

# Chapter 4

## Network Layer

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

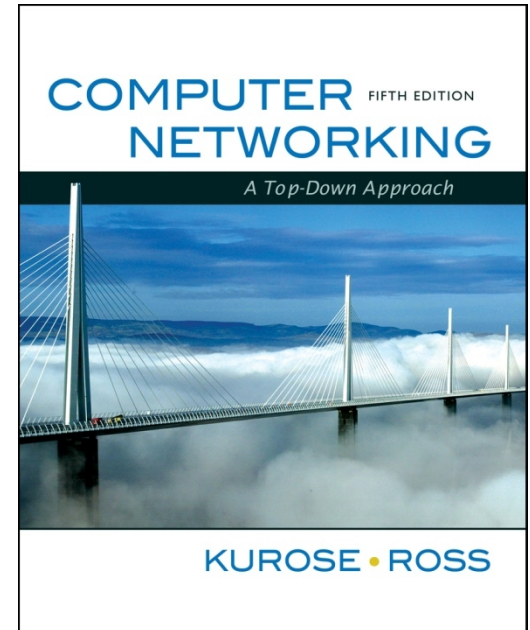
We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- ❑ 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

All material copyright 1996-2009

J.F Kurose and K.W. Ross, All Rights Reserved



*Computer Networking:  
A Top Down Approach  
5<sup>th</sup> edition.*

*Jim Kurose, Keith Ross  
Addison-Wesley, April  
2009.*

# Chapter 4: Network Layer

## Chapter goals:

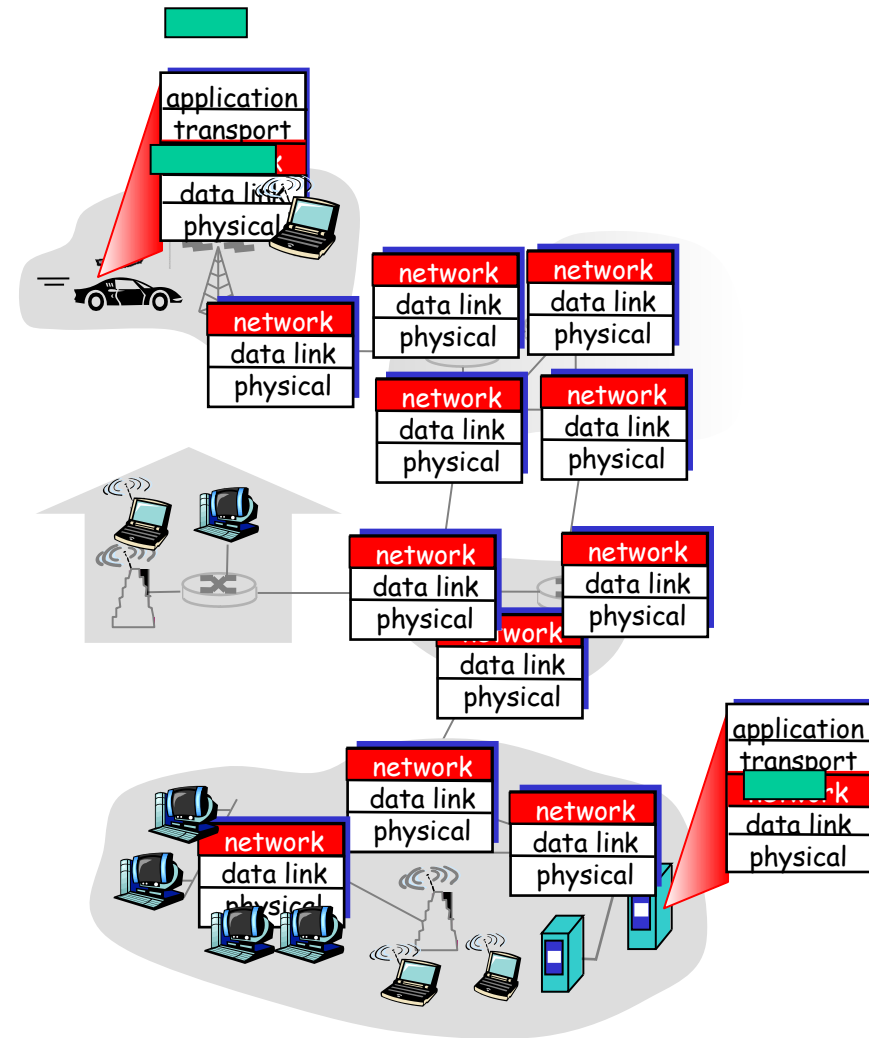
- ❑ understand principles behind network layer services:
  - network layer service models
  - forwarding versus routing
  - how a router works
  - routing (path selection)
  - dealing with scale
  - advanced topics: IPv6, mobility
- ❑ instantiation, implementation in the Internet

# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link-state
- 4.3 Distance Vector
- 4.3 What's inside a router
  - Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 4.7 Broadcast and multicast routing

# Network layer

- ❑ transport segment from sending to receiving host
- ❑ on sending side encapsulates segments into datagrams
- ❑ on rcving side, delivers segments to transport layer
- ❑ network layer protocols in *every* host, router
- ❑ router examines header fields in all IP datagrams passing through it



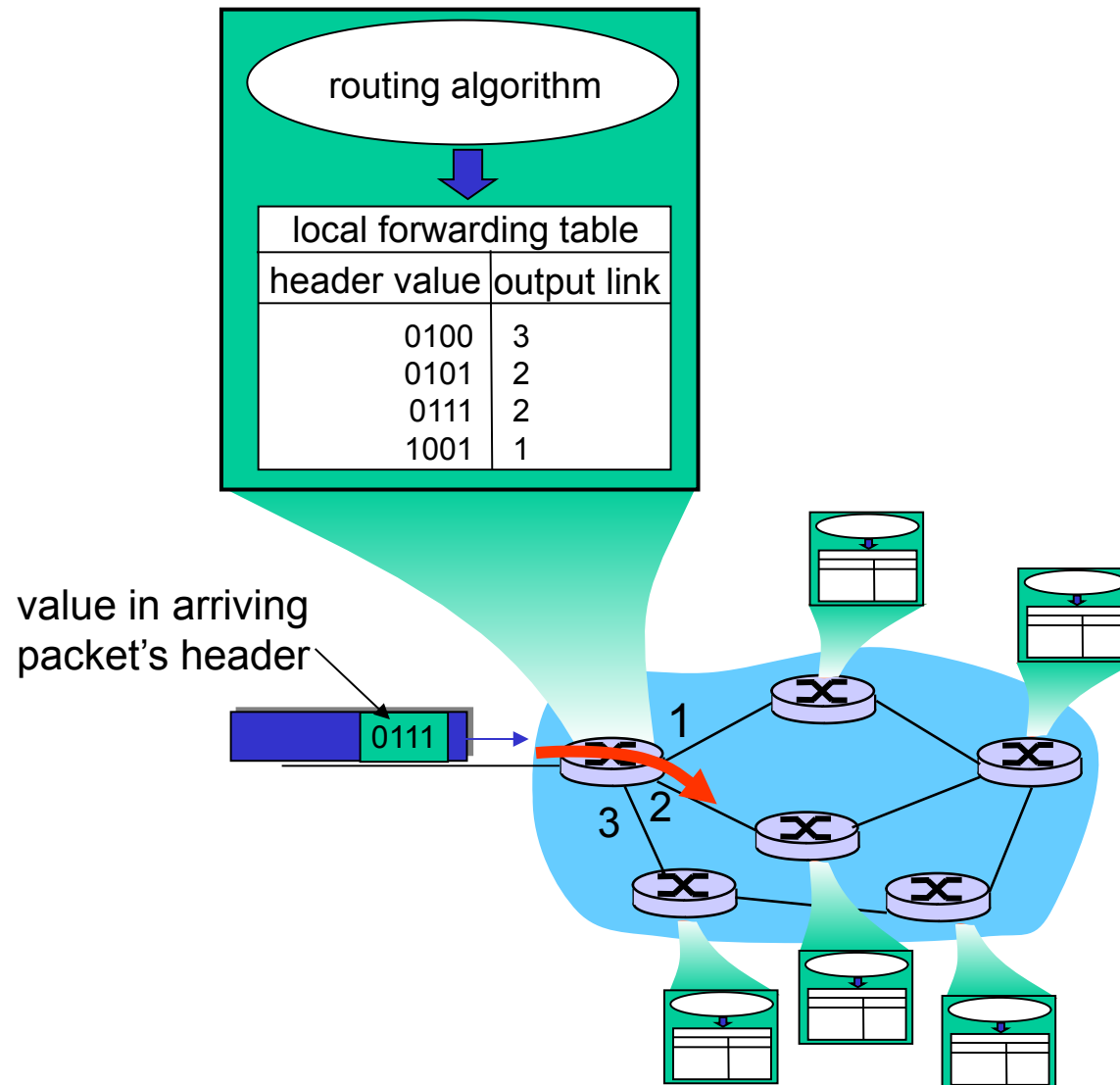
# Two Key Network-Layer Functions

- ❑ *forwarding*: move packets from router's input to appropriate router output
- ❑ *routing*: determine route taken by packets from source to dest.
  - *routing algorithms*

## analogy:

- ❑ *routing*: process of planning trip from source to dest
- ❑ *forwarding*: process of getting through single interchange

# Interplay between routing and forwarding



# Connection setup

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

# 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 rate	yes	yes	yes	no congestion
ATM	VBR	guaranteed rate	yes	yes	yes	no congestion
ATM	ABR	guaranteed minimum	no	yes	no	yes
ATM	UBR	none	no	yes	no	no

# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link state
- 4.2.1 Virtual circuit and datagram networks
- 4.2.2 Distance Vector
- 4.3 What's inside a router
  - Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 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 network core

# Virtual circuits

“source-to-dest path behaves much like telephone circuit”

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

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

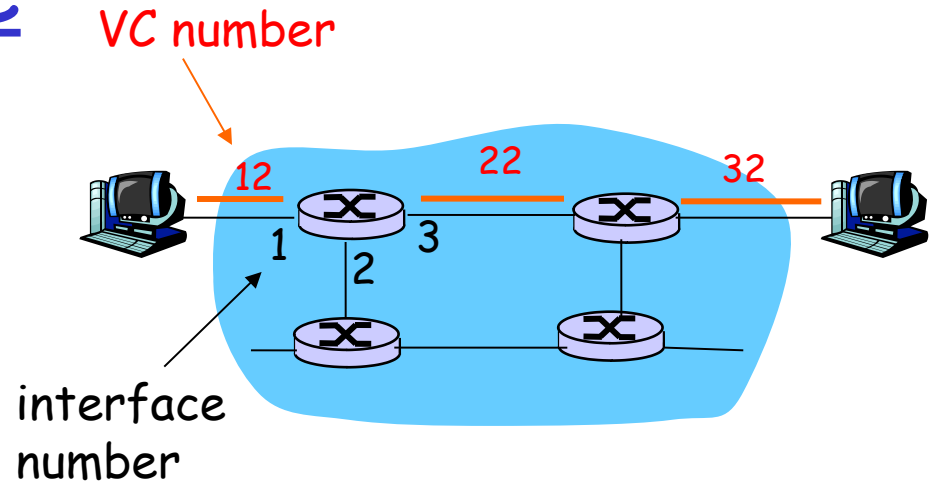
# VC implementation

a VC consists of:

1. path from source to destination
  2. VC numbers, one number for each link along path
  3. entries in forwarding tables in routers along path
- ❑ packet belonging to VC carries VC number (rather than dest address)
  - ❑ VC number can be changed on each link.
    - New VC number comes from forwarding table

# Forwarding table

Forwarding table in  
northwest router:

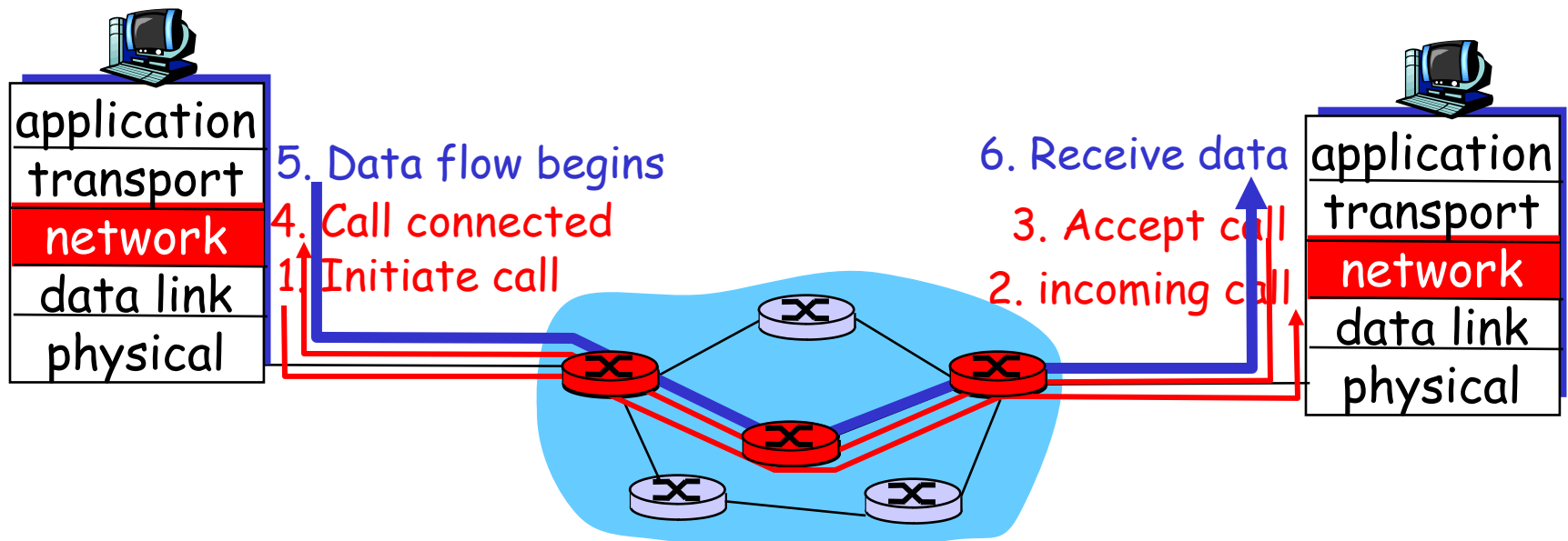


Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	3	22
2	63	1	18
3	7	2	17
1	97	3	87
...	...	...	...

**Routers maintain connection state information!**

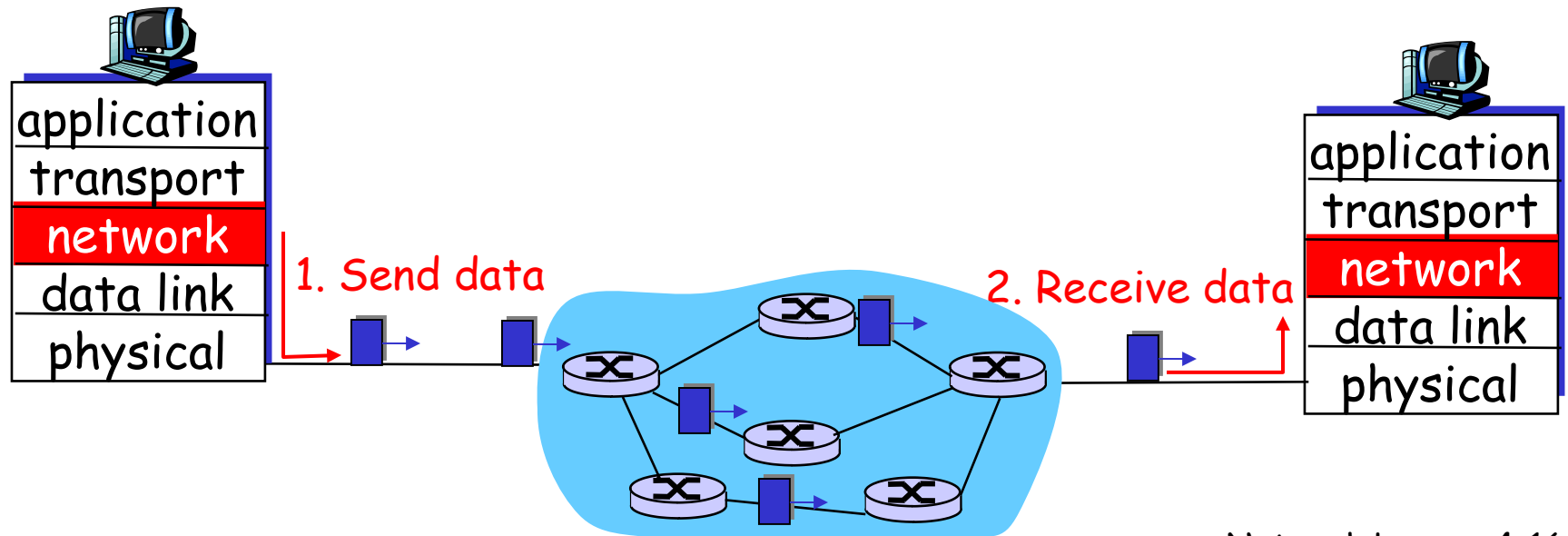
# Virtual circuits: signaling protocols

- used to setup, maintain teardown VC
- used in ATM, frame-relay, X.25
- not used in today's Internet



# 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





# Forwarding table

4 billion  
possible entries

<u>Destination Address Range</u>	<u>Link Interface</u>
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 11111111	0
11001000 00010111 00011000 00000000 through 11001000 00010111 00011000 11111111	1
11001000 00010111 00011001 00000000 through 11001000 00010111 00011111 11111111	2
otherwise	3

# Longest prefix matching

<u>Prefix Match</u>	<u>Link Interface</u>
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 (datagram)

- ❑ 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”
- ❑ many link types
  - different characteristics
  - uniform service difficult

## ATM (VC)

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

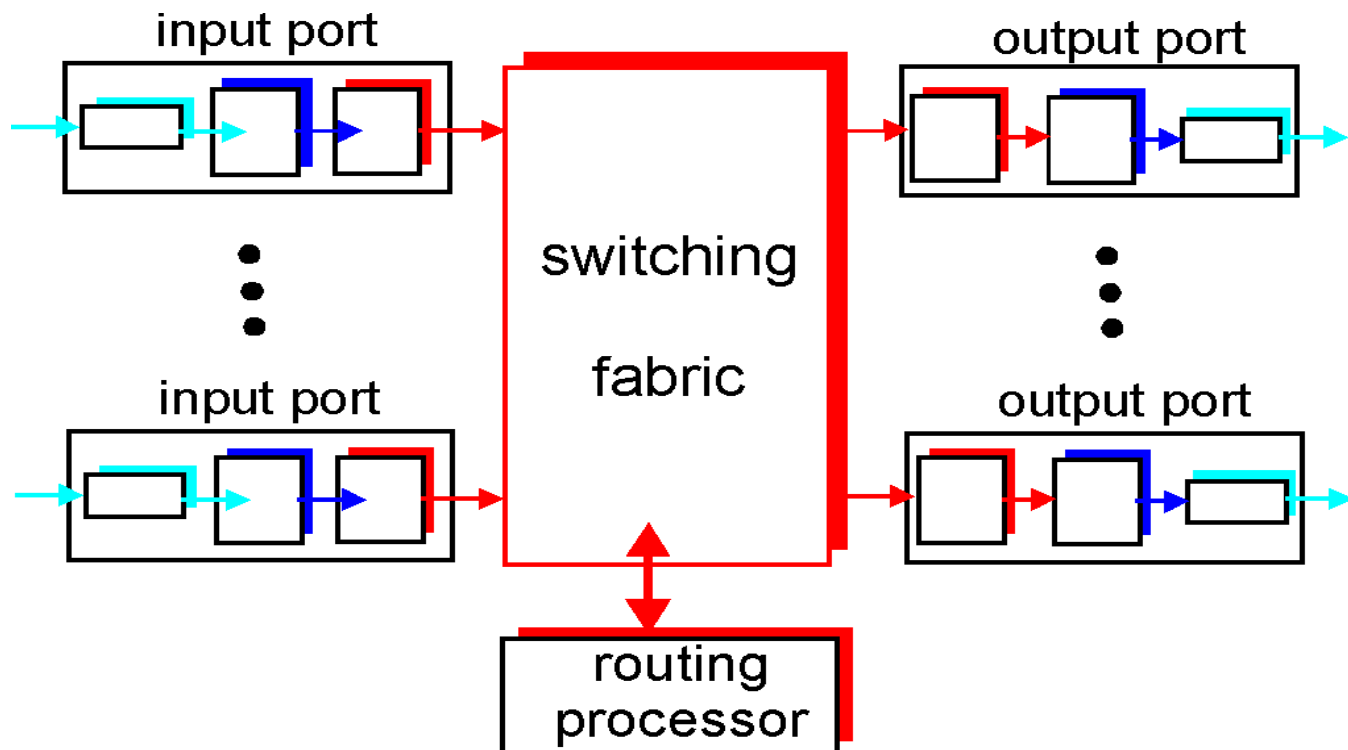
# Chapter 4: Network Layer

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

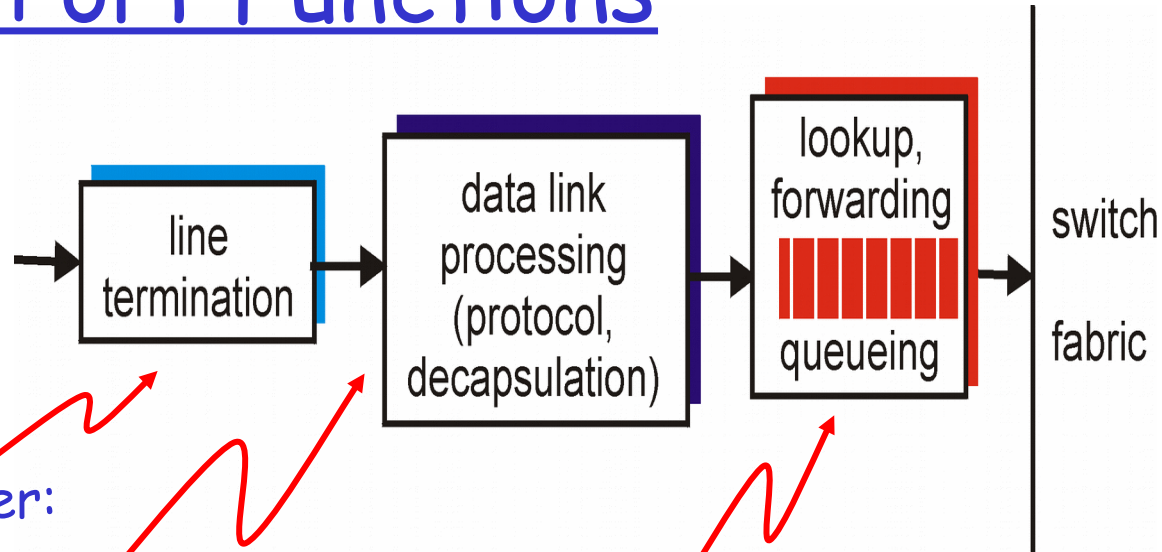
# Router Architecture Overview

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

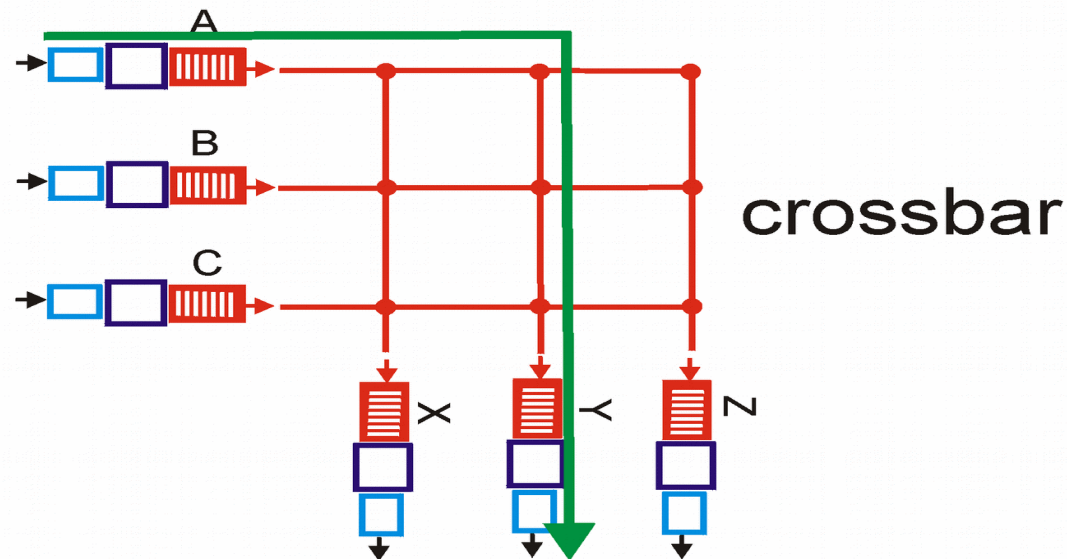
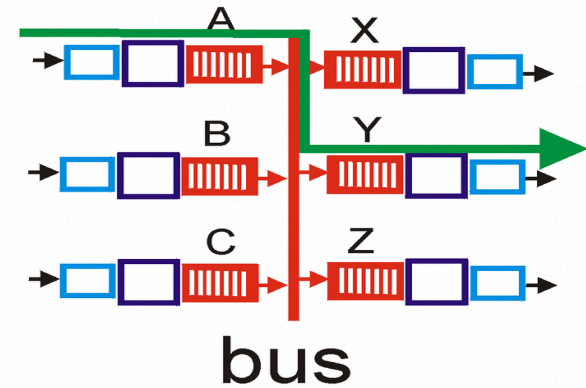
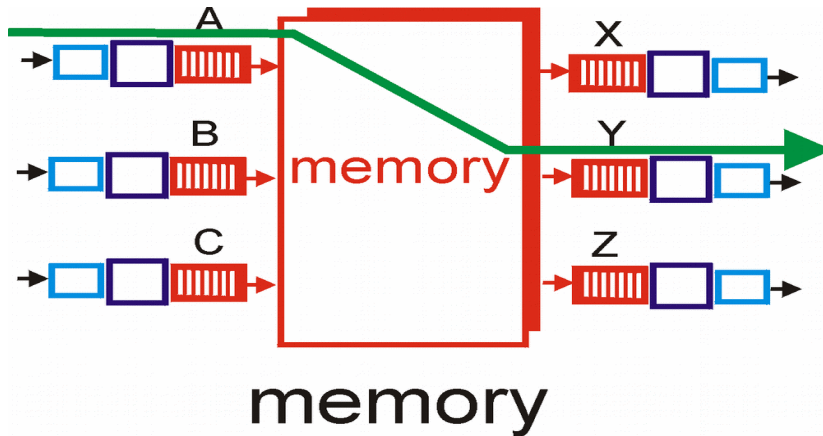
## Decentralized switching:

- 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

# Forwarding in high speed links

- ❑ Requires input processing at line speed
- ❑ Linear search through a large forwarding table is impossible
- ❑ Possible solutions
  - Tree data structure
  - Content Addressable Memories
  - Use caching for recently accessed entries

# Three types of switching fabrics

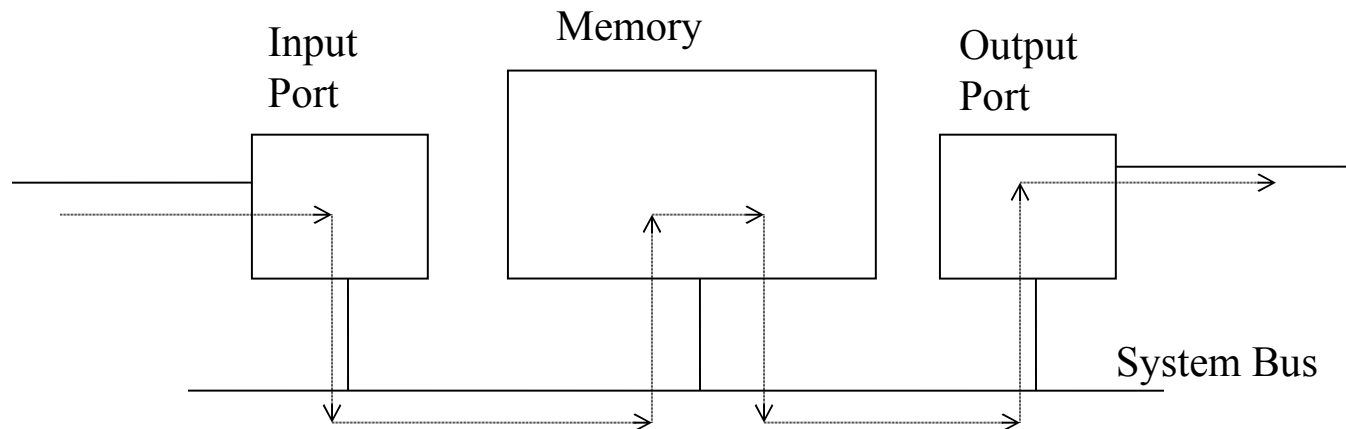




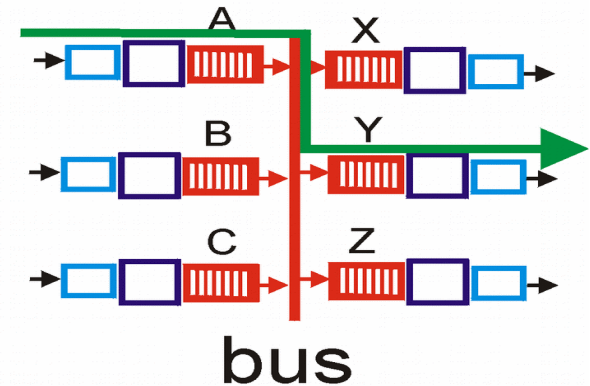
# Switching Via Memory

## 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)



# Switching Via a Bus

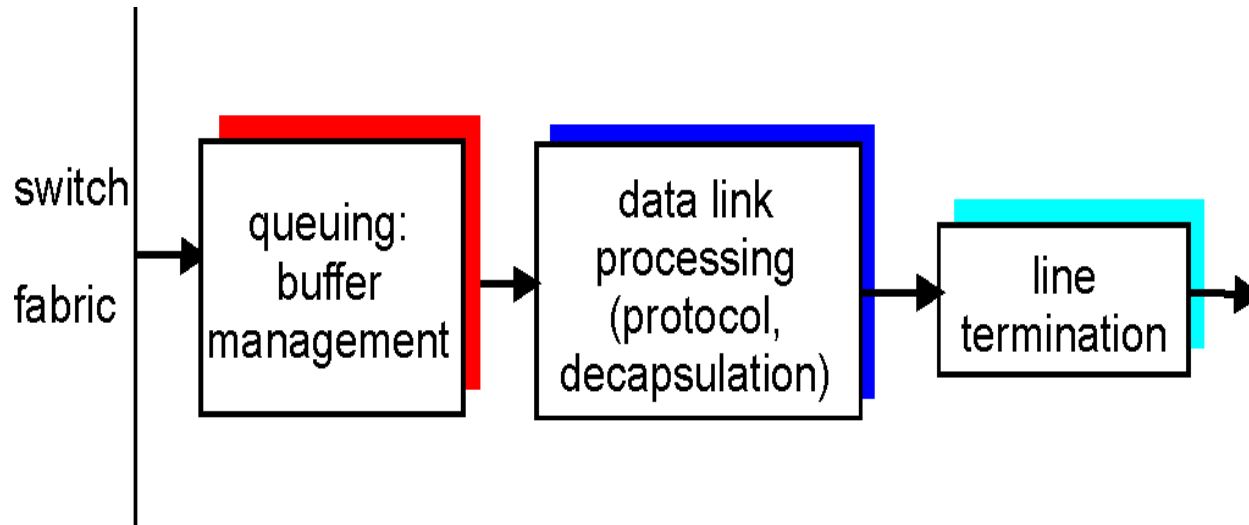


- ❑ datagram from input port memory to output port memory via a shared bus
- ❑ **bus contention:** switching speed limited by bus bandwidth
- ❑ 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers

# Switching Via An Interconnection Network

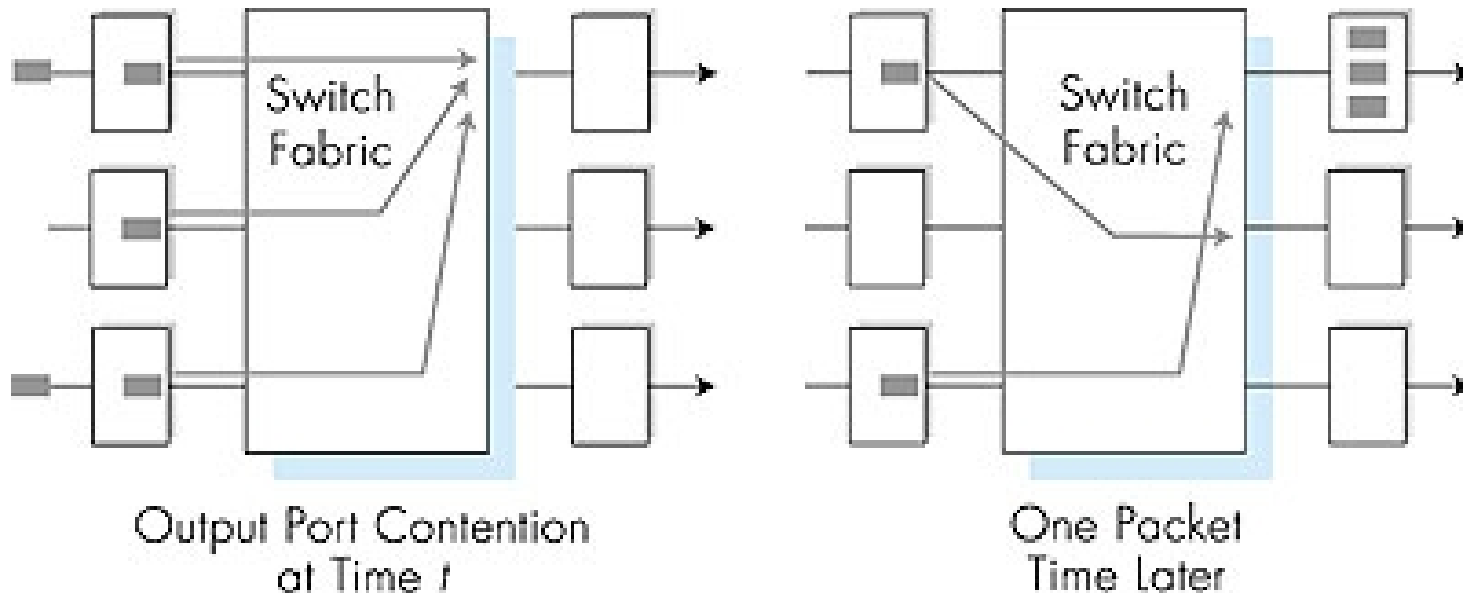
- ❑ overcome bus bandwidth limitations
- ❑ Banyan networks, other interconnection nets initially developed to connect processors in multiprocessor
- ❑ advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- ❑ Cisco 12000: switches 60 Gbps through the interconnection 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



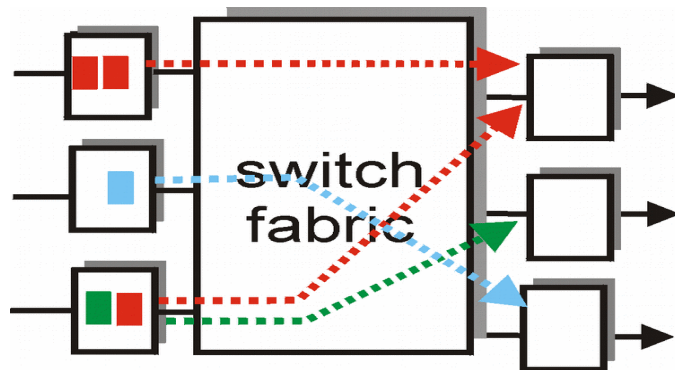
- buffering when arrival rate via switch exceeds output line speed
- *queueing (delay) and loss due to output port buffer overflow!*

# How much buffering?

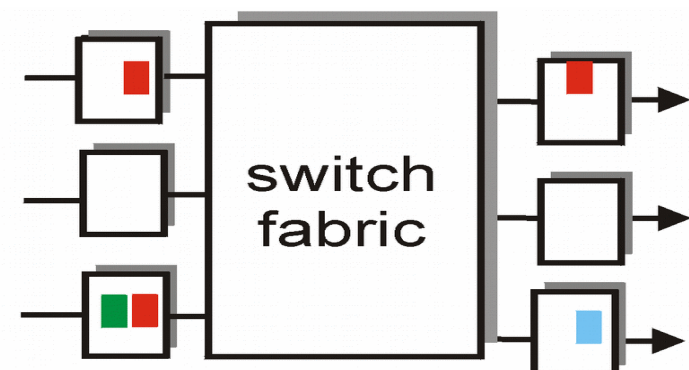
- ❑ RFC 3439 rule of thumb: average buffering equal to “typical” RTT (say 250 msec) times link capacity  $C$ 
  - e.g.,  $C = 10$  Gps link: 2.5 Gbit buffer
- ❑ Recent recommendation: with  $N$  flows, buffering equal to  $\frac{RTT \cdot C}{\sqrt{N}}$

# 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

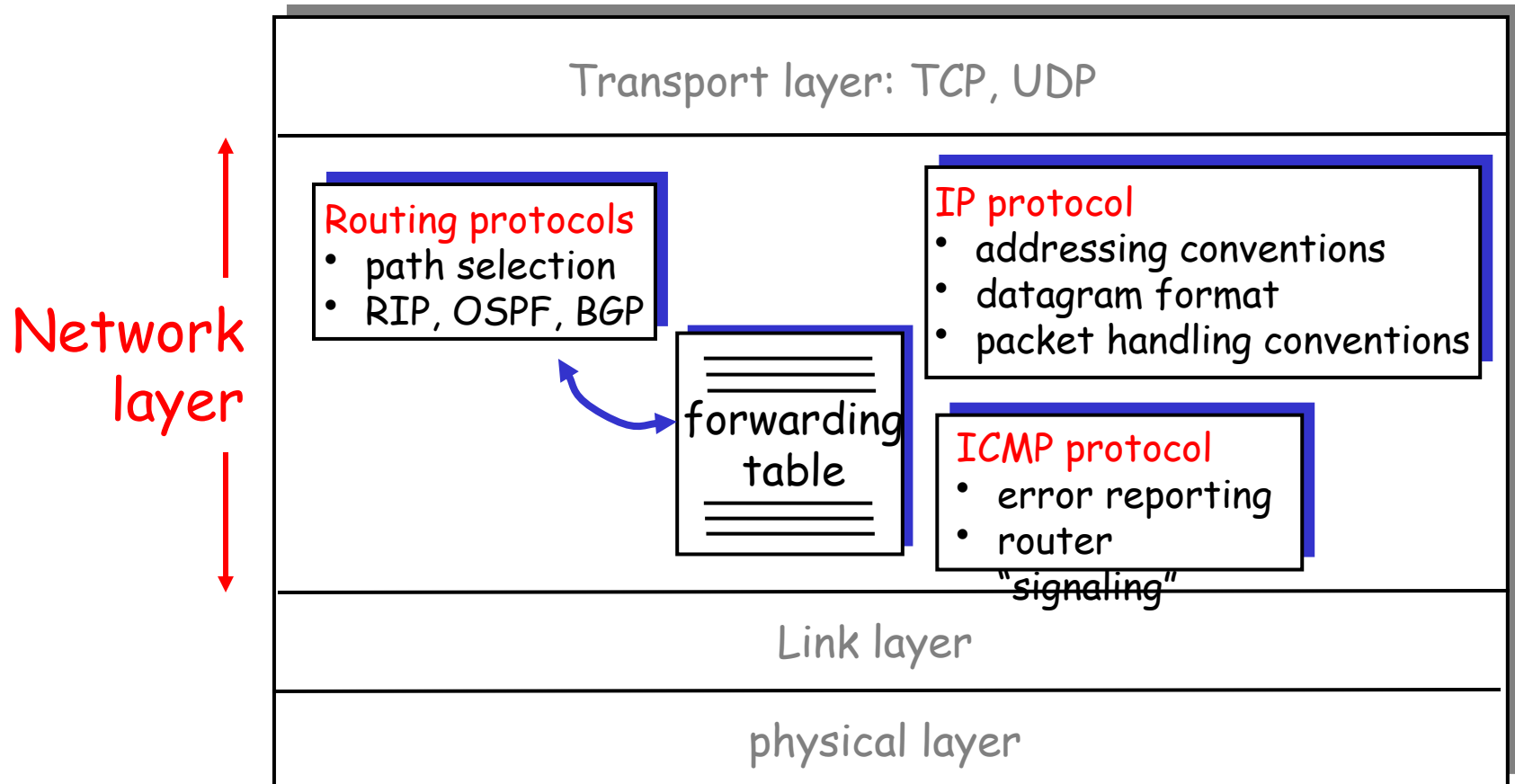
# Chapter 4: Network Layer

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



# The Internet Network layer

Host, router network layer functions:



# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link-state routing algorithms
- 4.3 Distance Vector
- 4.3.1 What's inside a router
  - Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 4.7 Broadcast and multicast routing

# IP datagram format

IP protocol version  
number

header length  
(bytes)

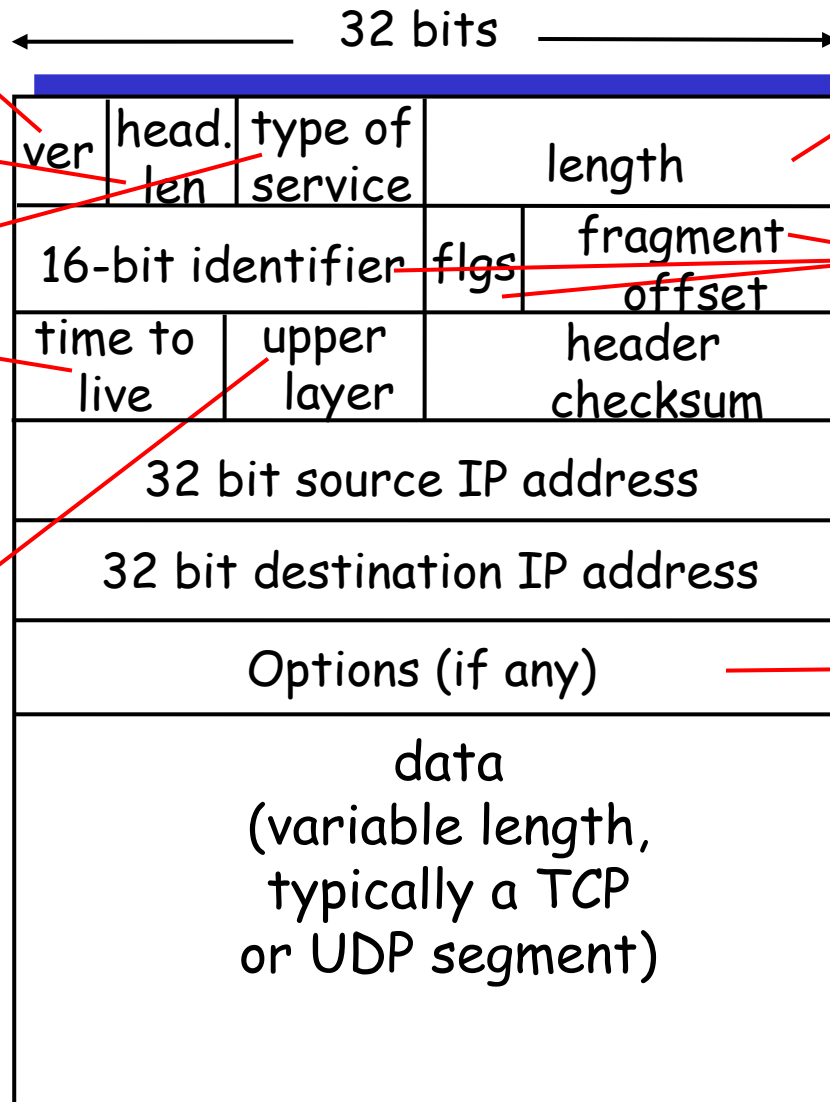
"type" of data

max number  
remaining hops  
(decremented at  
each router)

upper layer protocol  
to deliver payload to

how much overhead  
with TCP?

- ❑ 20 bytes of TCP
- ❑ 20 bytes of IP
- ❑ = 40 bytes + app layer overhead



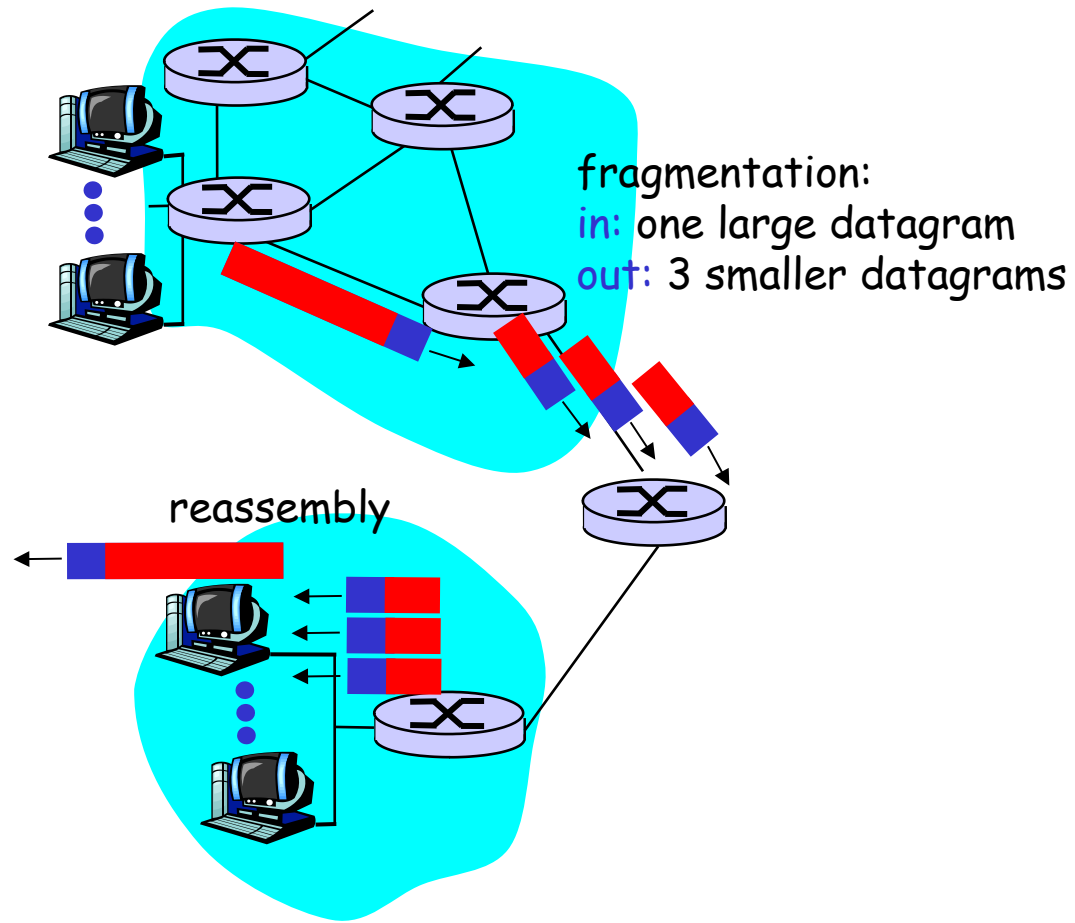
total datagram  
length (bytes)

for  
fragmentation/  
reassembly

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

# IP Fragmentation & Reassembly

- ❑ network links have MTU (max.transfer size) - largest possible link-level frame.
  - different link types, different MTUs
- ❑ large IP datagram divided ("fragmented") within net
  - one datagram becomes several datagrams
  - "reassembled" only at final destination
  - IP header bits used to identify, order related fragments



# IP Fragmentation and Reassembly

## Example

- ❑ 4000 byte datagram
- ❑ MTU = 1500 bytes

1480 bytes in  
data field

offset =  
 $1480/8$

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

One large datagram becomes  
several smaller datagrams

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

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

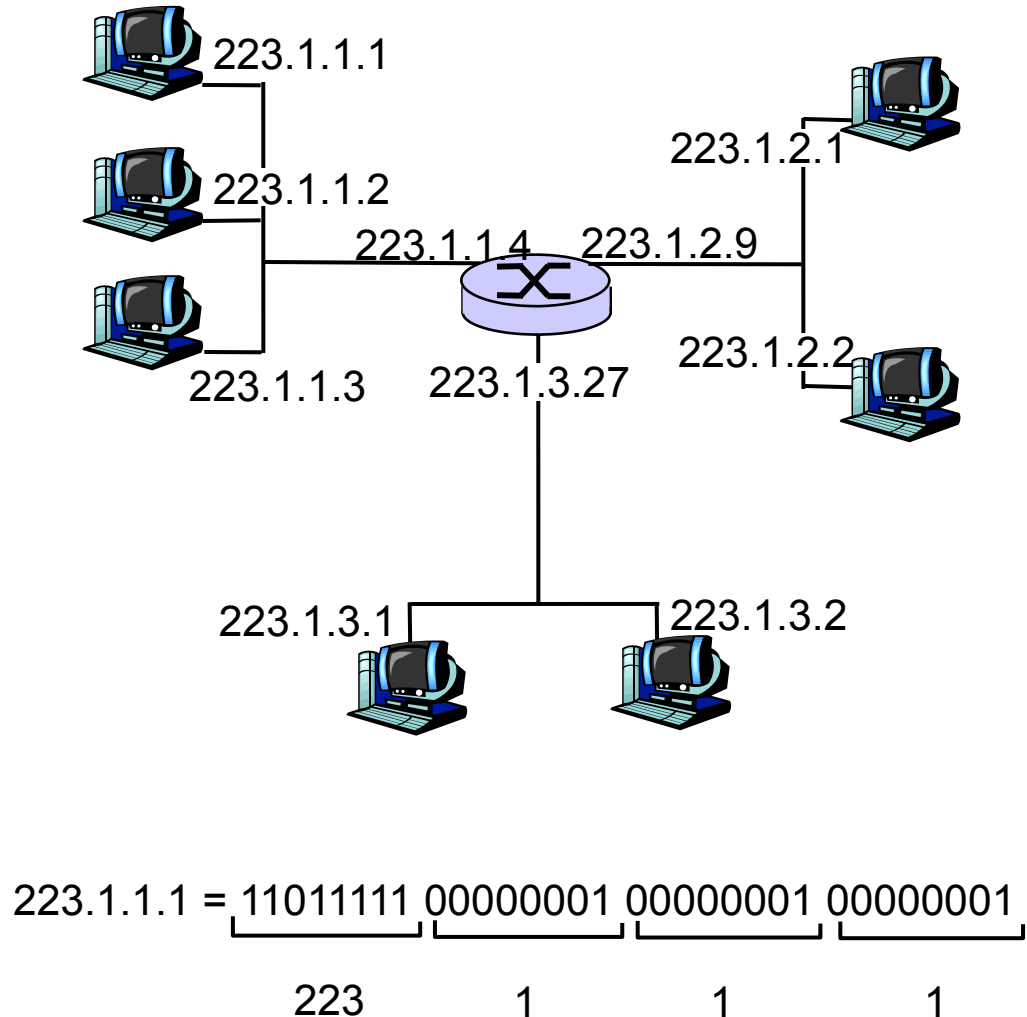
	length	ID	fragflag	offset	
	=1040	=x	=0	=370	

# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link-state
- 4.3 Distance Vector
- 4.3.1 What's inside a router
- 4.3.2 Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 4.7 Broadcast and multicast routing

# IP Addressing: introduction

- ❑ IP address: 32-bit identifier for host, router *interface*
- ❑ *interface*: connection between host/router and physical link
  - router's typically have multiple interfaces
  - host typically has one interface
  - IP addresses associated with each interface



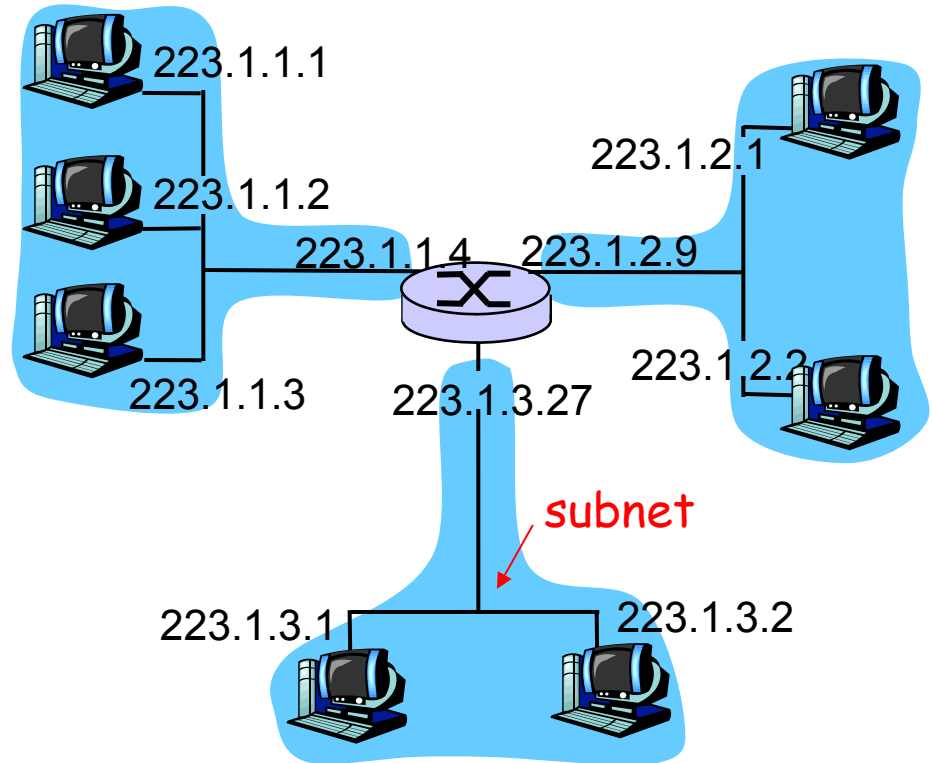
# Subnets

## ❑ 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



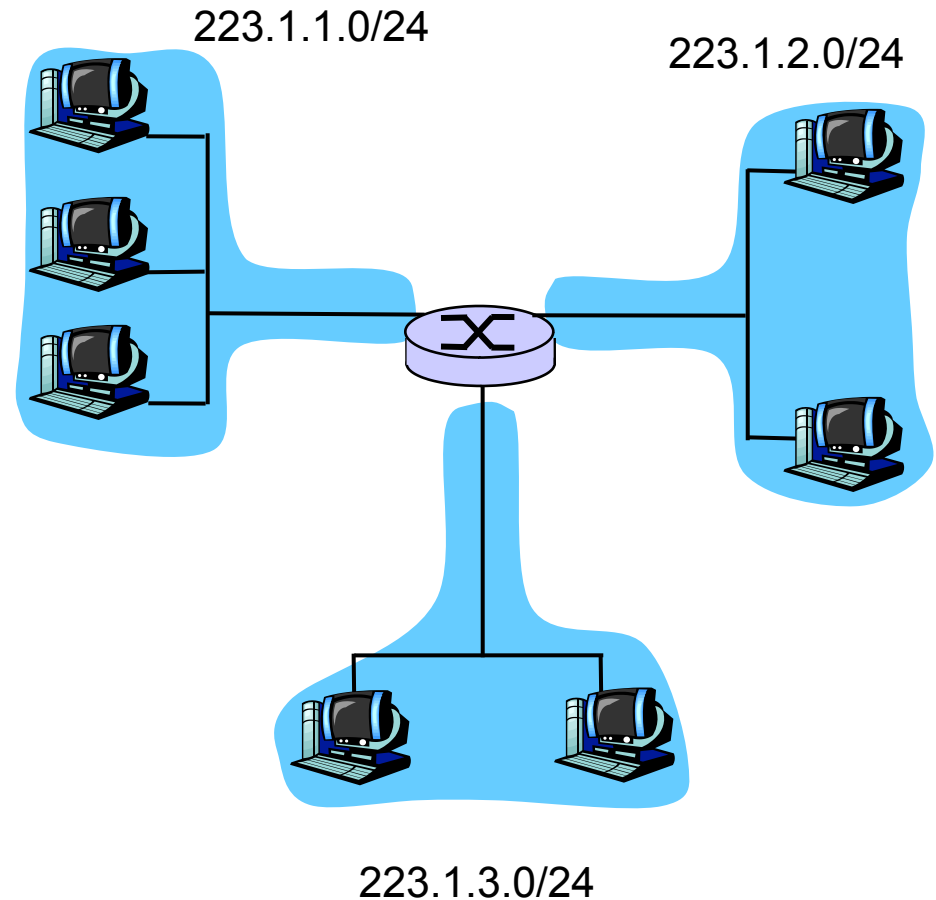
network consisting of 3 subnets



# Subnets

## Recipe

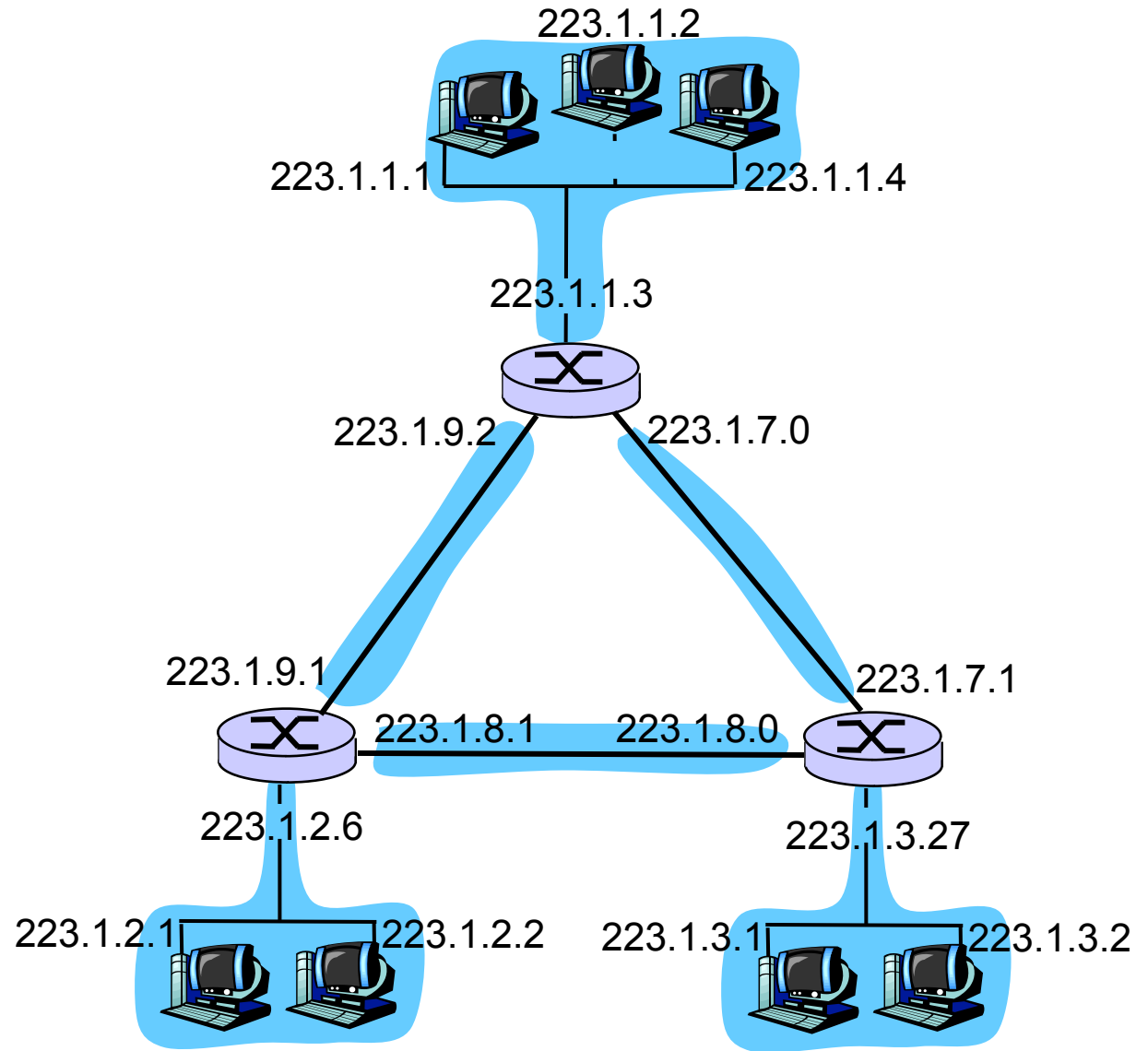
- To determine the subnets, detach each interface from its host or router, creating islands of isolated networks. Each isolated network is called a **subnet**.



Subnet mask: /24

# Subnets

How many?



# 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



200.23.16.0/23

# IP addresses: how to get one?

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

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

# DHCP: Dynamic Host Configuration Protocol

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

Can renew its lease on address in use

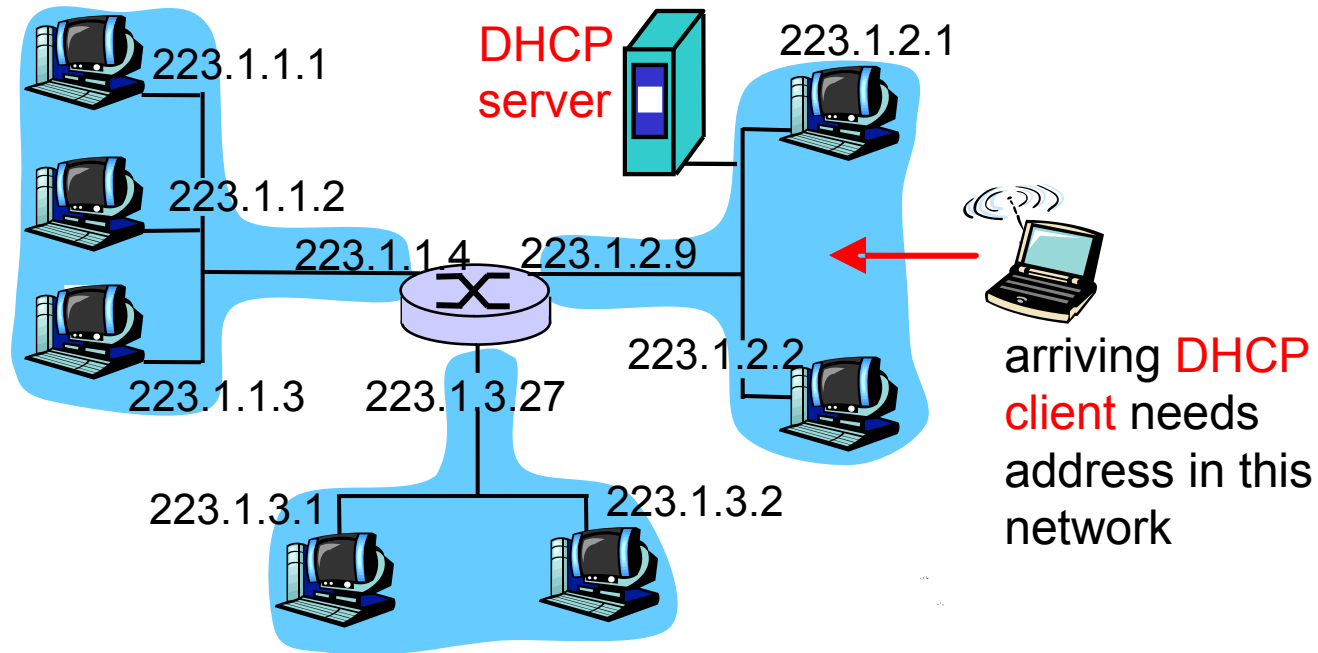
Allows reuse of addresses (only hold address while connected and "on")

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

DHCP overview:

- host broadcasts "DHCP discover" msg [optional]
- DHCP server responds with "DHCP offer" msg [optional]
- host requests IP address: "DHCP request" msg
- DHCP server sends address: "DHCP ack" msg

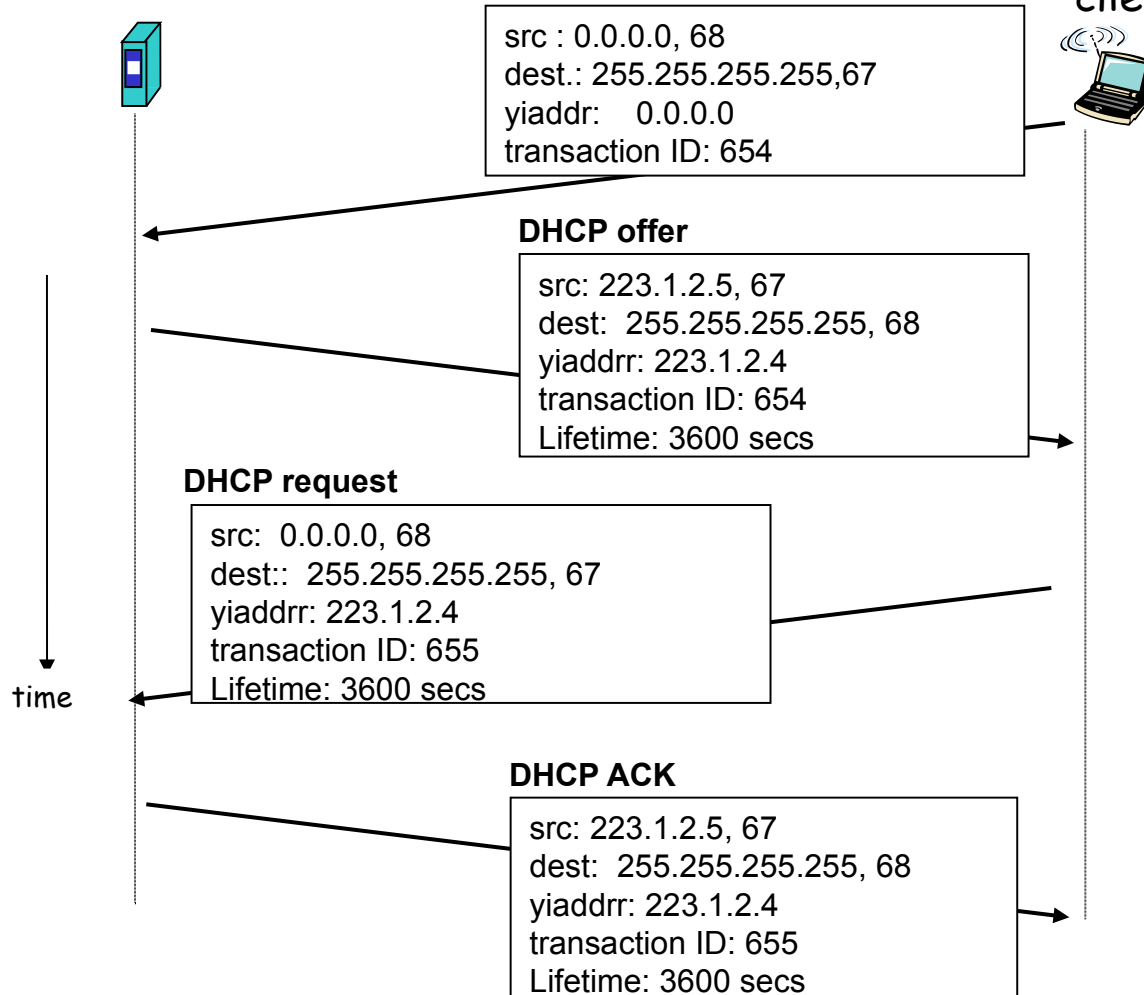
# DHCP client-server scenario



# DHCP client-server scenario

DHCP server: 223.1.2.5

arriving client



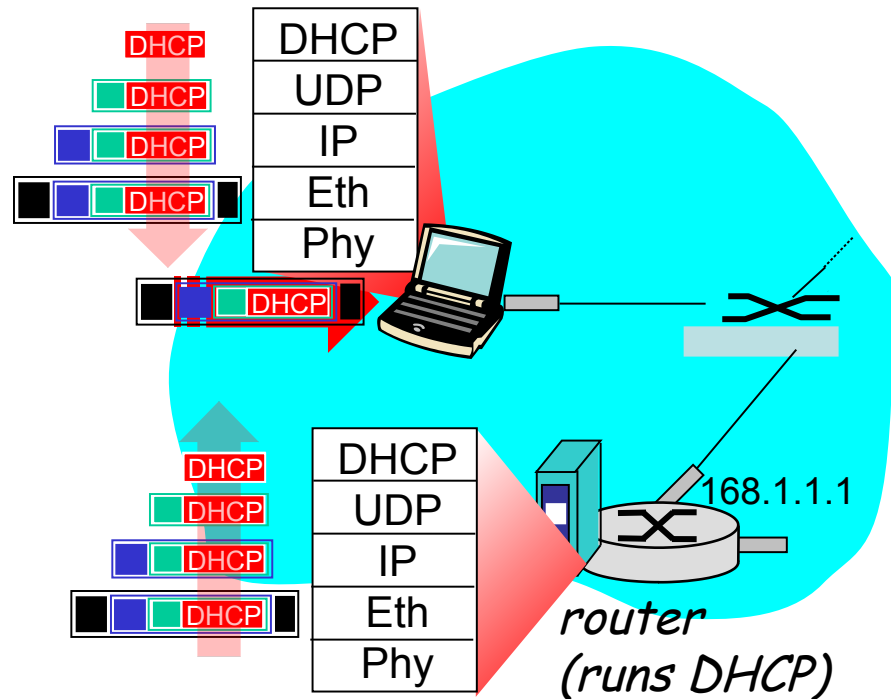
# DHCP: more than IP address

DHCP can return more than just allocated IP address on subnet:

- address of first-hop router for client
- name and IP address of DNS sever
- network mask (indicating network versus host portion of address)

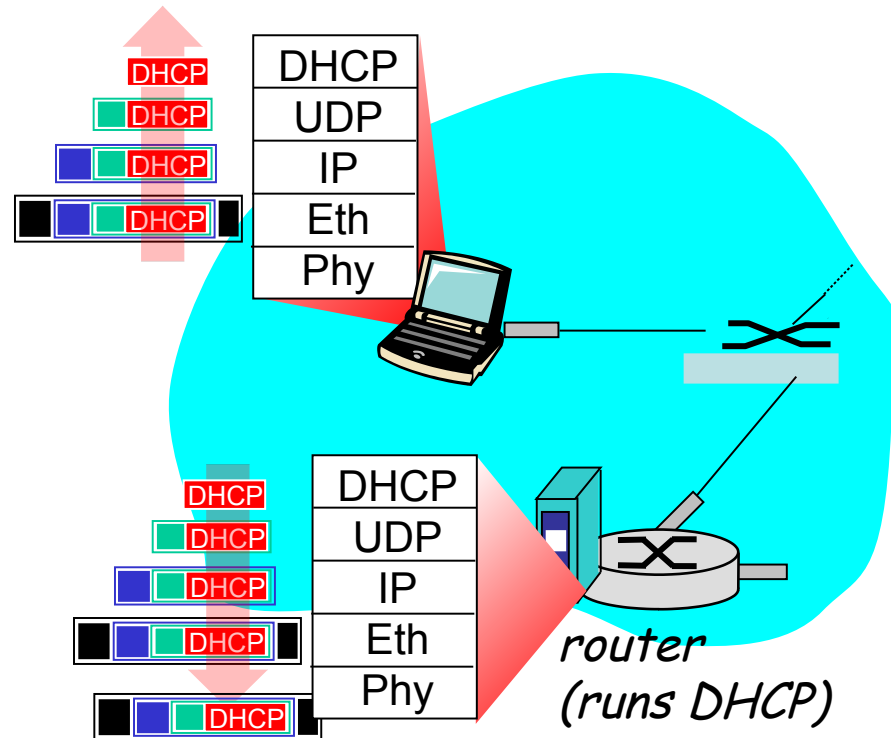


# DHCP: example



- connecting laptop needs its IP address, addr of first-hop router, addr of DNS server: use DHCP
- DHCP request encapsulated in UDP, encapsulated in IP, encapsulated in 802.1 Ethernet
- Ethernet frame broadcast (dest: FFFFFFFFFFFFFFFF) on LAN, received at router running DHCP server
- Ethernet demux'ed to IP demux'ed, UDP demux'ed to DHCP

# DHCP: example



- DCP server formulates DHCP ACK containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server
- encapsulation of DHCP server, frame forwarded to client, demux'ing up to DHCP at client
- client now knows its IP address, name and IP address of DSN server, IP address of its first-hop router

# DHCP: wireshark output (home LAN)

Message type: **Boot Request (1)**

Hardware type: Ethernet

Hardware address length: 6

Hops: 0

**Transaction ID: 0x6b3a11b7**

Seconds elapsed: 0

Bootp flags: 0x0000 (Unicast)

Client IP address: 0.0.0.0 (0.0.0.0)

Your (client) IP address: 0.0.0.0 (0.0.0.0)

Next server IP address: 0.0.0.0 (0.0.0.0)

Relay agent IP address: 0.0.0.0 (0.0.0.0)

**Client MAC address: Wistron\_23:68:8a (00:16:d3:23:68:8a)**

Server host name not given

Boot file name not given

Magic cookie: (OK)

Option: (t=53,l=1) **DHCP Message Type = DHCP Request**

Option: (61) Client identifier

Length: 7; Value: 010016D323688A;

Hardware type: Ethernet

Client MAC address: Wistron\_23:68:8a (00:16:d3:23:68:8a)

Option: (t=50,l=4) Requested IP Address = 192.168.1.101

Option: (t=12,l=5) Host Name = "nomad"

**Option: (55) Parameter Request List**

Length: 11; Value: 010F03062C2E2F1F21F92B

**1 = Subnet Mask; 15 = Domain Name**

**3 = Router; 6 = Domain Name Server**

44 = NetBIOS over TCP/IP Name Server

.....

request

Message type: **Boot Reply (2)**

Hardware type: Ethernet

Hardware address length: 6

Hops: 0

**Transaction ID: 0x6b3a11b7**

Seconds elapsed: 0

Bootp flags: 0x0000 (Unicast)

**Client IP address: 192.168.1.101 (192.168.1.101)**

Your (client) IP address: 0.0.0.0 (0.0.0.0)

**Next server IP address: 192.168.1.1 (192.168.1.1)**

Relay agent IP address: 0.0.0.0 (0.0.0.0)

Client MAC address: Wistron\_23:68:8a (00:16:d3:23:68:8a)

Server host name not given

Boot file name not given

Magic cookie: (OK)

**Option: (t=53,l=1) DHCP Message Type = DHCP ACK**

**Option: (t=54,l=4) Server Identifier = 192.168.1.1**

**Option: (t=1,l=4) Subnet Mask = 255.255.255.0**

**Option: (t=3,l=4) Router = 192.168.1.1**

**Option: (6) Domain Name Server**

**Length: 12; Value: 445747E2445749F244574092;**

**IP Address: 68.87.71.226;**

**IP Address: 68.87.73.242;**

**IP Address: 68.87.64.146**

**Option: (t=15,l=20) Domain Name = "hsd1.ma.comcast.net."**

reply

# IP addresses: how to get one?

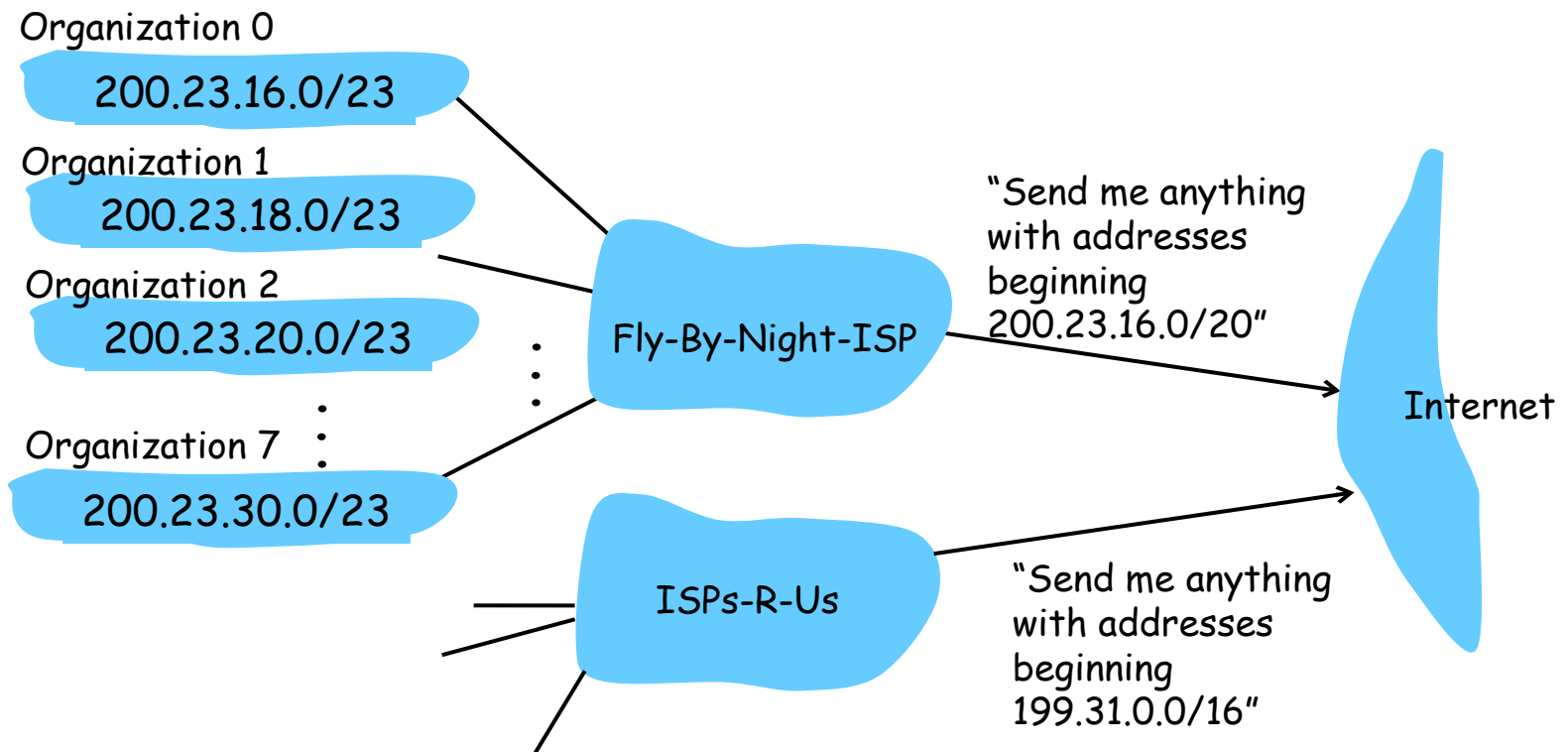
Q: How does *network* get subnet part of IP addr?

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

ISP's block	<u>11001000</u>	<u>00010111</u>	<u>00010000</u>	00000000	200.23.16.0/20
Organization 0	<u>11001000</u>	<u>00010111</u>	<u>00010000</u>	00000000	200.23.16.0/23
Organization 1	<u>11001000</u>	<u>00010111</u>	<u>00010010</u>	00000000	200.23.18.0/23
Organization 2	<u>11001000</u>	<u>00010111</u>	<u>00010100</u>	00000000	200.23.20.0/23
...	.....			....	....
Organization 7	<u>11001000</u>	<u>00010111</u>	<u>00011110</u>	00000000	200.23.30.0/23

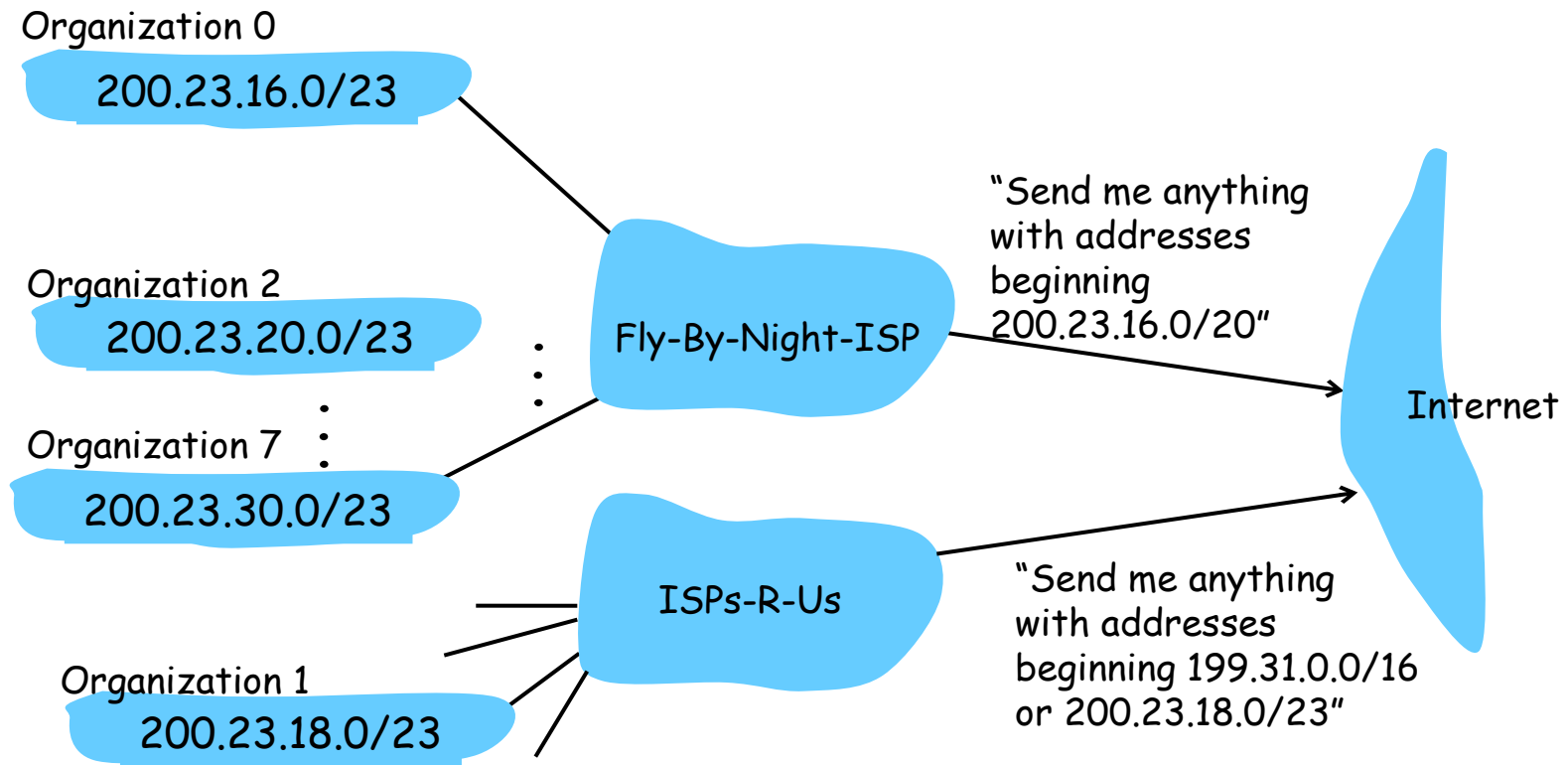
# Hierarchical addressing: route aggregation

Hierarchical addressing allows efficient advertisement of routing information:



# Hierarchical addressing: more specific routes

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



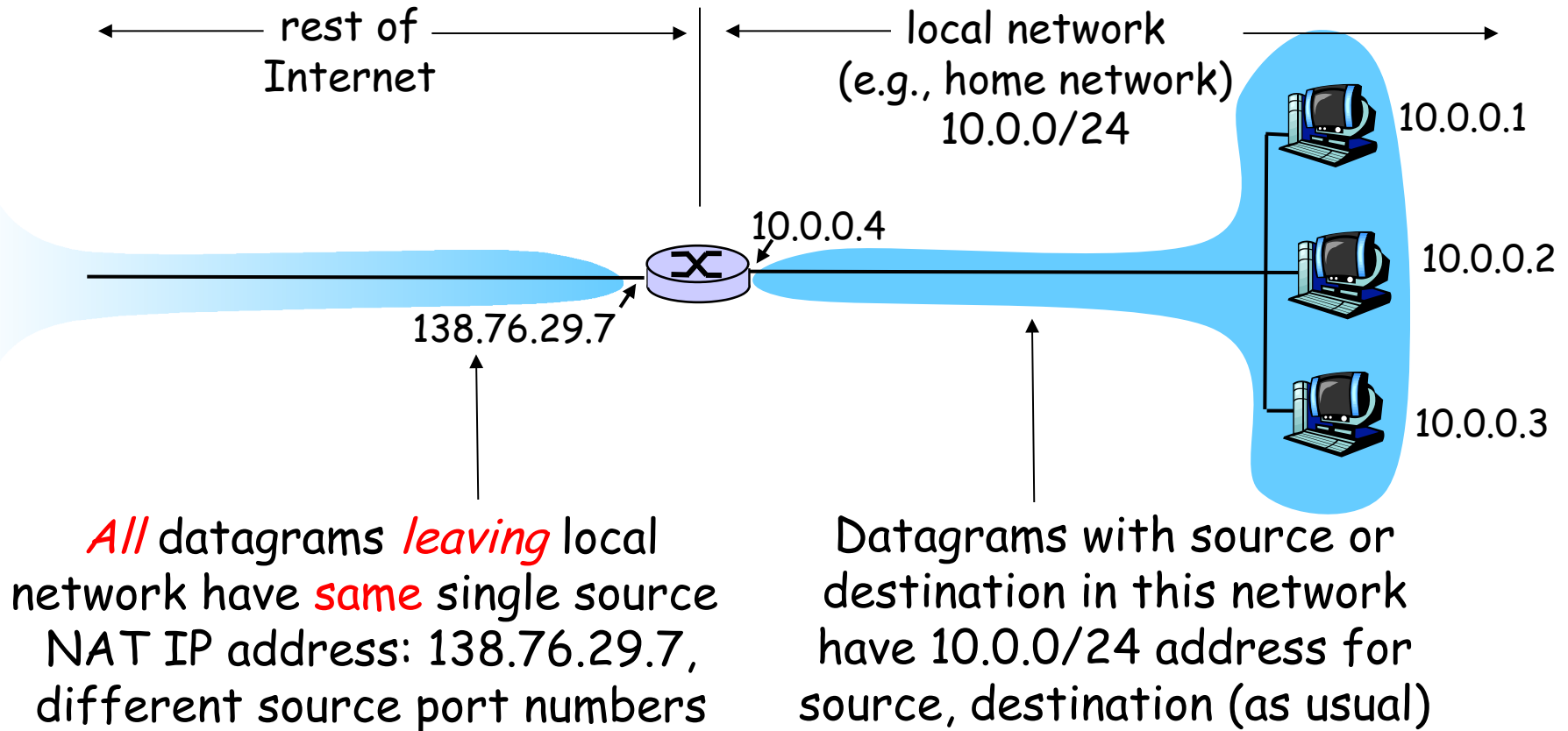
## IP addressing: the last word...

Q: How does an ISP get block of addresses?

A: **ICANN**: Internet **C**orporation for **A**ssigned  
**N**ames and **N**umbers

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

# NAT: Network Address Translation





# NAT: Network Address Translation

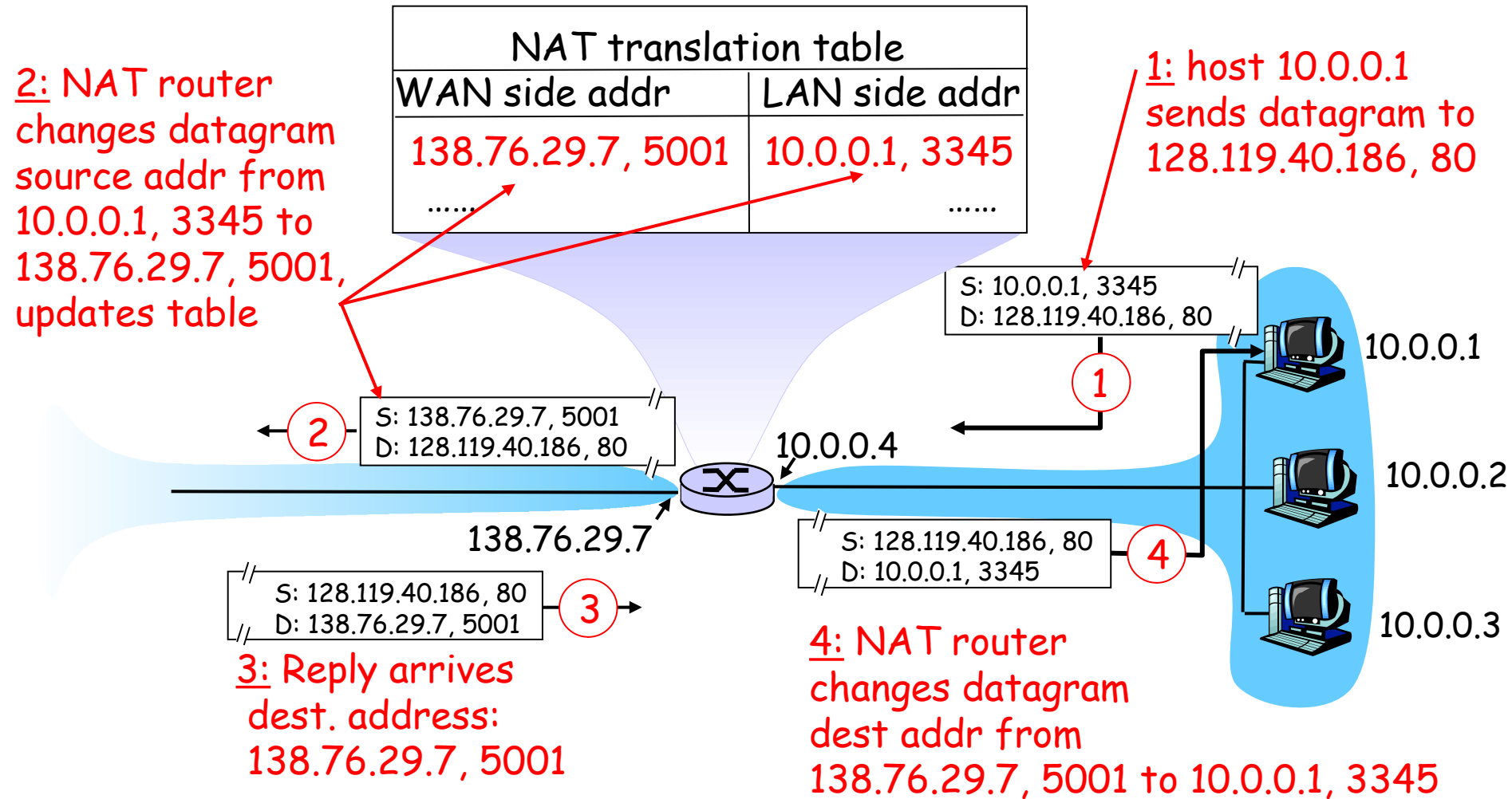
- ❑ **Motivation:** local network uses just one IP address as far as outside world is concerned:
  - range of addresses not needed from ISP: just one IP address for all devices
  - can change addresses of devices in local network without notifying outside world
  - can change ISP without changing addresses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

# NAT: Network Address Translation

**Implementation:** NAT router must:

- *outgoing datagrams: replace* (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #)  
... remote clients/servers will respond using (NAT IP address, new port #) as destination addr.
- *remember (in NAT translation table)* every (source IP address, port #) to (NAT IP address, new port #) translation pair
- *incoming datagrams: replace* (NAT IP address, new port #) in dest fields of every incoming datagram with corresponding (source IP address, port #) stored in NAT table

# NAT: Network Address Translation

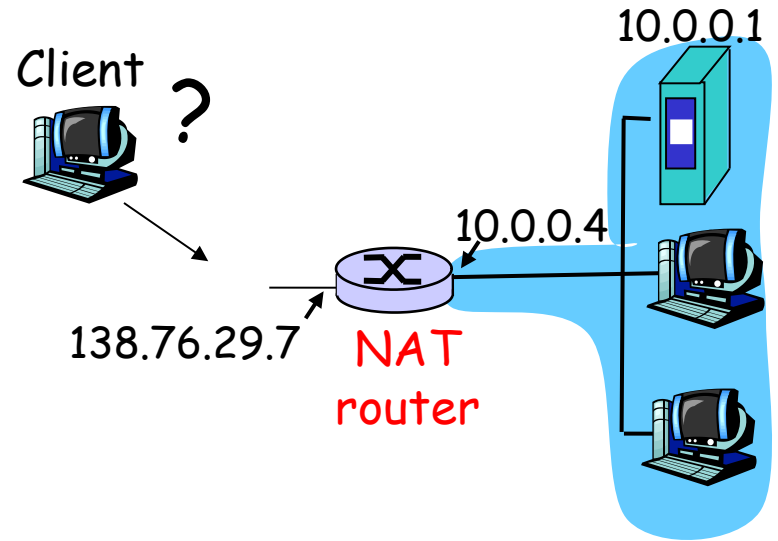


# NAT: Network Address Translation

- ❑ 16-bit port-number field:
  - 60,000 simultaneous connections with a single LAN-side address!
- ❑ NAT is controversial:
  - routers should only process up to layer 3
  - violates end-to-end argument
    - NAT possibility must be taken into account by app designers, eg, P2P applications
  - address shortage should instead be solved by IPv6

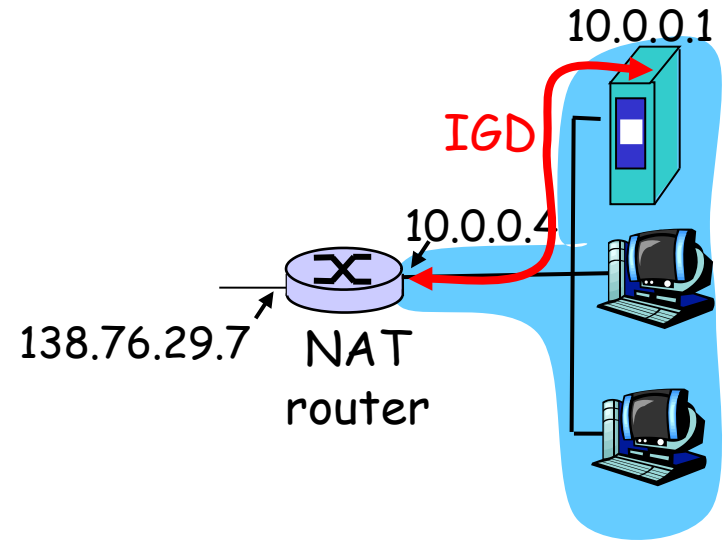
# NAT traversal problem

- ❑ client wants to connect to server with address 10.0.0.1
  - server address 10.0.0.1 local to LAN (client can't use it as destination addr)
  - only one externally visible NATted address: 138.76.29.7
- ❑ solution 1: statically configure NAT to forward incoming connection requests at given port to server
  - e.g., (123.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000



# NAT traversal problem

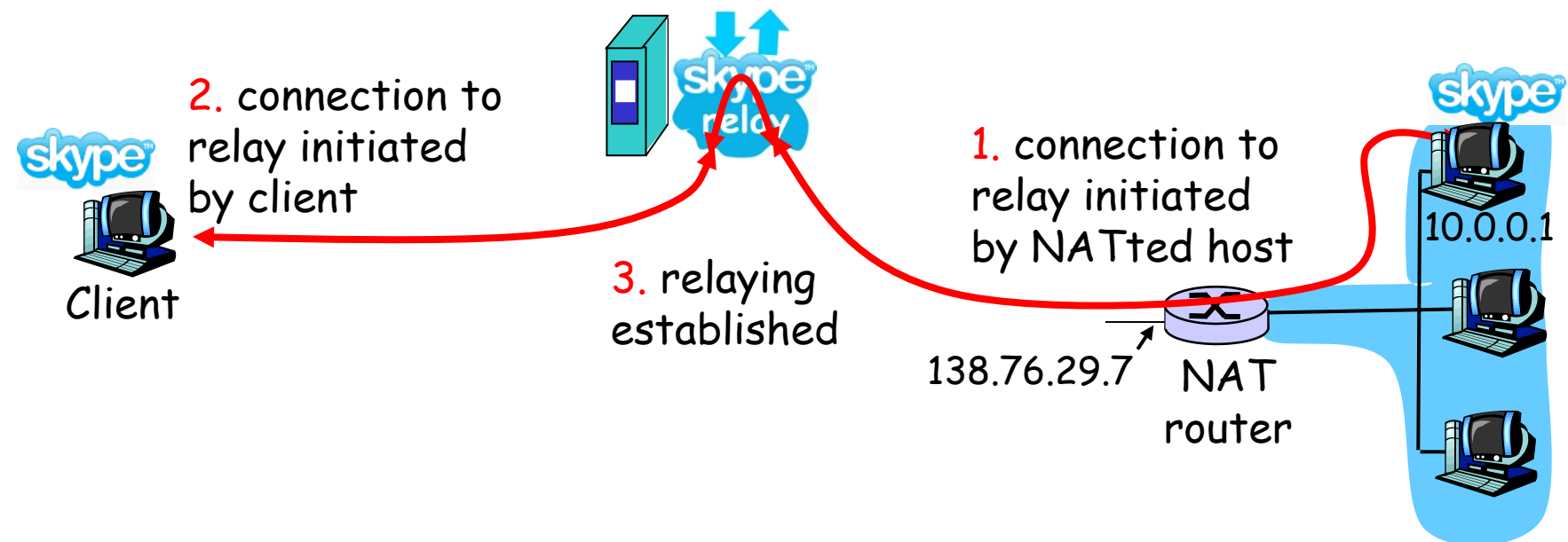
- solution 2: Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATted host to:
  - ❖ learn public IP address (138.76.29.7)
  - ❖ add/remove port mappings (with lease times)



i.e., automate static NAT port map configuration

# NAT traversal problem

- solution 3: relaying (used in Skype)
  - NATed client establishes connection to relay
  - External client connects to relay
  - relay bridges packets between to connections



# 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



# 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:
 

Type	Code	description
0	0	echo reply (ping)
3	1	dest. network unreachable
3	2	dest host unreachable
3	3	dest protocol unreachable
3	4	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
- **ICMP message:** type, code plus first 8 bytes of IP datagram causing error

# Traceroute and ICMP

- ❑ Source sends series of UDP segments to dest
  - First has TTL =1
  - Second has TTL=2, etc.
  - Unlikely port number
- ❑ When nth datagram arrives to nth router:
  - Router discards datagram
  - And sends to source an ICMP message (type 11, code 0)
  - Message includes name of router& IP address
- ❑ When ICMP message arrives, source calculates RTT
- ❑ Traceroute does this 3 times

## Stopping criterion

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

# 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.
- ❑ **Additional motivation:**
  - header format helps speed processing/forwarding
  - header changes to facilitate QoS

## **IPv6 datagram format:**

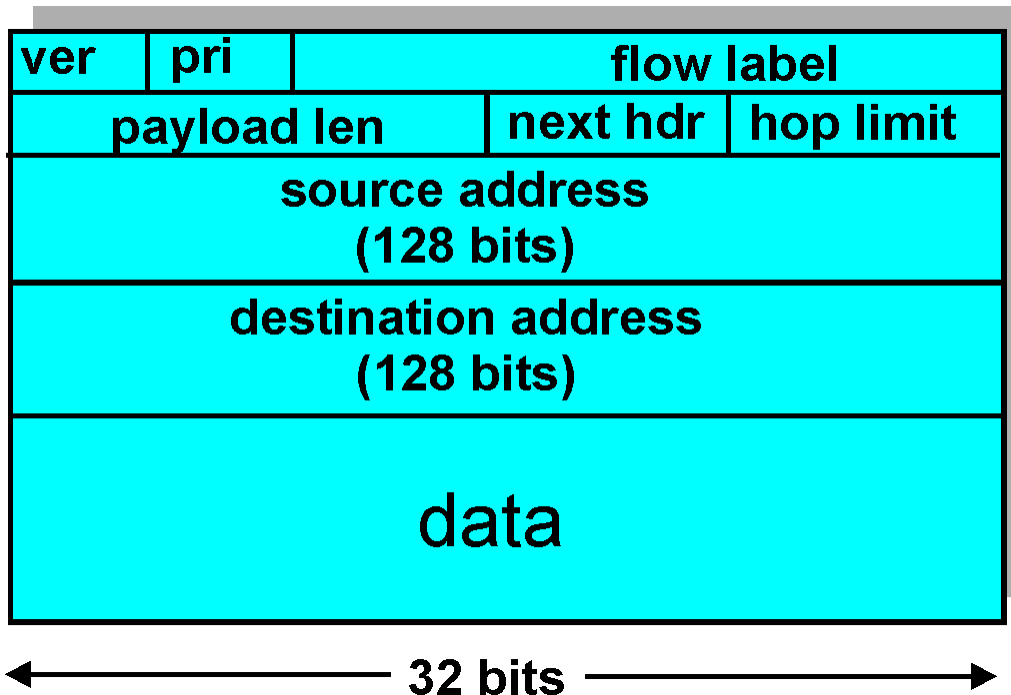
- fixed-length 40 byte header
- no fragmentation allowed

# IPv6 Header (Cont)

*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



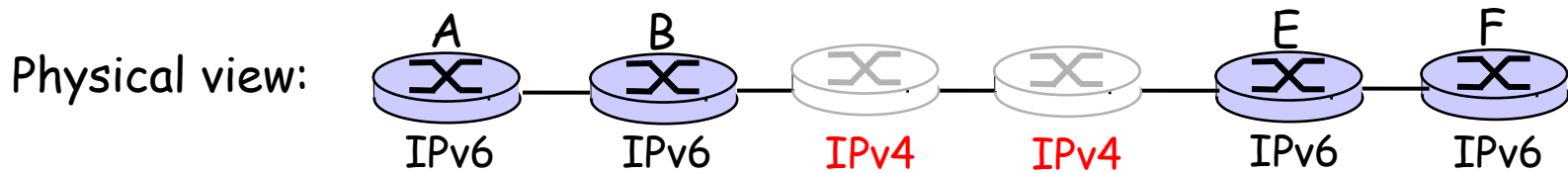
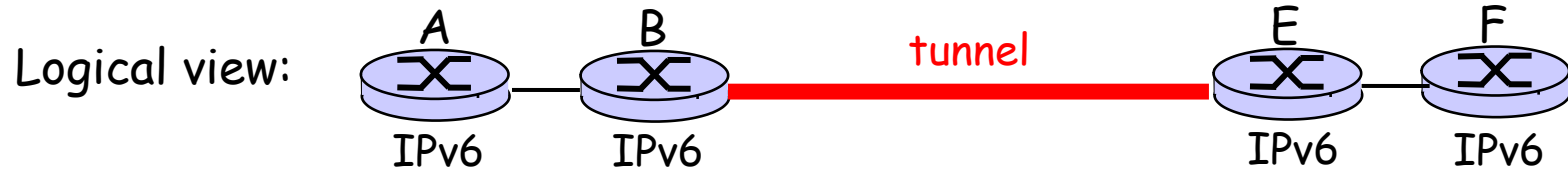
# Other Changes from IPv4

- ❑ *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

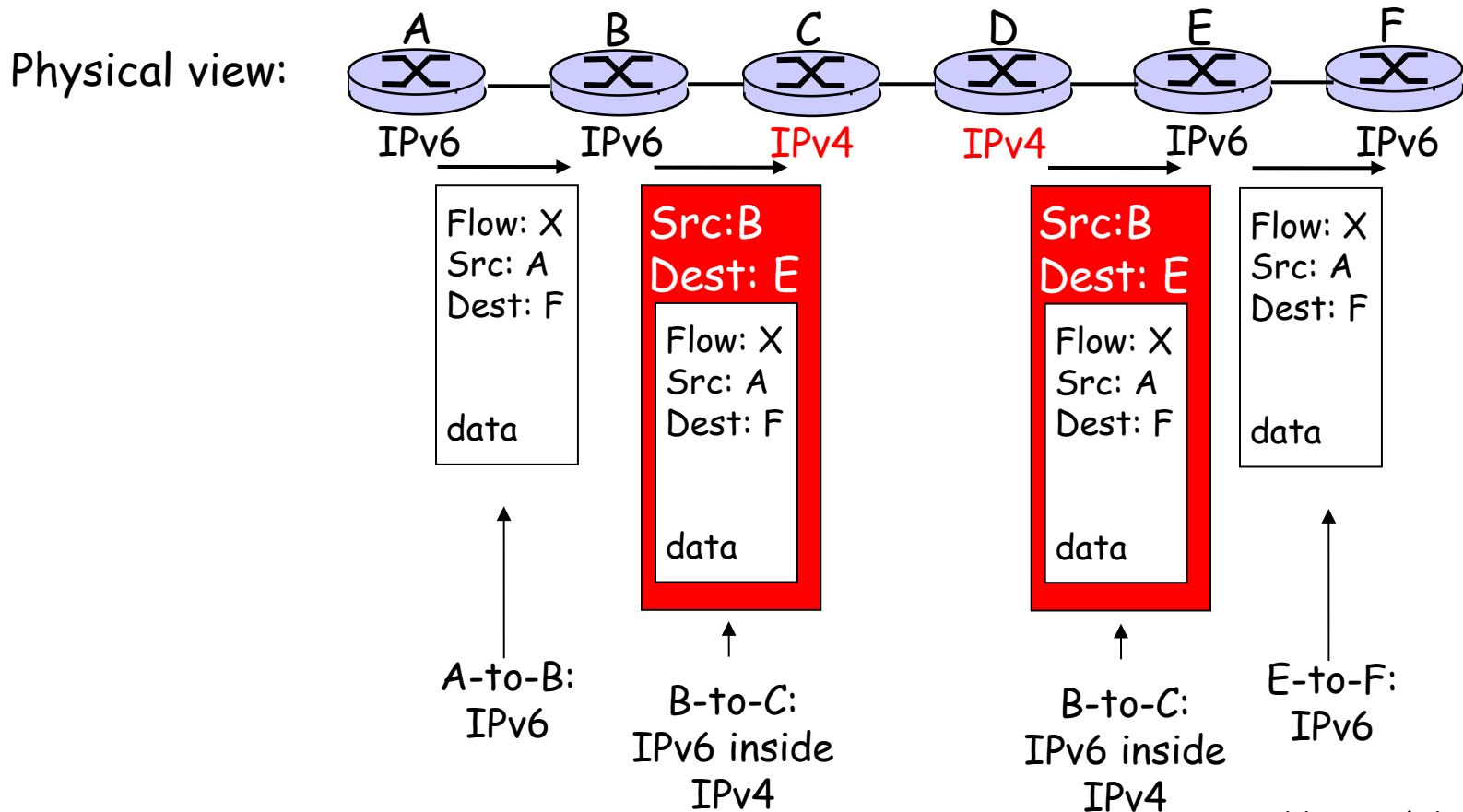
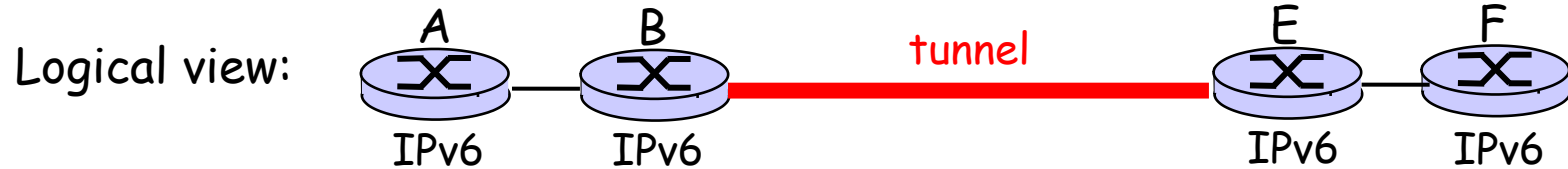
- ❑ Not all routers can be upgraded simultaneous
  - no “flag days”
  - How will the network operate with mixed IPv4 and IPv6 routers?
- ❑ *Tunneling*: IPv6 carried as payload in IPv4 datagram among IPv4 routers

# Tunneling





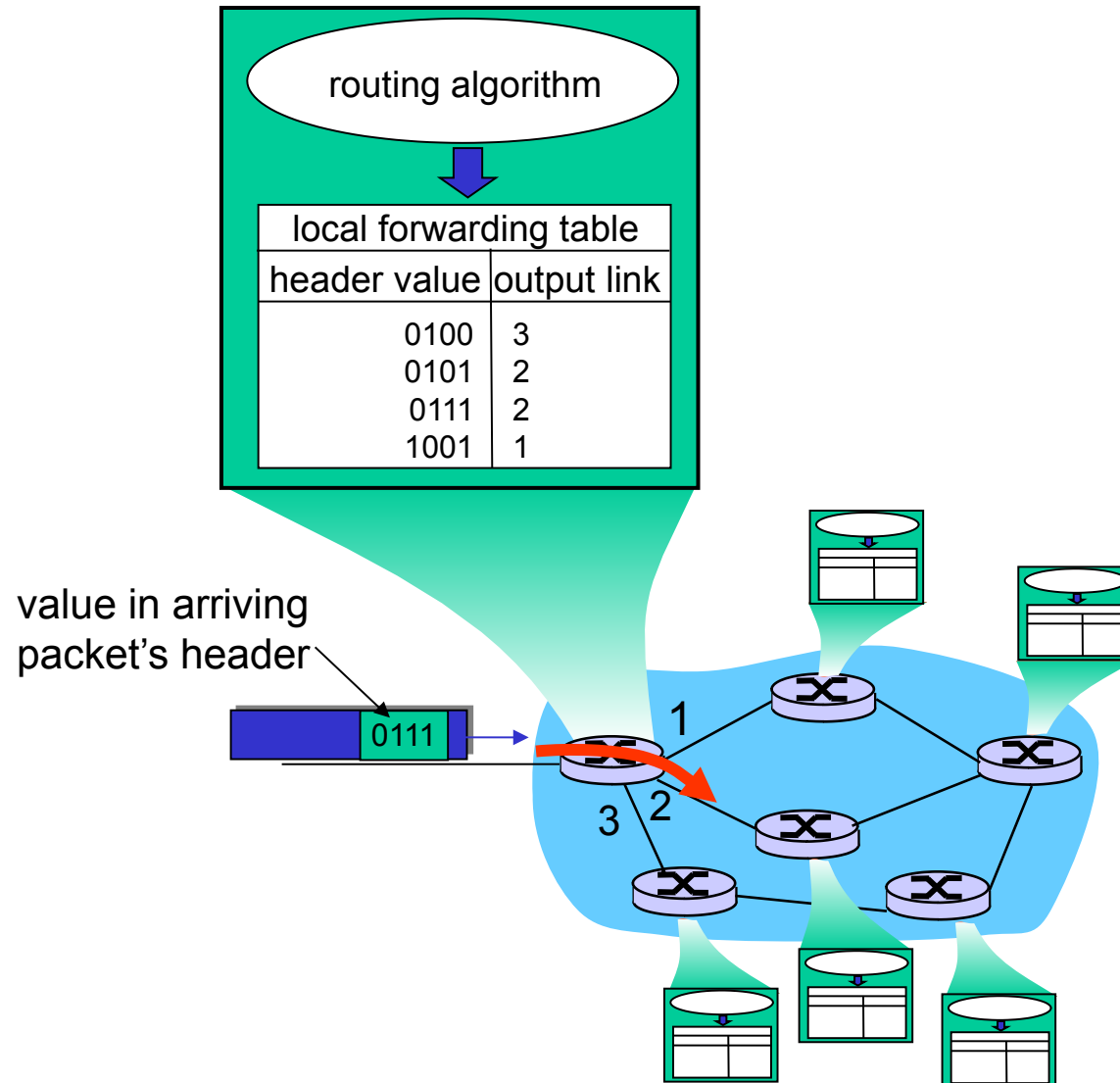
# Tunneling



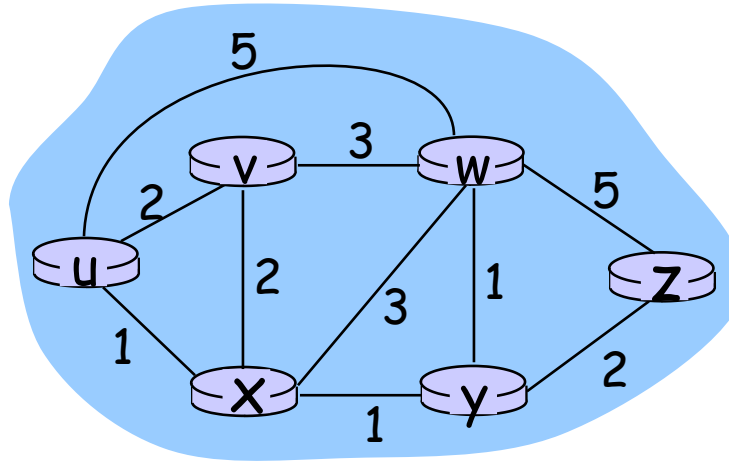
# Chapter 4: Network Layer

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

# Interplay between routing, forwarding



# Graph abstraction



Graph:  $G = (N, E)$

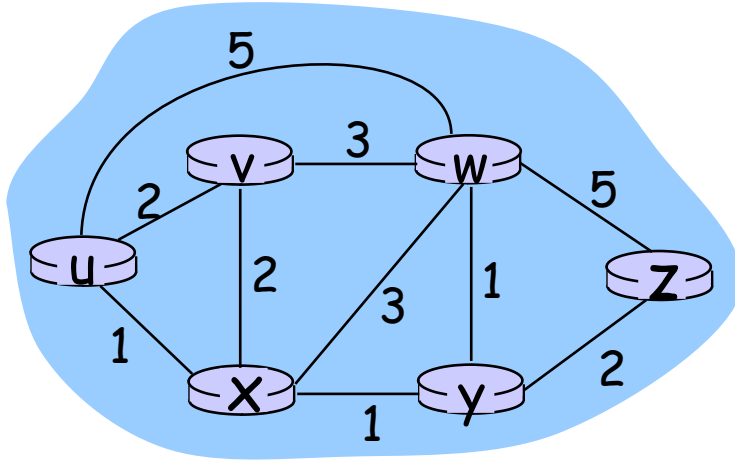
$N$  = set of routers =  $\{ u, v, w, x, y, z \}$

$E$  = 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



- $c(x, x') = \text{cost of link } (x, x')$ 
  - e.g.,  $c(w, z) = 5$
- cost could always be 1, or inversely related to bandwidth, or inversely related to congestion

Cost of path  $(x_1, x_2, x_3, \dots, x_p) = c(x_1, x_2) + c(x_2, x_3) + \dots + c(x_{p-1}, x_p)$

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

Routing algorithm: algorithm that finds least-cost path

# Routing Algorithm classification

## Global or decentralized information?

### Global:

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

### Decentralized:

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

## Static or dynamic?

### Static:

- ❑ routes change slowly over time

### Dynamic:

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

# Chapter 4: Network Layer

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

# A Link-State Routing Algorithm

## Dijkstra's algorithm

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



## Notation:

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

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

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

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

# Dijkstra's Algorithm

1 **Initialization:**

2  $N' = \{u\}$

3 for all nodes  $v$

4 if  $v$  adjacent to  $u$

5 then  $D(v) = c(u,v)$

6 else  $D(v) = \infty$

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  $D(v) = \min( D(v), D(w) + c(w,v) )$

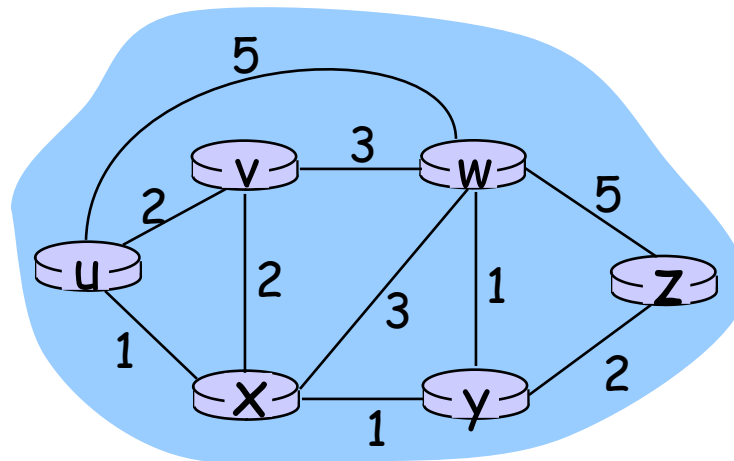
13 /\* new cost to  $v$  is either old cost to  $v$  or known

14 shortest path cost to  $w$  plus cost from  $w$  to  $v$  \*/

15 **until all nodes in  $N'$**

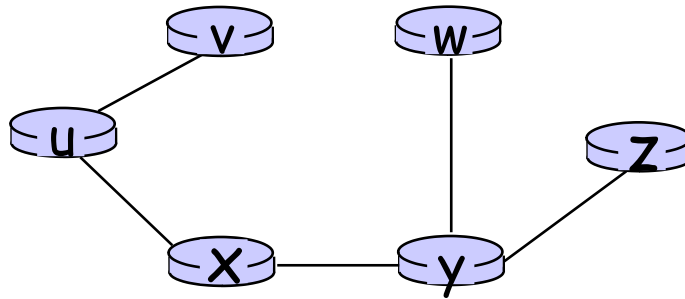
# Dijkstra's algorithm: example

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	$\infty$	$\infty$
1	ux	2,u	4,x		2,x	$\infty$
2	uxy	2,u	3,y			4,y
3	uxyv		3,y			4,y
4	uxyvw					4,y
5	uxyvwz					



# Dijkstra's algorithm: example (2)

Resulting shortest-path tree from u:



Resulting forwarding table in u:

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

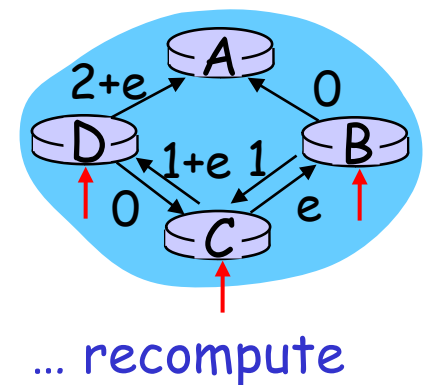
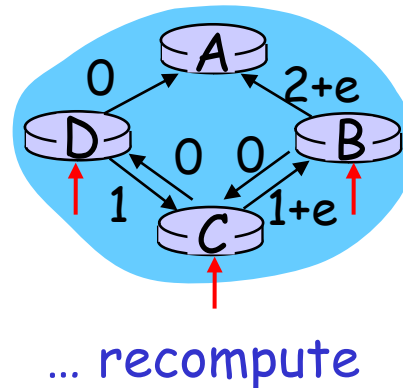
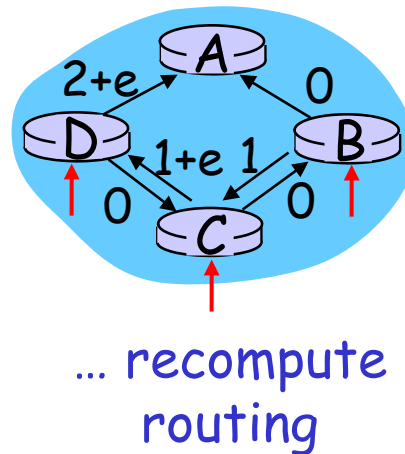
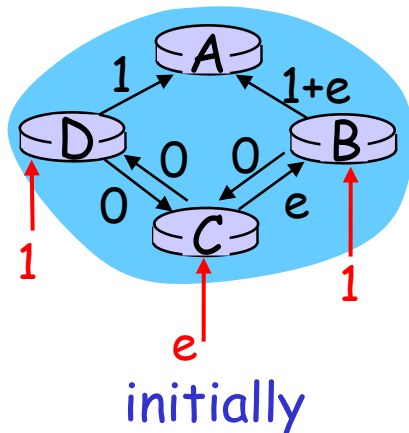
# Dijkstra's algorithm, discussion

**Algorithm complexity:**  $n$  nodes

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

**Oscillations possible:**

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



# Chapter 4: Network Layer

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

# Distance Vector Algorithm

## Bellman-Ford Equation (dynamic programming)

Define

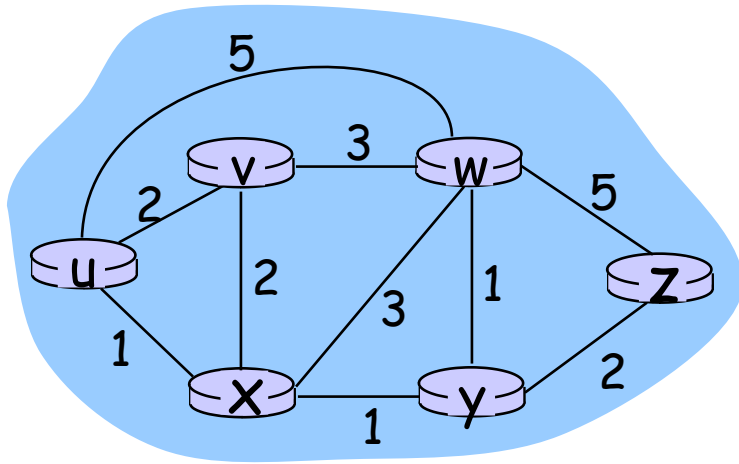
$d_x(y) :=$  cost of least-cost path from  $x$  to  $y$

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:

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

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



# Distance Vector Algorithm

- $D_x(y)$  = estimate of least cost from  $x$  to  $y$
- Node  $x$  knows cost to each neighbor  $v$ :  
 $c(x,v)$
- Node  $x$  maintains distance vector  $D_x = [D_x(y): y \in N]$
- 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:

- ❑ From time-to-time, each node sends its own distance vector estimate to neighbors
- ❑ Asynchronous
- ❑ When a node  $x$  receives new DV estimate from neighbor, it updates its own DV using B-F equation:

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

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

# Distance Vector Algorithm (5)

## Iterative, asynchronous:

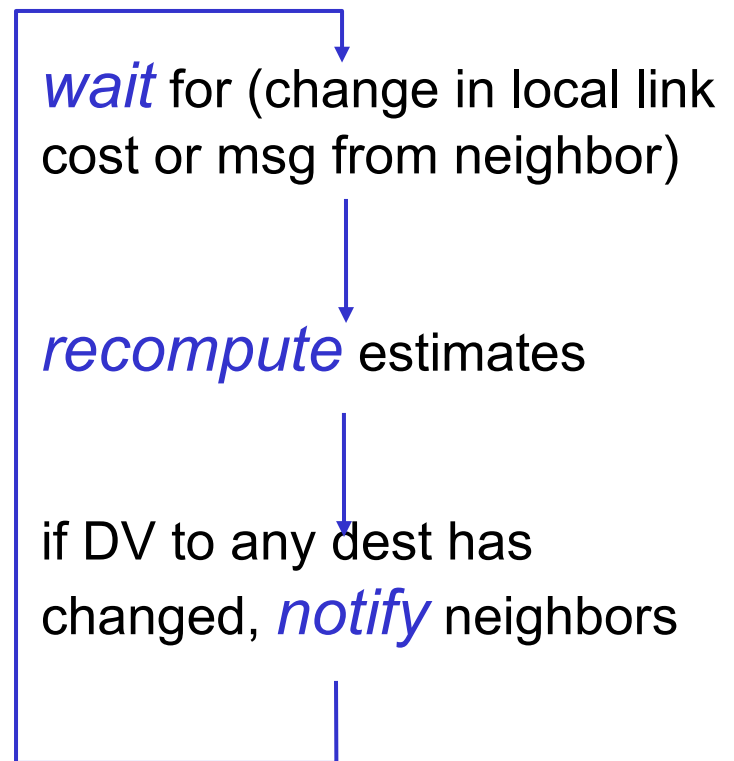
each local iteration caused by:

- local link cost change
- DV update message from neighbor

## Distributed:

- each node notifies neighbors *only* when its DV changes
  - neighbors then notify their neighbors if necessary

## Each node:



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

$$= \min\{2+0, 7+1\} = 2$$

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

$$= \min\{2+1, 7+0\} = 3$$

node x table

		cost to		
from		x	y	z
	x	0	2	7
	y	∞	∞	∞
	z	∞	∞	∞

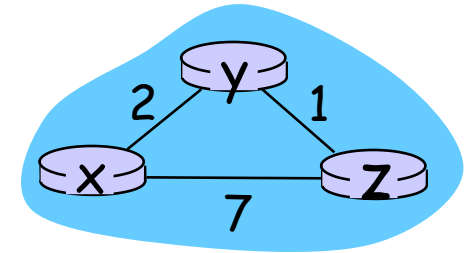
node y table

		cost to		
from		x	y	z
	x	∞	∞	∞
	y	2	0	1
	z	∞	∞	∞

node z table

		cost to		
from		x	y	z
	x	∞	∞	∞
	y	∞	∞	∞
	z	7	1	0

		cost to		
from		x	y	z
	x	0	2	3
	y	2	0	1
	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		
		x	y	z
from	x	0	2	7
	y	$\infty$	$\infty$	$\infty$
	z	$\infty$	$\infty$	$\infty$

node y table

		cost to		
from		x	y	z
	x	$\infty$	$\infty$	$\infty$
	y	2	0	1
	z	$\infty$	$\infty$	$\infty$

node z table

		cost to		
		x	y	z
from	x	$\infty$	$\infty$	$\infty$
	y	$\infty$	$\infty$	$\infty$
	z	7	1	0

		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	7	1	0

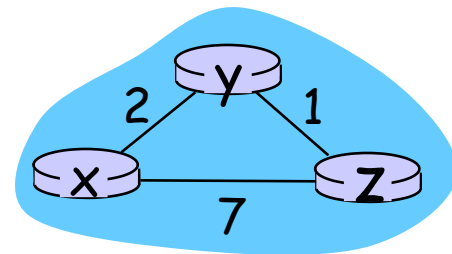
		cost to		
from		x	y	z
	x	0	2	7
	y	2	0	1
	z	7	1	0

		cost to		
		x	y	z
from	x	0	2	7
	y	2	0	1
	z	3	1	0

		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	3	1	0

		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	3	1	0

		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	3	1	0

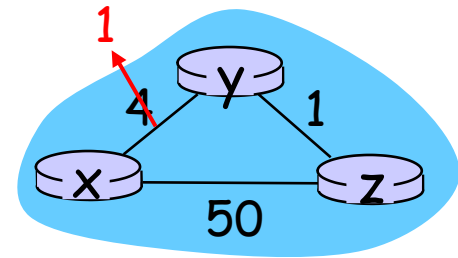


time →

# 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.

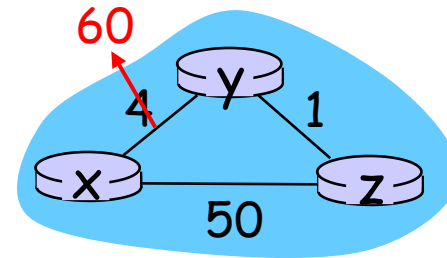
At time  $t_1$ ,  $z$  receives the update from  $y$  and updates its table. It computes a new least cost to  $x$  and sends its neighbors its DV.

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

## Link cost changes:

- ❑ good news travels fast
- ❑ bad news travels slow - "count to infinity" problem!
- ❑ 44 iterations before algorithm stabilizes: see text



## Poisoned reverse:

- ❑ If Z routes through Y to get to X :
  - Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)
- ❑ will this completely solve count to infinity problem?

# Comparison of LS and DV algorithms

## Message complexity

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

## Speed of Convergence

- LS:  $O(n^2)$  algorithm requires  $O(nE)$  msgs
  - may have oscillations
- DV: convergence time varies
  - may be routing loops
  - count-to-infinity problem

**Robustness:** what happens if router malfunctions?

## LS:

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

## DV:

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



# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link-state circuit and datagram networks
- 4.3 Distance Vector
- 4.3.1 What's inside a router
- 4.3.2 Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 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 dest's in routing tables!
- ❑ routing table exchange would swamp links!

**administrative autonomy**

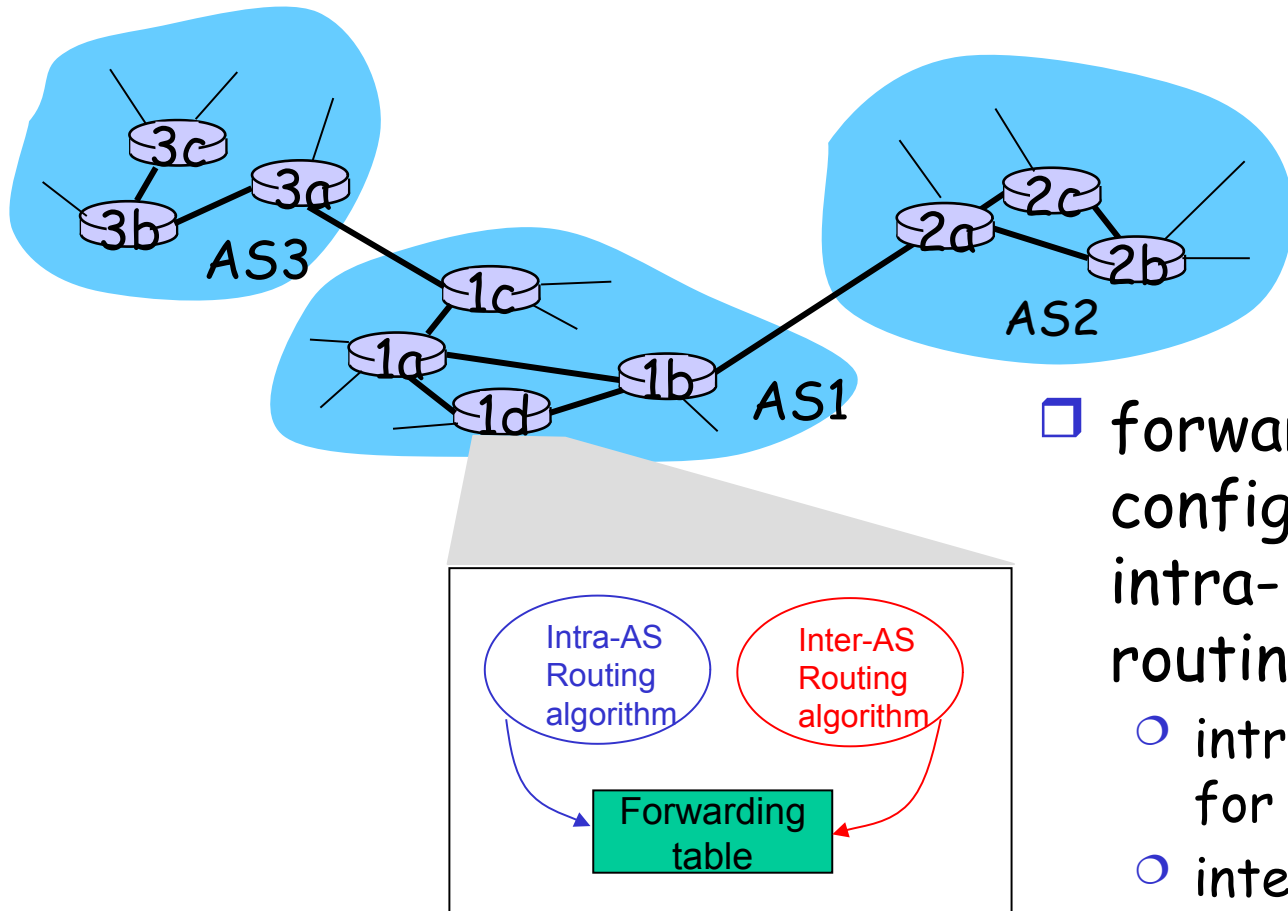
- ❑ internet = network of networks
- ❑ each network admin may want to control routing in its own network

# Hierarchical Routing

## Gateway router

- ❑ Direct link to router in another AS
- ❑ 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

# Interconnected ASes



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

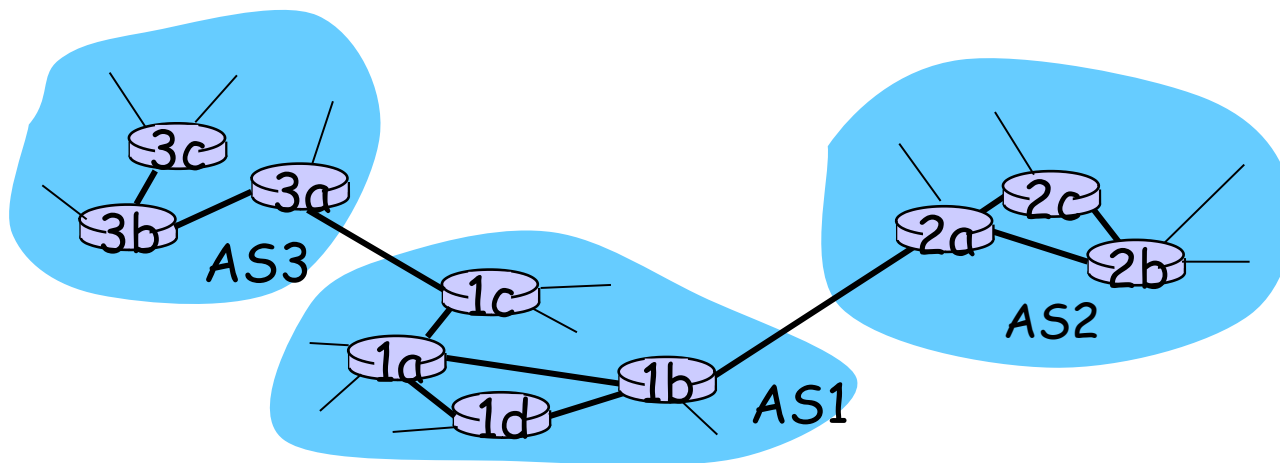
# Inter-AS tasks

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

## AS1 must:

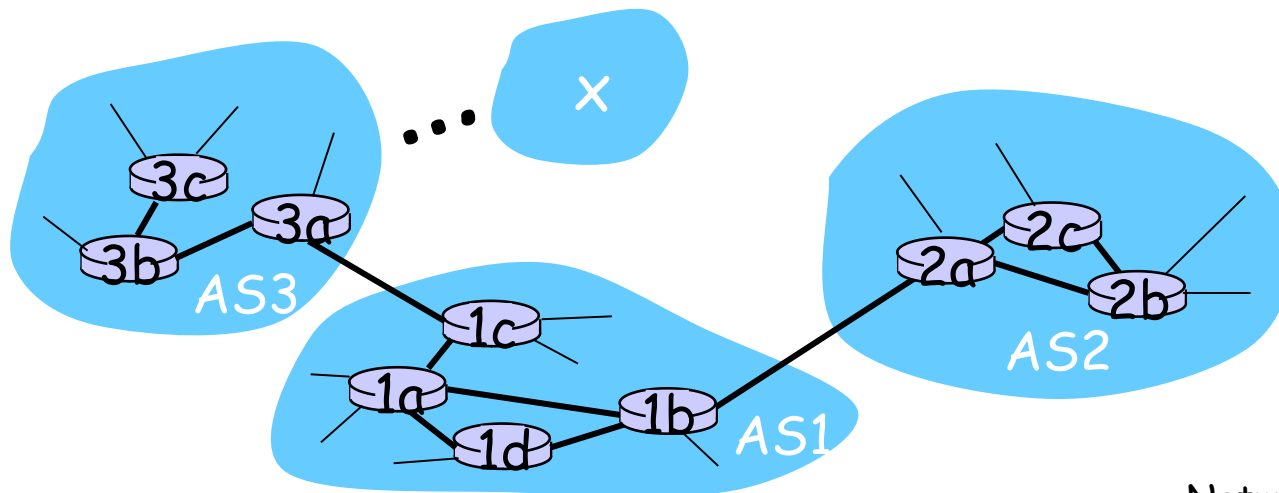
1. learn which destds are reachable through AS2, which through AS3
2. propagate this reachability info to all routers in AS1

**Job of inter-AS routing!**



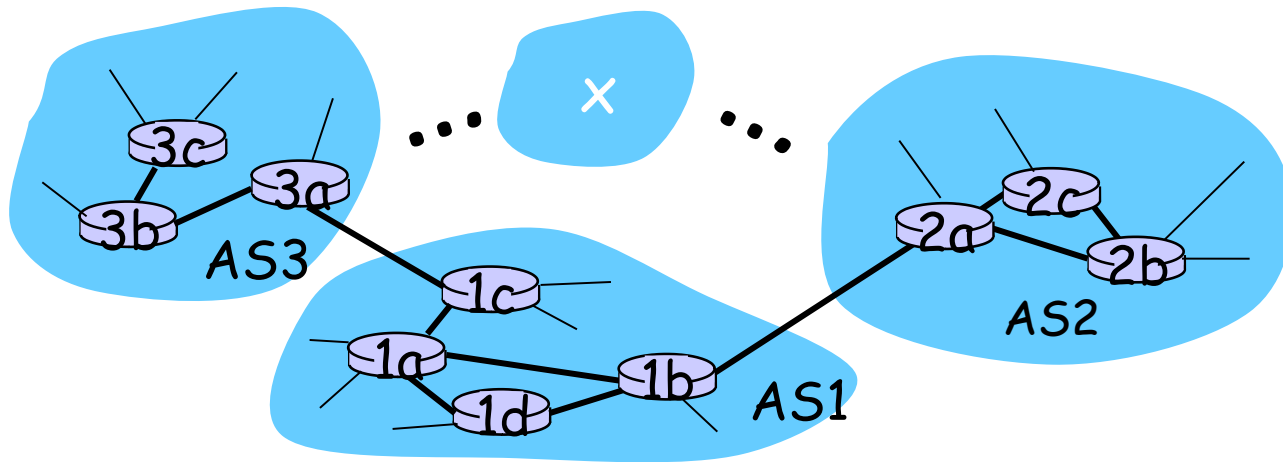
## Example: Setting forwarding table in router 1d

- suppose AS1 learns (via inter-AS protocol) that subnet  $x$  reachable via AS3 (gateway 1c) but not via AS2.
- inter-AS protocol propagates reachability info to all internal routers.
- router 1d determines from intra-AS routing info that its interface  $I$  is on the least cost path to 1c.
  - installs forwarding table entry  $(x, I)$



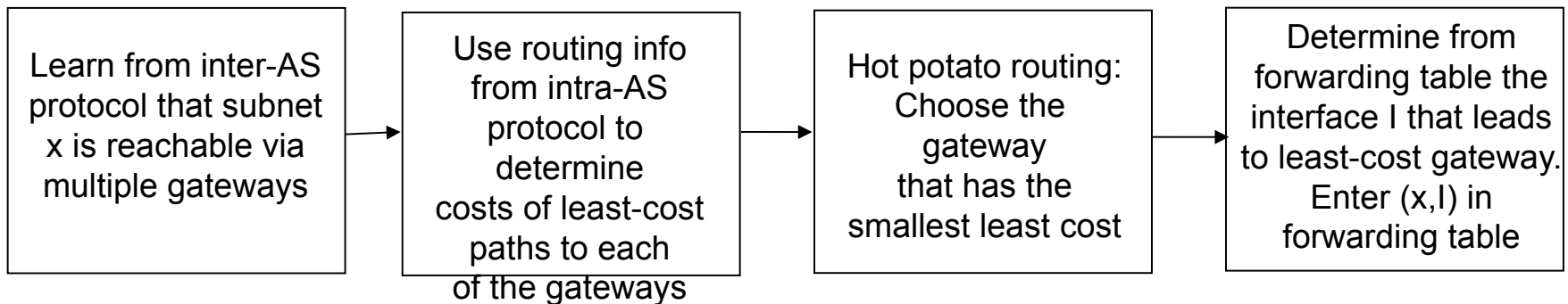
# Example: Choosing among multiple ASes

- now suppose AS1 learns from inter-AS protocol that subnet **x** is reachable from AS3 *and* from AS2.
- to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest **x**.
  - this is also job of inter-AS routing protocol!



# Example: Choosing among multiple ASes

- ❑ now suppose AS1 learns from inter-AS protocol that subnet **x** is reachable from AS3 *and* from AS2.
- ❑ to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest **x**.
  - this is also job of inter-AS routing protocol!
- ❑ **hot potato routing**: send packet towards closest of two routers.





# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link-state
- 4.3 Distance Vector
- 4.3 What's inside a router
- 4.4 Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 4.7 Broadcast and multicast routing

# Intra-AS Routing

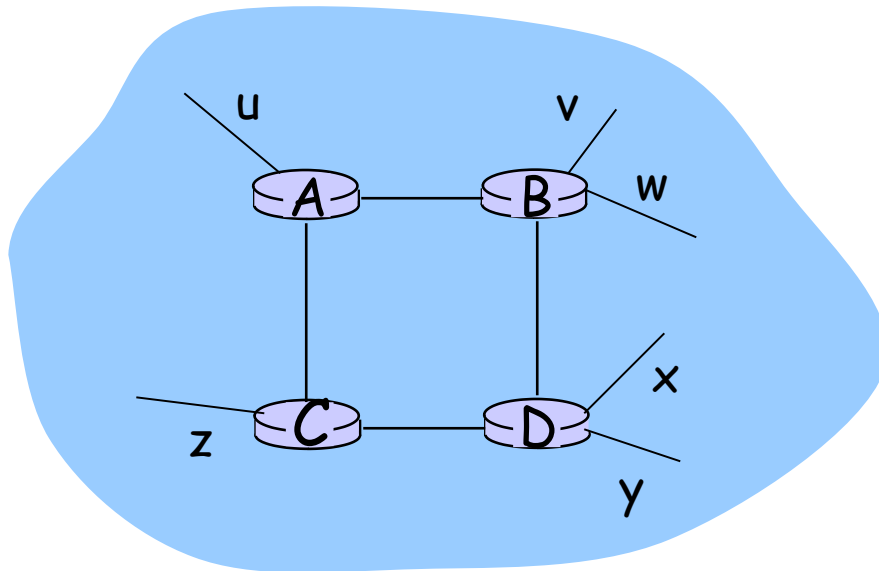
- ❑ also known as **Interior Gateway Protocols (IGP)**
- ❑ most common Intra-AS routing protocols:
  - RIP: Routing Information Protocol
  - OSPF: Open Shortest Path First
  - IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link-state
- 4.3 Distance Vector
- 4.3.1 What's inside a router
- 4.3.2 Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 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)



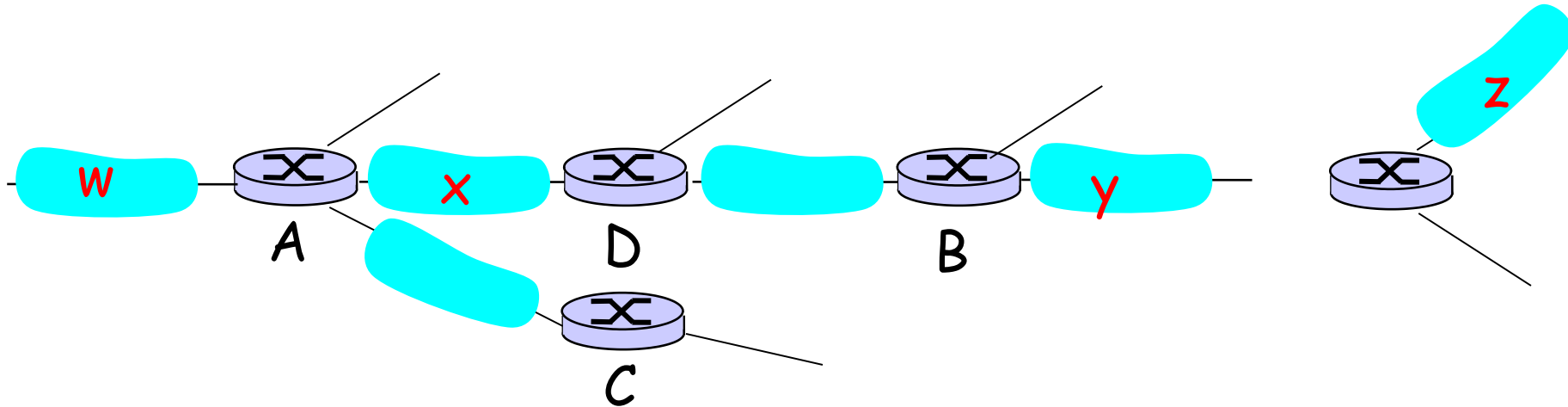
From router A to subnets:

<u>destination</u>	<u>hops</u>
u	1
v	2
w	2
x	3
y	3
z	2

# RIP advertisements

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

# RIP: Example



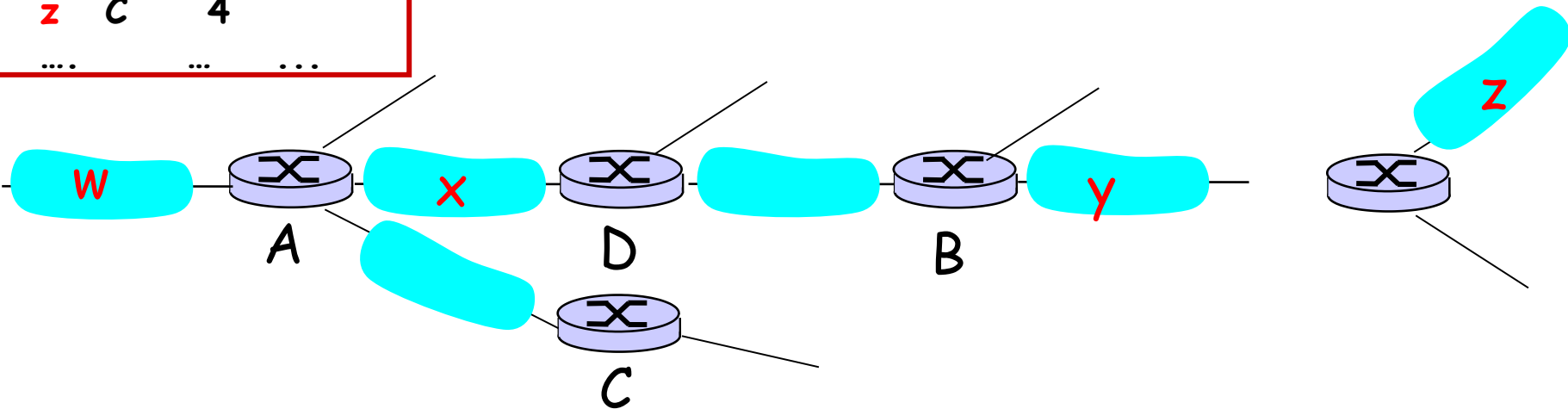
Destination Network	Next Router	Num. of hops to dest.
w	A	2
y	B	2
z	B	7
x	--	1
...	...	...

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	<del>7</del> 5	×
x	--	1	
...	...	...	

Routing/Forwarding table in D

Network Layer 4-111

# RIP: Link Failure and Recovery

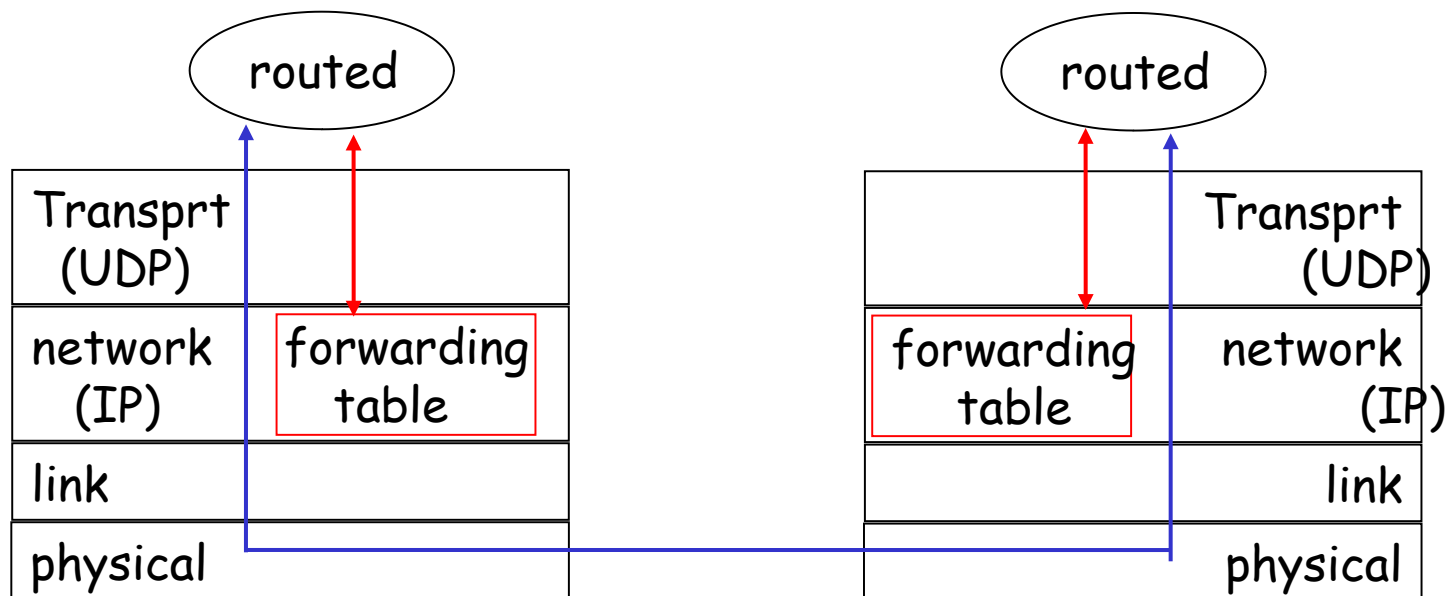
If no advertisement heard after 180 sec -->  
neighbor/link declared dead

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



# RIP Table processing

- ❑ RIP routing tables managed by **application-level** process called route-d (daemon)
- ❑ advertisements sent in UDP packets, periodically repeated



# Chapter 4: Network Layer

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

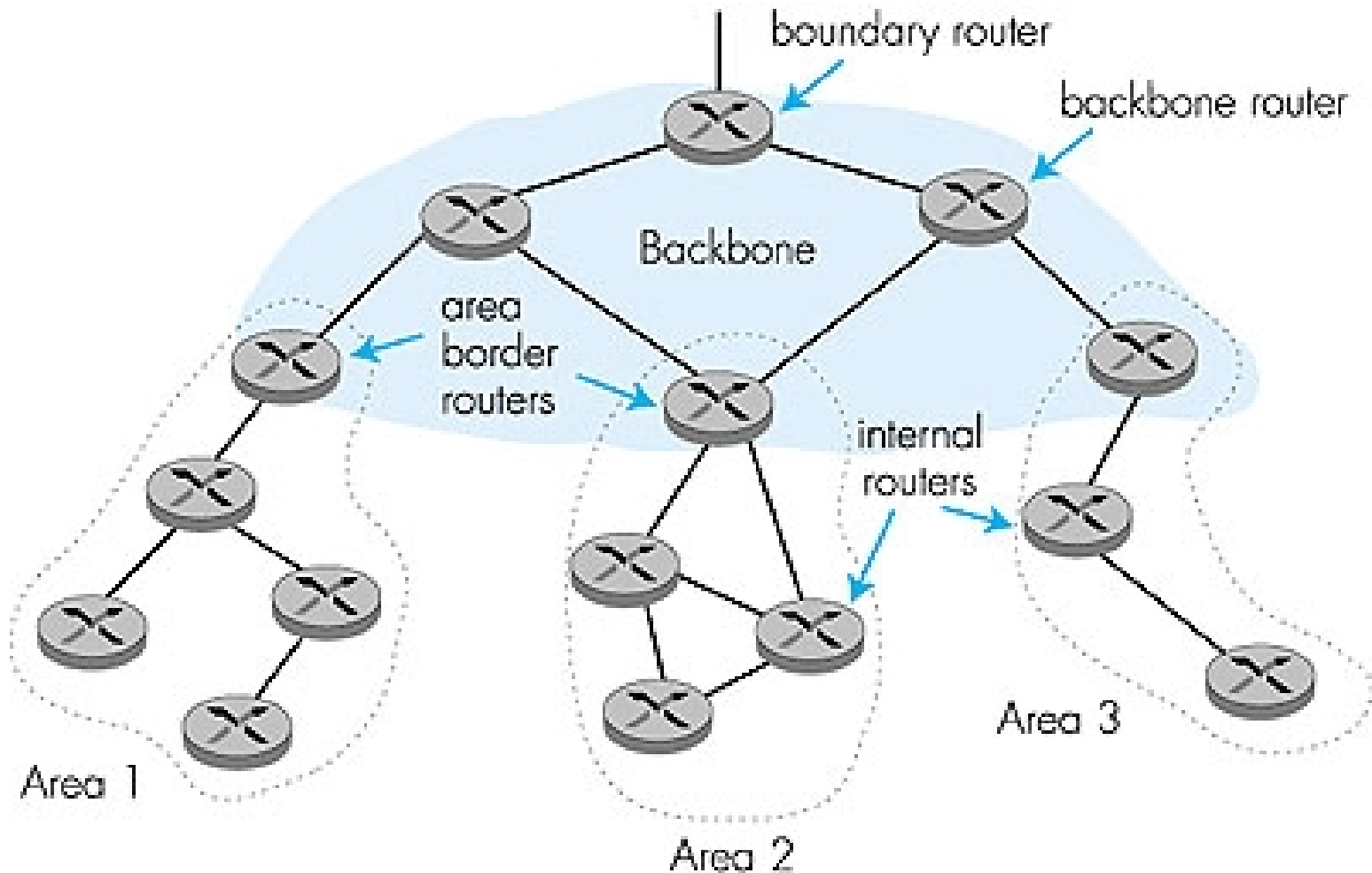
# OSPF (Open Shortest Path First)

- ❑ “open”: publicly available
- ❑ uses Link State algorithm
  - LS packet dissemination
  - topology map at each node
  - route computation using Dijkstra's algorithm
- ❑ OSPF advertisement carries one entry per neighbor router
- ❑ advertisements disseminated to **entire** AS (via flooding)
  - carried in OSPF messages directly over IP (rather than TCP or UDP)

# OSPF "advanced" features (not in RIP)

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

# Hierarchical OSPF



# Hierarchical OSPF

- ❑ **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.
- ❑ **area border routers:** “summarize” distances to nets in own area, advertise to other Area Border routers.
- ❑ **backbone routers:** run OSPF routing limited to backbone.
- ❑ **boundary routers:** connect to other AS's.

# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link-state
- 4.3 Distance Vector
- 4.3 What's inside a router
- 4.4 Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 4.7 Broadcast and multicast routing

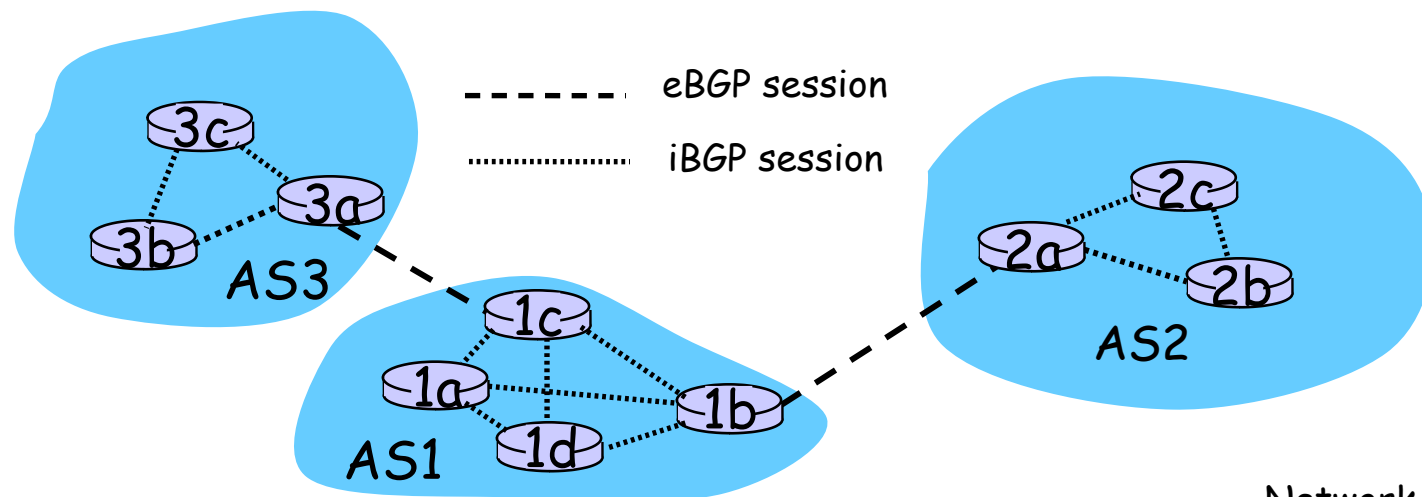
# Internet inter-AS routing: BGP

- ❑ **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 reachability information to all AS-internal routers.
  3. Determine “good” routes to subnets based on reachability information and policy.
- ❑ allows subnet to advertise its existence to rest of Internet: *“I am here”*



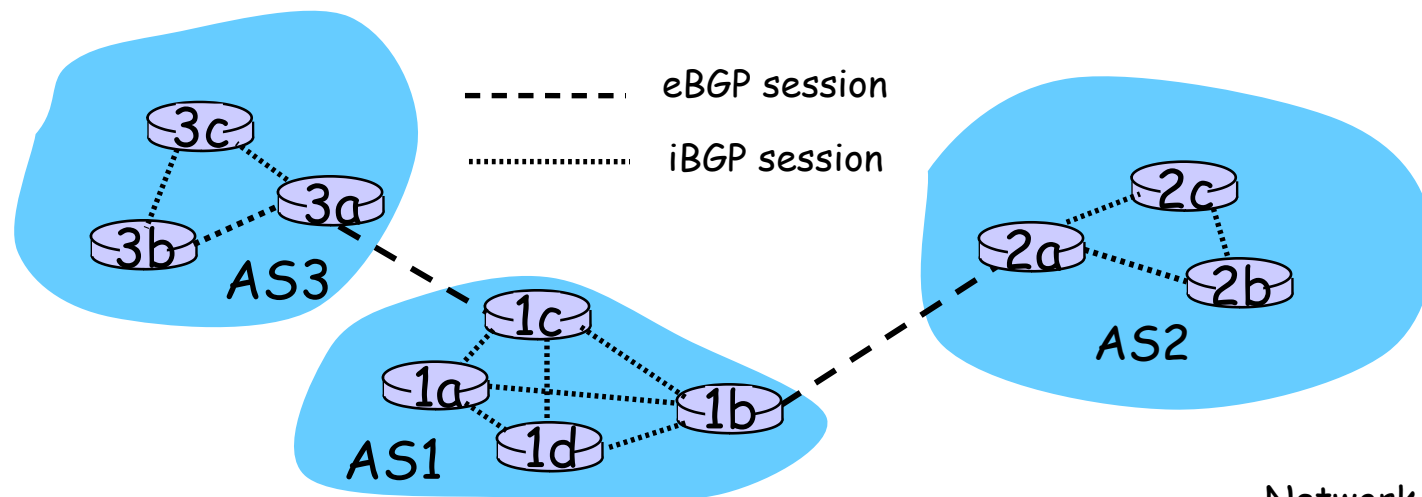
# BGP basics

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



# Distributing reachability info

- using eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
  - 1c can then use iBGP to distribute new prefix info to all routers in AS1
  - 1b can then re-advertise new reachability info to AS2 over 1b-to-2a eBGP session
- when router learns of new prefix, it creates entry for prefix in its forwarding table.



# Path attributes & BGP routes

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

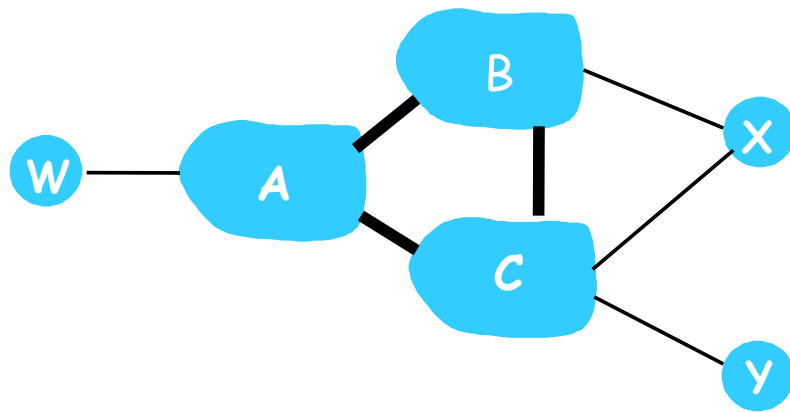
# BGP route selection



- ❑ router may learn about more than 1 route to some prefix. Router must select route.
- ❑ elimination rules:
  1. local preference value attribute: policy decision
  2. shortest AS-PATH
  3. closest NEXT-HOP router: hot potato routing
  4. additional criteria

# BGP messages

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

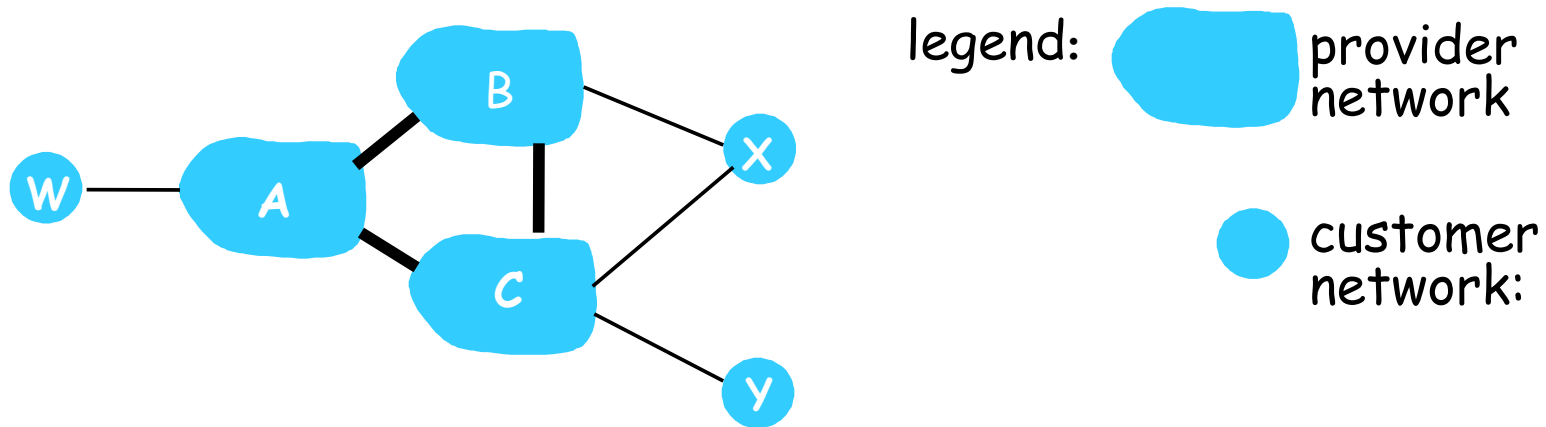
# BGP routing policy



legend:  provider network  
 customer network:

- A,B,C are **provider networks**
- X,W,Y are customer (of provider networks)
- X is **dual-homed**: attached to two networks
  - X does not want to route from B via X to C
  - .. so X will not advertise to B a route to C

## BGP routing policy (2)



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

# Why different Intra- and Inter-AS routing ?

## Policy:

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

## Scale:

- ❑ hierarchical routing saves table size, reduced update traffic

## Performance:

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

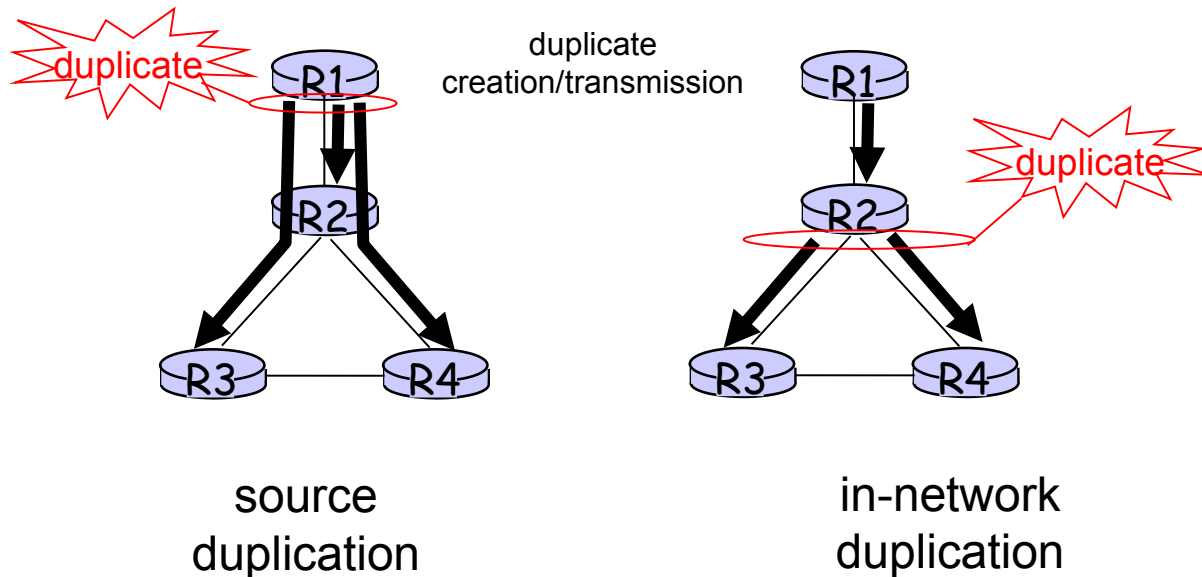


# Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Link-state
- 4.3 Distance Vector
- 4.3 What's inside a router
- 4.4 Hierarchical routing
- 4.4 IP: Internet Protocol
- 4.6 Routing in the Internet
  - Datagram format
  - RIP
  - IPv4 addressing
  - OSPF
  - ICMP
  - BGP
  - IPv6
- 4.7 Broadcast and multicast routing

# Broadcast Routing

- ❑ deliver packets from source to all other nodes
- ❑ source duplication is inefficient:



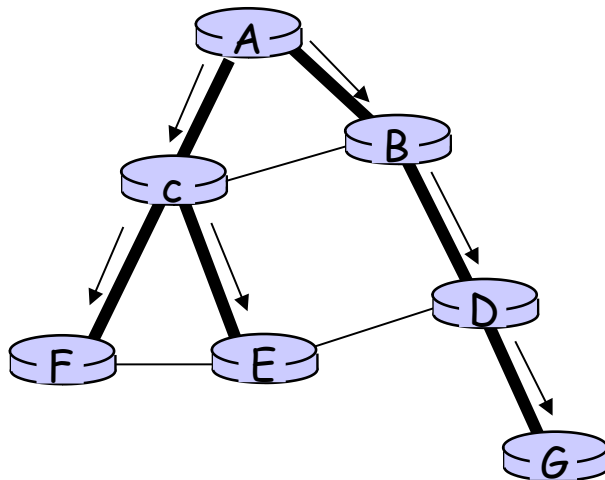
- ❑ source duplication: how does source determine recipient addresses?

# In-network duplication

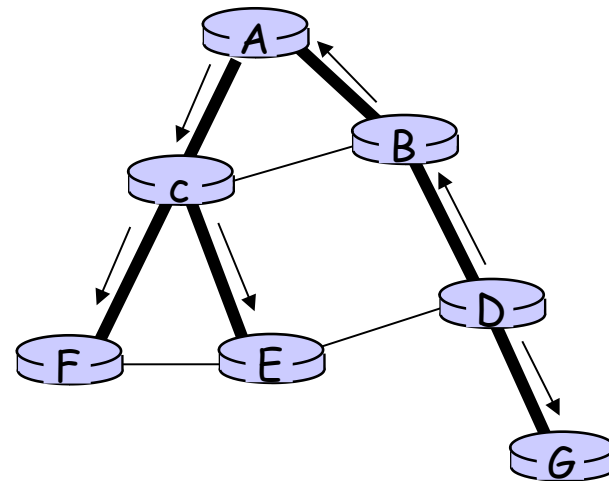
- ❑ flooding: when node receives brdcst pkt, sends copy to all neighbors
  - Problems: cycles & broadcast storm
- ❑ controlled flooding: node only brdcsts pkt if it hasn't brdcst same packet before
  - Node keeps track of pkt ids already brdcsted
  - Or reverse path forwarding (RPF): only forward pkt if it arrived on shortest path between node and source
- ❑ spanning tree
  - No redundant packets received by any node

# Spanning Tree

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



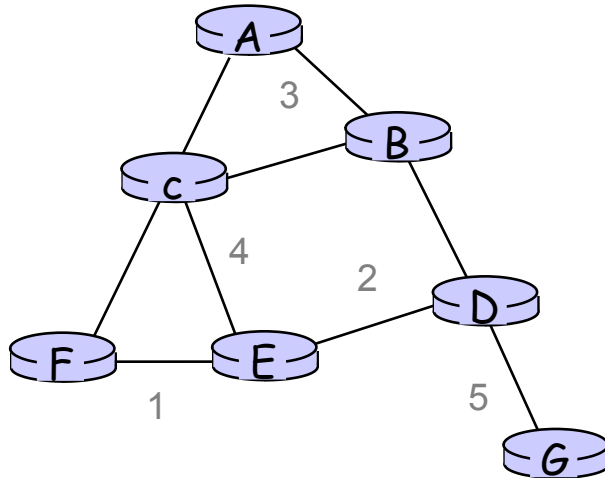
(a) Broadcast initiated at A



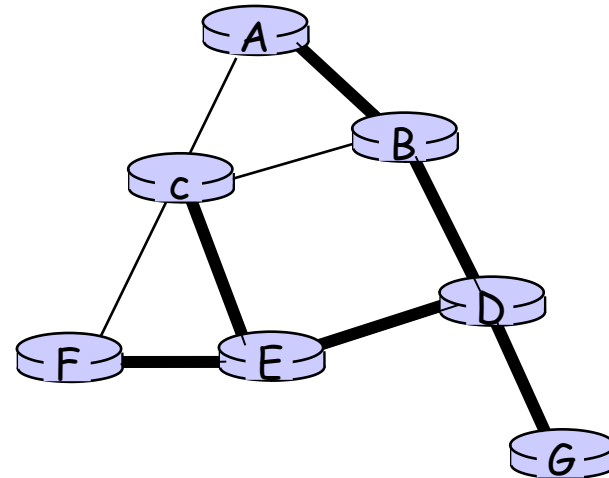
(b) Broadcast initiated at D

# Spanning Tree: Creation

- Center node
- Each node sends unicast join message to center node
  - Message forwarded until it arrives at a node already belonging to spanning tree



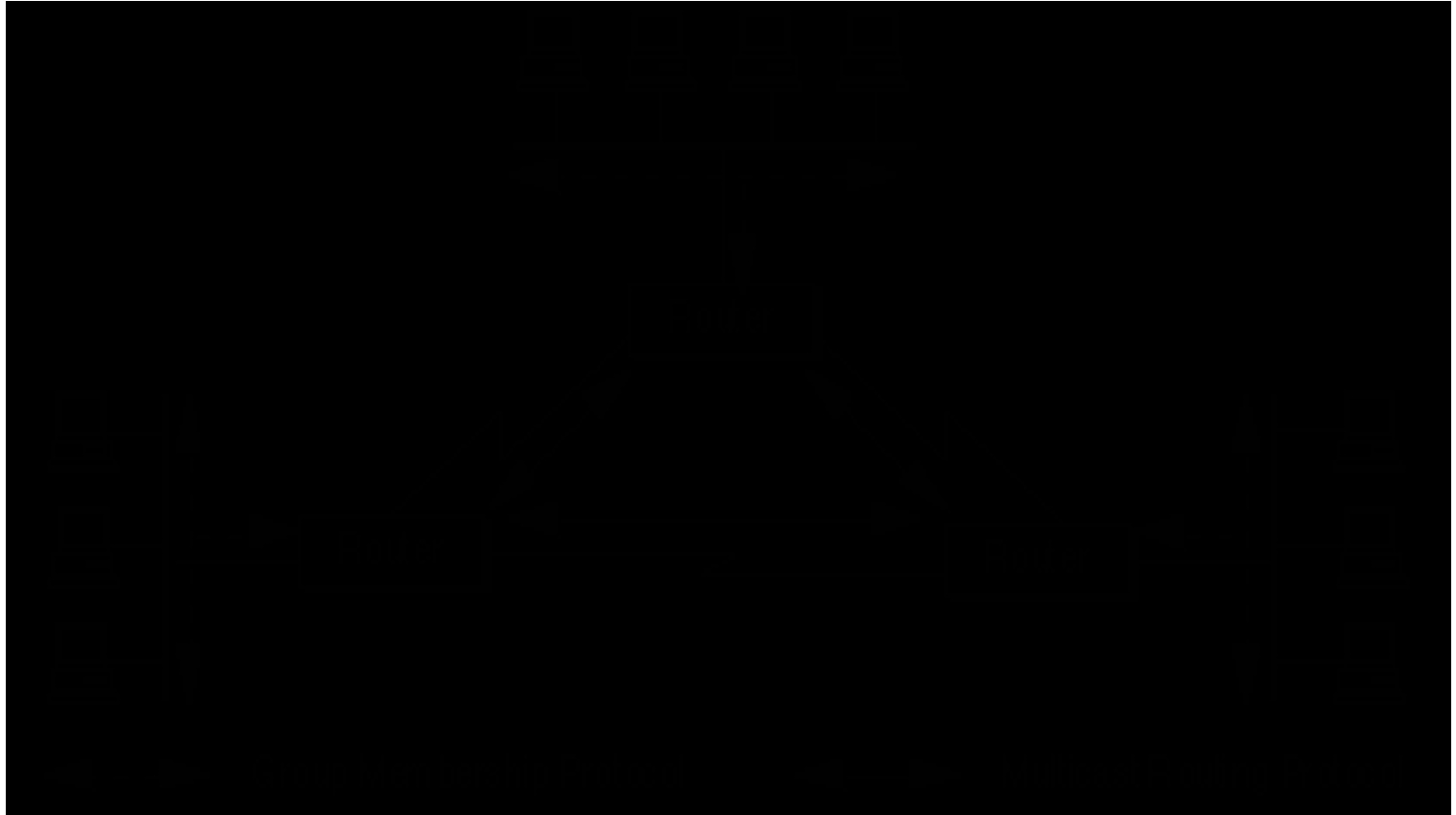
(a) Stepwise construction of spanning tree



(b) Constructed spanning tree

# Multicasting

# Multicast Forwarding



# Time To Live (TTL)

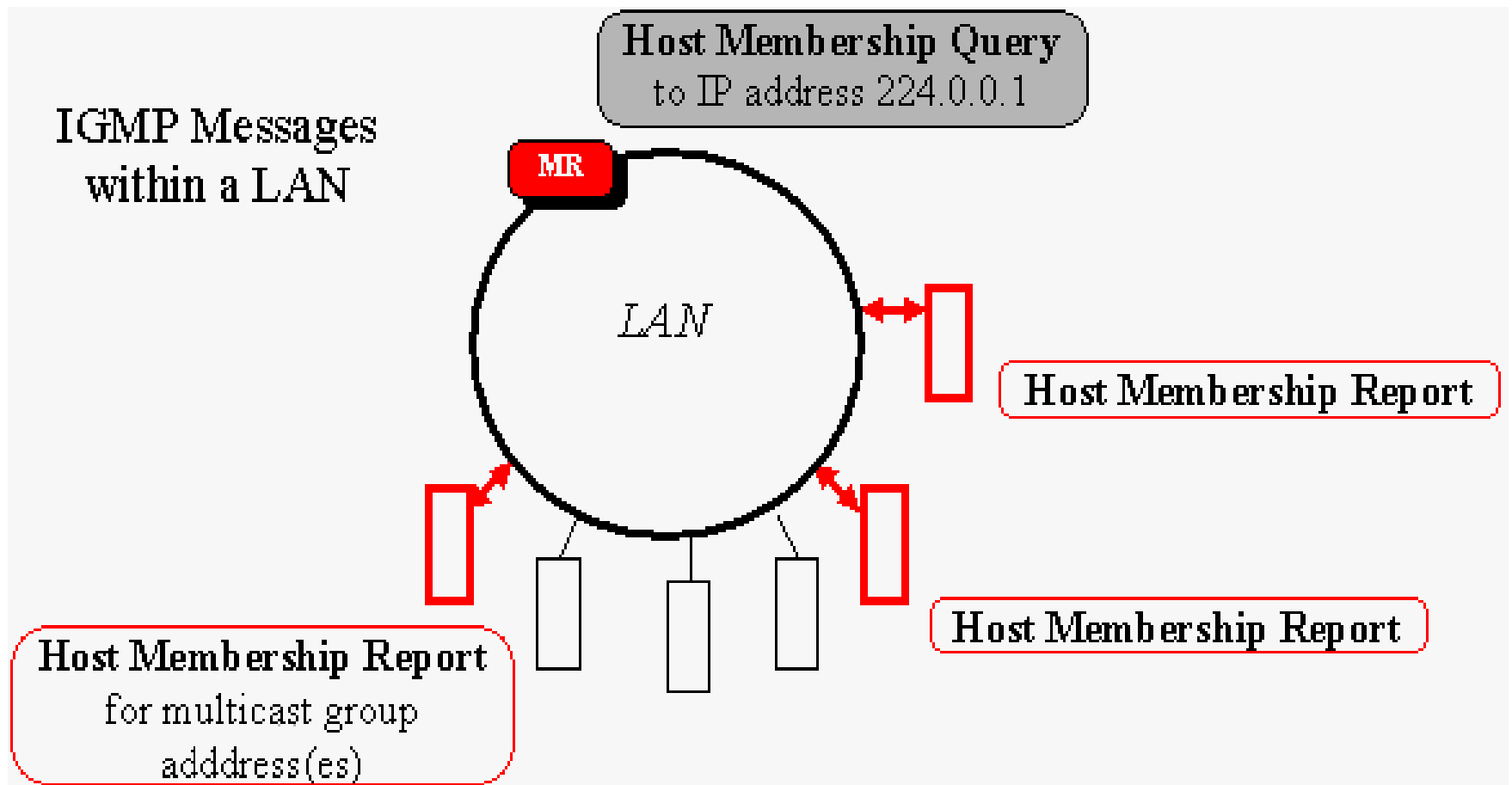
- ❑ Scope-limiting parameter for IP Multicast datagrams
- ❑ Controls the number of hops that a IP Multicast packet is allowed to propagate
- ❑  $TTL = 1$ : local network multicast
- ❑  $TTL > 1$ : Multicast router(s) attached to the local network forward IP Multicast datagrams



# Internet Group Management Protocol (IGMP)

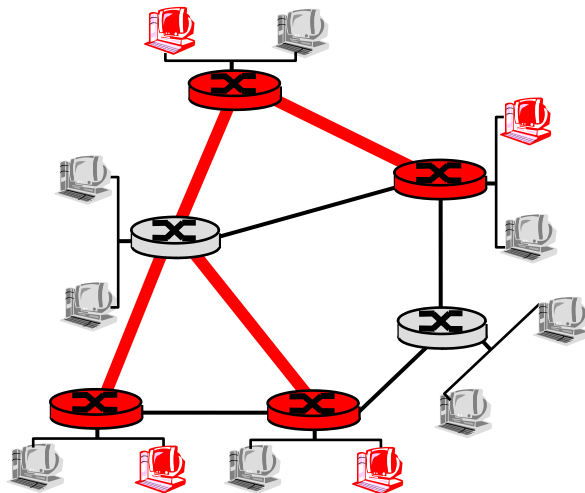
- ❑ Used by mrouters to learn about Multicast Group Memberships on their directly attached subnets
  - the existence of at least one member/group
- ❑ Implemented over IP
- ❑ Designated Router
  - Each network has one Querier
  - All router begin as Queriers
  - Mrouter with the lowest IP address chosen

# IGMP Messages

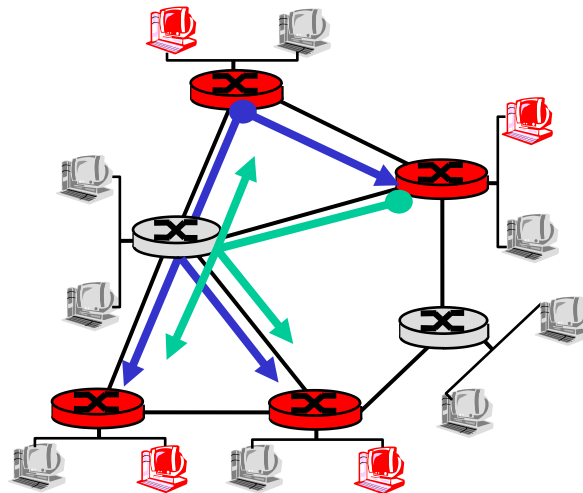


# Multicast Routing: Problem Statement

- **Goal:** find a tree (or trees) connecting routers having local mcast group members
  - **tree:** not all paths between routers used
  - **source-based:** different tree from each sender to rcvrs
  - **shared-tree:** same tree used by all group members



Shared tree



Source-based trees

# Approaches for building mcast trees

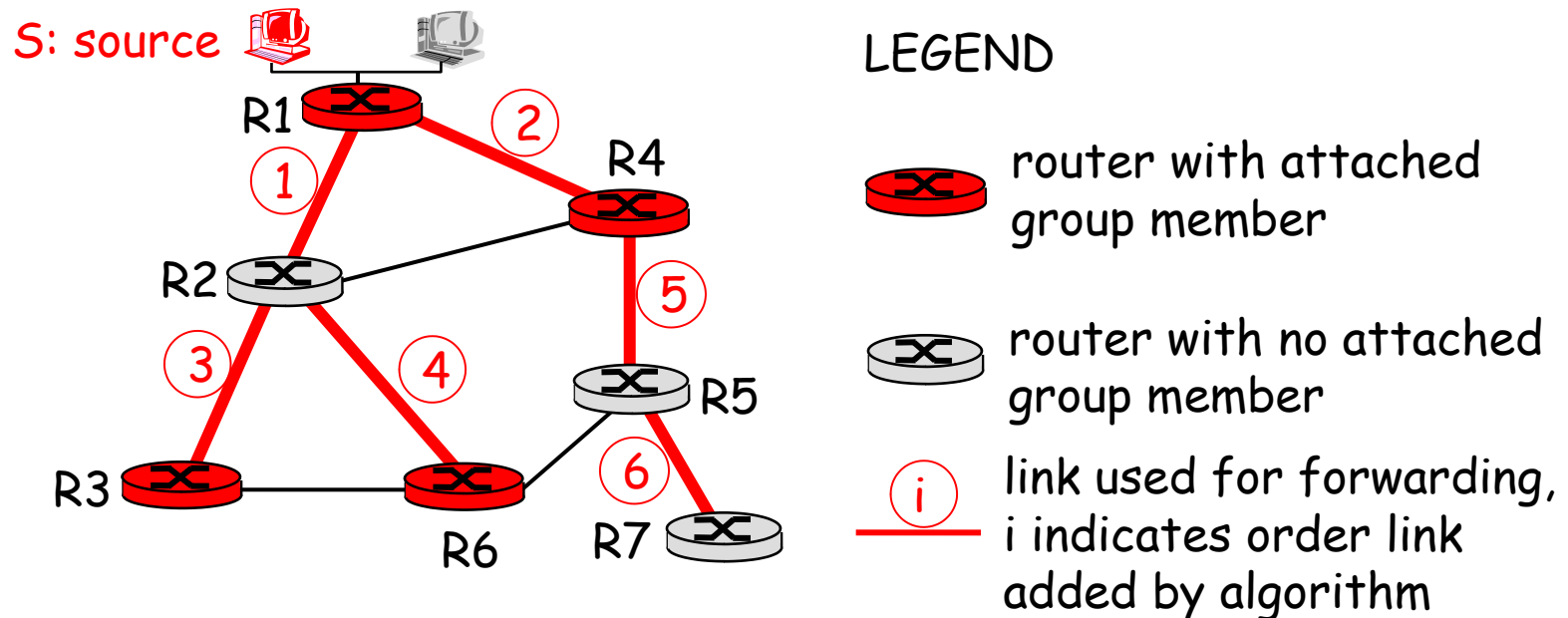
## Approaches:

- **source-based tree:** one tree per source
  - shortest path trees
  - reverse path forwarding
- **group-shared tree:** group uses one tree
  - minimal spanning (Steiner)
  - center-based trees

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

# Shortest Path Tree

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

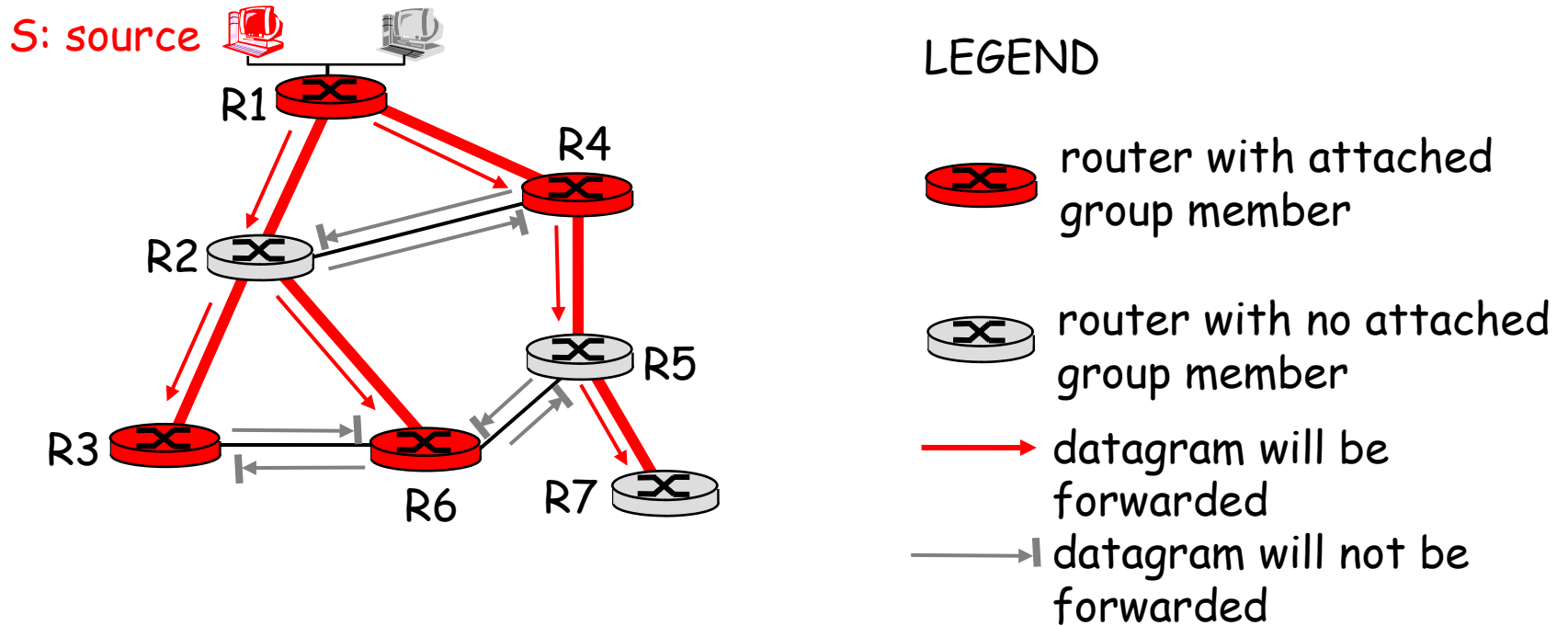


# Reverse Path Forwarding

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

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

# Reverse Path Forwarding: example

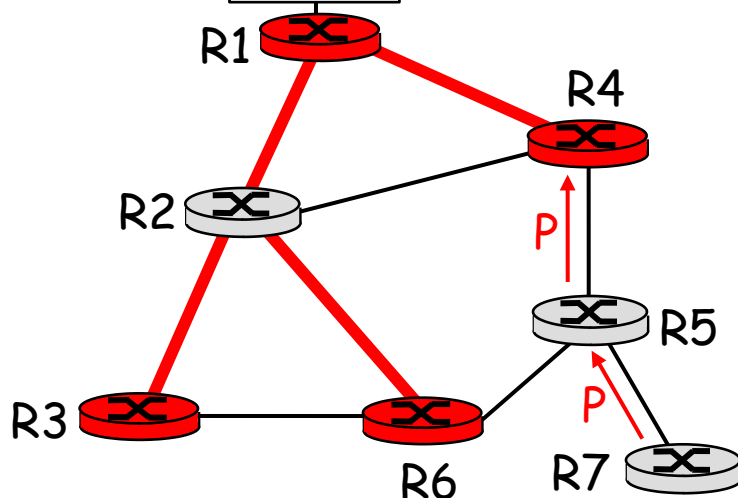


- result is a source-specific *reverse* SPT
  - may be a bad choice with asymmetric links

# Reverse Path Forwarding: pruning

- forwarding tree contains subtrees with no mcast group members
  - no need to forward datagrams down subtree
  - “prune” msgs sent upstream by router with no downstream group members

S: source



## LEGEND



router with attached group member



router with no attached group member



prune message



links with multicast forwarding



# Shared-Tree: Steiner Tree

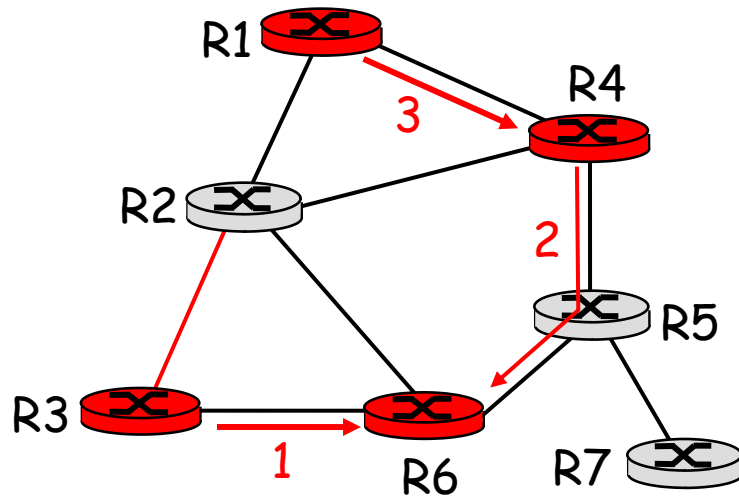
- ❑ **Steiner Tree:** minimum cost tree connecting all routers with attached group members
- ❑ problem is NP-complete
- ❑ excellent heuristics exists
- ❑ not used in practice:
  - computational complexity
  - information about entire network needed
  - monolithic: rerun whenever a router needs to join/leave

# Center-based trees




- ❑ single delivery tree shared by all
- ❑ one router identified as "*center*" of tree
- ❑ to join:
  - edge router sends unicast *join-msg* addressed to center router
  - *join-msg* "processed" by intermediate routers and forwarded towards center
  - *join-msg* either hits existing tree branch for this center, or arrives at center
  - path taken by *join-msg* becomes new branch of tree for this router

# Center-based trees: an example

Suppose R6 chosen as center:



## LEGEND

-  router with attached group member
-  router with no attached group member
-  path order in which join messages generated

# Internet Multicasting Routing: DVMRP

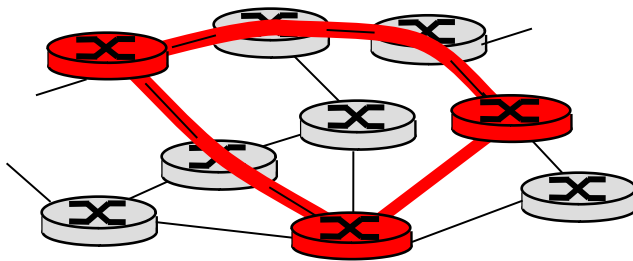
- **DVMRP**: distance vector multicast routing protocol, RFC1075
- *flood and prune*: reverse path forwarding, source-based tree
  - RPF tree based on DVMRP's own routing tables constructed by communicating DVMRP routers
  - no assumptions about underlying unicast
  - initial datagram to mcast group flooded everywhere via RPF
  - routers not wanting group: send upstream prune msgs

# DVMRP: continued...

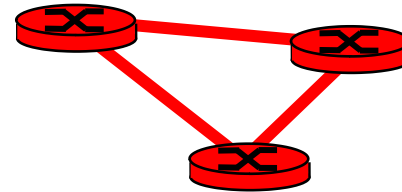
- ❑ *soft state*: DVMRP router periodically (1 min.) "forgets" branches are pruned:
  - mcast data again flows down unpruned branch
  - downstream router: reprune or else continue to receive data
- ❑ routers can quickly regraft to tree
  - following IGMP join at leaf
- ❑ odds and ends
  - commonly implemented in commercial routers
  - Mbone routing done using DVMRP

# Tunneling

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



physical topology



logical topology

- ❑ mcast datagram encapsulated inside "normal" (non-multicast-addressed) datagram
- ❑ normal IP datagram sent thru "tunnel" via regular IP unicast to receiving mcast router
- ❑ receiving mcast router unencapsulates to get mcast datagram

# PIM: Protocol Independent Multicast

- ❑ not dependent on any specific underlying unicast routing algorithm (works with all)
- ❑ two different multicast distribution scenarios :

## Dense:

- ❑ group members densely packed, in "close" proximity.
- ❑ bandwidth more plentiful

## Sparse:

- ❑ # networks with group members small wrt # interconnected networks
- ❑ group members "widely dispersed"
- ❑ bandwidth not plentiful

# Consequences of Sparse-Dense Dichotomy:

## Dense

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

## Sparse:

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



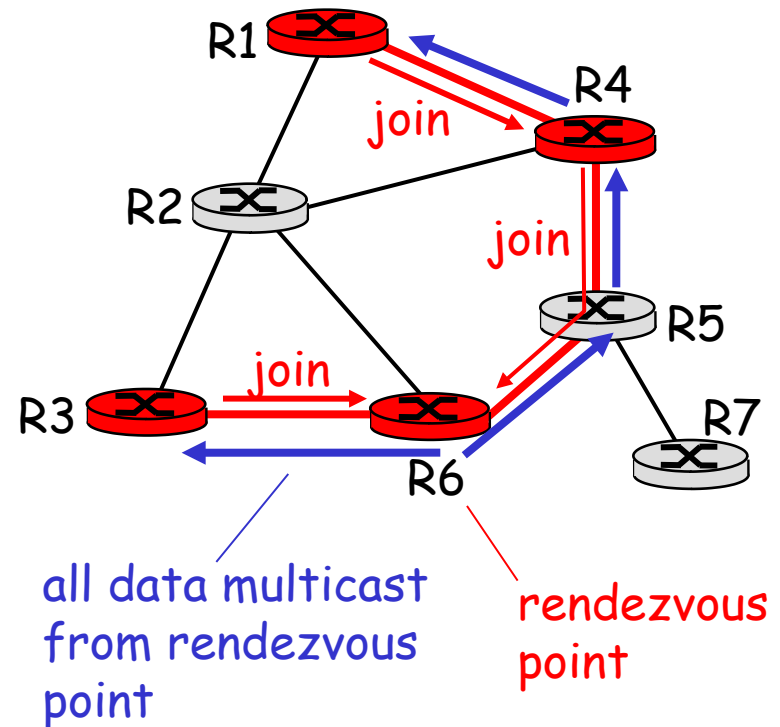
# PIM- Dense Mode

**flood-and-prune RPF**, similar to DVMRP but

- ❑ underlying unicast protocol provides RPF info for incoming datagram
- ❑ less complicated (less efficient) downstream flood than DVMRP reduces reliance on underlying routing algorithm
- ❑ has protocol mechanism for router to detect it is a leaf-node router

# PIM - Sparse Mode

- ❑ center-based approach
- ❑ router sends *join* msg to rendezvous point (RP)
  - intermediate routers update state and forward *join*
- ❑ after joining via RP, router can switch to source-specific tree
  - increased performance: less concentration, shorter paths



# PIM - Sparse Mode

sender(s):

- ❑ unicast data to RP, which distributes down RP-rooted tree
- ❑ RP can extend mcast tree upstream to source
- ❑ RP can send *stop* msg if no attached receivers
  - "no one is listening!"

