcast

Spanning tree – a tree that contains all nodes in a graph

Minimum spanning tree – a spanning tree whose cost is the minimum among all the spanning trees of a graph

 Broadcast along a spanning tree



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

# Construction of Spanning-tree

- Many algorithms have been developed
- Center-based approach
  - Select a center node (rendezvous or core)
  - Each node unicasts tree-join message to the center node



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

**(b) Constructed spanning tree**

# Multicast Routing: Problem Statement

- **_Goal:_** find a tree (or trees) connecting routers having local multicast group members
  - _tree:_ not all paths between routers used
  - _source-based:_ different tree from each sender to receivers
  - _shared-tree:_ same tree used by all group members



Shared tree               Source-based trees

# Approaches for building multicast 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

- multicast 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* (multicast datagram received on incoming
    link on shortest path back to sender)
    *then* flood datagram onto all outgoing links
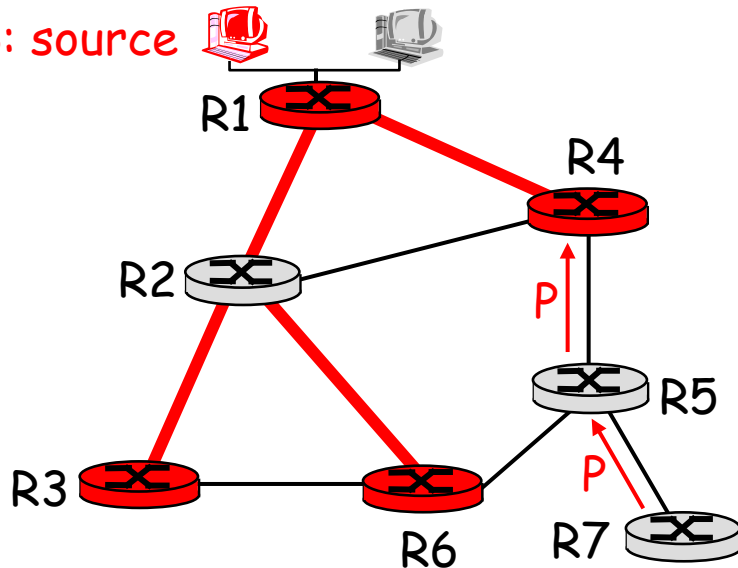  *else* ignore datagram

# 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

# Reverse Path Forwarding: pruning

- forwarding tree contains subtrees with no multicast group members
  - no need to forward datagrams down subtree
  - "prune" messages sent upstream by router with no downstream group members

S: source

LEGEND

router with attached group member

router with no attached group member

P → prune message

links with multicast forwarding

R1

R2

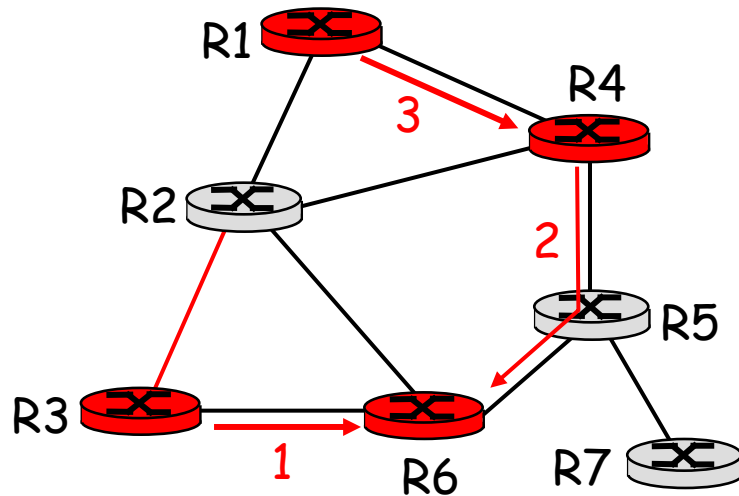R3

R4

R5

R6

R7

P

P

# 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-message* addressed to center router
  - *join-message* "processed" by intermediate routers and forwarded towards center
  - *join-message* either hits existing tree branch for this center, or arrives at center
  - path taken by *join-message* 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
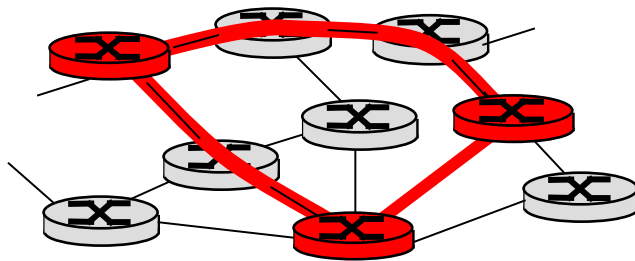
# Internet Multicasting Routing: DVMRP

- 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 multicast group flooded everywhere via RPF

  - routers not wanting group: send upstream prune messages
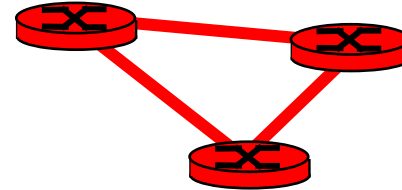
# DVMRP: continued...

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

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

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

❑ receiving multicast router decapsulates to get multicast 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:
- ❑ # of routers with group members is small wrt total # of routers
- ❑ 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 of multicast tree (e.g., RPF)
- bandwidth and non-group-router processing *profligate*

## Sparse:

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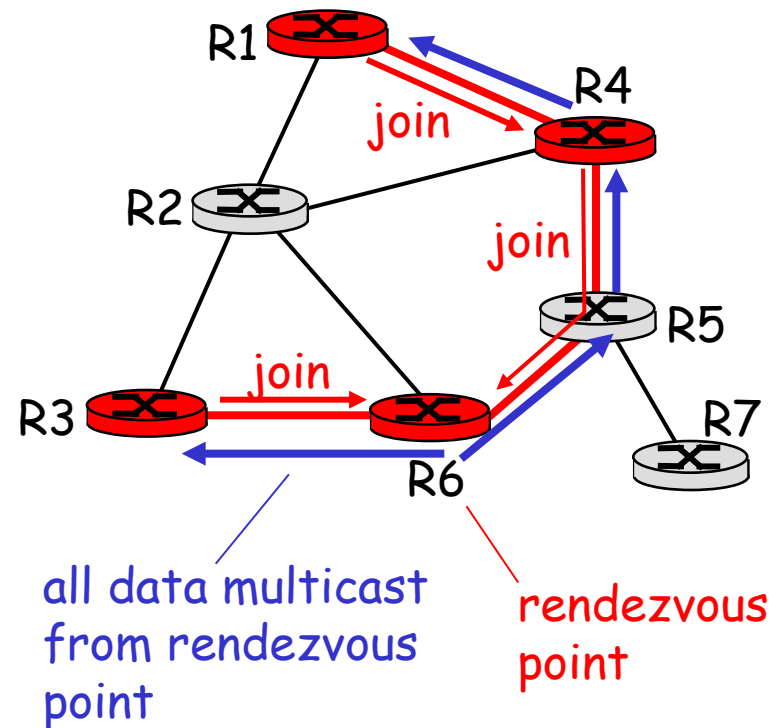- bandwidth and non-group-router processing *conservative*

# PIM- Dense Mode

flood-and-prune RPF, similar to DVMRP but

- ❑ underlying unicast protocol provides RPF information for incoming datagram

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

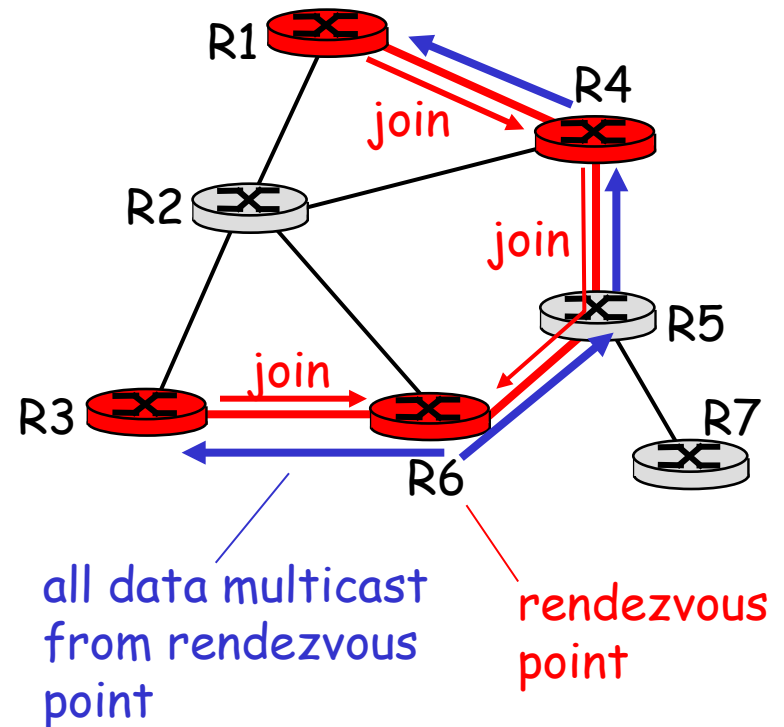- ❑ has protocol mechanism for router to detect it is a leaf-node router

# PIM - Sparse Mode

 center-based approach

 router sends *join* message 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



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 multicast tree upstream to source
- RP can send *stop* message to the source if no attached receivers
  - "no one is listening!"

R1

R4

join

R2

join

R5

join

R3

join

R7

R6

all data multicast from rendezvous point

rendezvous point

# Network Layer: summary

What we've covered:

- 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

Next stop:
the Data
link layer!