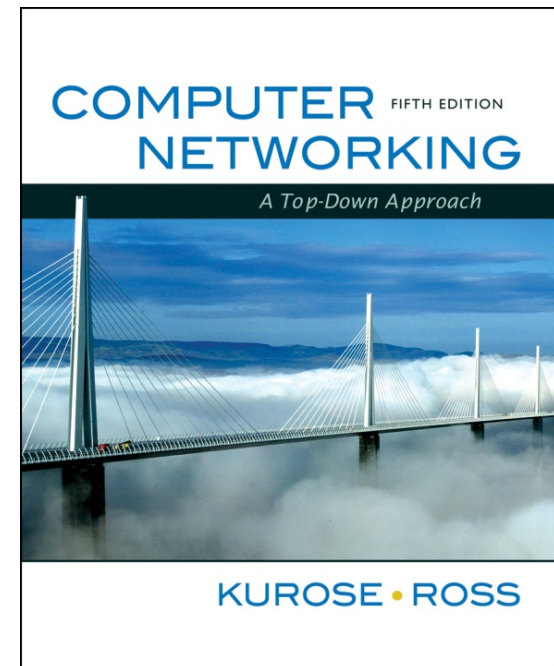


Chapter 4

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 1996-2012
J.F Kurose and K.W. Ross, All Rights Reserved

*Computer Networking:
A Top Down Approach ,
6th edition.
Jim Kurose, Keith Ross
Addison-Wesley, March
2012.*

Chapter 4: Network Layer

Chapter goals:

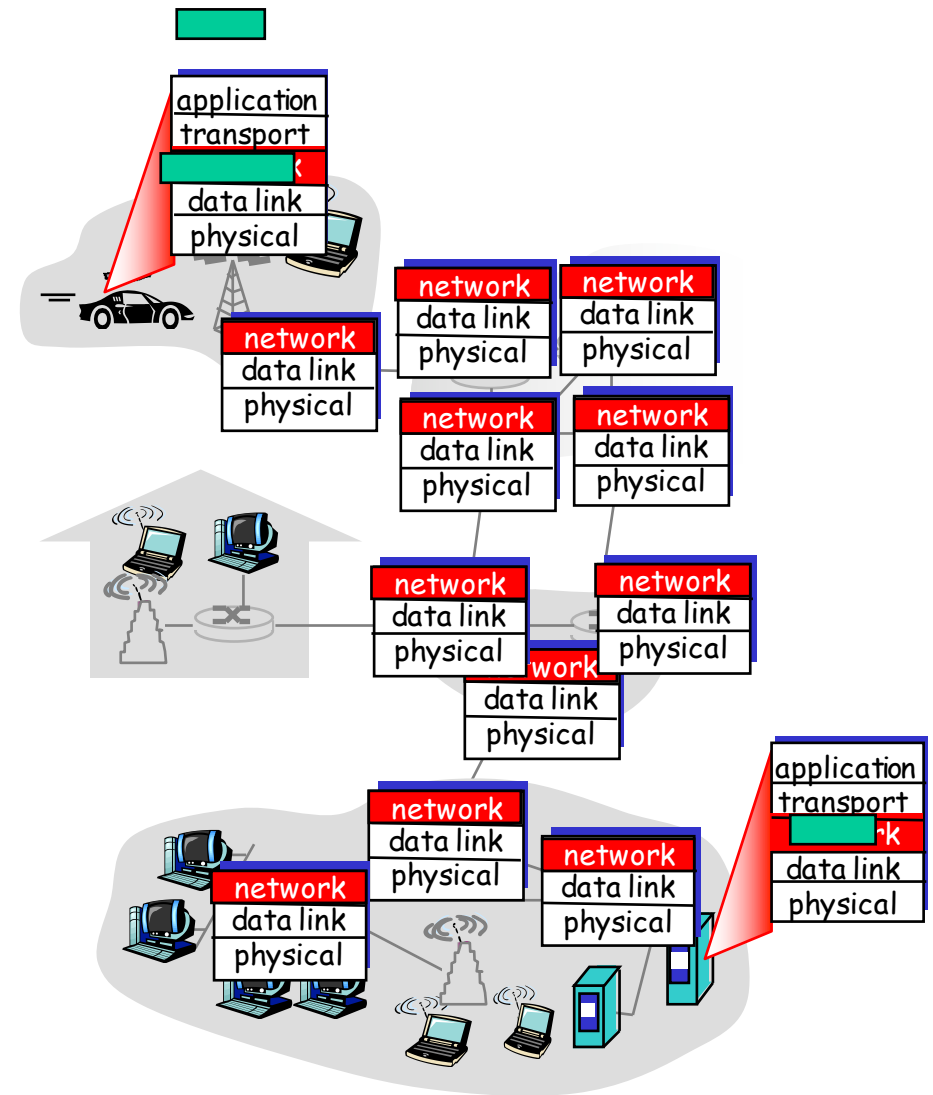
- ❑ understand principles behind network layer services:
 - network layer service models
 - forwarding versus routing
 - how a router works
 - routing (path selection)
 - dealing with scale
 - advanced topics: IPv6, mobility
- ❑ instantiation, implementation in the Internet

Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 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

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



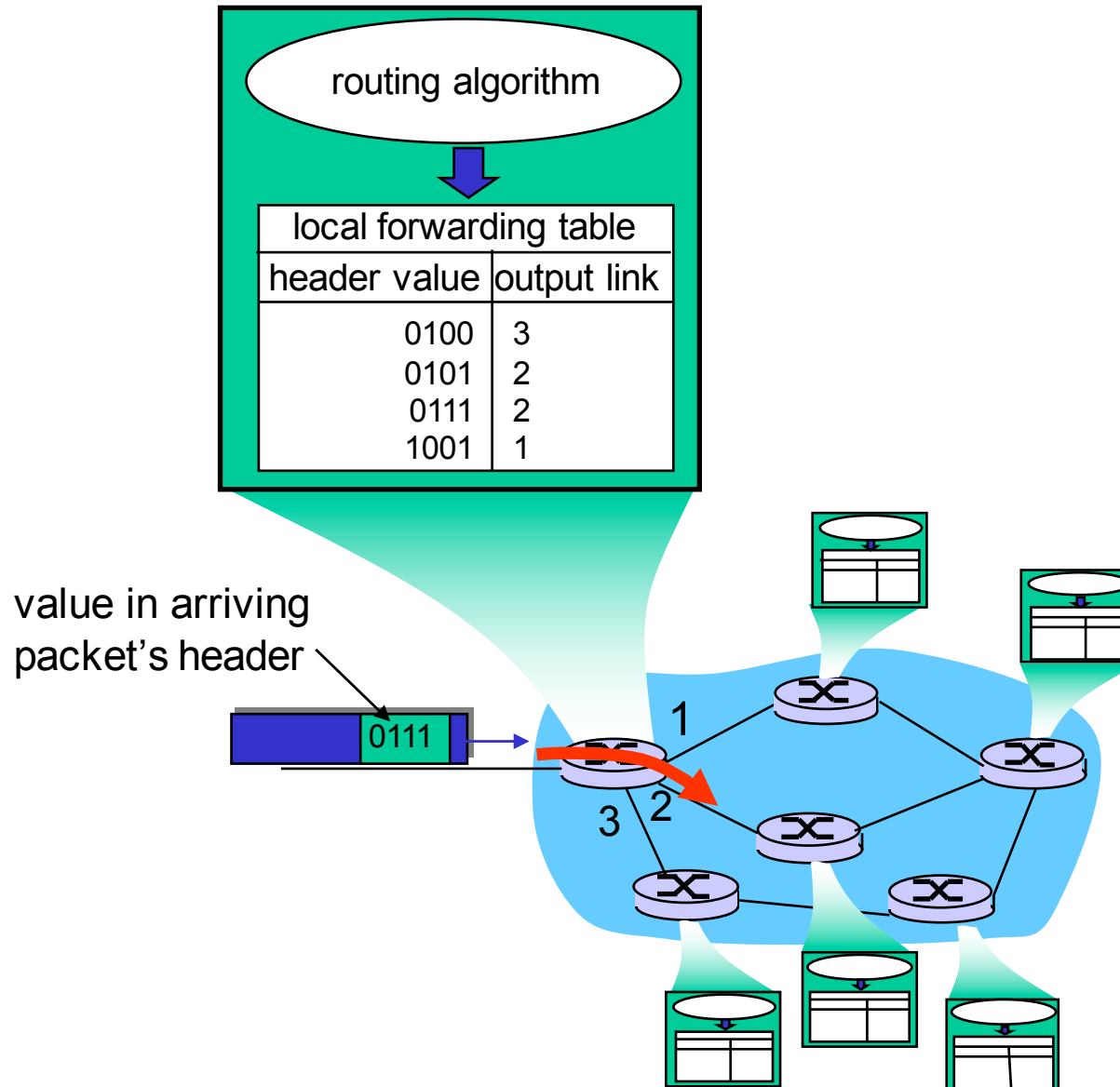
Two Key Network-Layer Functions

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

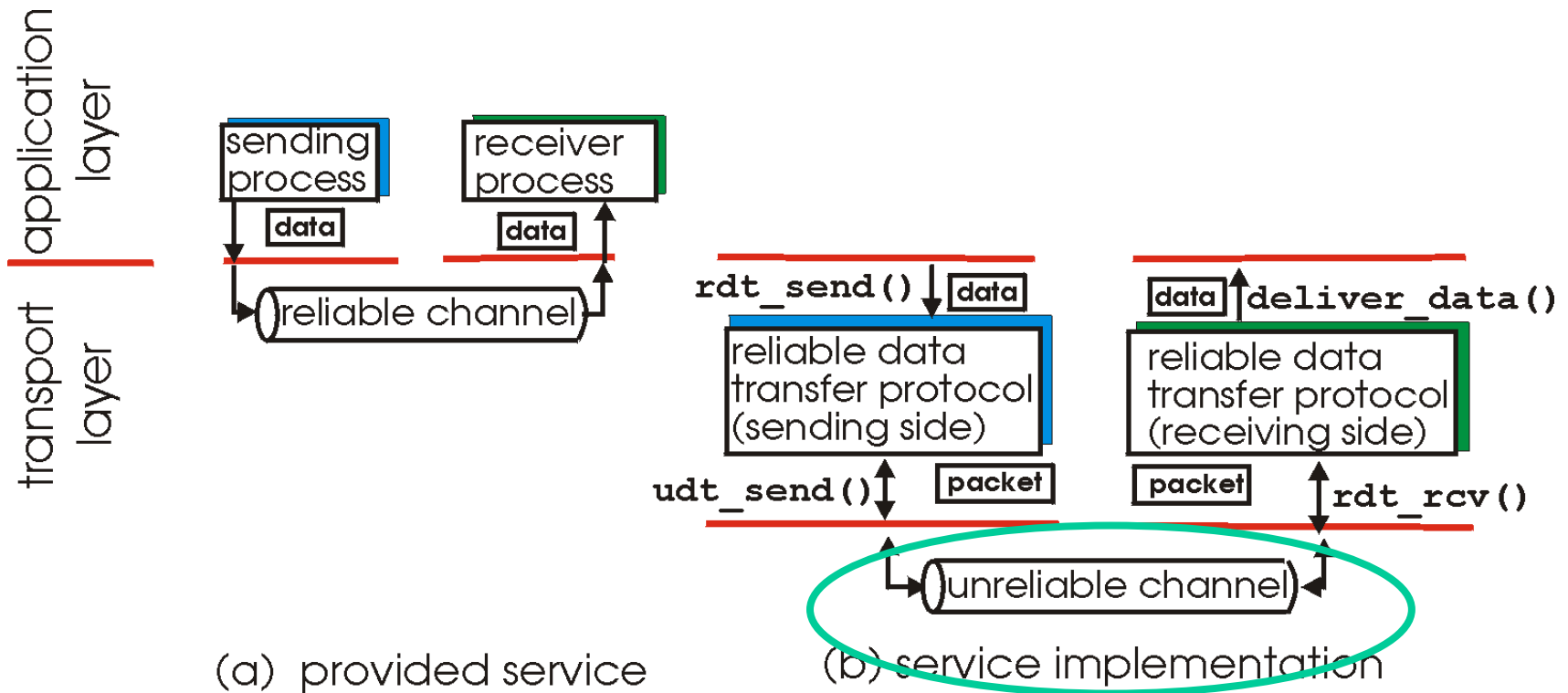
analogy:

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

Interplay between routing and forwarding



Remember This?



(a) provided service

(b) service implementation

Connectionless service

- ❑ best effort:
 - TCP/IP (aka. the Internet)
- ❑ two end hosts *and* intervening routers do not set up any 'connection specific' states before forwarding datagrams

Connection oriented service

- ❑ 3rd 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 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 (not process-to-process)
 - **no choice:** network provides one or the other (not both)
 - **implementation:** in network core (not at the end systems)

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)

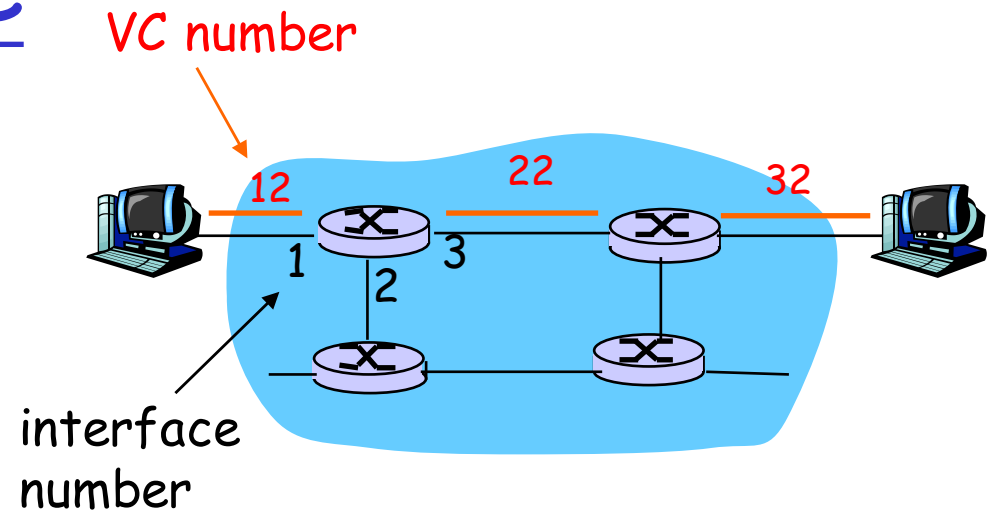
VC implementation

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:

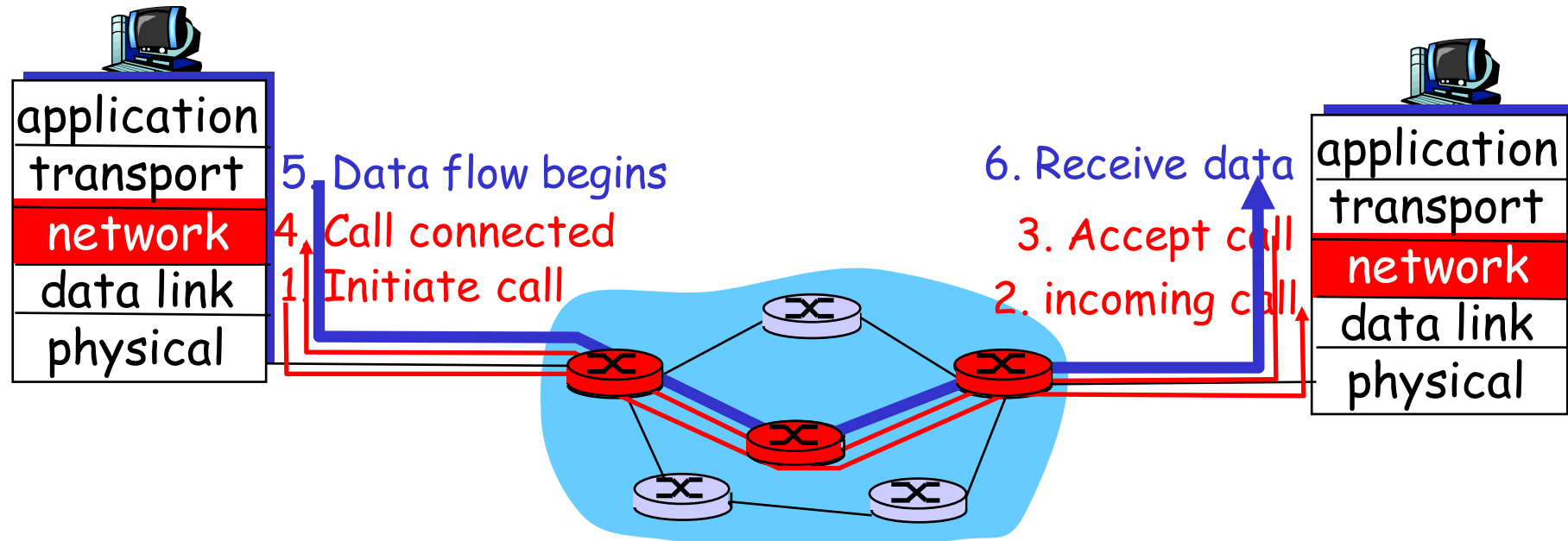


Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	3	22
2	63	1	18
3	7	2	17
1	97	3	87
...

Routers maintain connection state information!

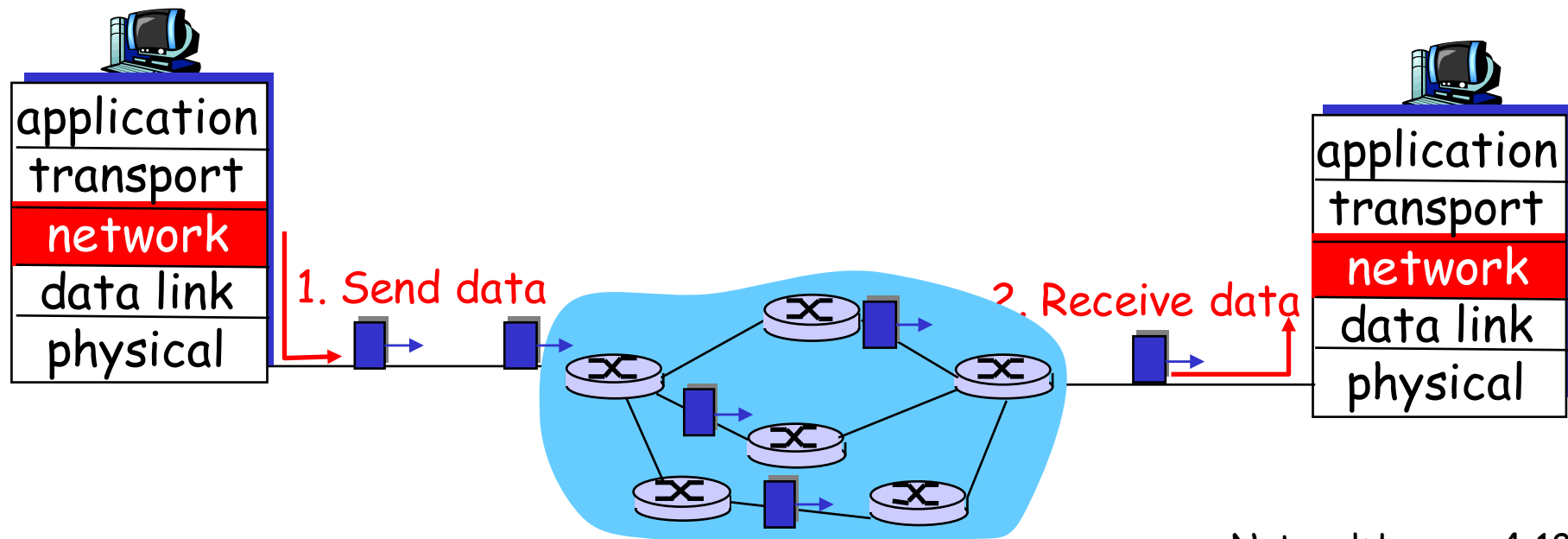
Virtual circuits: signaling protocols

- ❑ 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>
11001000000101110001000000000000 through 11001000000101110001011111111111	0
11001000000101110001100000000000 through 11001000000101110001100011111111	1
11001000000101110001100100000000 through 11001000000101110001111111111111	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?

Double Quiz Time!

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

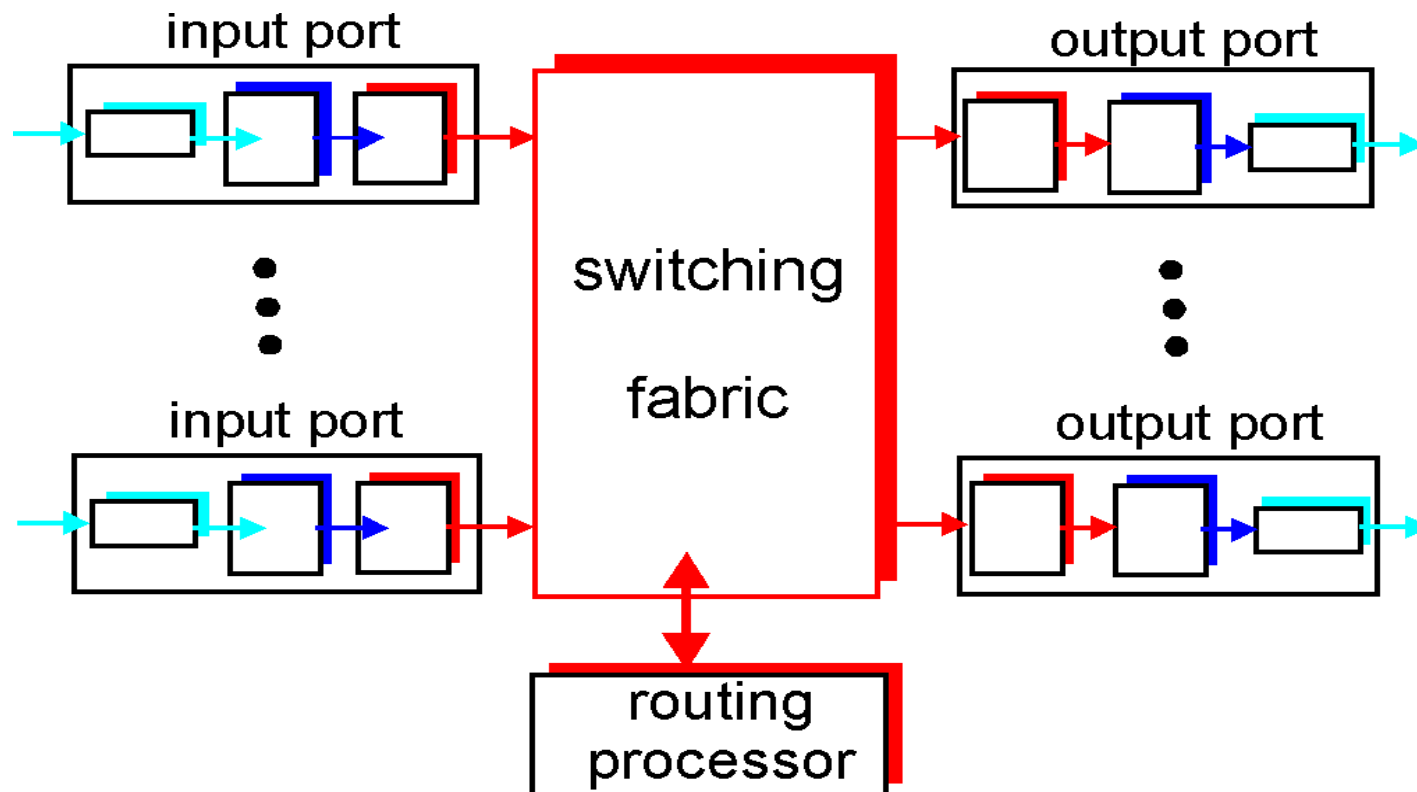
Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 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

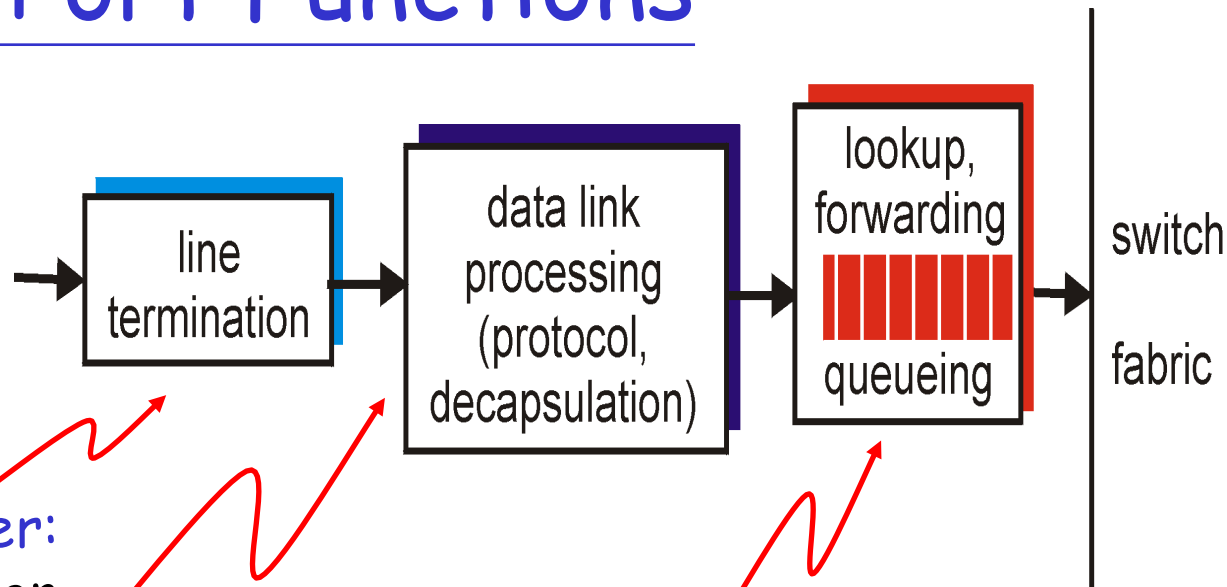
Router Architecture Overview

Two key router functions:

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



Input Port Functions



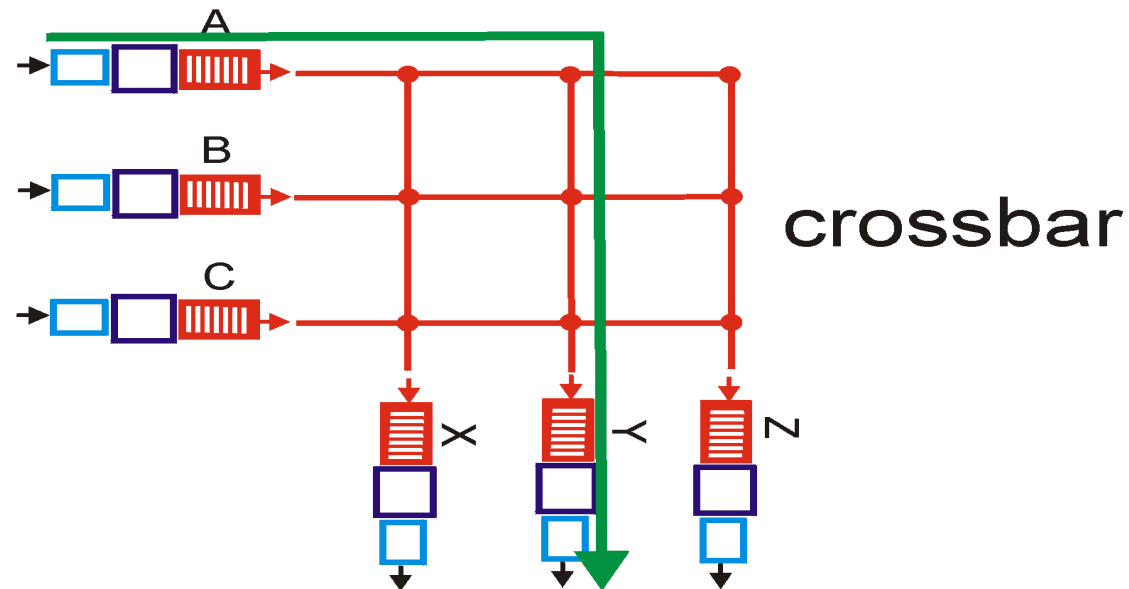
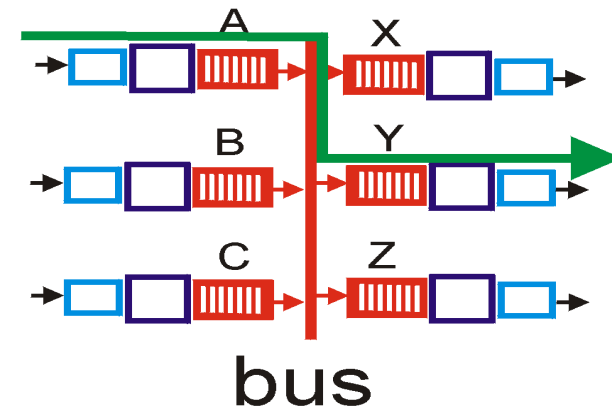
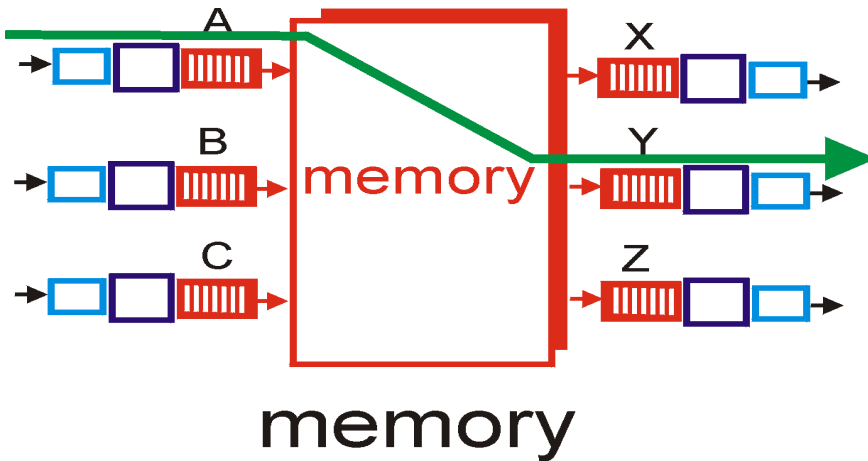
Physical layer:
bit-level reception

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

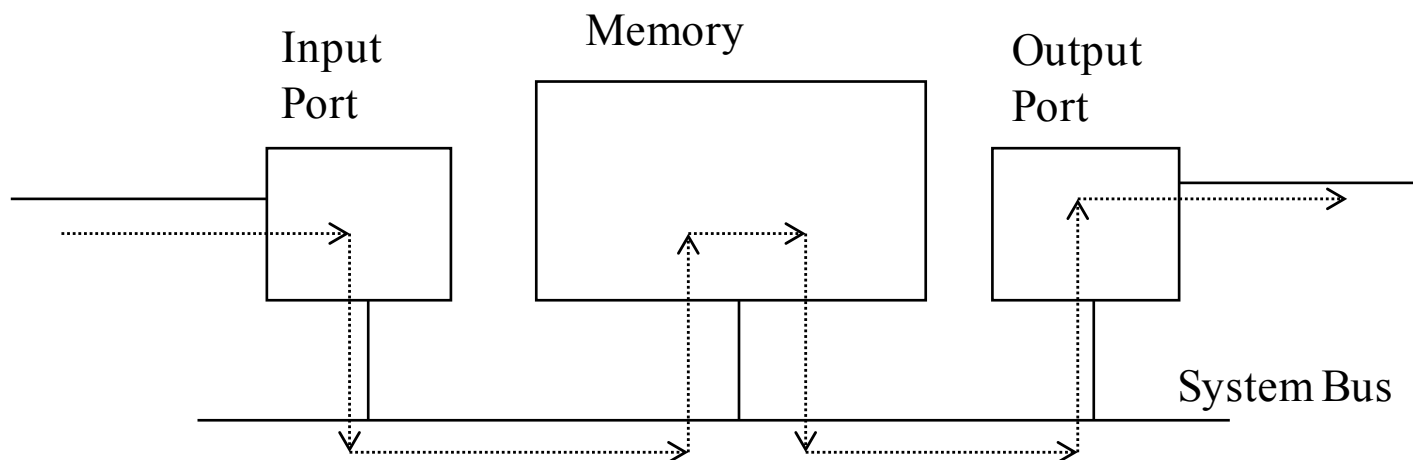
Three types of switching fabrics



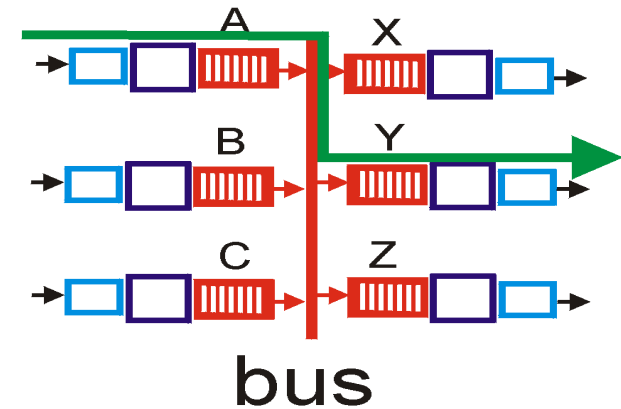
Switching Via Memory

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

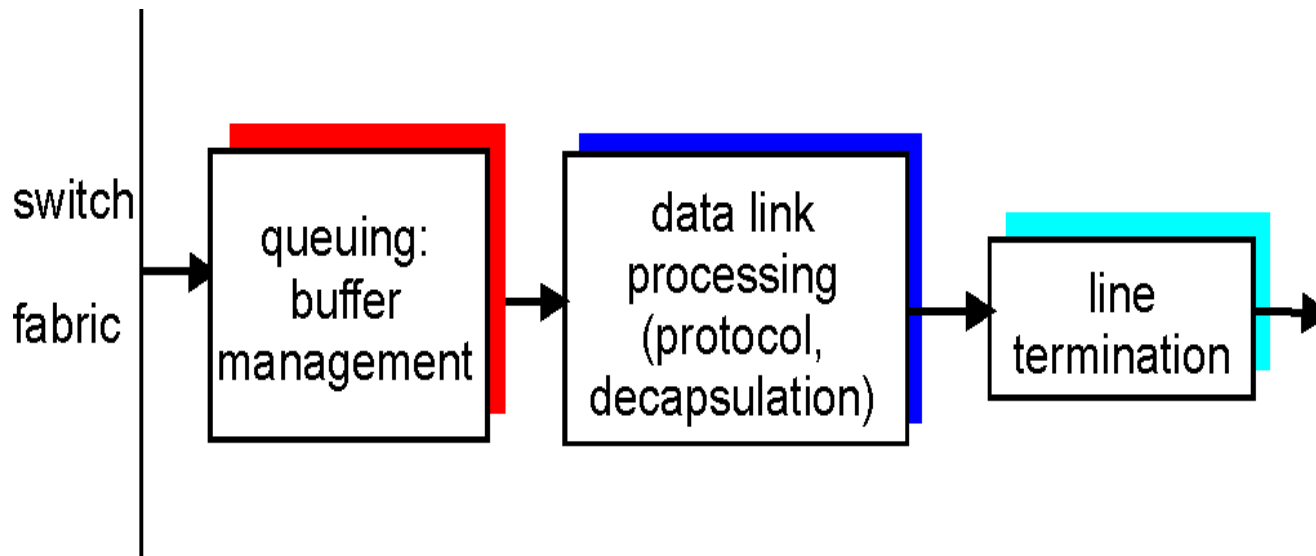


- ❑ datagram from input port memory to output port memory via a shared bus
- ❑ **bus contention:** switching speed limited by bus bandwidth
- ❑ 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers

Switching Via An Interconnection Network

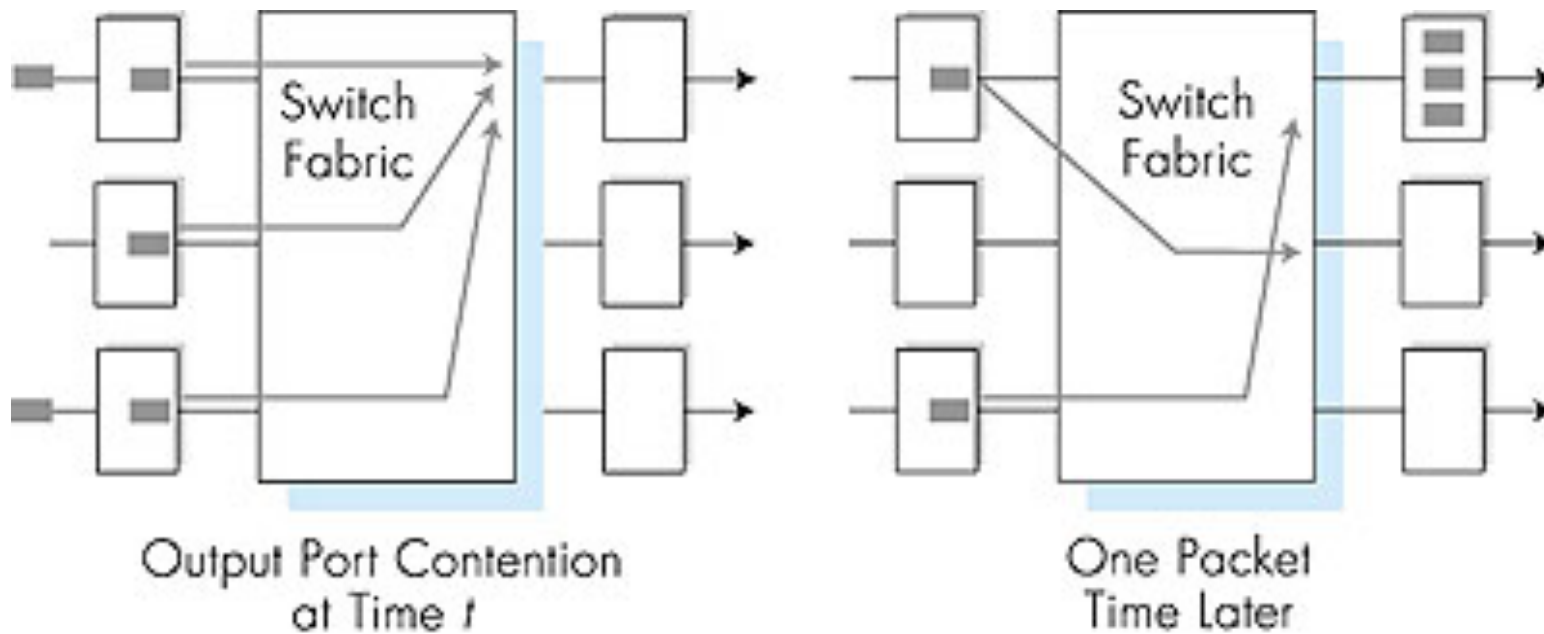
- ❑ overcome bus bandwidth limitations
- ❑ Banyan networks, other interconnection nets initially developed to connect processors in multiprocessor
- ❑ advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- ❑ Cisco 12000: switches 60 Gbps through the interconnection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



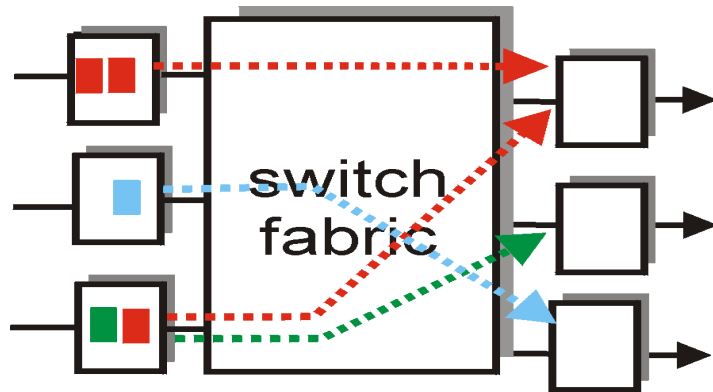
- ❑ buffering when arrival rate via switch exceeds output line speed
- ❑ *queueing (delay) and loss due to output port buffer overflow!*

How much buffering?

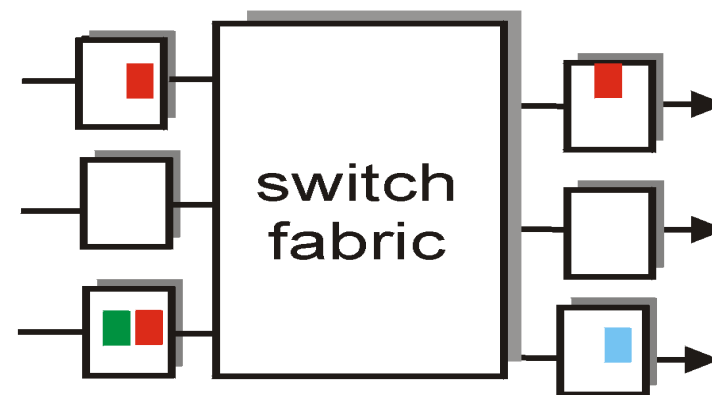
- ❑ RFC 3439 rule of thumb: average buffering equal to "typical" RTT (say 250 msec) times link capacity C
 - e.g., $C = 10$ Gps link: 2.5 Gbit buffer
- ❑ Recent recommendation: with N flows, buffering equal to $\frac{RTT \cdot C}{\sqrt{N}}$

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



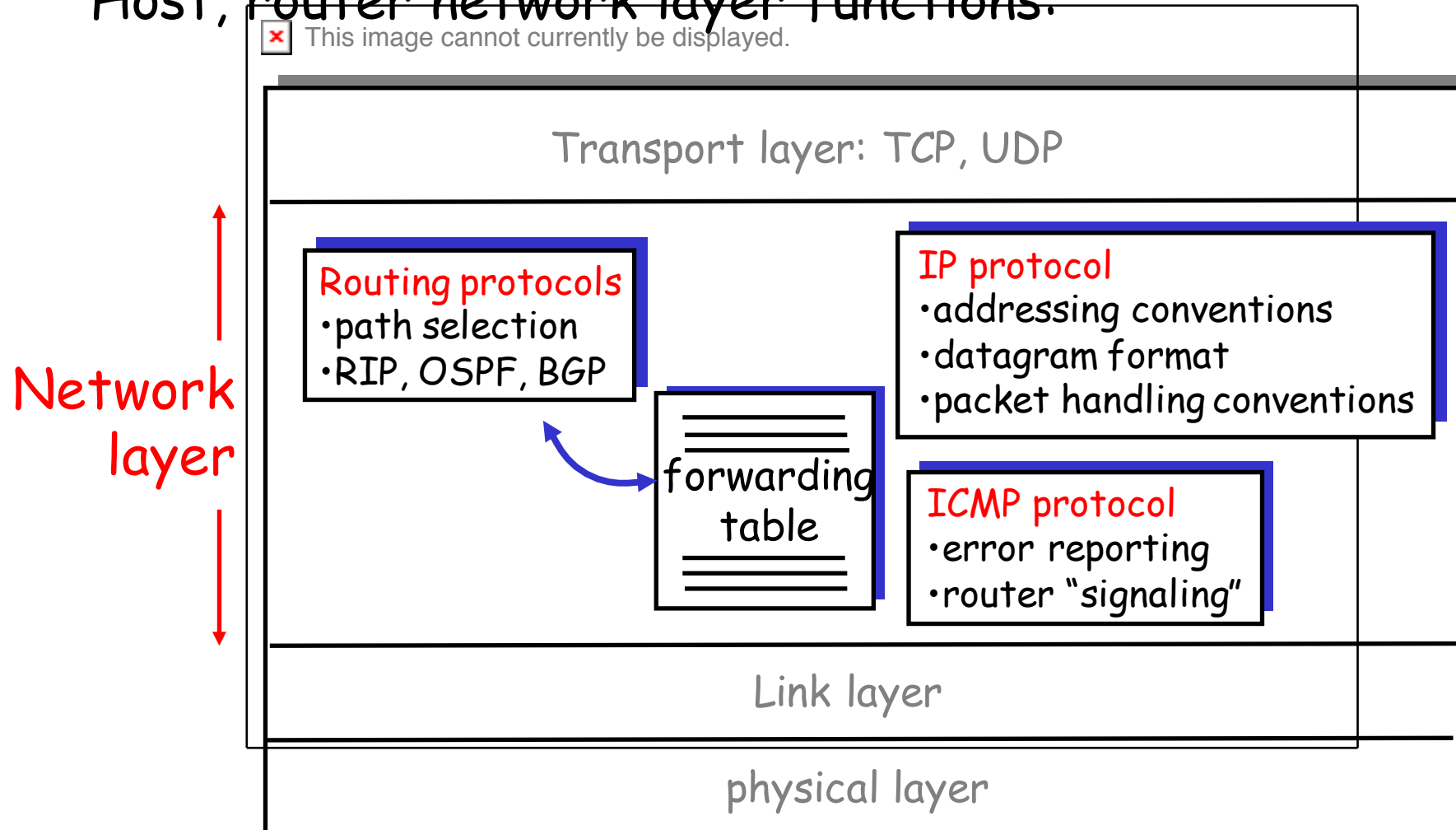
green packet experiences HOL blocking

Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 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:



Chapter 4: Network Layer

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

IP datagram format

IP protocol version number

header length (bytes)

"type" of data

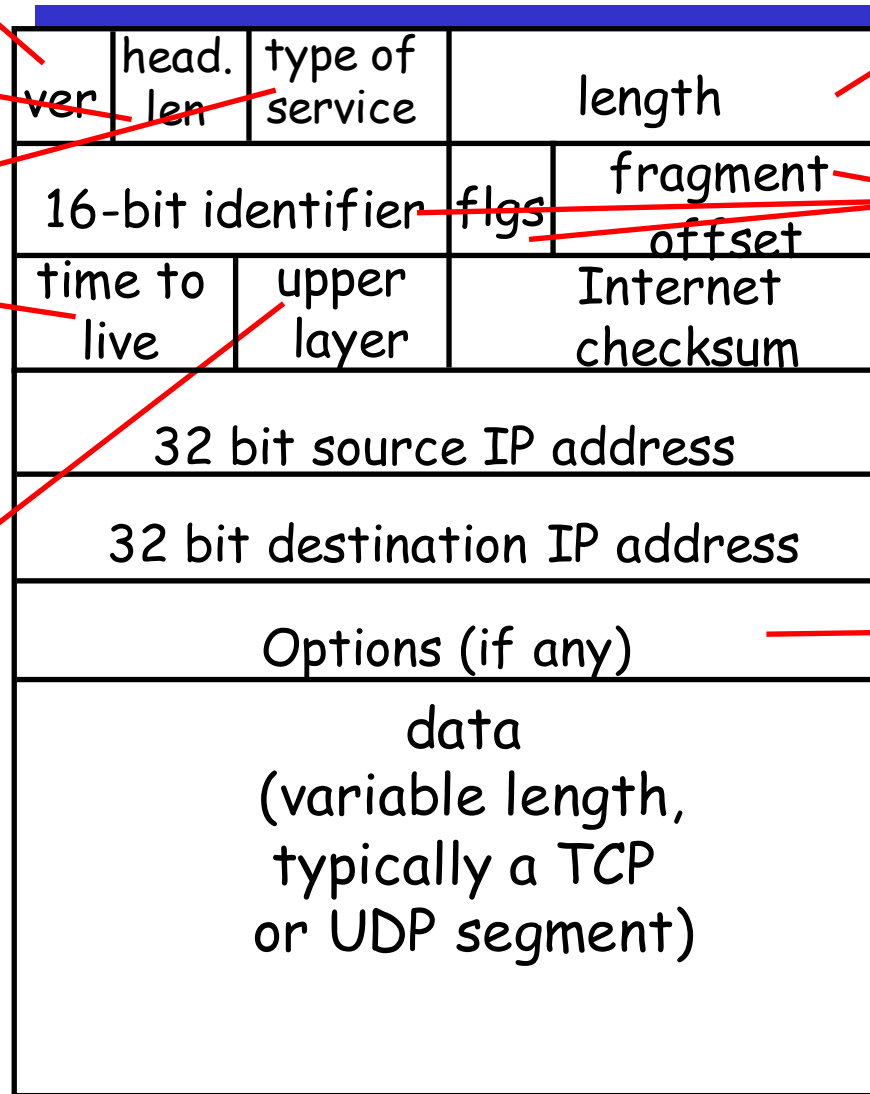
max number remaining hops (decremented at each router)

upper layer protocol to deliver payload to

how much overhead with TCP?

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

32 bits



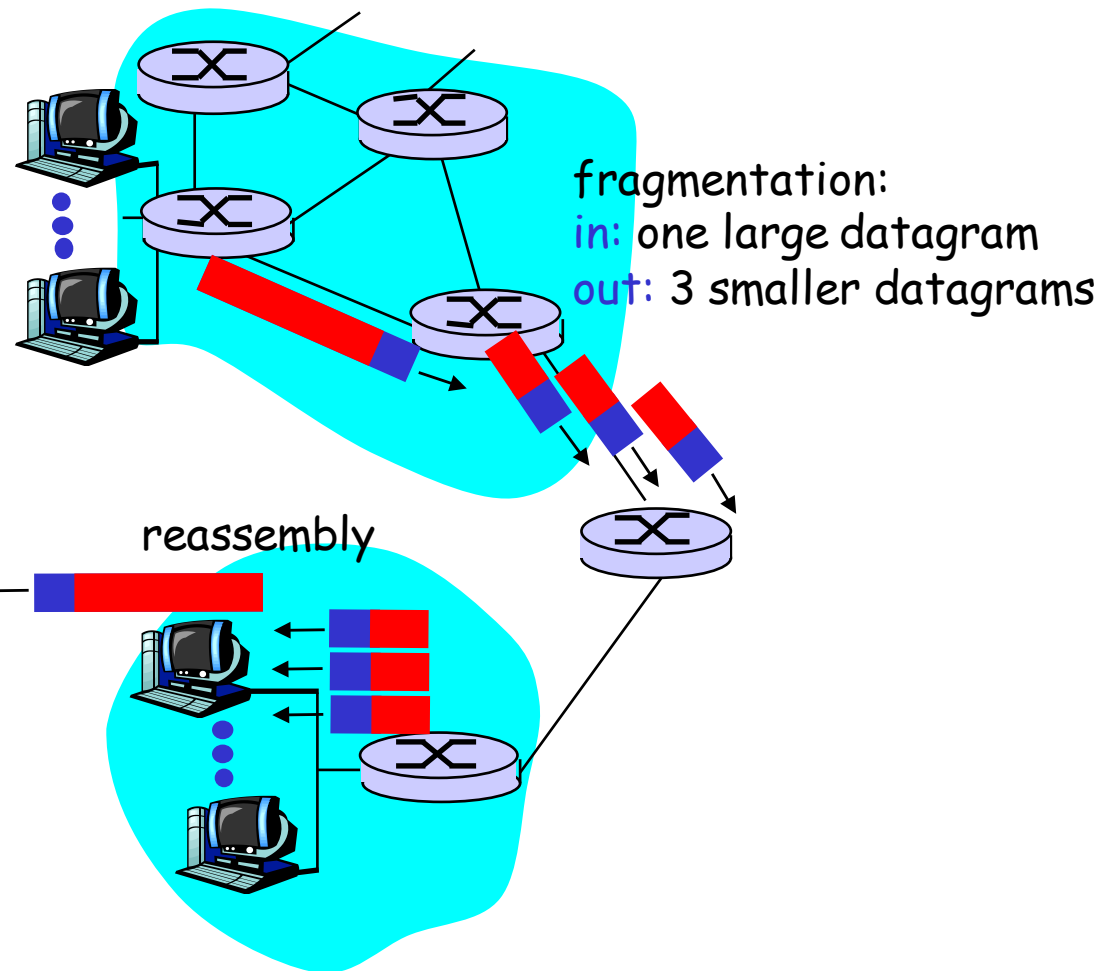
total datagram length (bytes)

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

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

One large datagram becomes several smaller datagrams

1480 bytes in data field

offset = $1480/8$

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

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

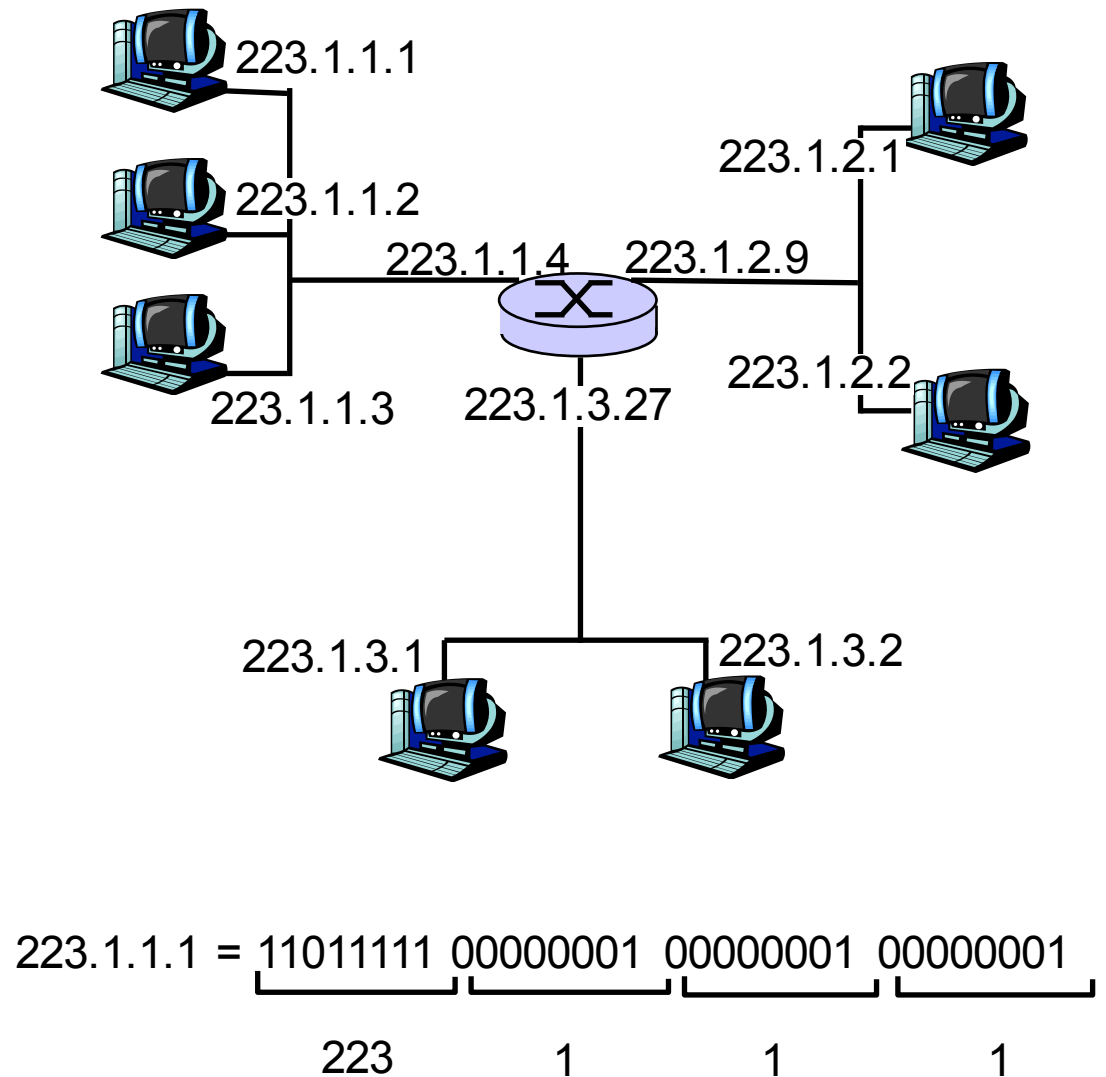
length	ID	fragflag	offset
=1040	=x	=0	=370

Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 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

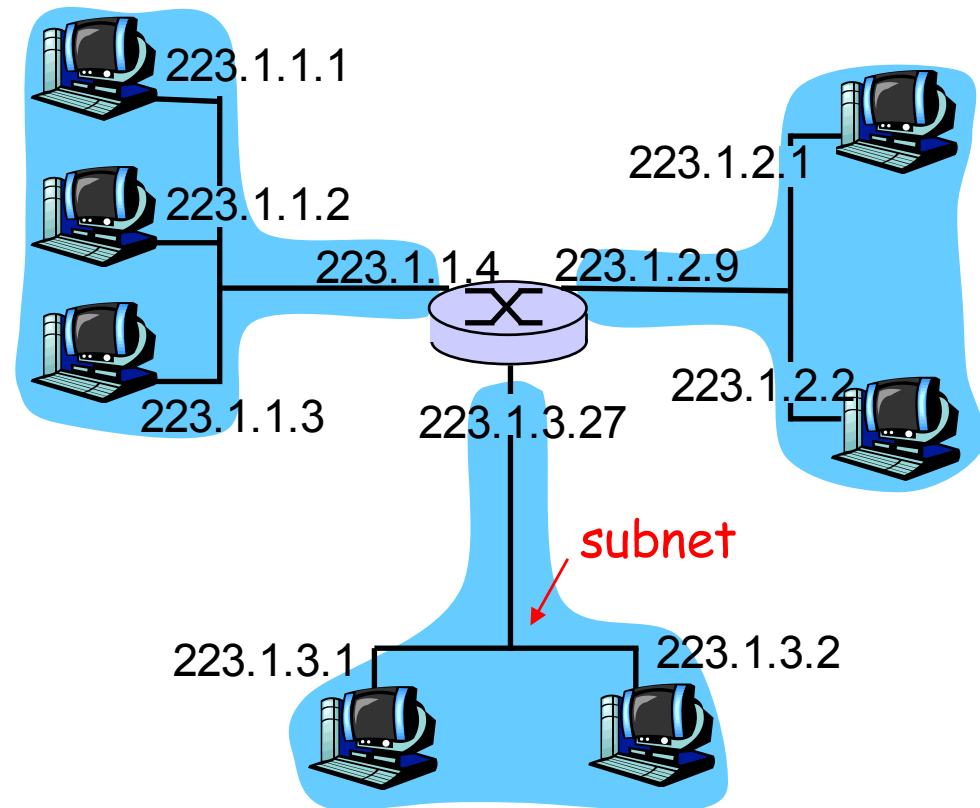
IP Addressing: introduction

- ❑ IP address: 32-bit identifier for host, router *interface*
- ❑ *interface*: connection between host/router and physical link
 - routers typically have multiple interfaces
 - Hosts typically have one but may have multiple interfaces
 - 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

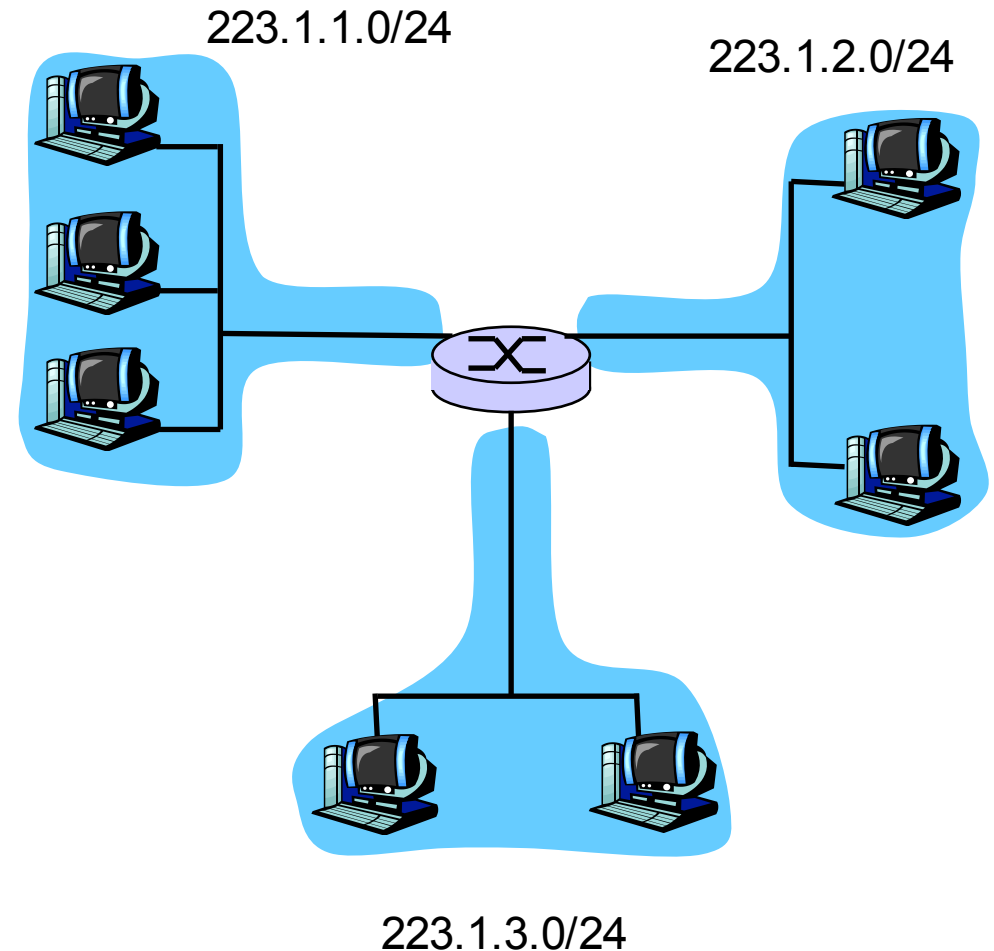


network consisting of 3 subnets

Subnets

Recipe

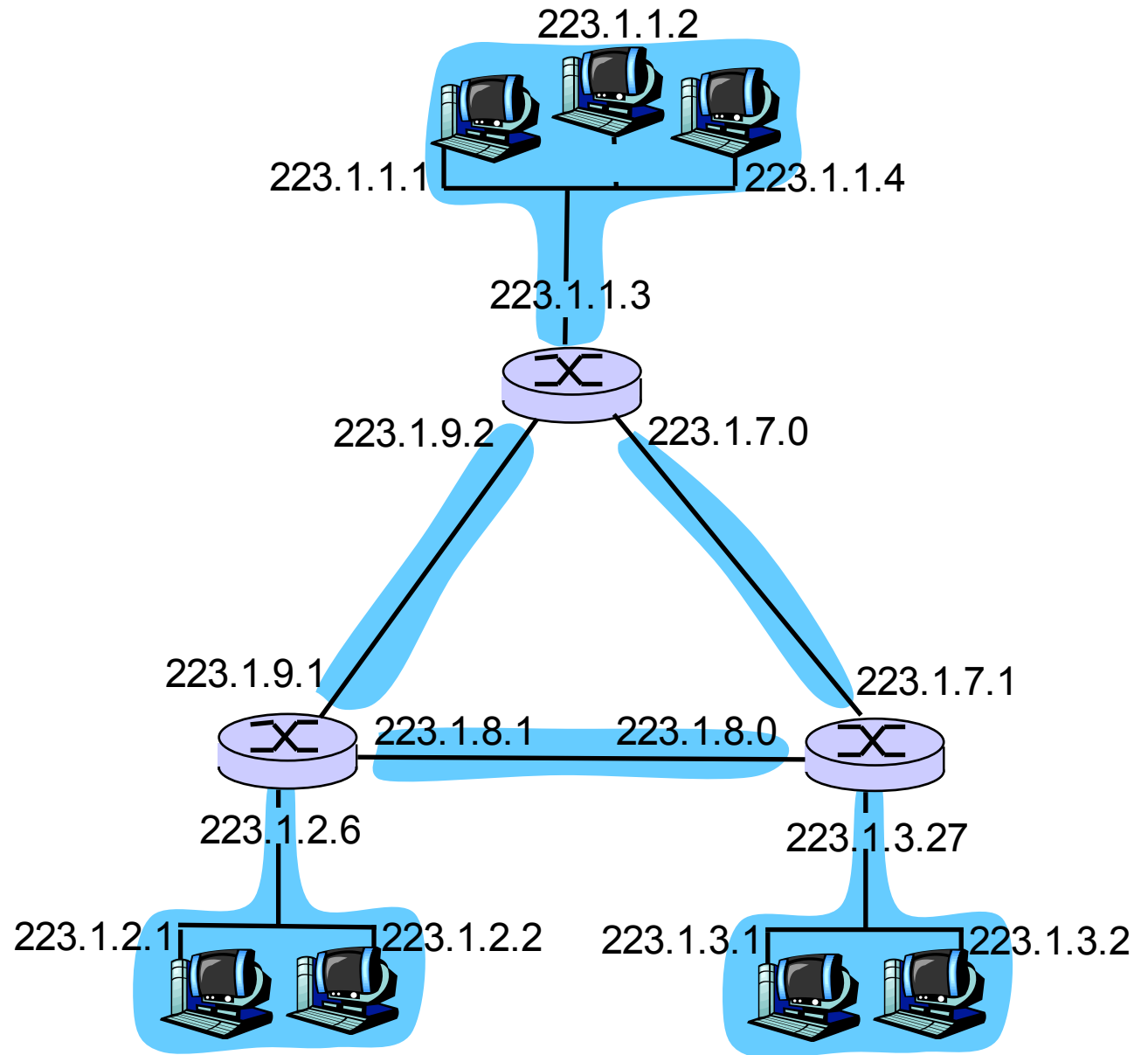
- ❑ To determine the subnets, detach each interface from its host or router, creating islands of isolated networks. Each isolated network is called a **subnet**.



Subnet mask: /24

Subnets

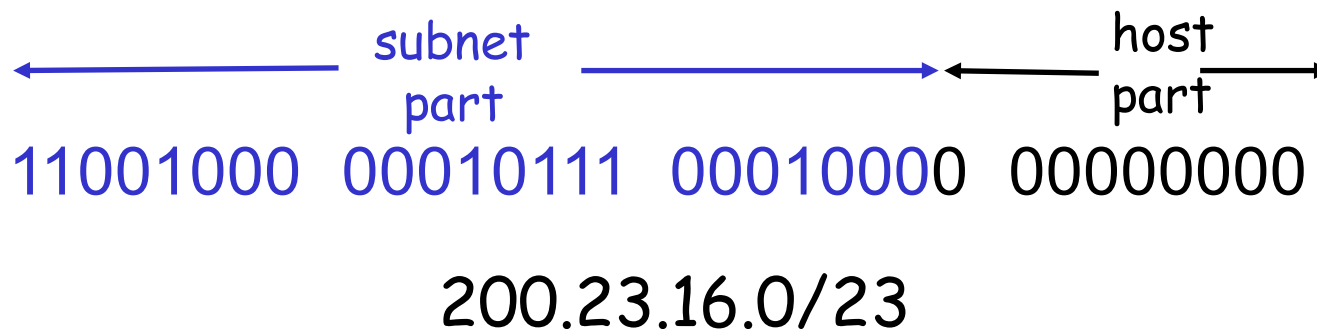
How many?



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



Quiz Time!

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 address server
 - "plug-and-play"

DHCP: Dynamic Host Configuration Protocol

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

Can renew its lease on address in use

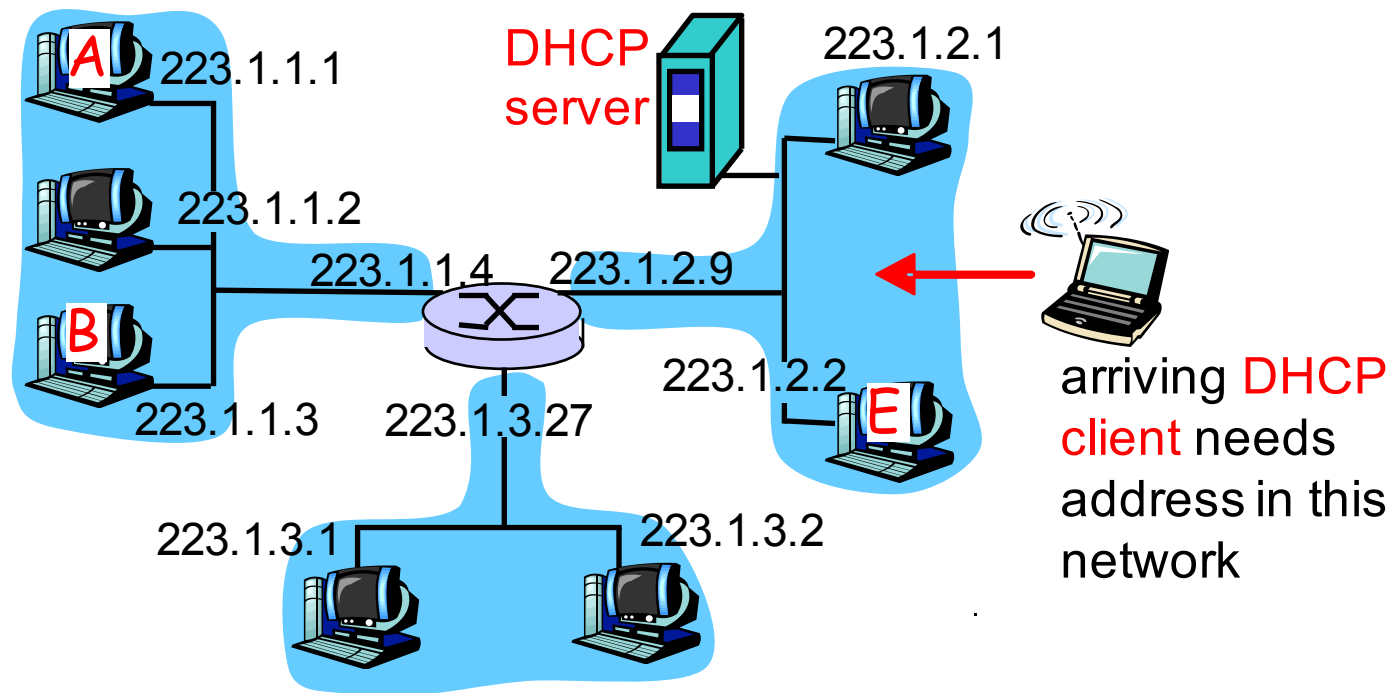
Allows reuse of addresses (only hold address while connected)

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

DHCP overview:

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

DHCP client-server scenario



DHCP client-server scenario

DHCP server: 223.1.2.5

DHCP discover

src : 0.0.0.0, 68
dest.: 255.255.255.255, 67
yiaddr: 0.0.0.0
transaction ID: 654

arriving
client



DHCP offer

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

DHCP request

src: 0.0.0.0, 68
dest.: 255.255.255.255, 67
yiaddr: 223.1.2.4
transaction ID: 655
Lifetime: 3600 secs

DHCP ACK

src: 223.1.2.5, 67
dest: 255.255.255.255, 68
yiaddr: 223.1.2.4
transaction ID: 655
Lifetime: 3600 secs

time

Quiz Time!

IP addresses: how to get one?

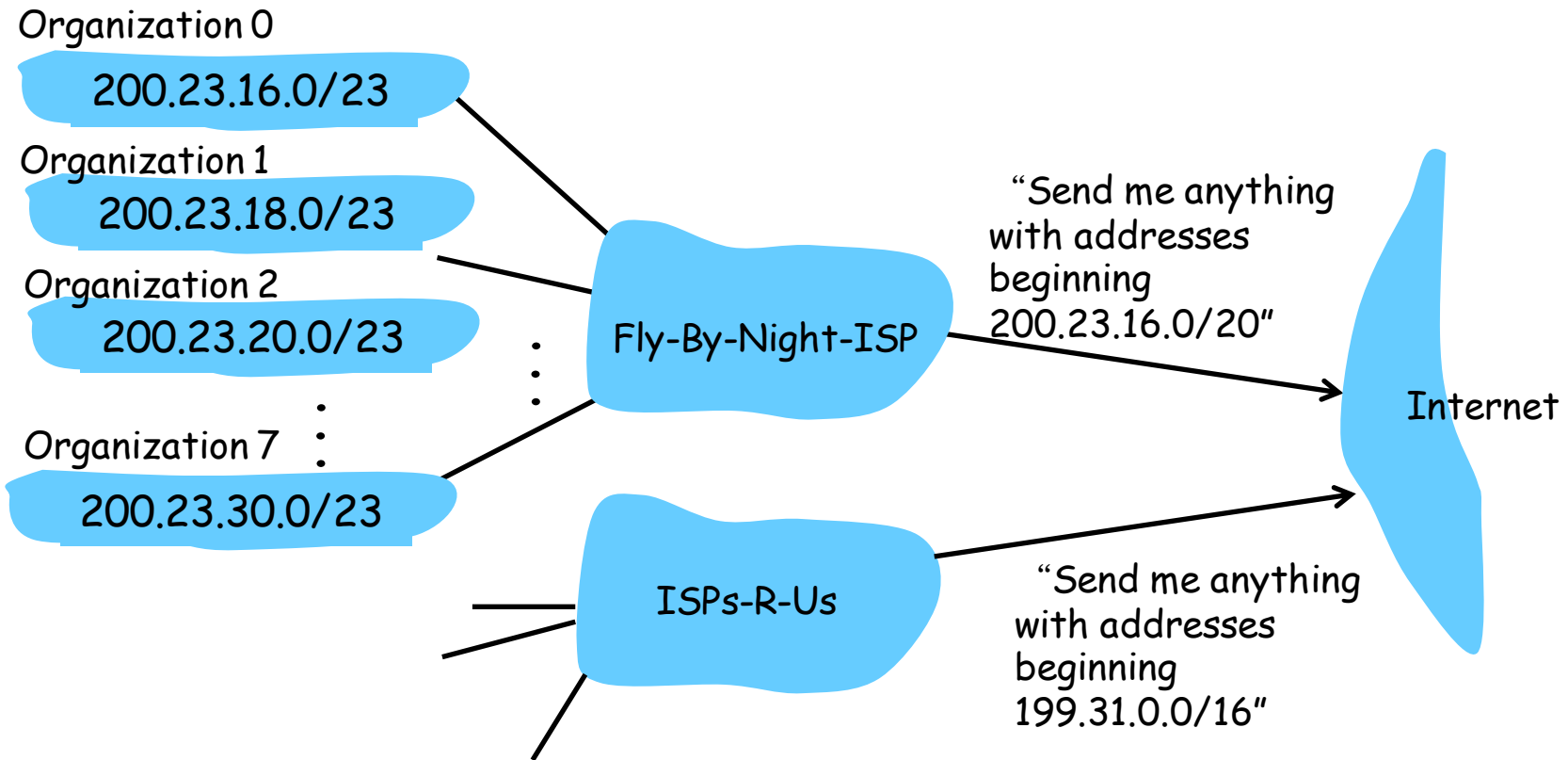
Q: How does *network* get network part of IP addr?

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

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

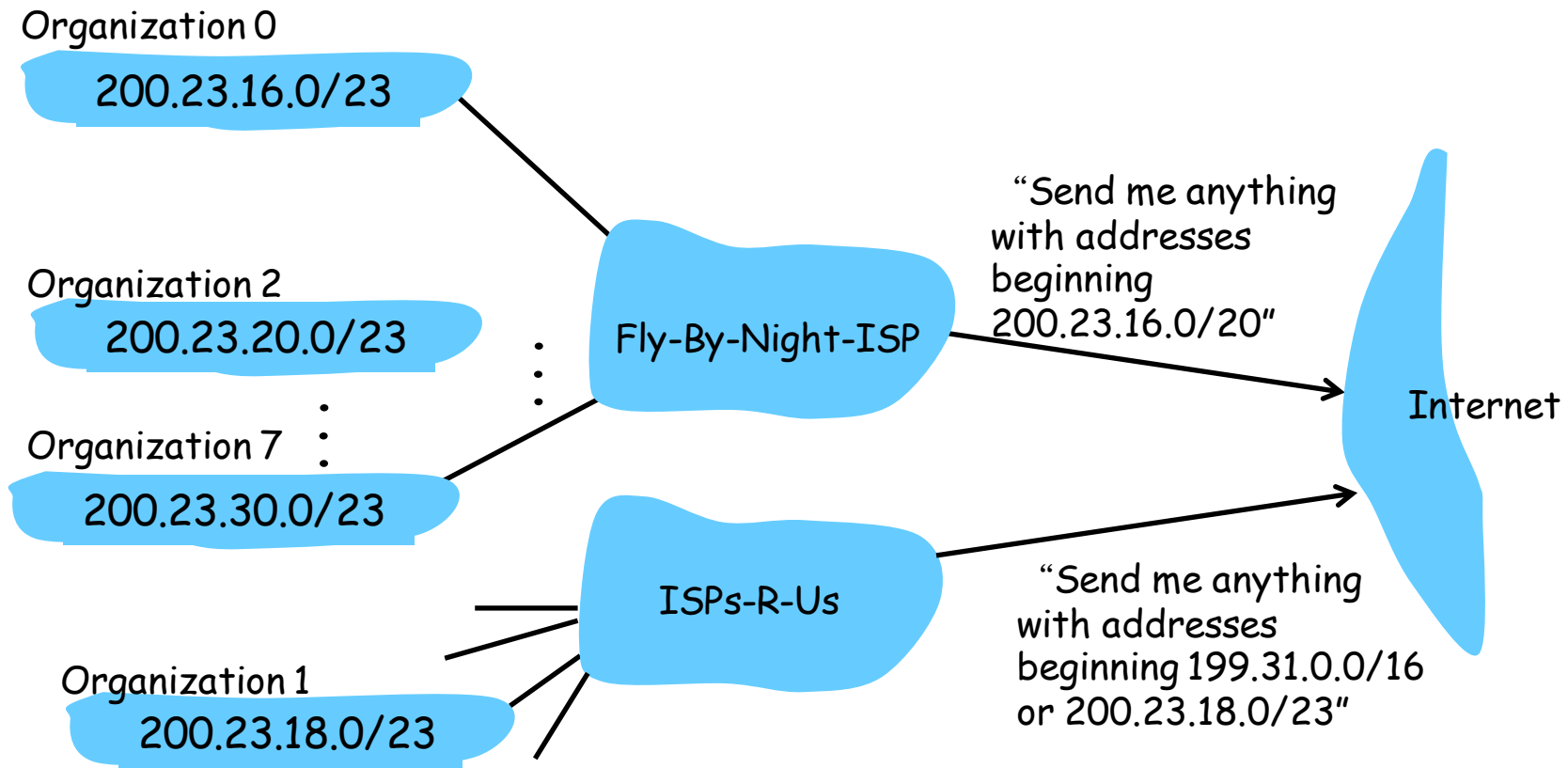
Hierarchical addressing: route aggregation

Hierarchical addressing allows efficient advertisement of routing information:



Hierarchical addressing: more specific routes

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



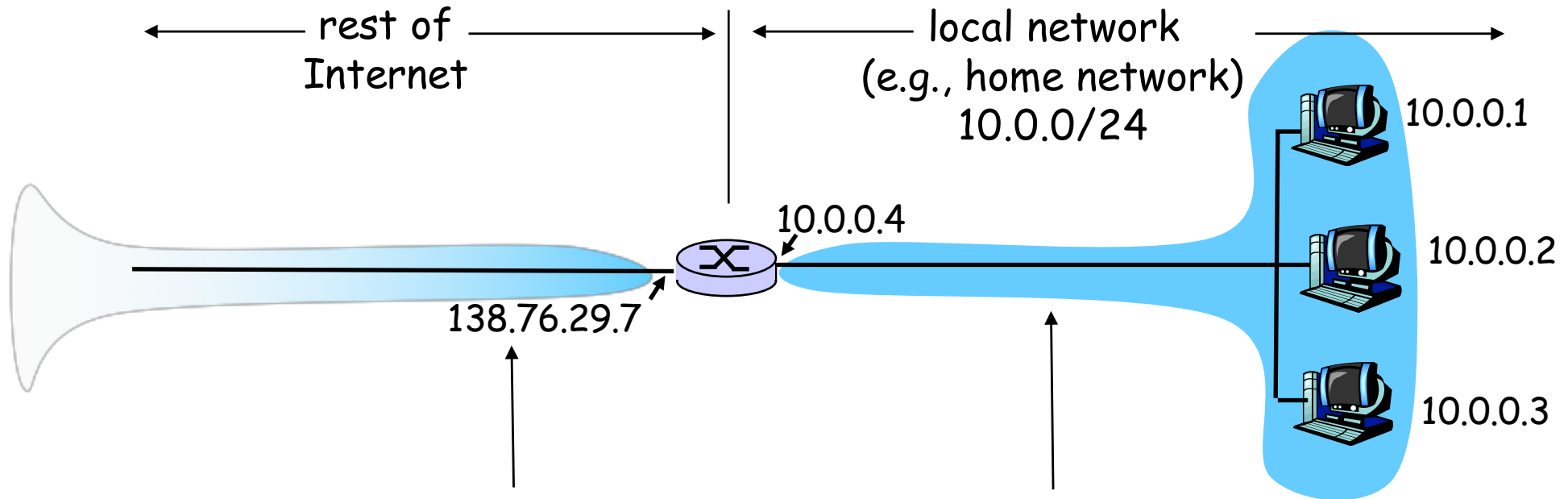
IP addressing: the last word...

Q: How does an ISP get block of addresses?

A: **ICANN**: Internet Corporation for Assigned Names and Numbers

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

NAT: Network Address Translation



All datagrams *leaving* local network have *same* single source NAT IP address: 138.76.29.7, different source port numbers

Datagrams with source or destination in this network have 10.0.0/24 address for source, destination (as usual)

NAT: Network Address Translation

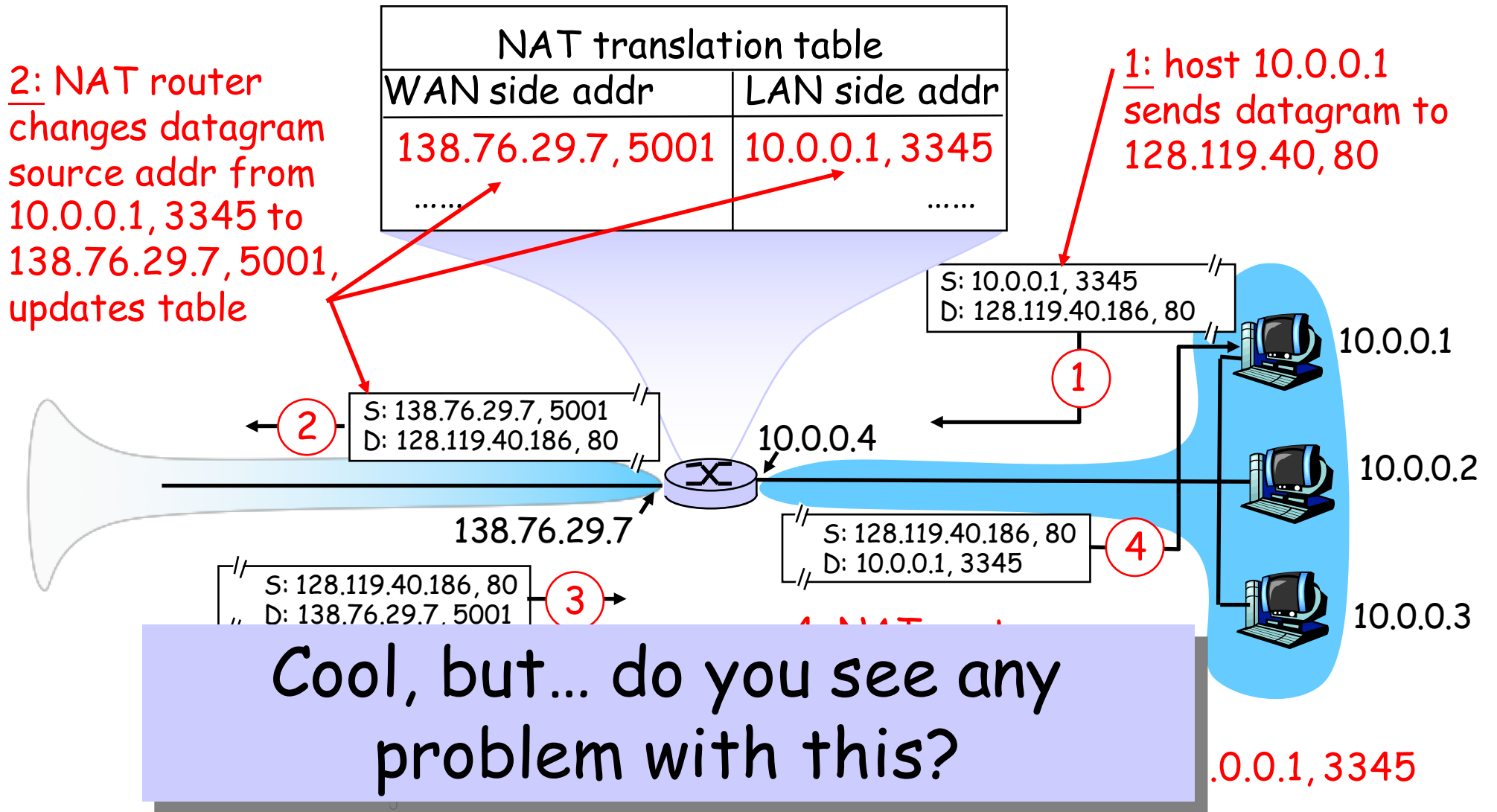
- **Motivation:** local network uses just one IP address as far as outside world is concerned:
 - no need to be allocated range of addresses from ISP:
 - just one IP address is used for all devices
 - can change addresses of devices in local network without notifying outside world
 - can change ISP without 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/gateway 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

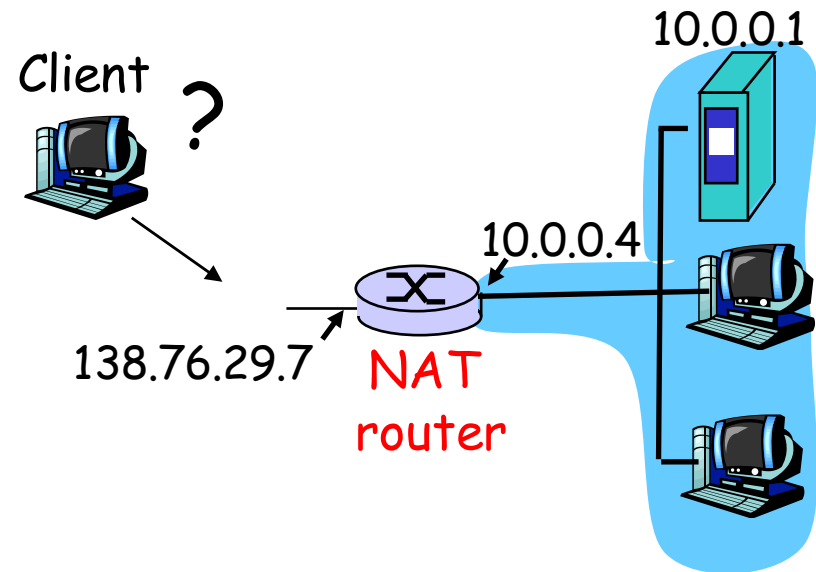


NAT: Network Address Translation

- ❑ 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 layer transparency 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 want 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., (138.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000

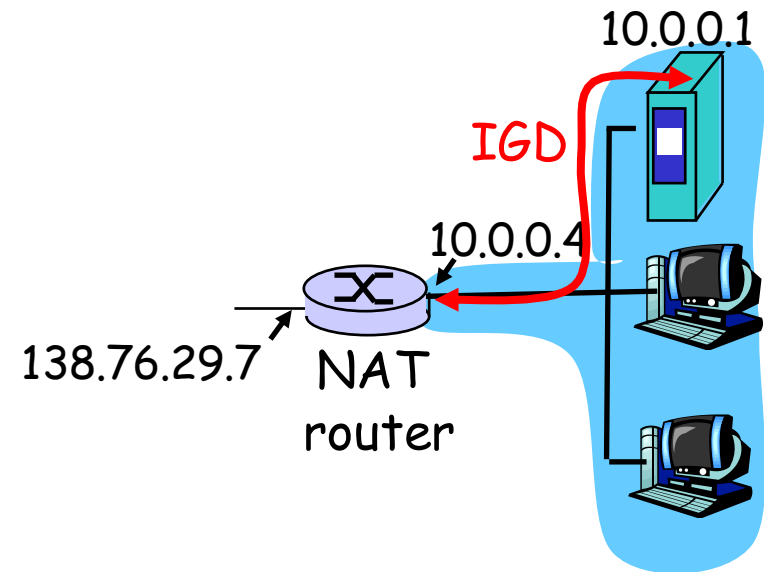


NAT traversal problem

□ solution 2: Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATted host to:

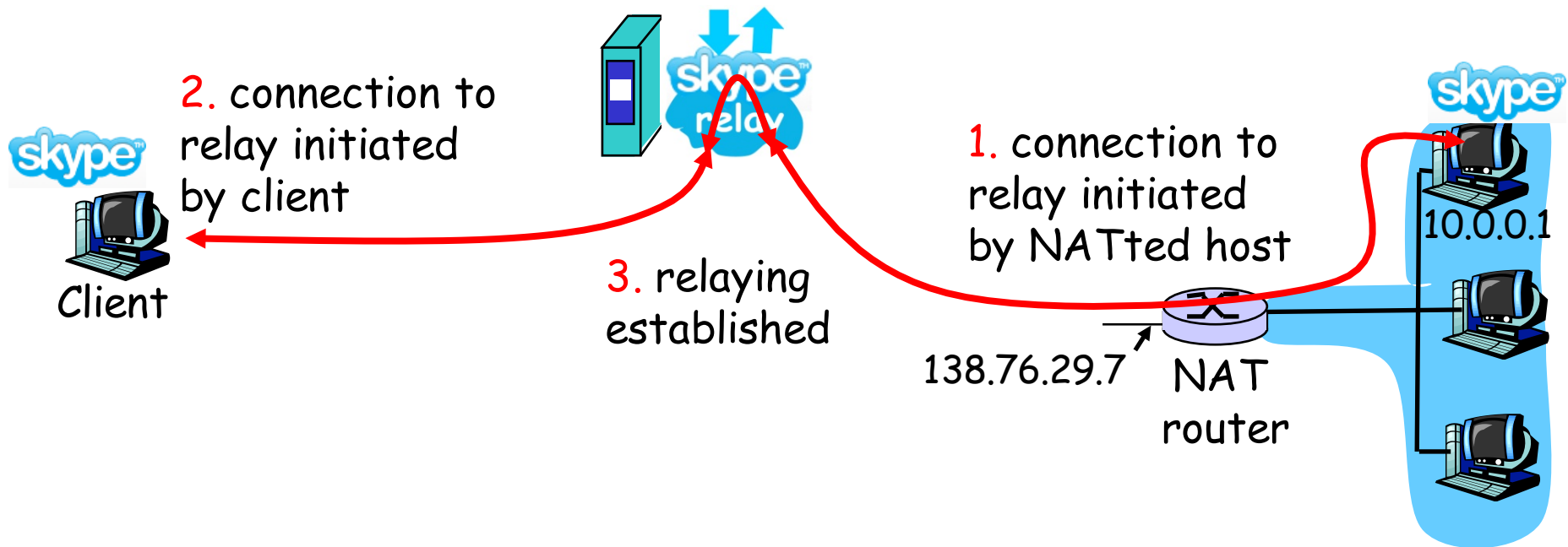
- ❖ learn public IP address (138.76.29.7)
- ❖ enumerate existing port mappings
- ❖ add/remove port mappings (with lease times)

i.e., automate static NAT port map configuration



NAT traversal problem

- solution 3: relaying (used in Skype)
 - NATed server 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
 - 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 msgs carried in IP datagrams
- ❑ **ICMP message:** type, code plus first 8 bytes of IP datagram causing error

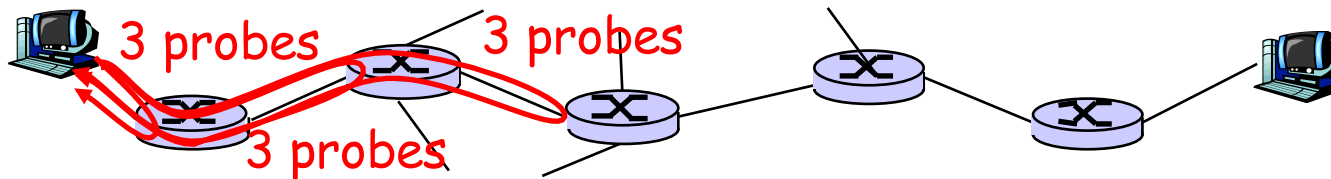
<u>Type</u>	<u>Code</u>	<u>description</u>
0	0	echo reply (ping)
3	0	dest. network unreachable
3	1	dest host unreachable
3	2	dest protocol unreachable
3	3	dest port unreachable
3	6	dest network unknown
3	7	dest host unknown
4	0	source quench (congestion control - not used)
8	0	echo request (ping)
9	0	route advertisement
10	0	router discovery
11	0	TTL expired
12	0	bad IP header

Traceroute and ICMP

- ❑ Source sends series of UDP segments to dest
 - First has TTL =1
 - Second has TTL=2, etc.
 - Unlikely port number
 - ❑ When nth datagram arrives to nth router:
 - Router discards datagram
 - And sends to source an ICMP message (type 11, code 0)
 - Message includes name of router& IP address
 - ❑ When ICMP message arrives, source calculates RTT
 - ❑ Traceroute does this 3 times
- Stopping criterion
- ❑ UDP segment eventually arrives at destination host
 - ❑ Destination returns ICMP "host unreachable" packet (type 3, code 3)
 - ❑ When source gets this ICMP, stops.

"Real" Internet delays and routes

- ❑ What do "real" Internet delay & loss look like?
- ❑ **Traceroute program**: provides delay measurement from source to router along end-end Internet path towards destination. For all i :
 - sends three packets that will reach router i on path towards destination
 - router i will return packets to sender
 - sender times interval between transmission and reply.



"Real" Internet delays and routes

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

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

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

trans-oceanic link

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

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

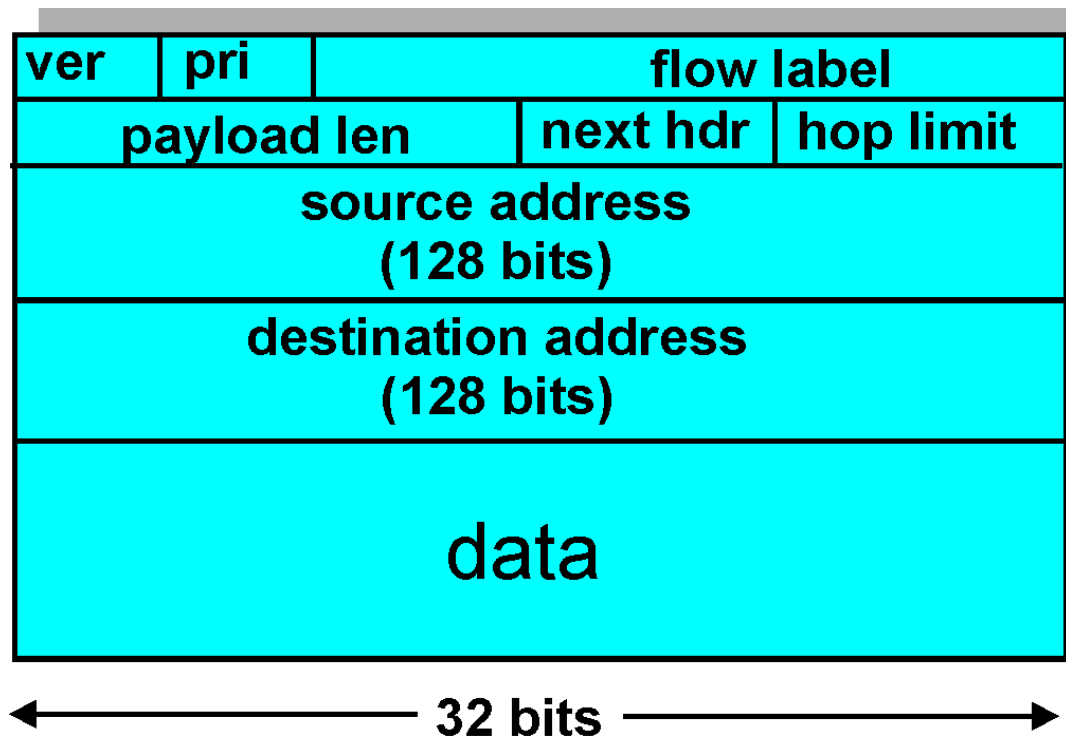
- ❑ **Initial motivation:** 32-bit address space soon to be completely allocated
- ❑ **Additional motivation:**
 - header format helps speed processing/forwarding
 - header changes to facilitate QoS
- ❑ **IPv6 datagram format:**
 - fixed-length 40 byte header
 - no fragmentation allowed

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



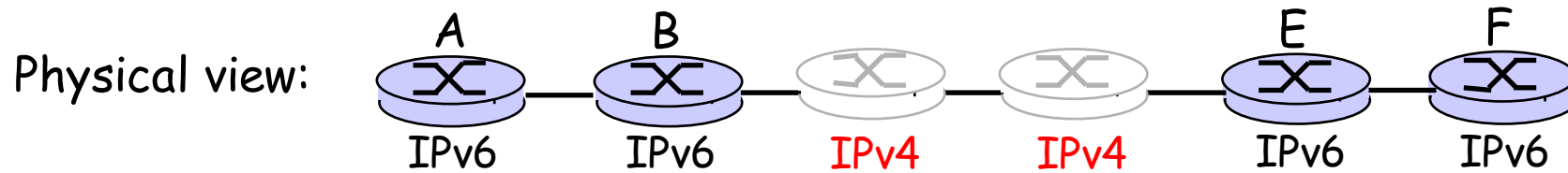
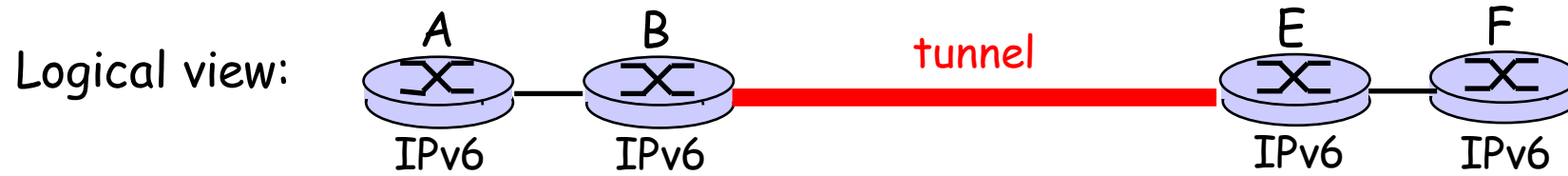
Other Changes from IPv4

- ❑ *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

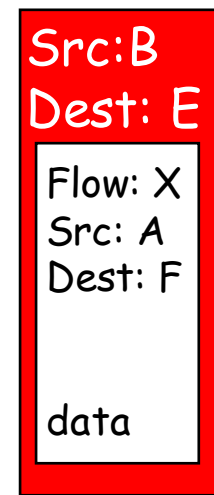
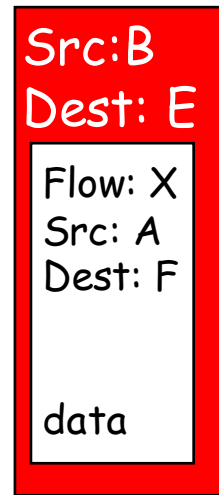
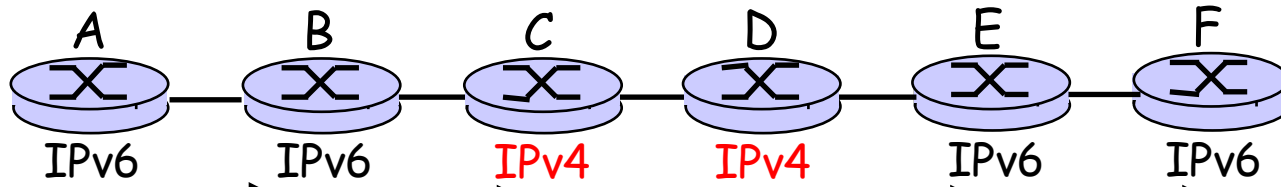


Tunneling

Logical view:



Physical view:



A-to-B:
IPv6

B-to-C:
IPv6 inside
IPv4

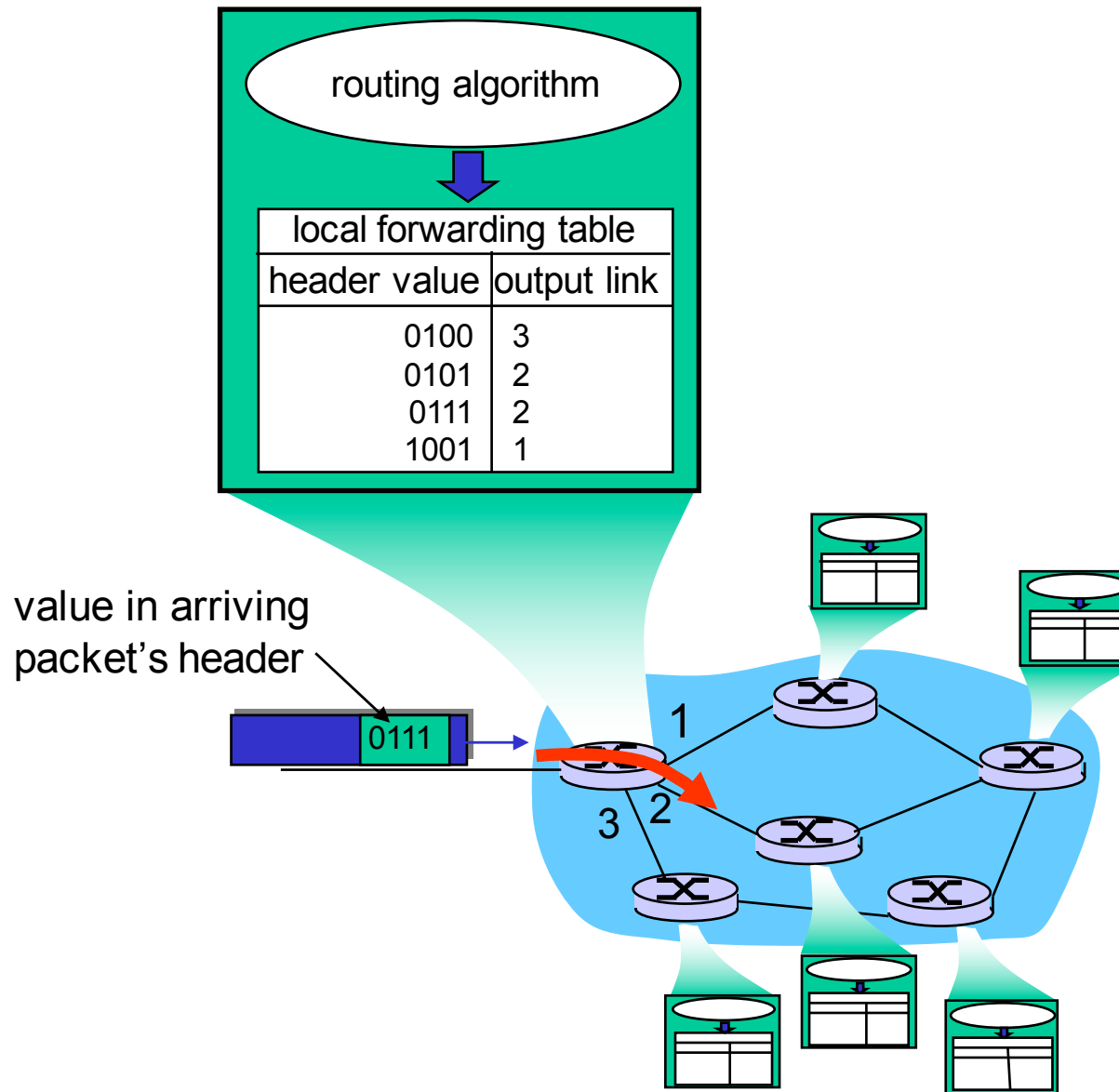
B-to-C:
IPv6 inside
IPv4

E-to-F:
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

Interplay between routing, forwarding



Getting a datagram from source to dest.

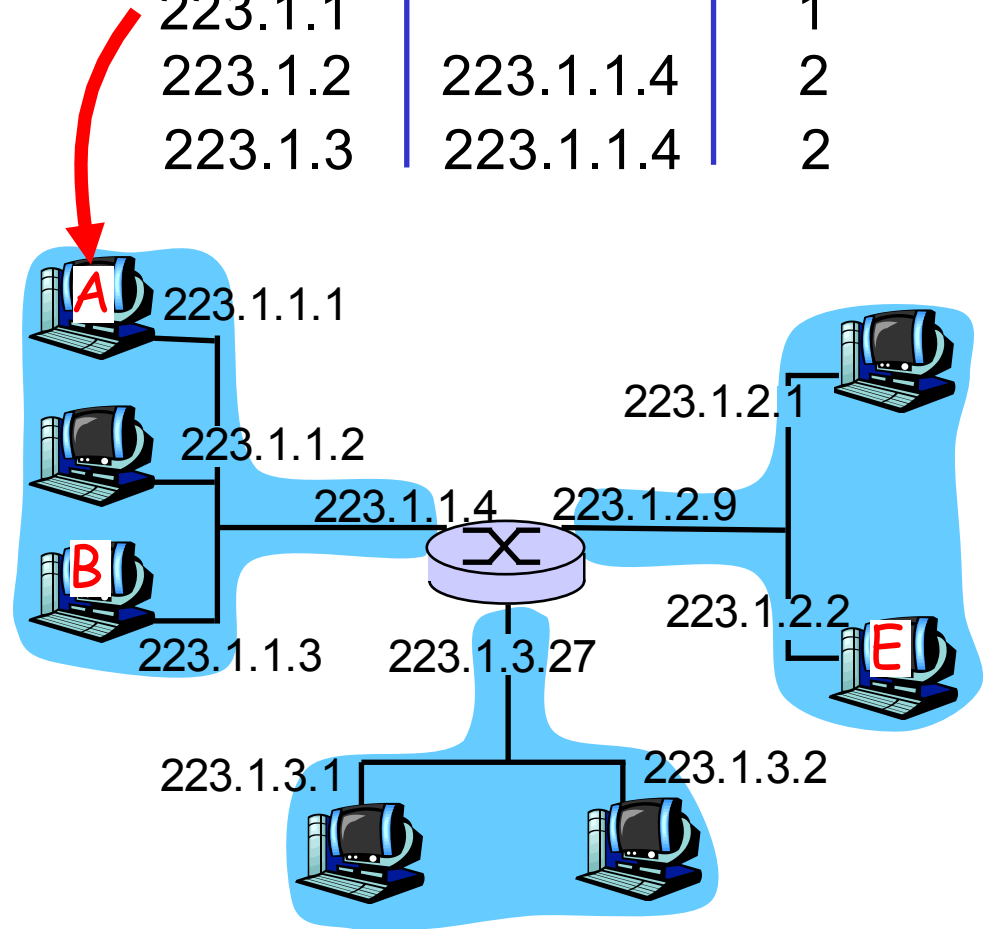
IP datagram:

misc fields	source IP addr	dest IP addr	data
-------------	----------------	--------------	------

- datagram remains **unchanged**, as it travels source to destination
- addr fields of interest here

forwarding table in A

Dest. Net.	next router	Nhops
223.1.1		1
223.1.2	223.1.1.4	2
223.1.3	223.1.1.4	2



Getting a datagram from source to dest.

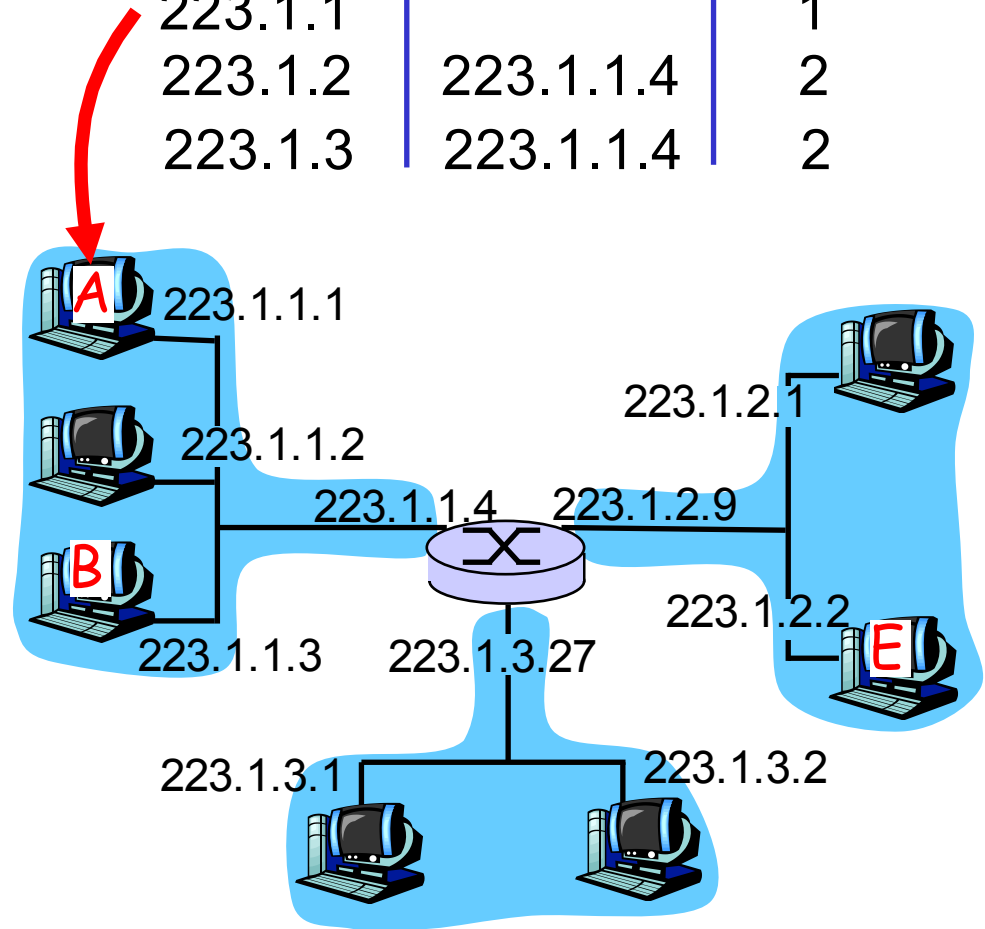
misc fields	223.1.1.1	223.1.1.3	data
-------------	-----------	-----------	------

Starting at A, send IP datagram addressed to B:

- ❑ look up net. address of B in forwarding table
- ❑ find B is on same net. as A
- ❑ link layer will send datagram directly to B inside link-layer frame
 - B and A are directly connected

forwarding table in A

Dest. Net.	next router	Nhops
223.1.1		1
223.1.2	223.1.1.4	2
223.1.3	223.1.1.4	2



Getting a datagram from source to dest.

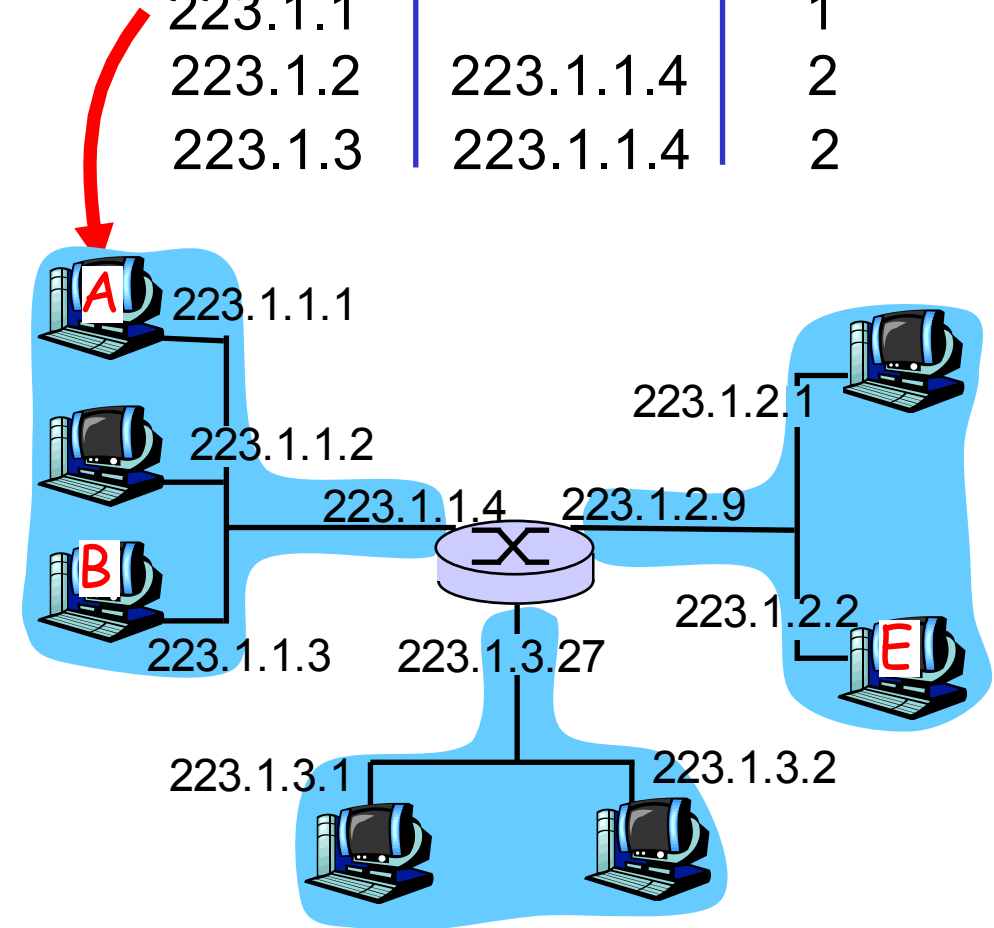
misc fields	223.1.1.1	223.1.2.2	data
-------------	-----------	-----------	------

Starting at A, dest. E:

- ❑ look up network address of E in forwarding table
- ❑ E on *different* network
 - A, E not directly attached
- ❑ routing table: next hop router to E is 223.1.1.4
- ❑ link layer sends datagram to router 223.1.1.4 inside link-layer frame
- ❑ datagram arrives at 223.1.1.4
- ❑ continued.....

forwarding table in A

Dest. Net.	next router	Nhops
223.1.1		1
223.1.2	223.1.1.4	2
223.1.3	223.1.1.4	2



Getting a datagram from source to dest.

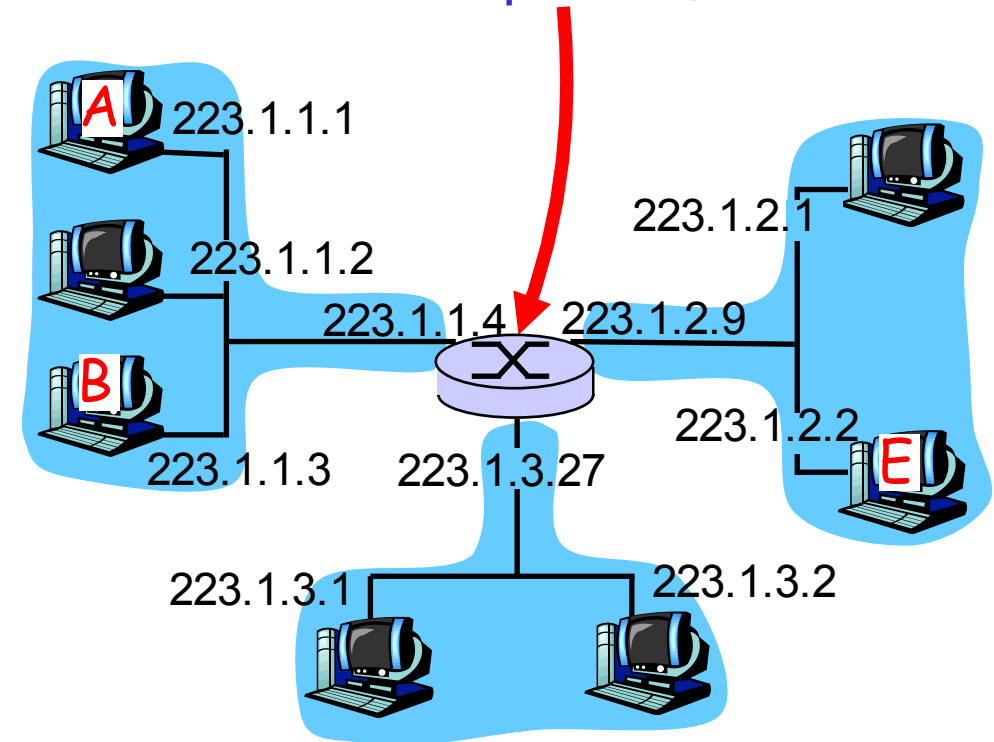
misc fields	223.1.1.1	223.1.2.2	data
-------------	-----------	-----------	------

Arriving at 223.1.4,
destined for 223.1.2.2

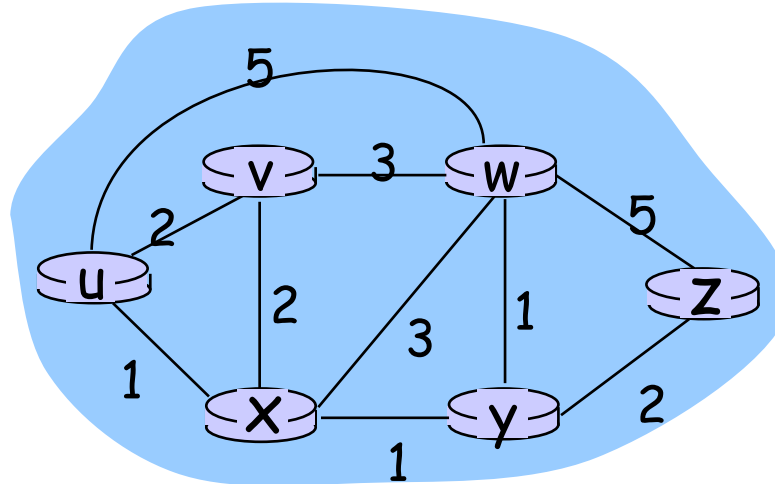
- ❑ look up network address of E in router's forwarding table
- ❑ E on *same* network as router's interface 223.1.2.9
 - router, E directly attached
- ❑ link layer sends datagram to 223.1.2.2 inside link-layer frame via interface 223.1.2.9
- ❑ datagram arrives at 223.1.2.2!!! (hooray!)

forwarding table in router

Dest. Net	router	Nhops	interface
223.1.1	-	1	223.1.1.4
223.1.2	-	1	223.1.2.9
223.1.3	-	1	223.1.3.27



Graph abstraction



Graph: $G = (N,E)$

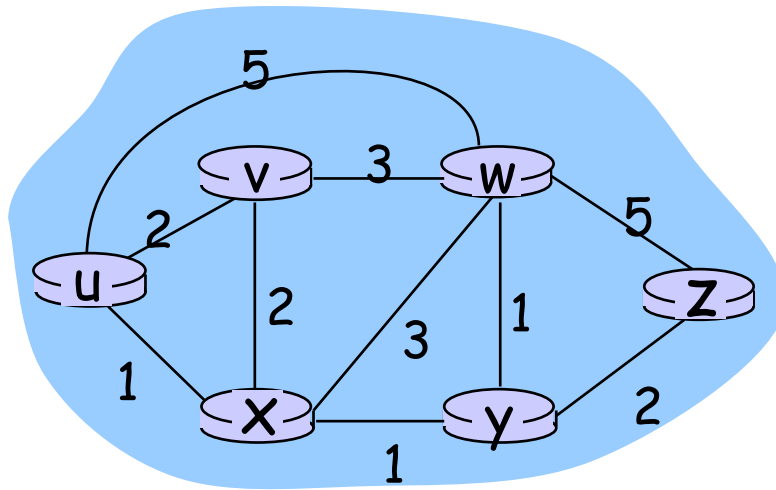
$N = \text{set of routers} = \{ u, v, w, x, y, z \}$

$E = \text{set of links} = \{ (u,v), (u,x), (u,w), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z) \}$

Remark: Graph abstraction is useful in other network contexts

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

Graph abstraction: costs



- $c(x,x')$ = cost of link (x,x')

- e.g., $c(w,z) = 5$

- cost could always be 1, or inversely related to bandwidth, or inversely related to congestion

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

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 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 Idea of Link State Routing

- ❑ Given the topology (graph)
- ❑ Compute the least cost path (the shortest path)
 - From each possible source
 - To each possible destination

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 each node ("source") to all other nodes
 - gives forwarding table for that node
- ❑ iterative: after k iterations, know least cost path to k dest.'s

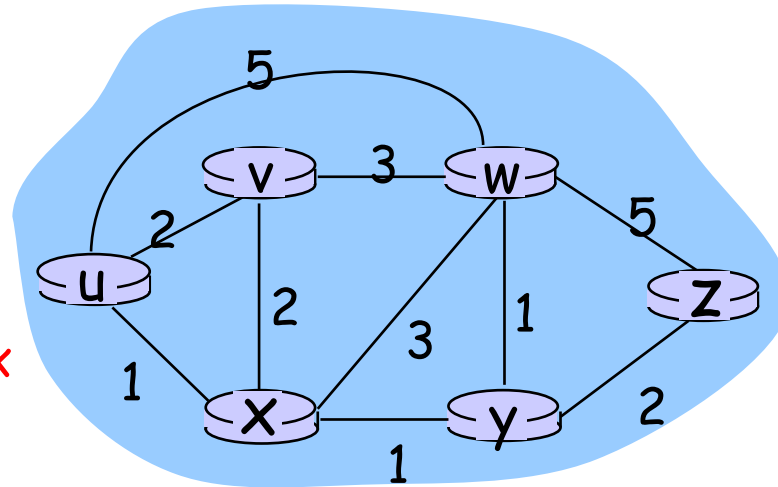
Notation:

- ❑ $c(a,b)$: link cost from node a to b ; $= \infty$ if not direct neighbors
- ❑ $D(d)$: current value of cost of path from source to dest. d
- ❑ $p(d)$: predecessor node along path from source to d
- ❑ N' : set of nodes whose least cost path definitively known

Dijkstra's Algorithm

1 **Initialization:**

- 2 Travel from the source a ^u
 3 for all nodes b
 4 if b adjacent to a
 5 then $D(b) = c(a,b)$ ^{v, w, x}
 6 else $D(b) = \text{infinity}$ ^{y, z}
 7

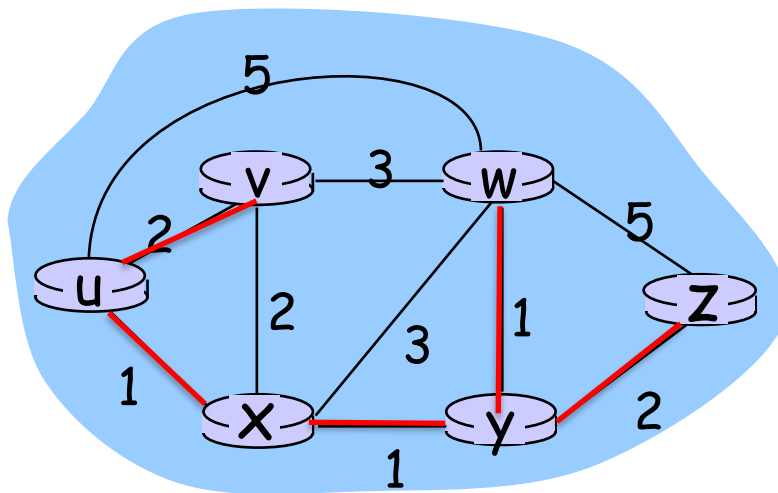


8 **Loop**

- 9 find another node c adjacent to the traveled nodes to travel ^x
 10 such that $D(c)$ is a minimum
 11 for all d adjacent to c and not yet traveled ^{u, v, w, y} v, w, y
 12 $D(d) = \min(D(d), D(c) + c(c,d))$ ← check if there's a better $D(d)$
 13 /* new cost to d is the better one between the old cost to d
 14 or the minimum path cost to c plus the c-d link cost */
 15 **until all nodes are traveled**

Dijkstra's algorithm: example

Step	Travel Set	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	infinity	infinity
→ 1	ux	2,u	4,x		2,x	infinity
→ 2	uxy	2,u	3,y			4,y
→ 3	uxyv		3,y			4,y
→ 4	uxyvw					4,y
5	uxyvwz					



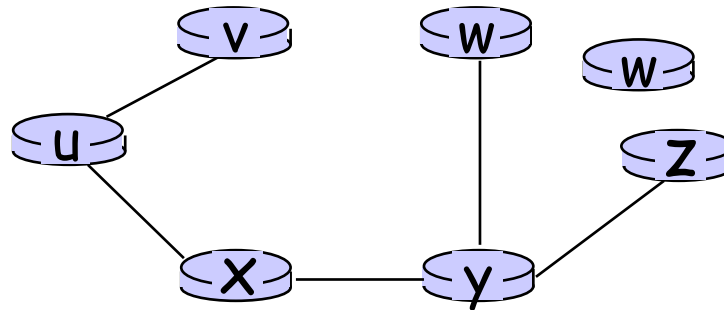
OK the min cost to w is 3
But what's the path from u to w?

How do computers figure it out from the table?

From w, get the previous node y
From y, get the previous node x
From x, get the previous node u

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)

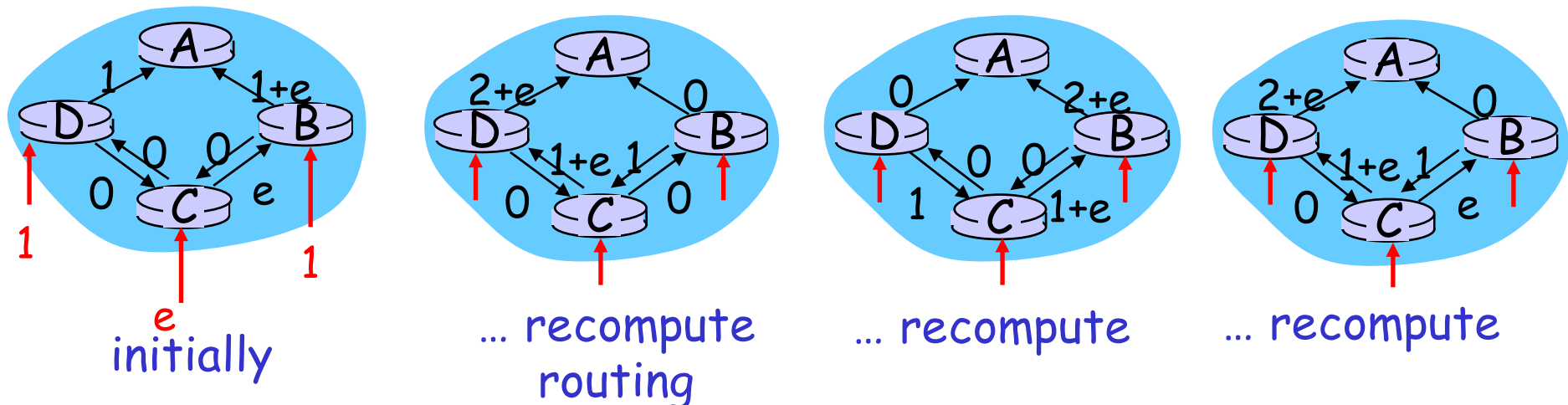
Dijkstra's algorithm, discussion

Algorithm complexity: n nodes

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

Oscillations possible:

- ❑ e.g., link cost = amount of carried traffic
- ❑ initially, link cost = 0
- ❑ traffic coming from B, C, D to A



LS Routing Summary

- ❑ net topology, link costs known to all nodes
 - accomplished via "link state broadcast"
 - all nodes have the entire topology info
- ❑ computes least cost paths from one node ('source') to all other nodes
 - gives **routing table** for that node

Quiz Time

LS Routing Summary

- ❑ net topology, link costs known to all nodes
 - accomplished via "link state broadcast"
 - all nodes have the entire topology info
- ❑ computes least cost paths from one node ('source') to all other nodes
 - gives **routing table** for that node
- ❑ Do you see any problems?
 - LS broadcast: consumes bandwidth
 - Topology info: occupies memory space

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

Can you do without knowing the Topology?

Yes, I tell my neighbors.
You tell yours

Quiz Time!

How does a router know the
best route without knowing the
topology?

Check which neighbor is the closest
(based on what the neighbors have told)

Distance Vector Algorithm (1)

Bellman-Ford Equation (dynamic programming)

Define

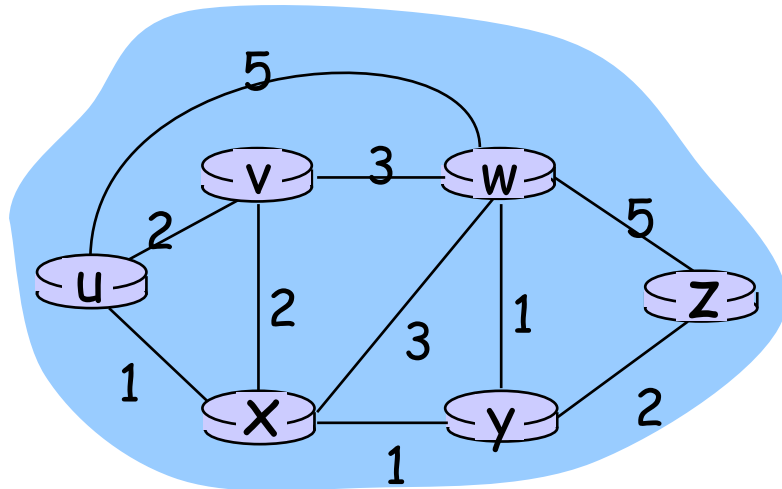
$d_x(y) :=$ cost of least-cost path from x to y

Then

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

where min is taken over all neighbors 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 (3)

- $d_x(y)$ = 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:

- Each node periodically sends its own distance vector estimate to neighbors
- 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*

Distance Vector Algorithm (5)

Iterative, asynchronous:

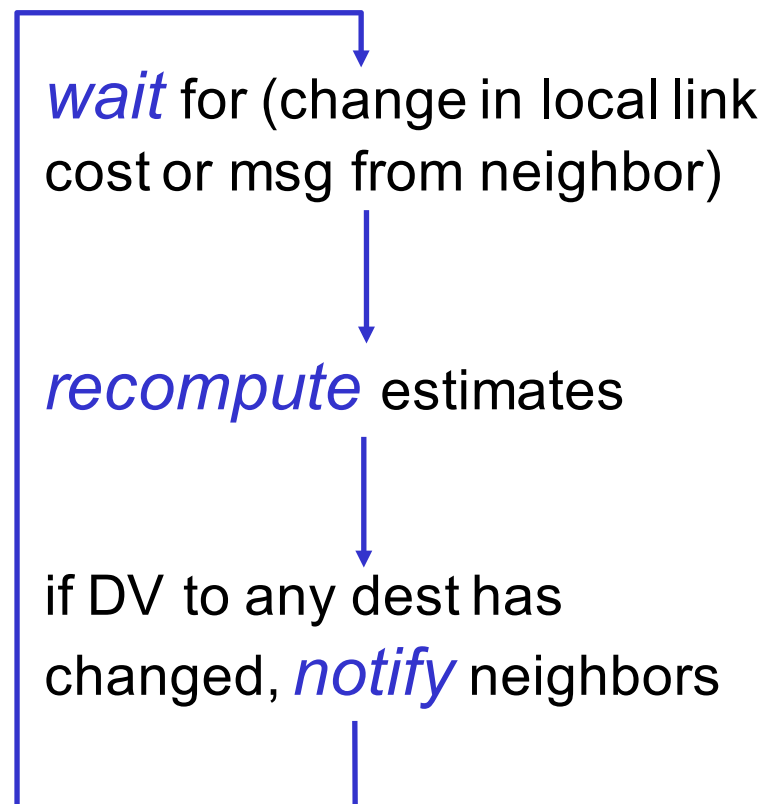
each local iteration caused by:

- ❑ 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:



Distance Vector Algorithm:

At each node, x :

- 1 **Initialization:**
- 2 for all destination y :
- 3 $d_x(y) = c(x,y)$
- 4 send D_x to all neighbor

Distance Vector Algorithm (cont.):

```
5 loop
6   wait (until I see a link cost change to a neighbor
7     or until I receive update from neighbor)
8
9   for each destination y
10     $d_x(y) = \min \{c(x,v)+d_v(y)\}$  /* for all v, the neighbors*/
11
12   if  $D_x$  changed, send  $D_x$  to all neighbors
13
14 forever
```

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

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

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

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

node x table

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

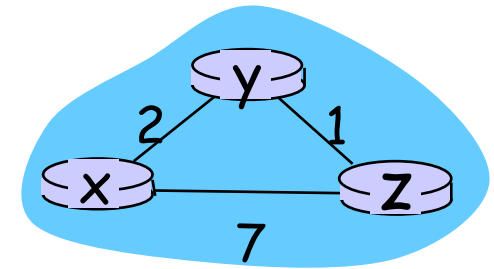
		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	7	1	0

node y table

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

node z table

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



.....→ time

Quiz 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	∞	∞	∞
	z	∞	∞	∞

node y table

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

node z table

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

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

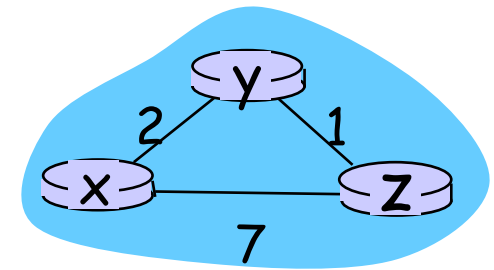
		cost to		
		x	y	z
from	x	0	2	7
	y	2	0	1
	z	7	1	0

		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

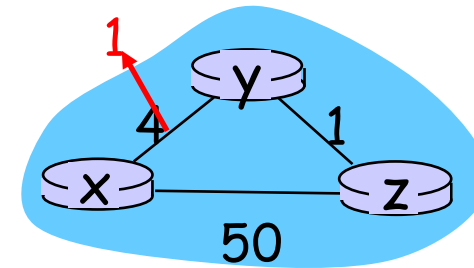


time →

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.

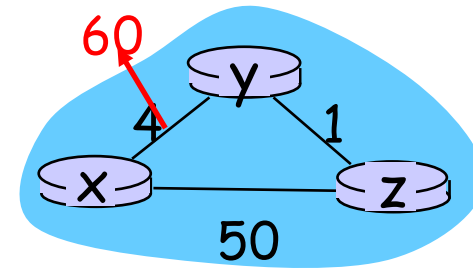
At time t_1 , z receives the update from y and updates its table. It computes a new least cost to x and sends its neighbors its DV.

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 [quiz!](#)



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 each
- ❑ DV: exchange between neighbors only

Speed of Convergence

- ❑ LS: $O(n^2)$ algorithm requires $O(nE)$ msgs
 - may have oscillations
- ❑ DV: convergence time varies
 - may have 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 dest's in routing tables!
- ❑ routing table exchange would swamp links!

administrative autonomy

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

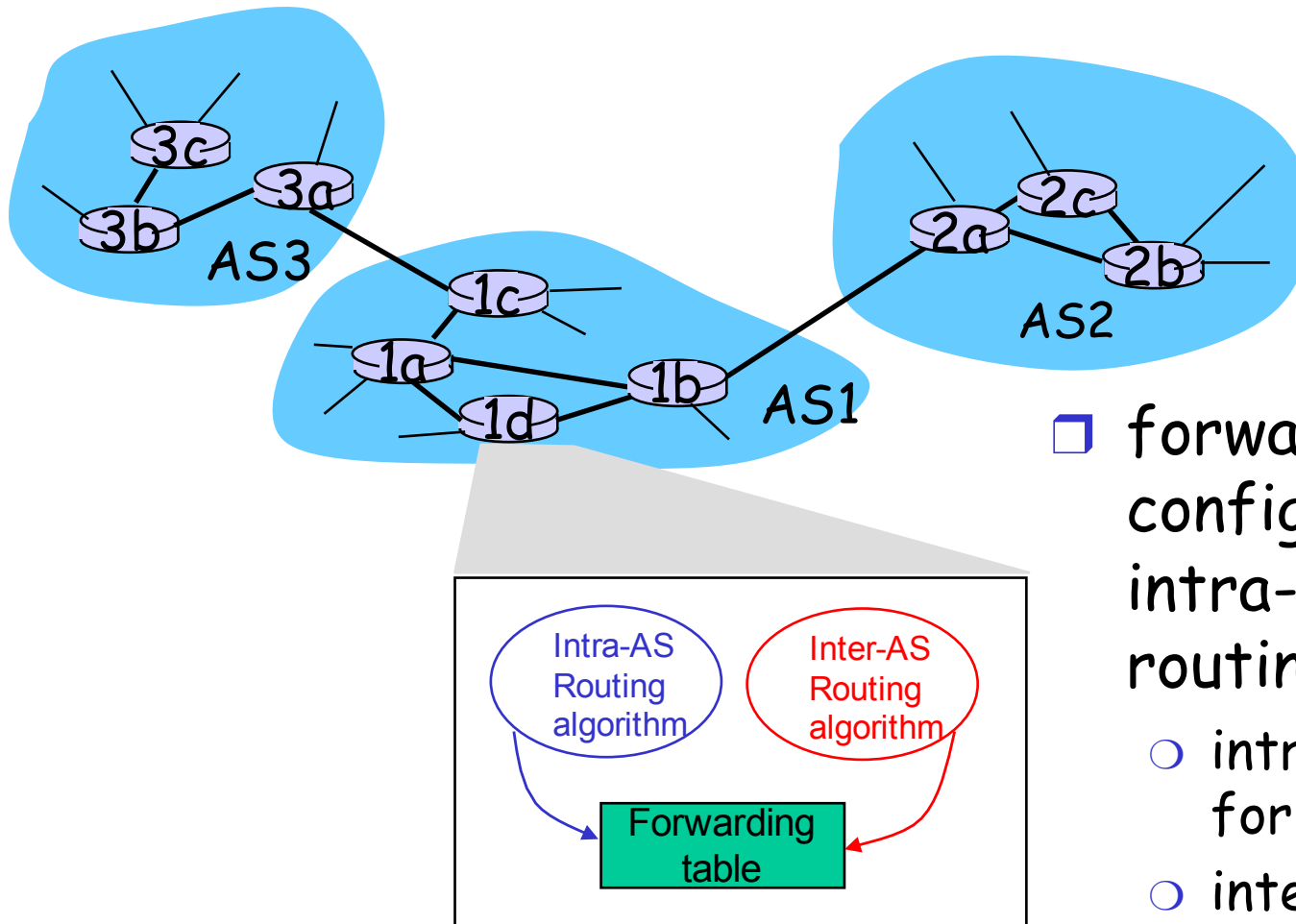
Hierarchical Routing

- ❑ aggregate routers into regions, "autonomous systems" (AS)
- ❑ routers in same AS run same routing protocol
 - "intra-AS" routing protocol
 - routers in different AS can run different intra-AS routing protocol

gateway routers

- ❑ special routers in AS
- ❑ run intra-AS routing protocol with all other routers in AS
- ❑ also responsible for routing to destinations outside AS
 - run *inter-AS routing* protocol with other gateway routers

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

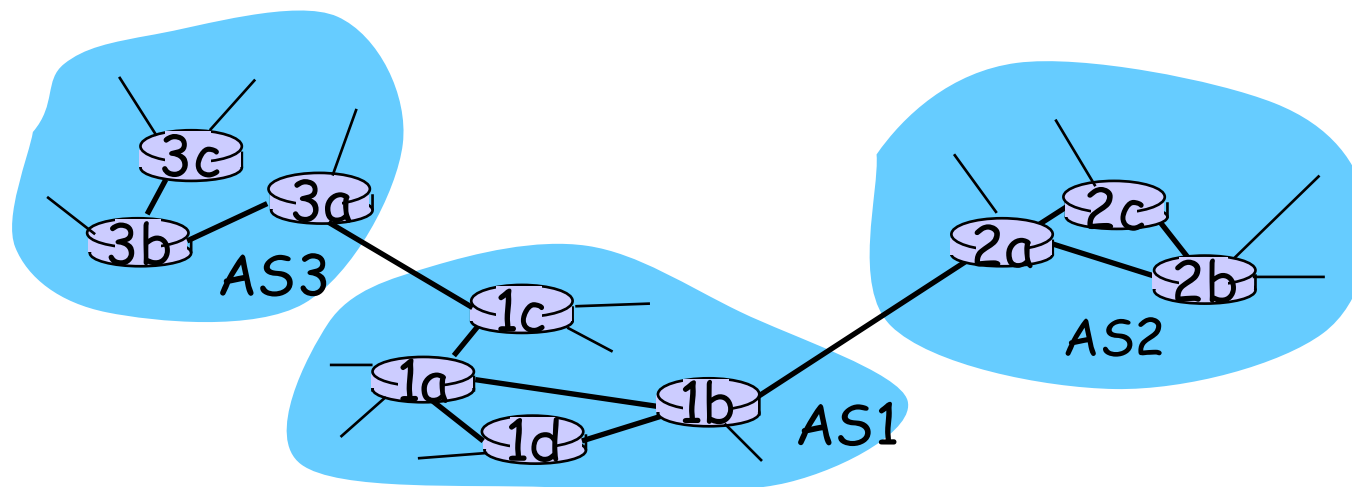
Inter-AS tasks

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

AS1 must:

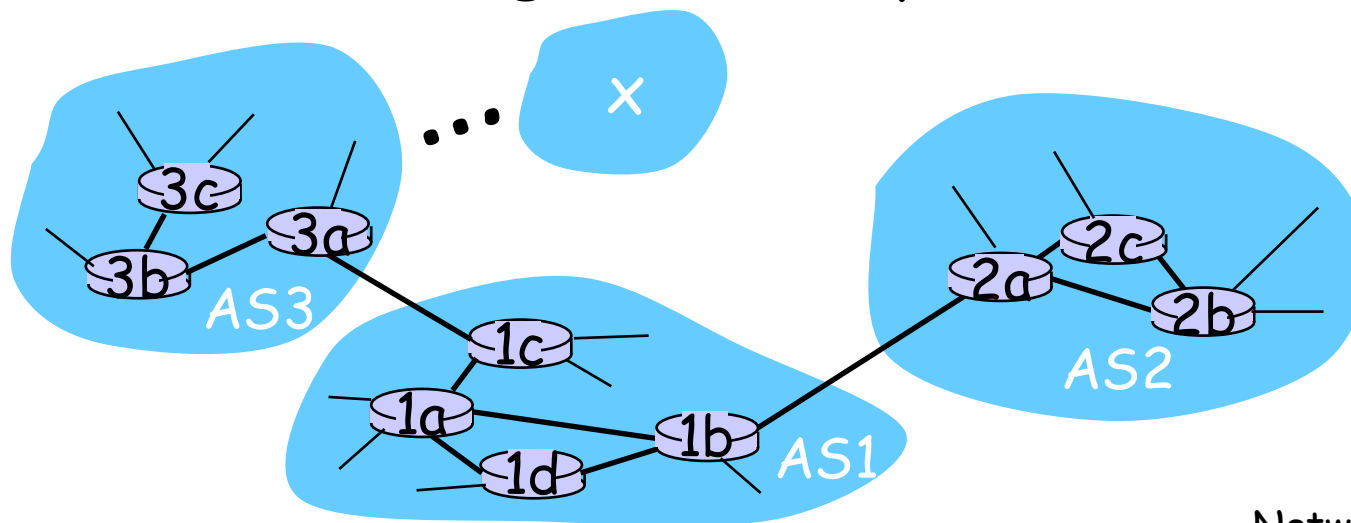
1. learn which dests reachable through AS2, which through AS3
2. propagate this reachability info to all routers in AS1

Job of inter-AS routing!



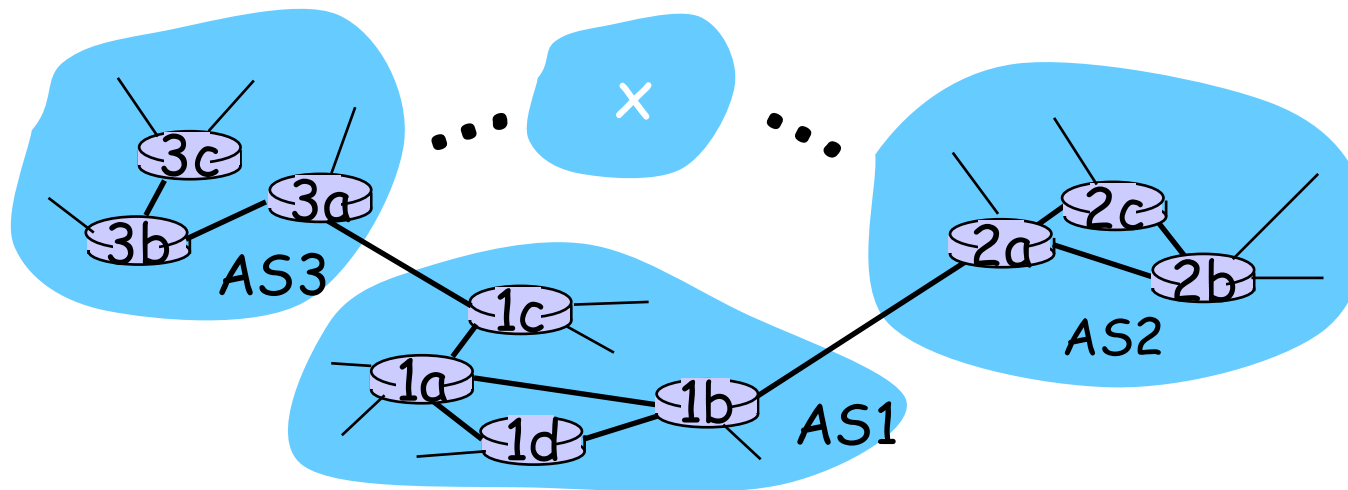
Example: Setting forwarding table in router 1d

- 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)



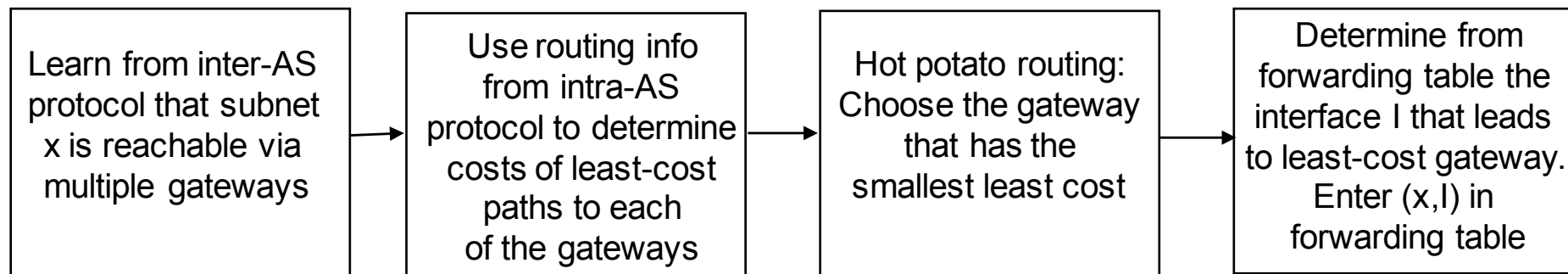
Example: Choosing among multiple ASes

- 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 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

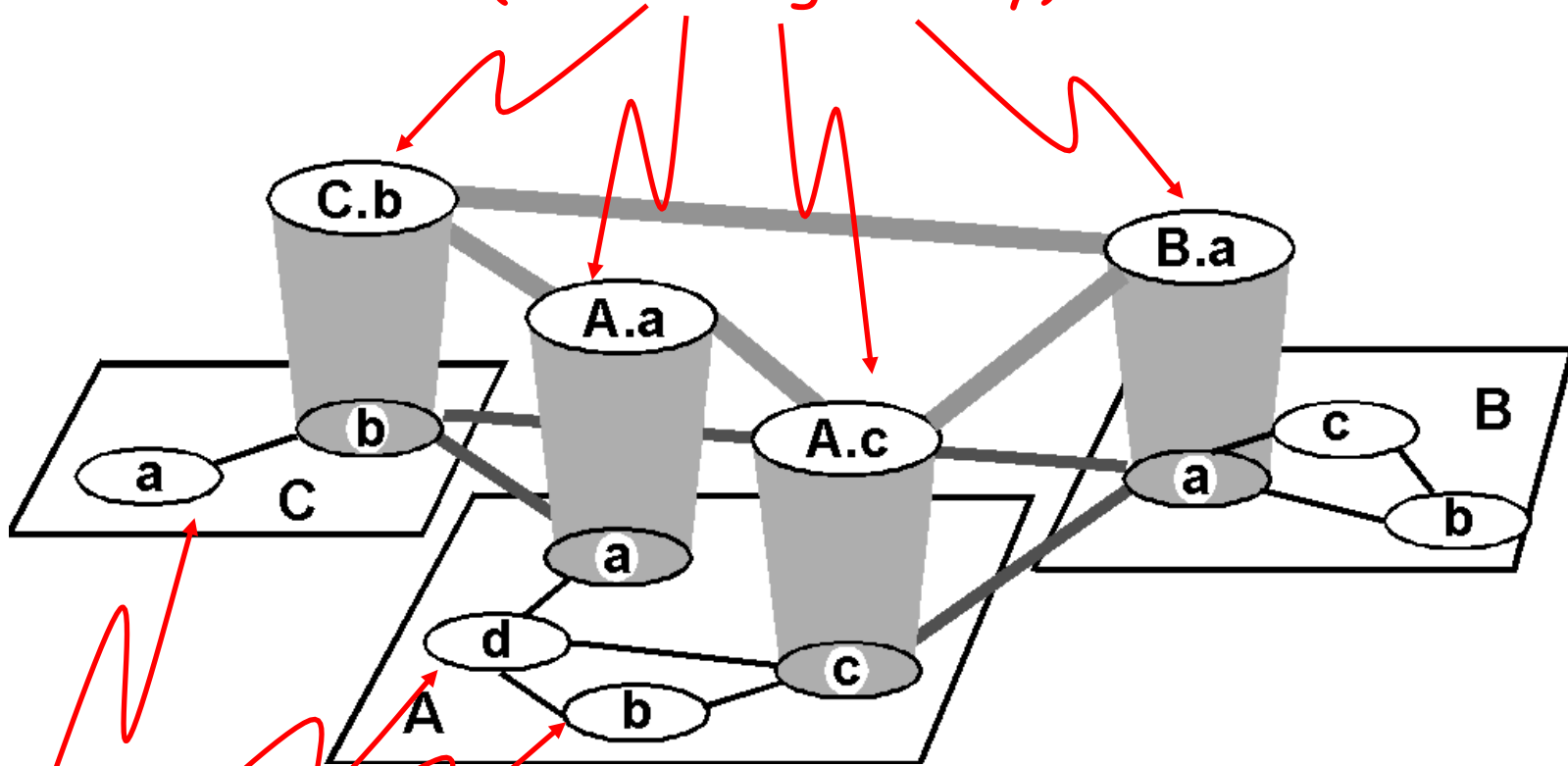
Routing in the Internet

- ❑ The Global Internet consists of **Autonomous Systems (AS)** interconnected with each other:
 - **Stub AS**: small corporation: one connection to other AS's
 - **Multihomed AS**: large corporation (no transit): multiple connections to other AS's
 - **Transit AS**: provider, hooking many AS's together

- ❑ Two-level routing:
 - **Intra-AS**: administrator responsible for choice of routing algorithm within network
 - **Inter-AS**: unique standard for inter-AS routing: BGP

Internet AS Hierarchy

Inter-AS border (exterior gateway) routers



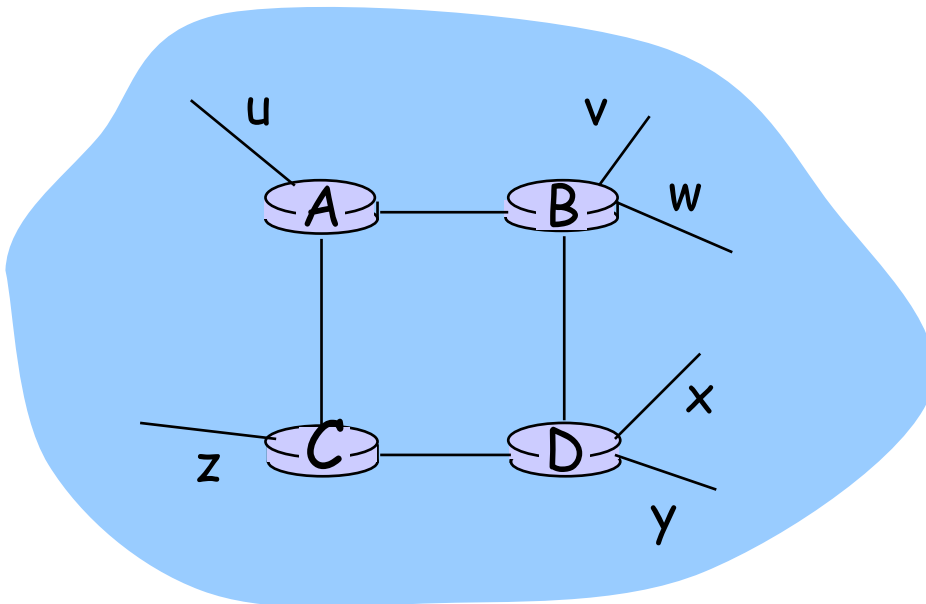
Intra-AS interior routers

Intra-AS Routing

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

RIP (Routing Information Protocol)

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



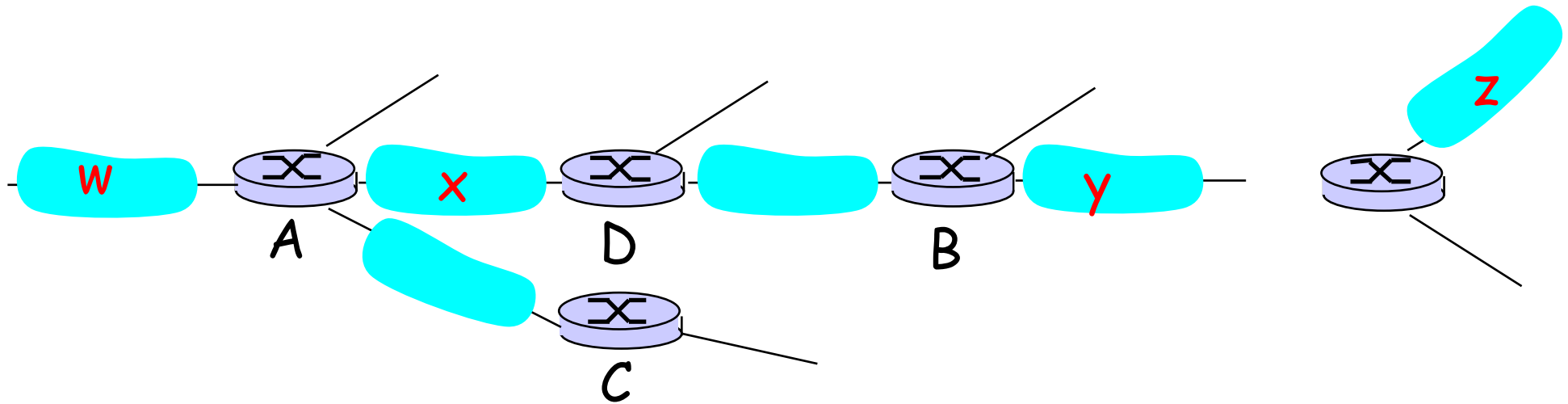
From router A to subsets:

<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 nets within AS

RIP: Example



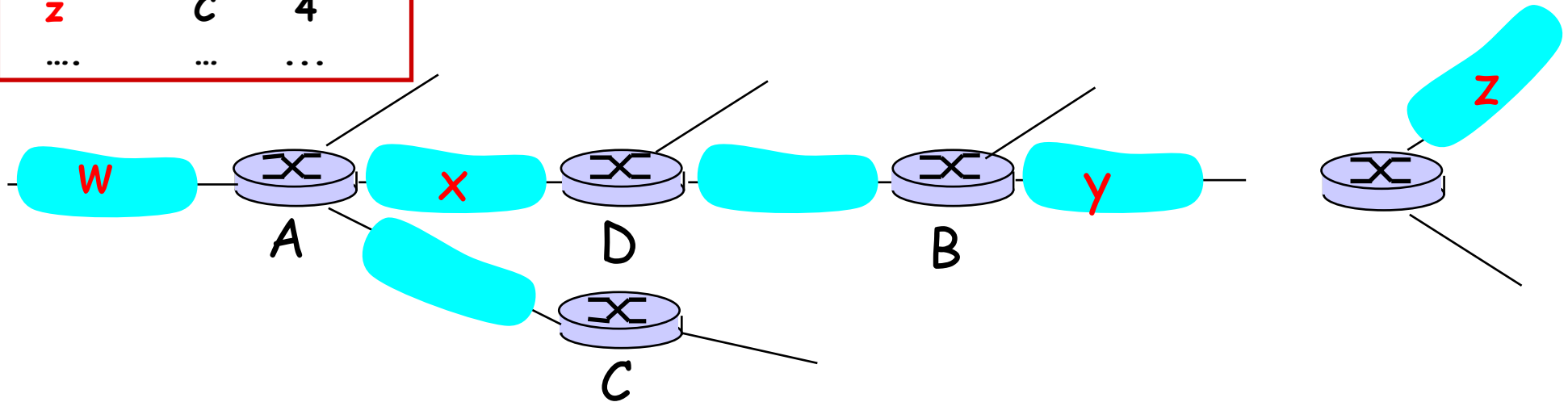
Destination Network	Next Router	Num. of hops to dest.
W	A	2
Y	B	2
Z	B	7
X	--	1
....

Routing table in D

RIP: Example

Dest	Next	hops
W	-	-
X	-	-
Z	C	4
...

Advertisement from A to D



Destination Network	Next Router	Num. of hops to dest.
W	A	2
Y	B	2
Z	B A	7 5
X	--	1
...

Routing table in D

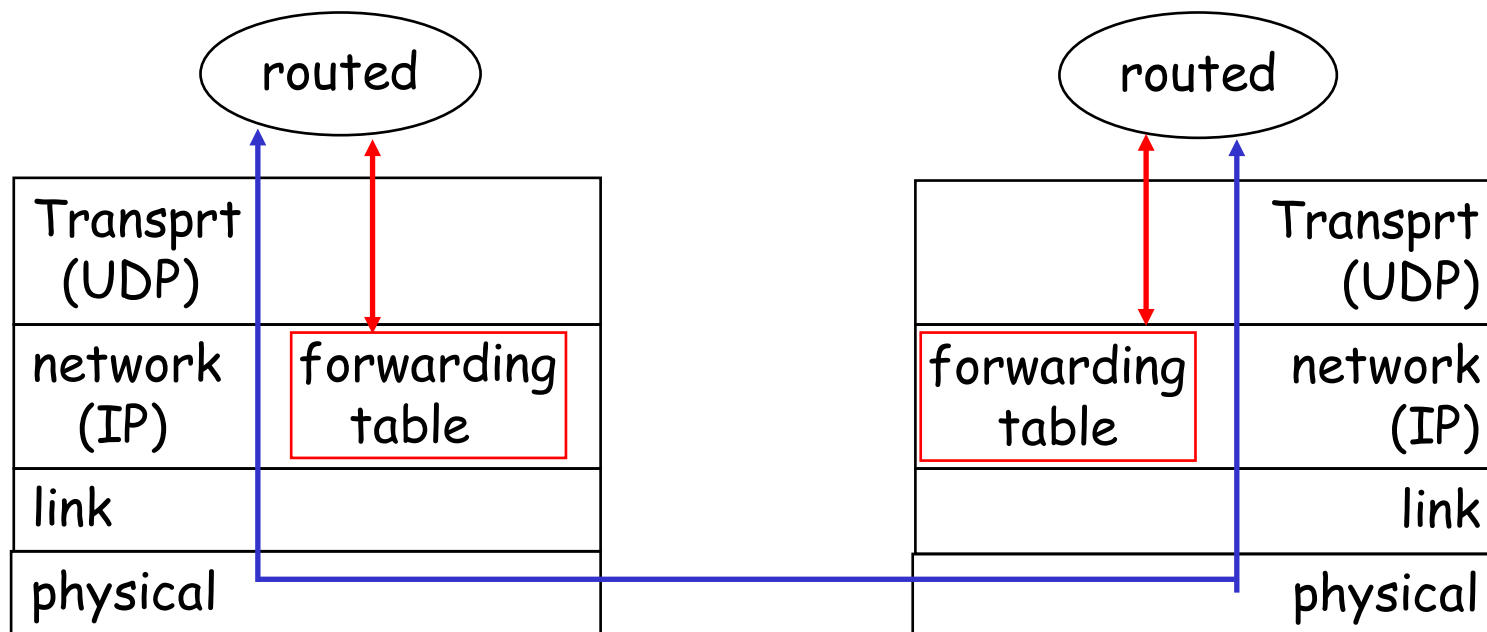
RIP: Link Failure and Recovery

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

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

RIP Table processing

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



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

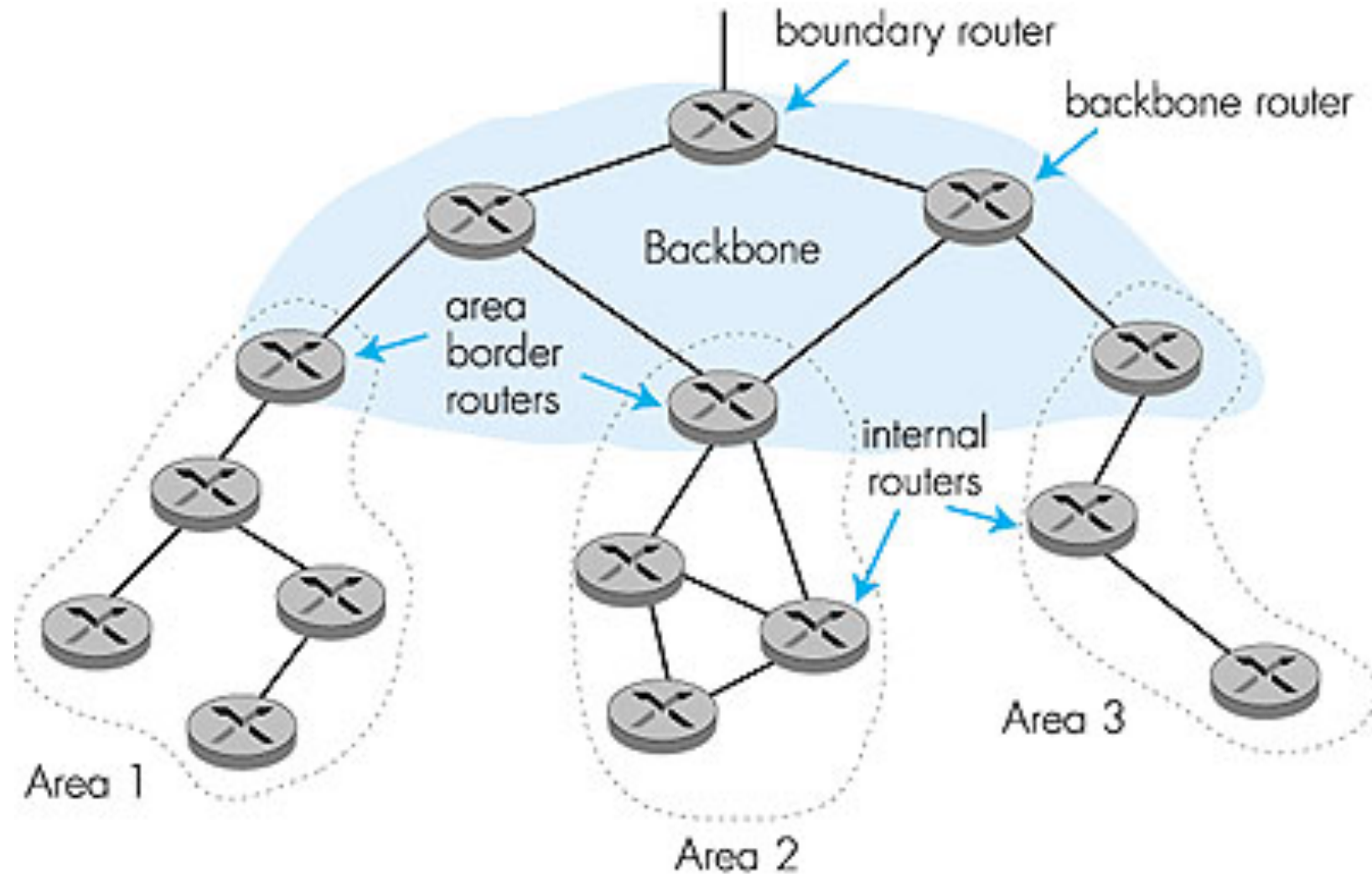
OSPF (Open Shortest Path First)

- ❑ “open”: publicly available
- ❑ 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



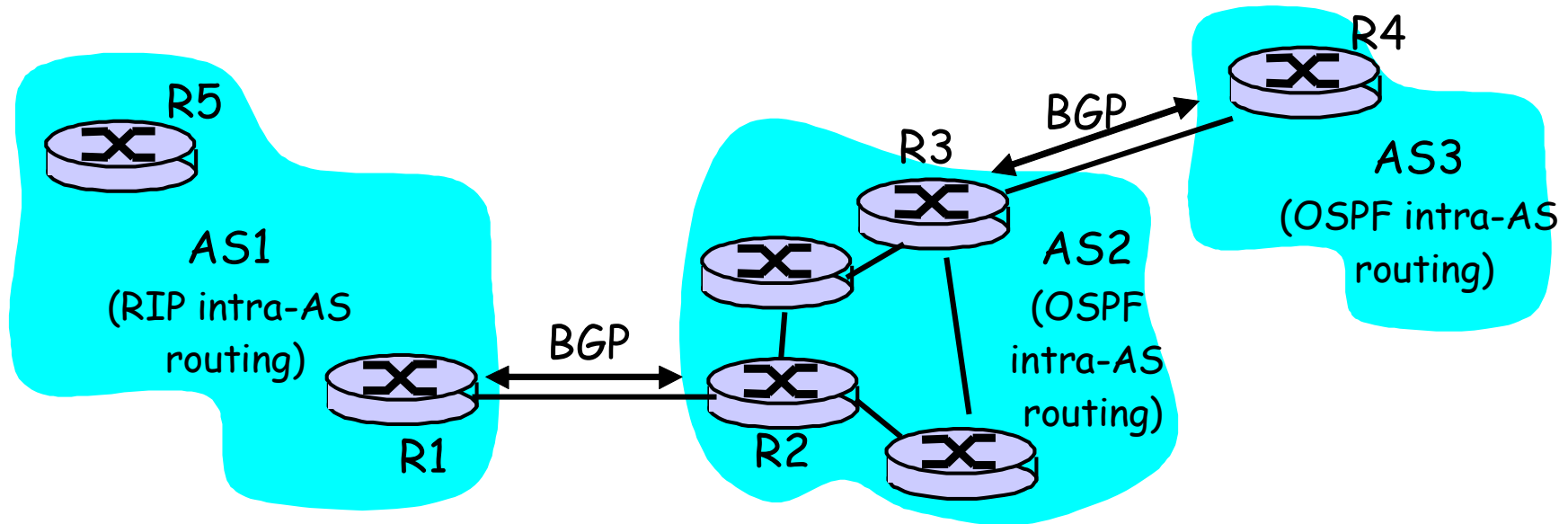
Hierarchical OSPF

- ❑ **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 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

Inter-AS routing in the Internet: BGP

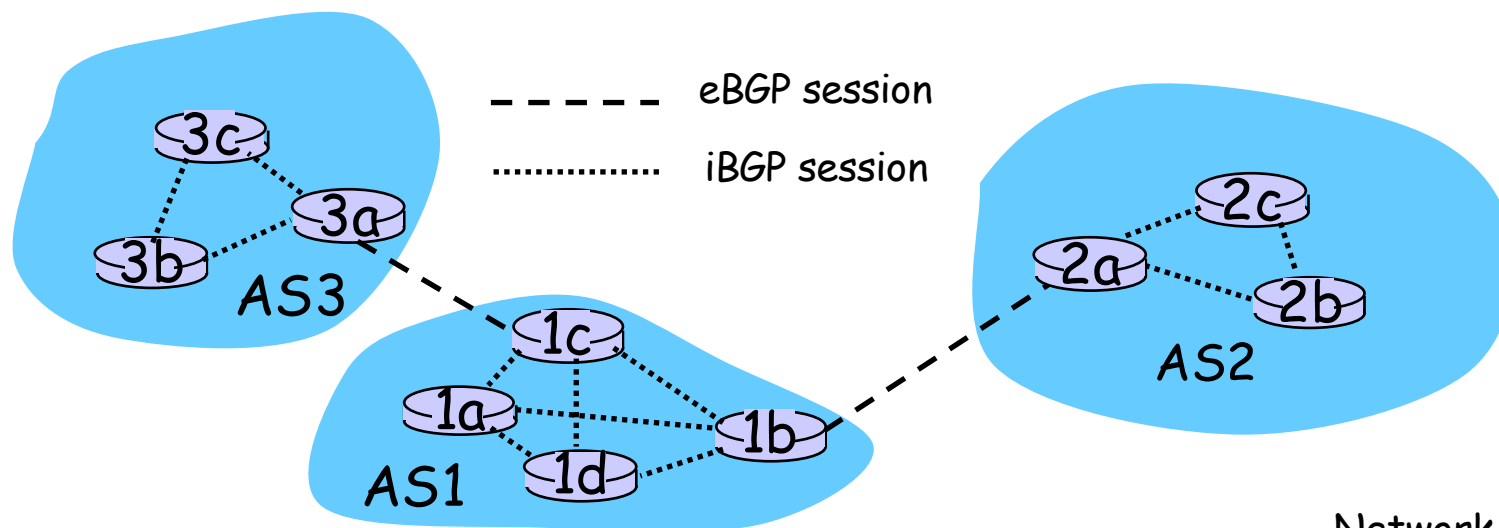


Internet inter-AS routing: BGP

- ❑ **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"*

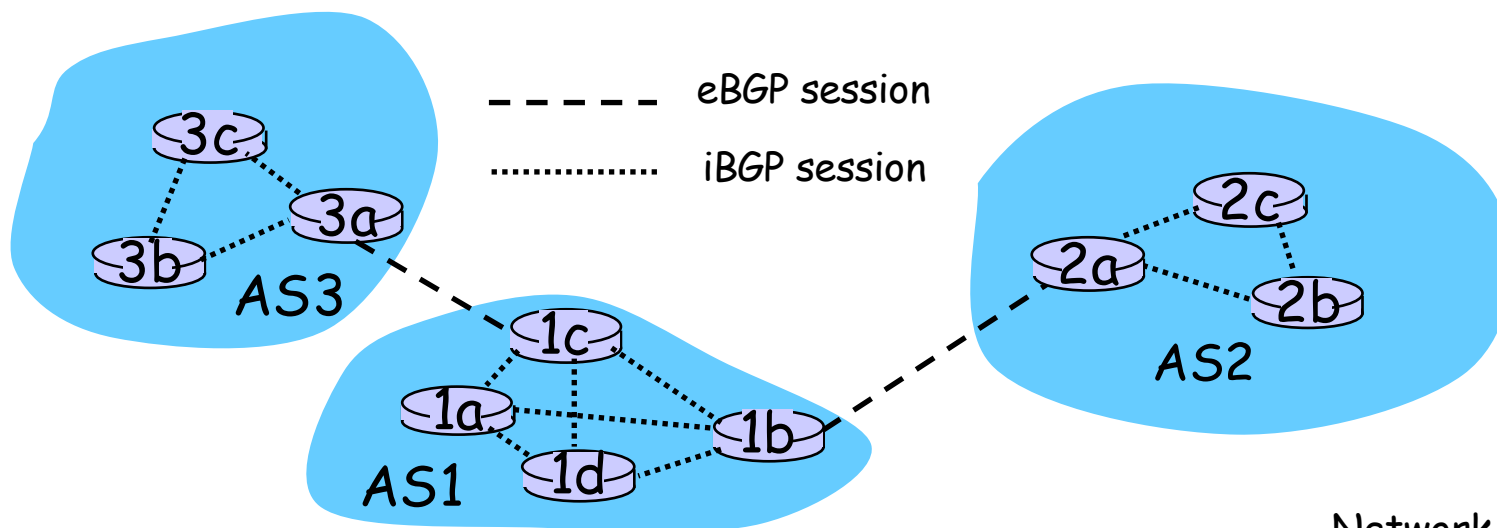
BGP basics

- 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 prefix to AS1:
 - AS2 **promises** it will forward any addresses 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, creates entry for prefix in its forwarding table.



Path attributes & BGP routes

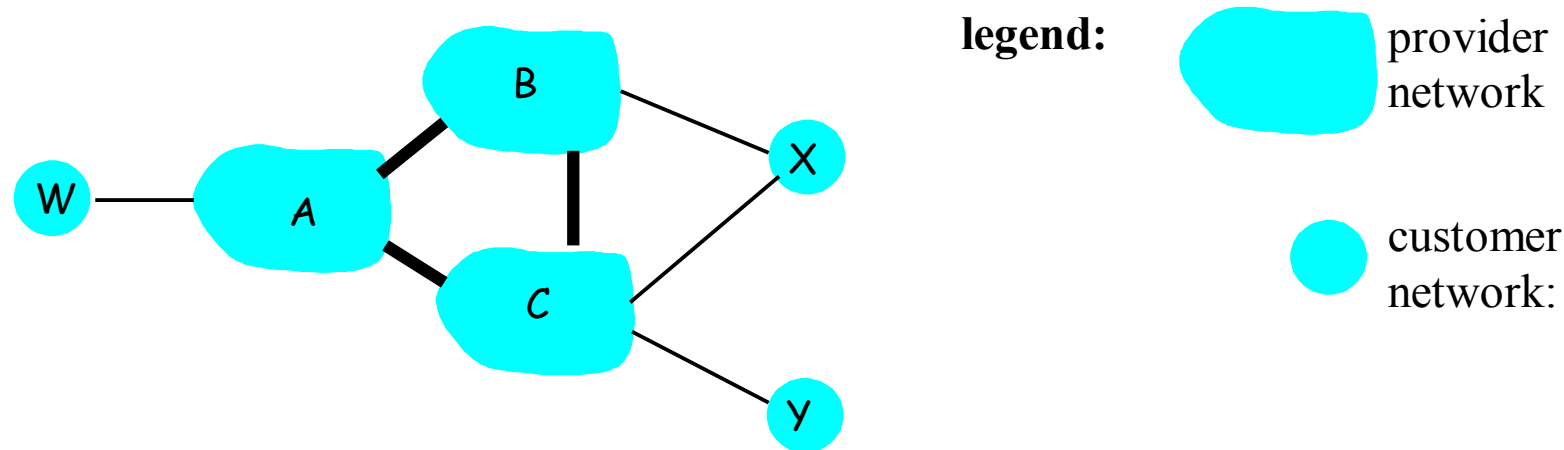
- ❑ advertised prefix includes BGP attributes.
 - prefix + attributes = "route"
- ❑ two important attributes:
 - **AS-PATH**: contains ASs through which prefix advertisement has passed: e.g, AS 67, AS 17
 - **NEXT-HOP**: indicates specific internal-AS router to next-hop AS. (may be multiple links from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

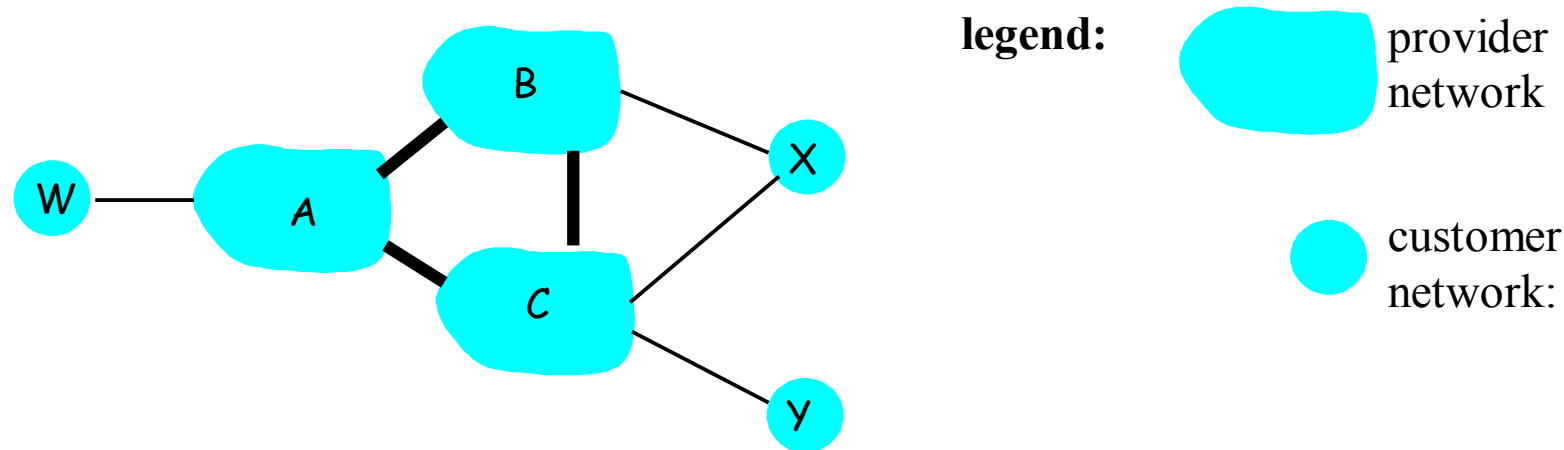
Quiz Time!

BGP: controlling who routes to you



- 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: controlling who routes to you



- ❑ 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!

Triple Quiz Time!

BGP operation

Q: What does a BGP router do?

- ❑ Receiving and filtering route advertisements from directly attached neighbor(s).
- ❑ Route selection.
 - To route to destination X, which path (of several advertised) will be taken?
- ❑ Sending route advertisements to neighbors.

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 messages; also used to close connection

Why different Intra- and Inter-AS routing ?

Policy:

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

Scale:

- ❑ hierarchical routing saves table size, reduced update traffic

Performance:

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

Chapter 4: summary

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