

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:

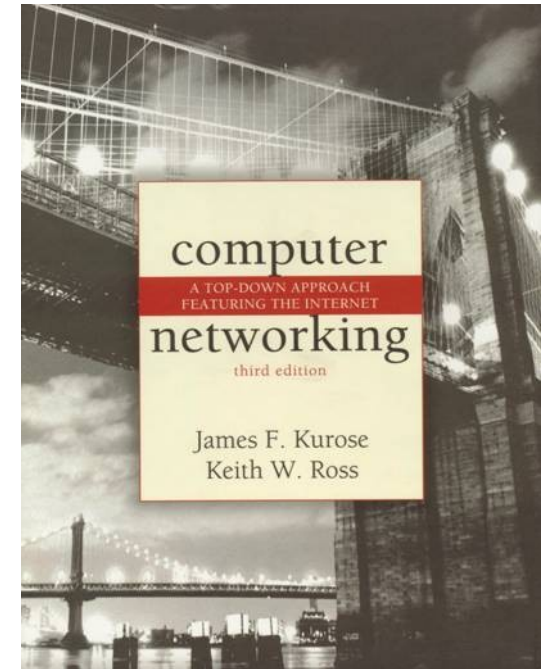
If you use these slides (e.g., in a class) in substantially unaltered form, that you mention their source (after all, we'd like people to use our book!)

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 1996-2004

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



*Computer Networking: A Top
Down Approach Featuring the
Internet,
3rd edition.*

*Jim Kurose, Keith Ross
Addison-Wesley, July 2004.*

Chapter 4: Network Layer

Chapter goals:

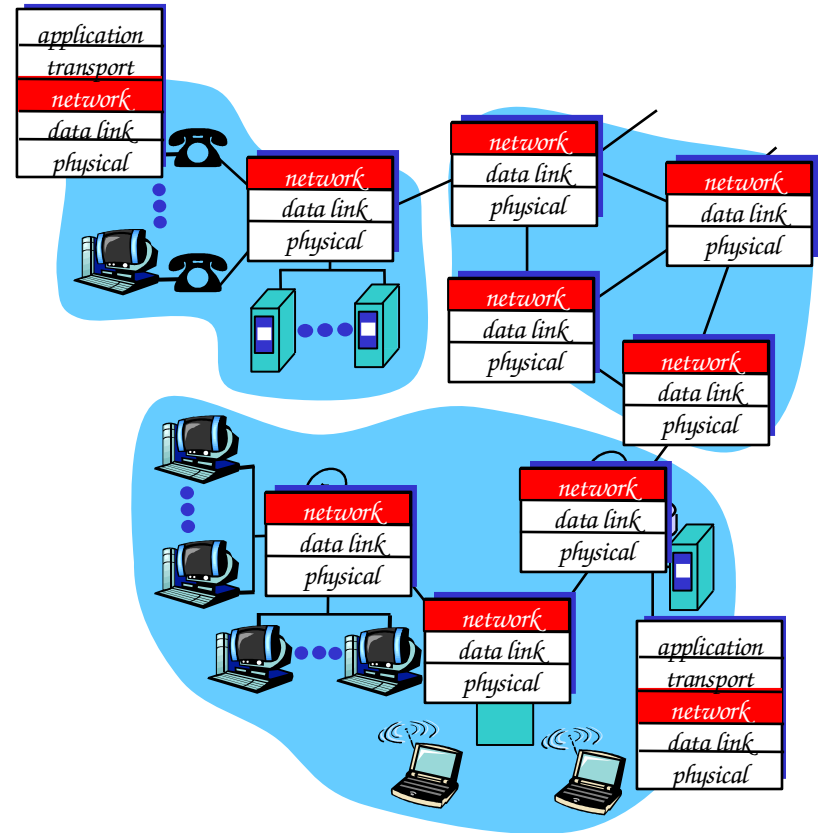
- *understand principles behind network layer services:*
 - *routing (path selection)*
 - *dealing with scale*
 - *how a router works*
 - *advanced topics: IPv6, mobility*
- *instantiation and implementation in the Internet*

Chapter 4: Network Layer

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



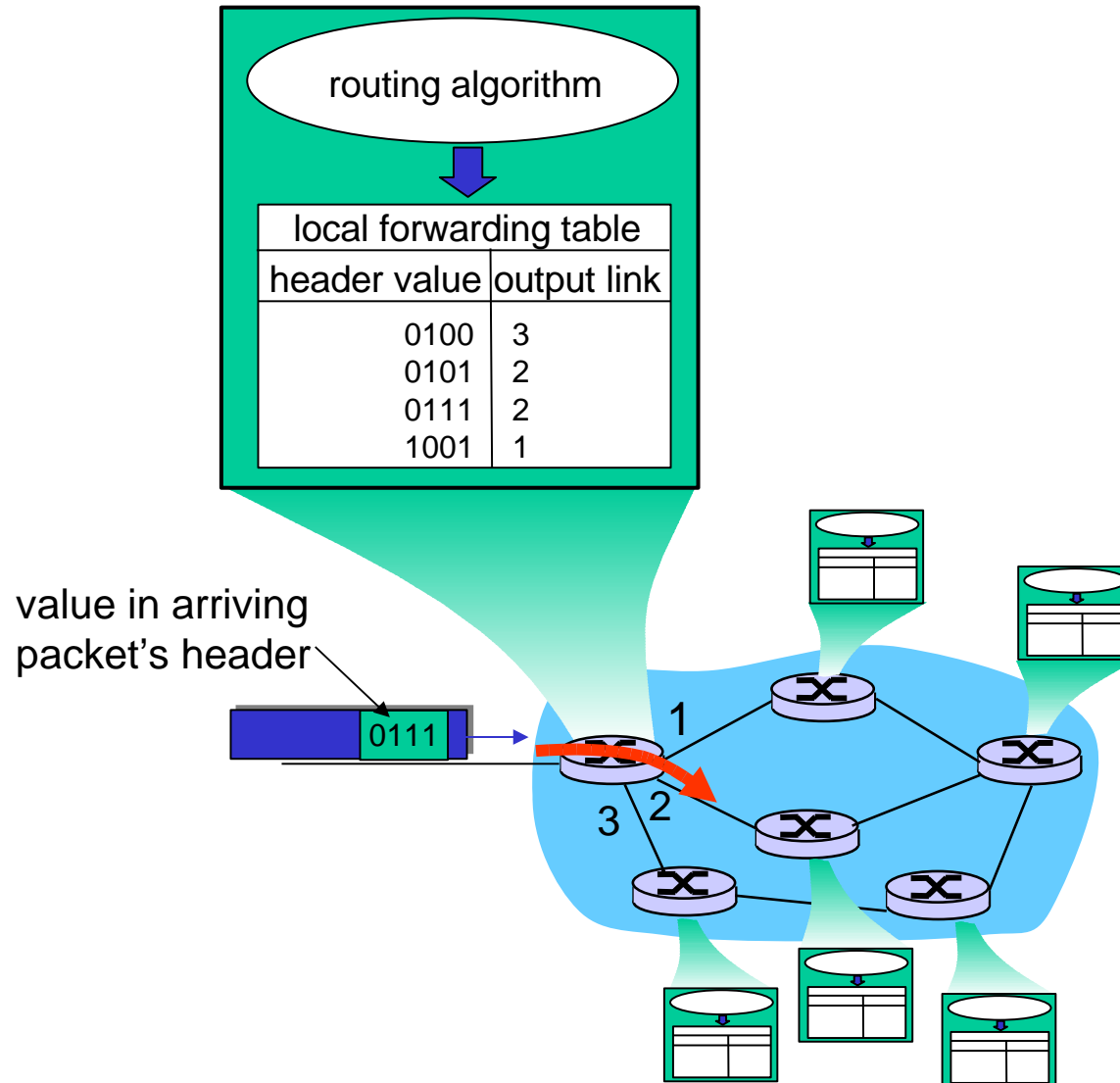
Key Network-Layer Functions

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

analogy:

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

Interplay between routing and forwarding



Connection setup

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

Network service model

*Q: What **service model** for “channel” transporting datagrams from sender to rcvr?*

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*
 - *No choice: network provides one or the other*
 - *Implementation: in the core*

Virtual circuits

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

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

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

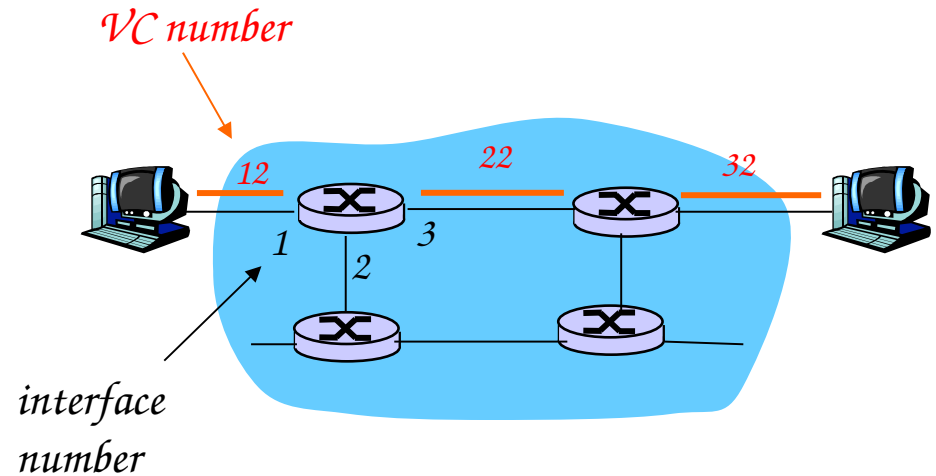
VC implementation

A VC consists of:

- 1. Path from source to destination*
 - 2. VC numbers, one number for each link along path*
 - 3. Entries in forwarding tables in routers along path*
- ☐ *Packet belonging to VC carries a VC number.*
 - ☐ *VC number must be changed on each link.*
 - ☐ *New VC number comes from forwarding table*

Forwarding table

Forwarding table in
northwest router:

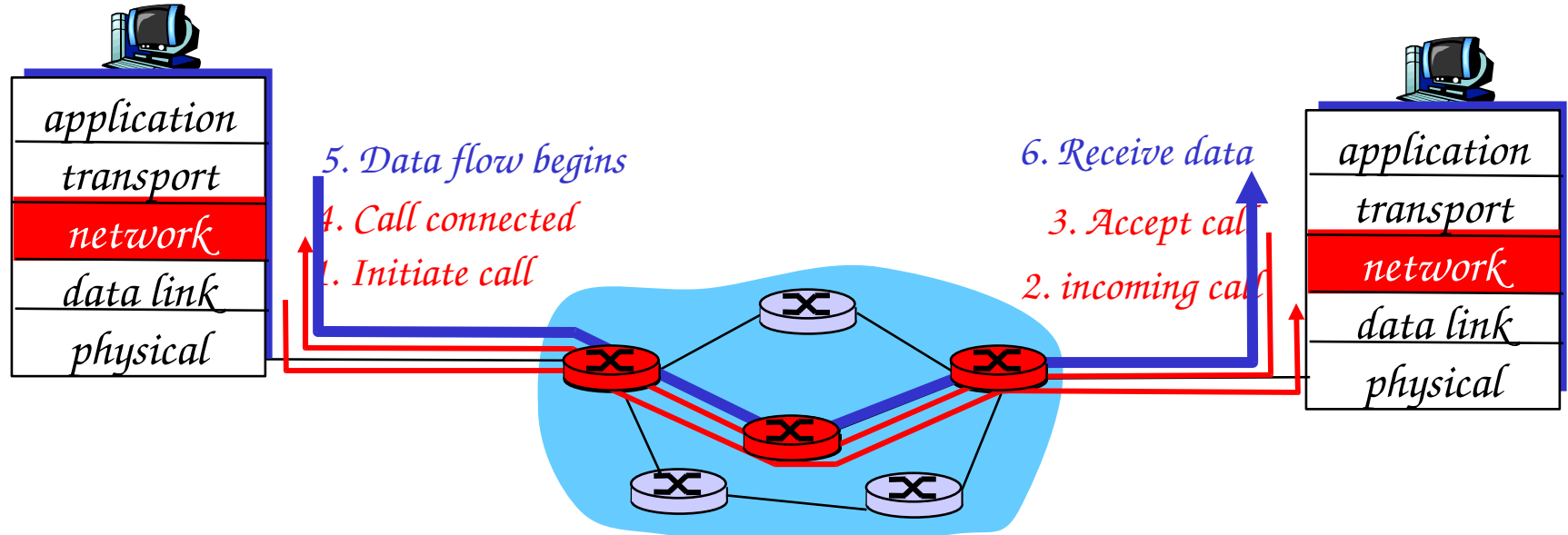


Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	3	22
2	63	1	18
3	7	2	17
1	97	3	87
...

Routers maintain connection state information!

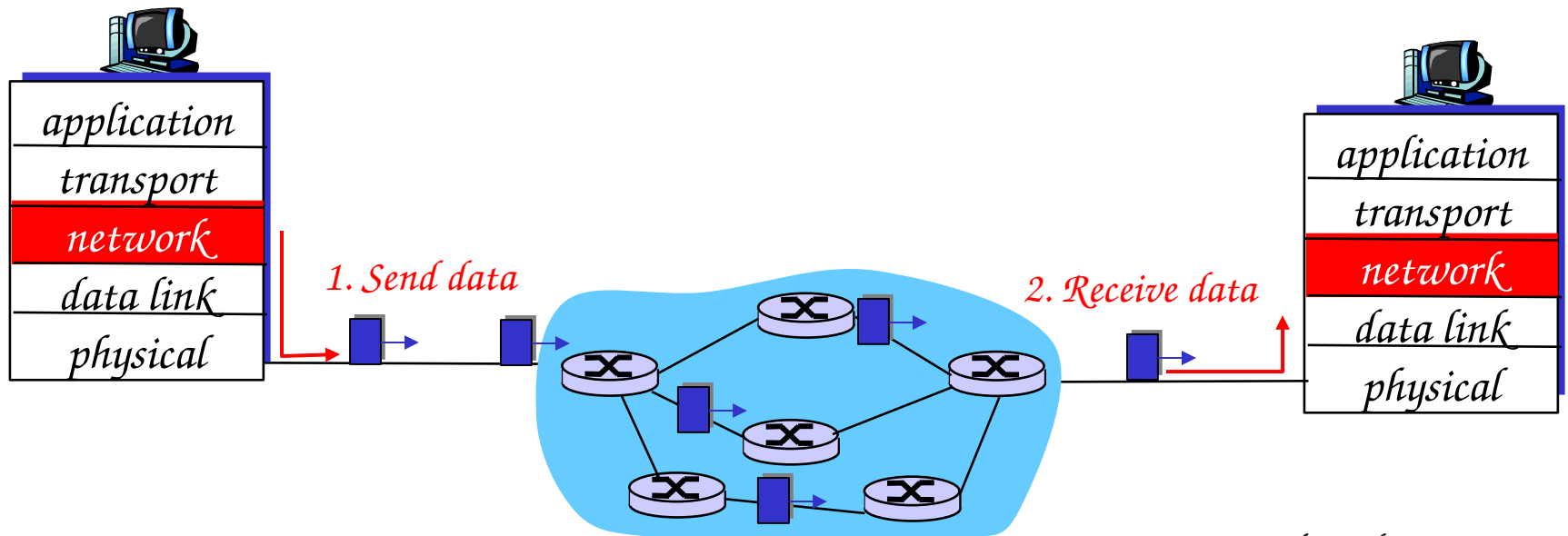
Virtual circuits: signaling protocols

- 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

- *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

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

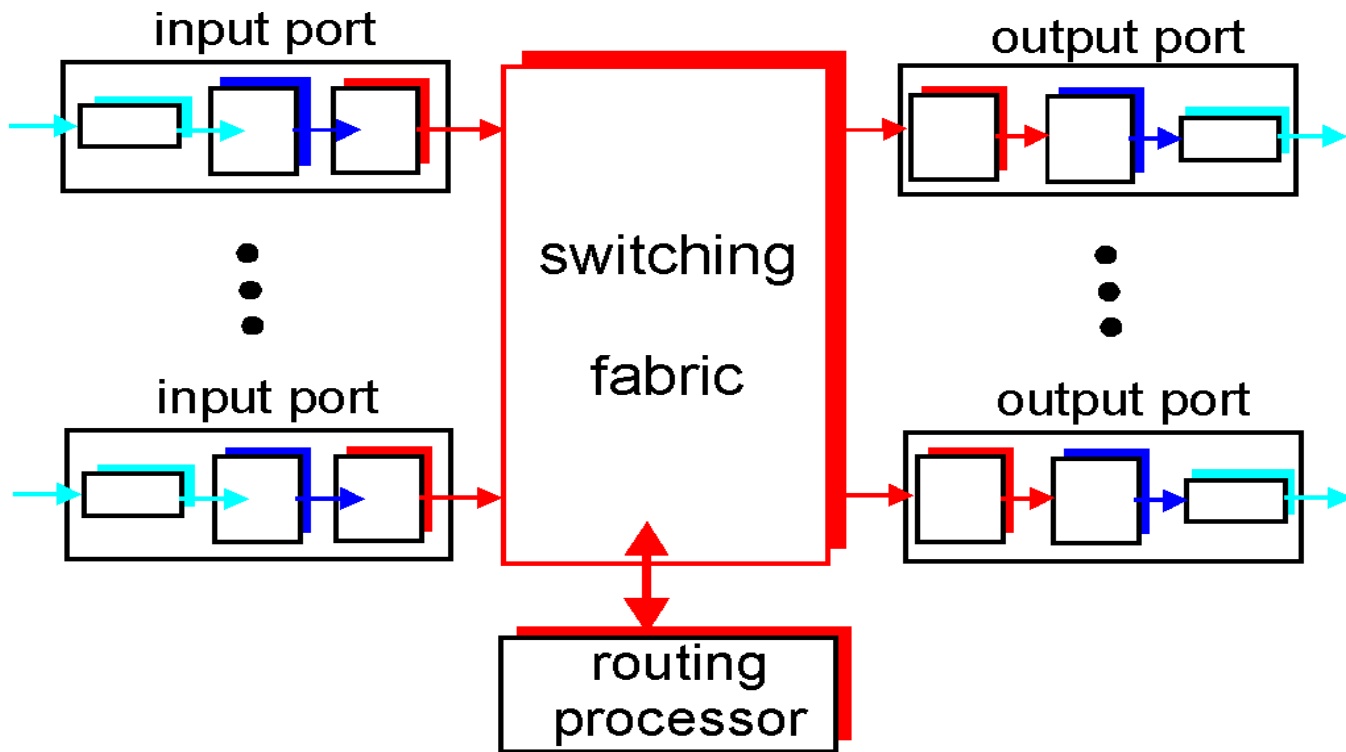
Chapter 4: Network Layer

- 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

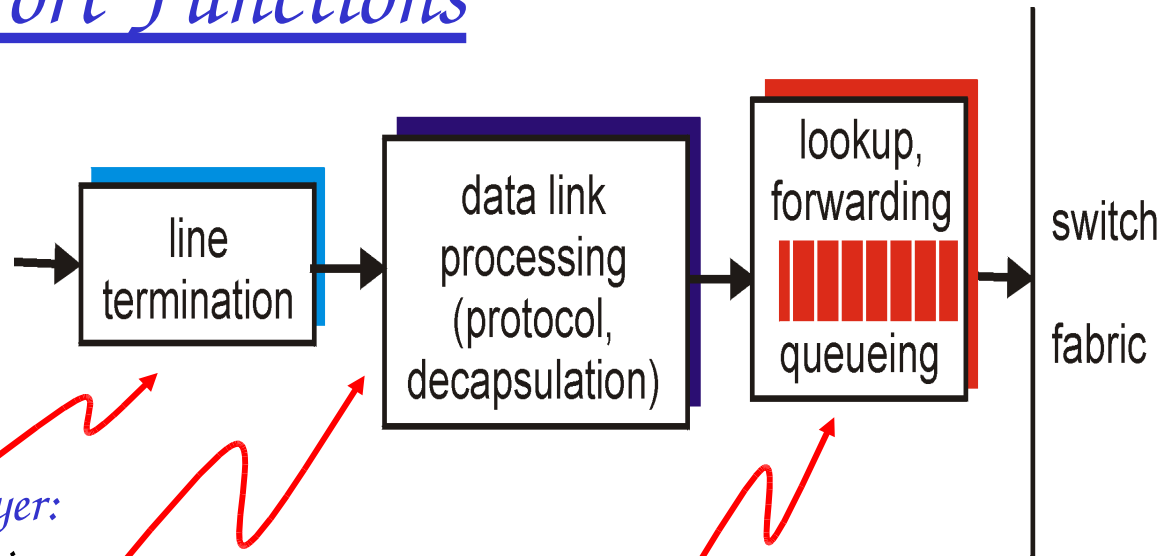
Router Architecture Overview

Two key router functions:

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



Input Port Functions



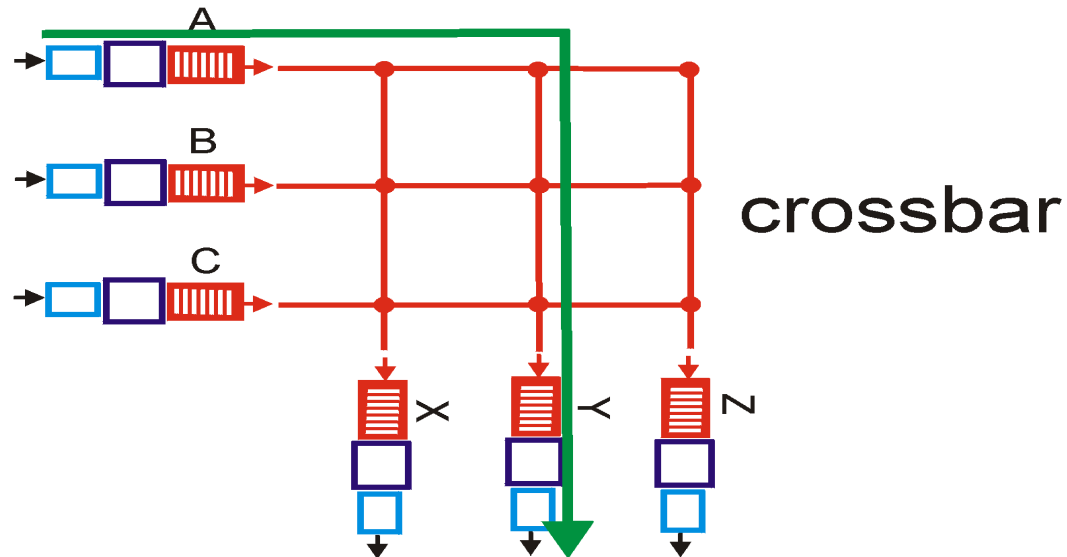
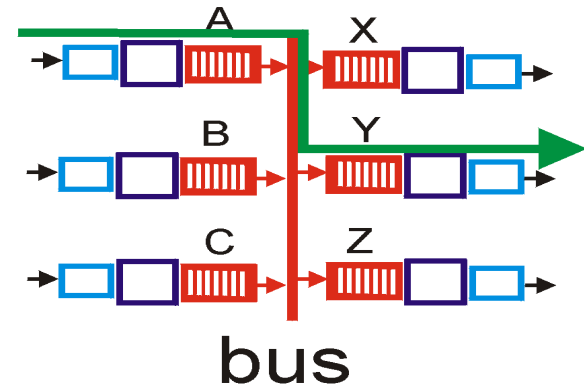
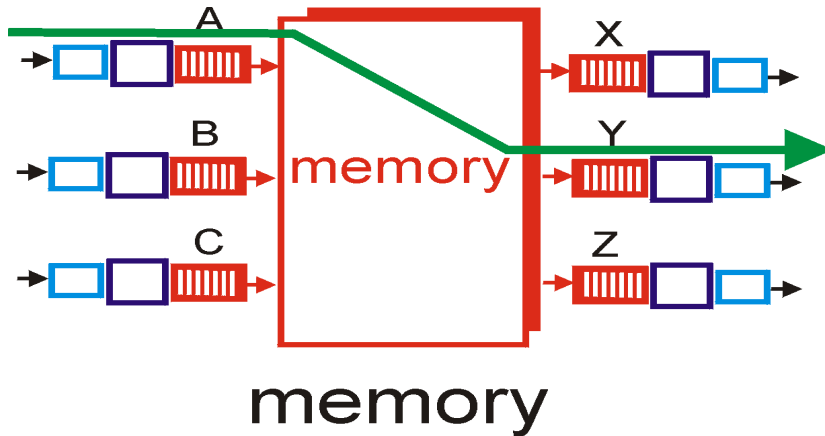
*Physical layer:
bit-level reception*

*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

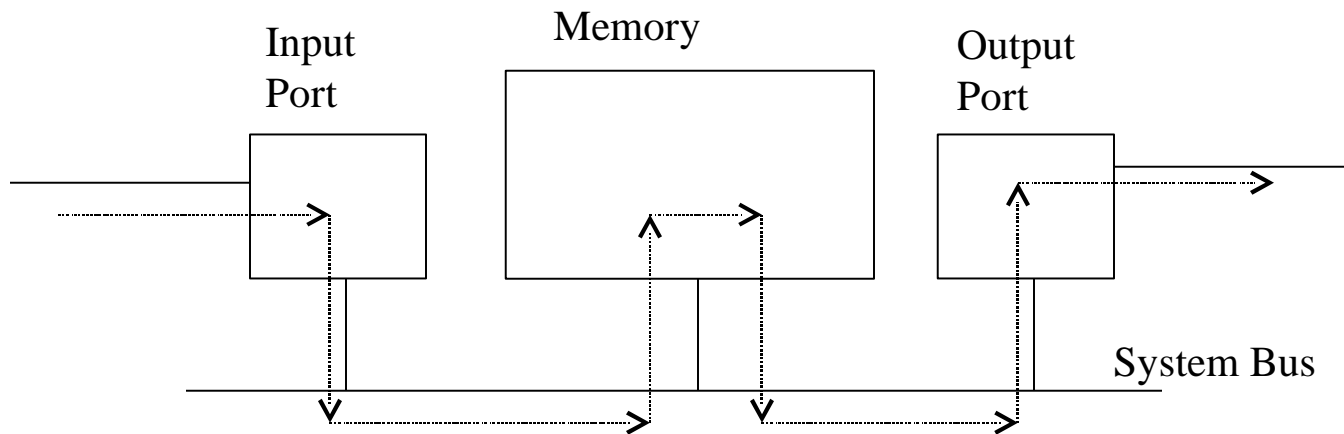
Three types of switching fabrics



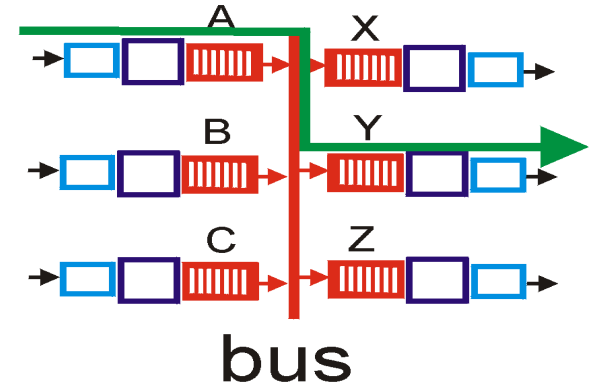
Switching Via Memory

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

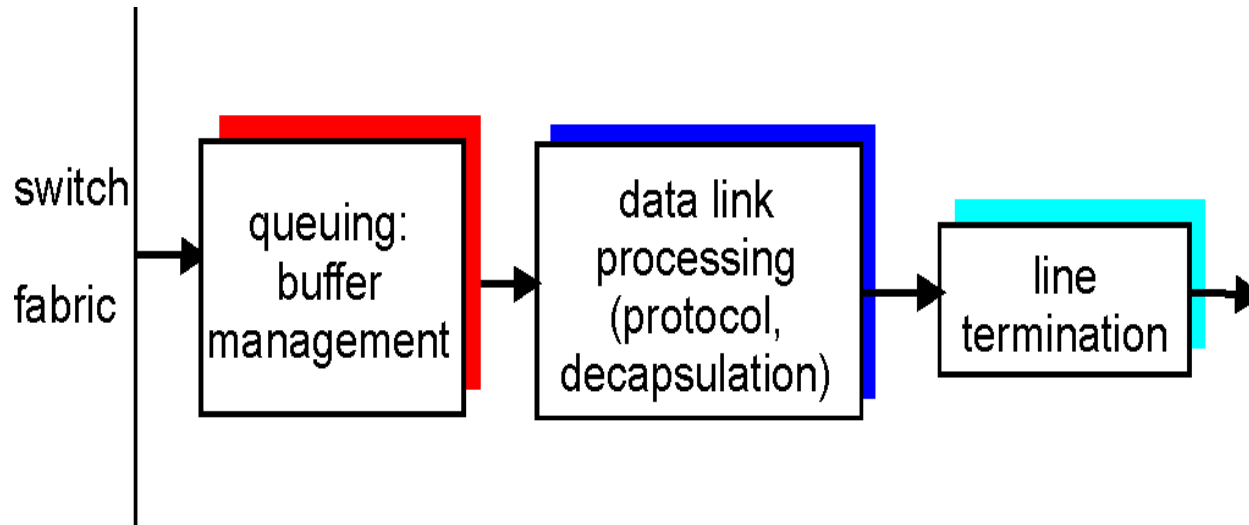


- ❑ *datagram from input port memory to output port memory via a shared bus*
- ❑ **bus contention:** *switching speed limited by bus bandwidth*
- ❑ *1 Gbps bus, Cisco 1900: sufficient speed for access and enterprise routers (not regional or backbone)*

Switching Via An Interconnection Network

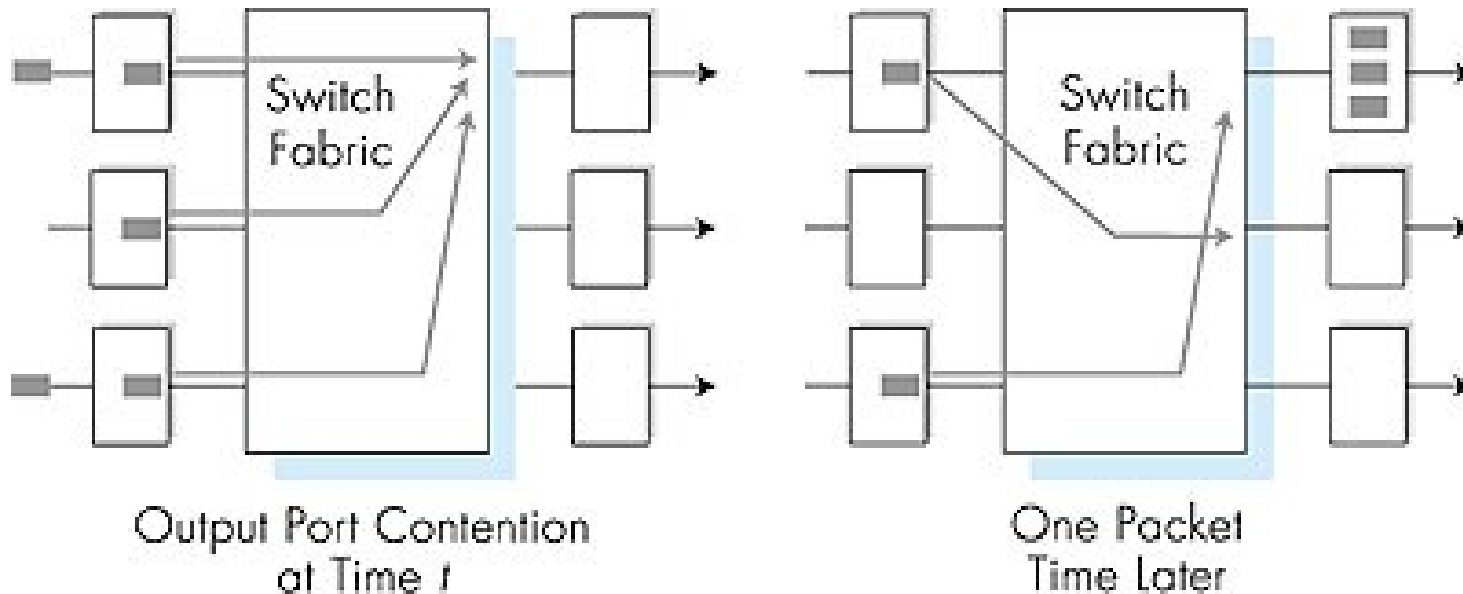
- ❑ *overcome bus bandwidth limitations*
- ❑ *Banyan networks, other interconnection nets initially developed to connect processors in multiprocessor*
- ❑ *Advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.*
- ❑ *Cisco 12000: switches Gbps through the interconnection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

Output Ports



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

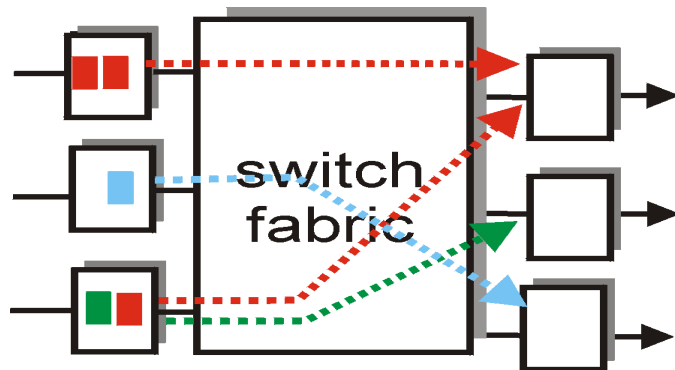
Output port queueing



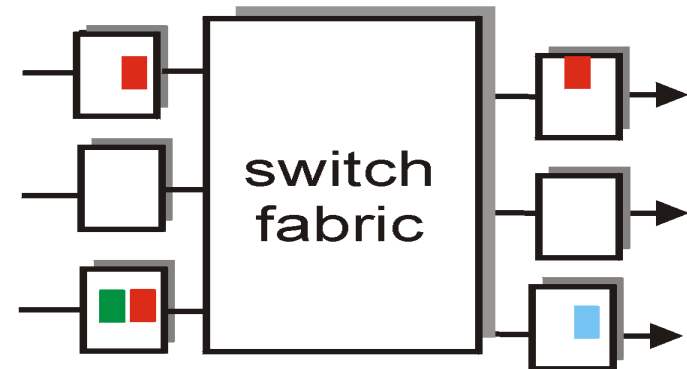
- *buffering when arrival rate via switch exceeds output line speed*
- *queueing (delay) and loss due to output port buffer overflow!*

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



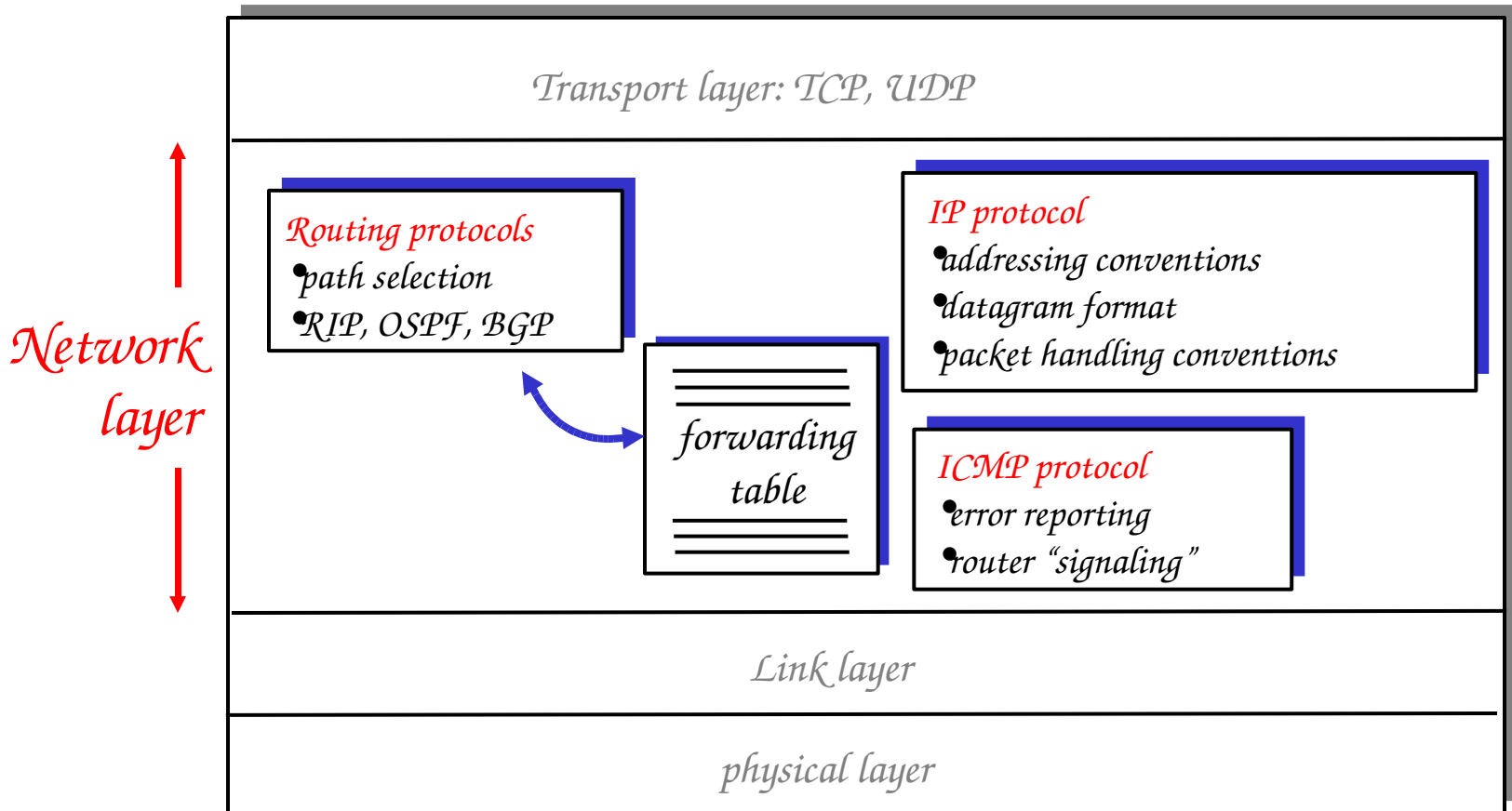
green packet
experiences HOL blocking

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

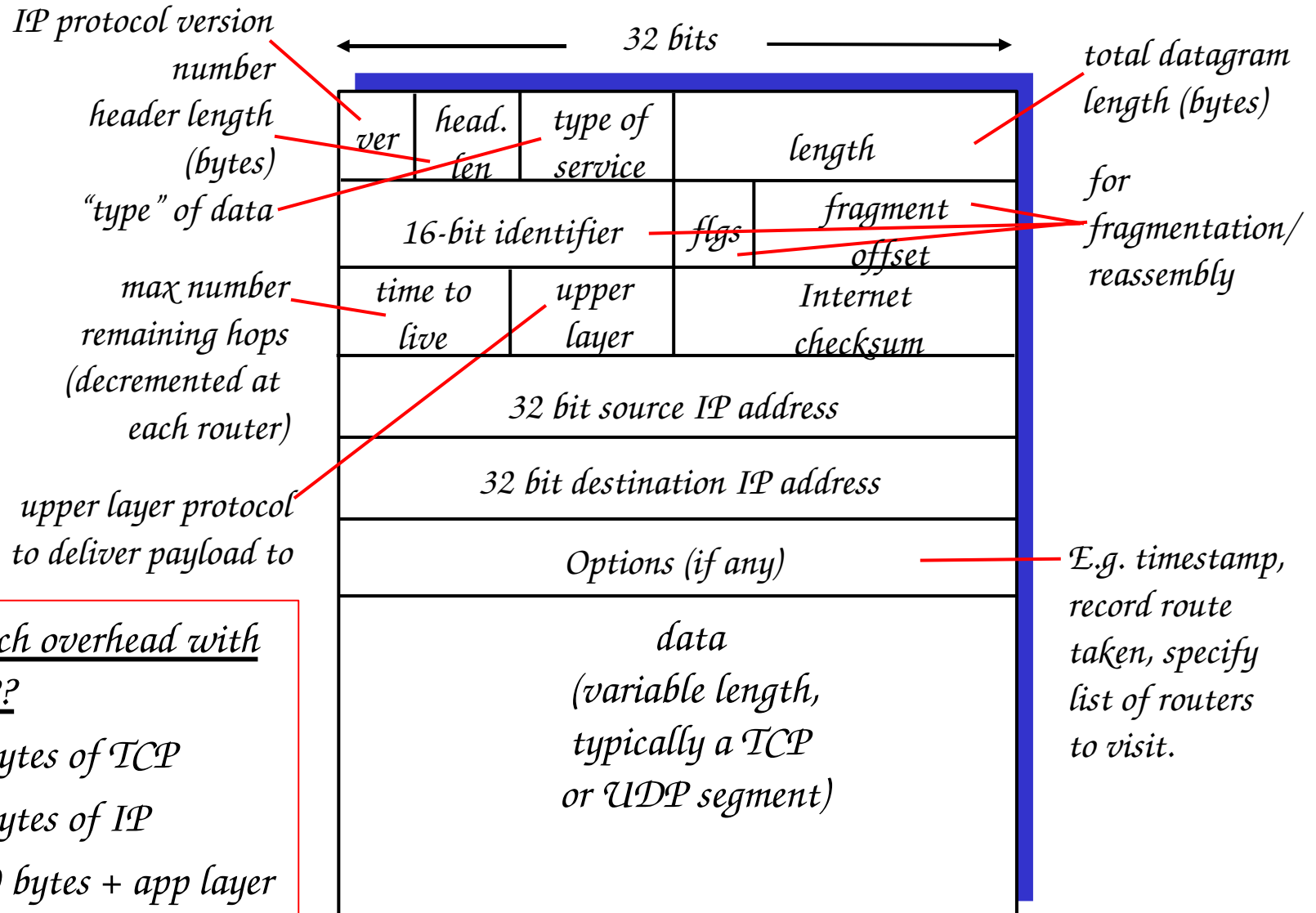
Host, router network layer functions:



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

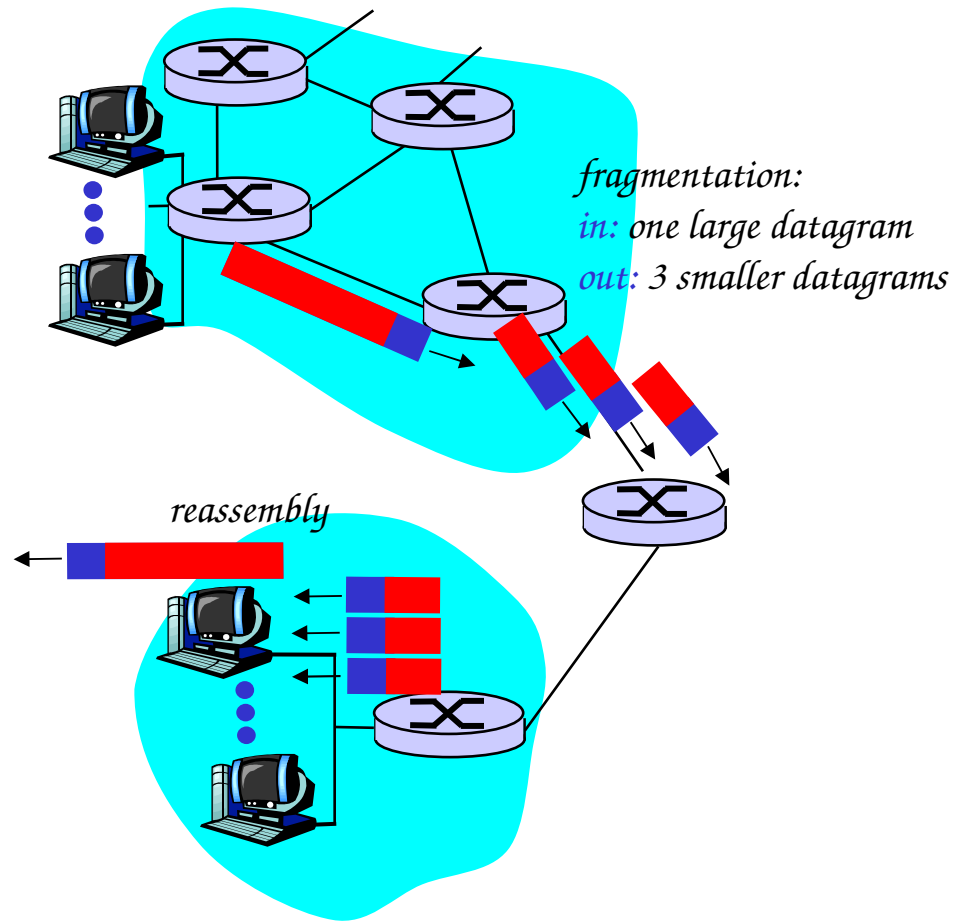


how much overhead with TCP?

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

IP Fragmentation & Reassembly

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



IP Fragmentation and Reassembly

Example

- 4000 byte datagram
- MTU = 1500 bytes

	<i>length</i> =4000	<i>ID</i> =χ	<i>fragflag</i> =0	<i>offset</i> =0	
--	------------------------	-----------------	-----------------------	---------------------	--

*One large datagram becomes
several smaller datagrams*

1480 bytes in
data field

offset =
1480/8

	<i>length</i> =1500	<i>ID</i> =χ	<i>fragflag</i> =1	<i>offset</i> =0	
--	------------------------	-----------------	-----------------------	---------------------	--

	<i>length</i> =1500	<i>ID</i> =χ	<i>fragflag</i> =1	<i>offset</i> =185	
--	------------------------	-----------------	-----------------------	-----------------------	--

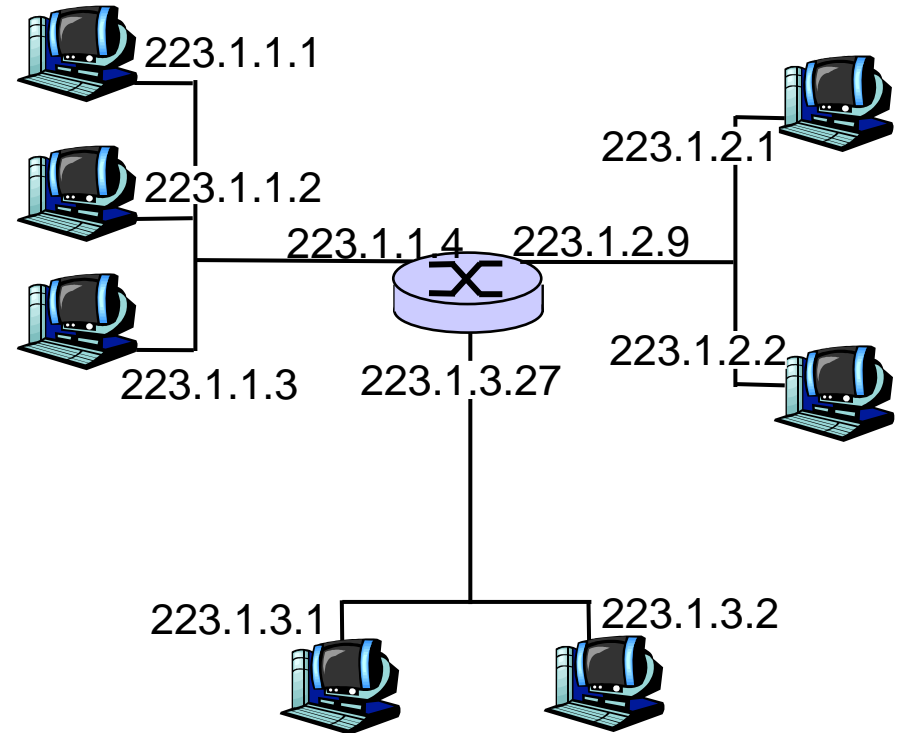
	<i>length</i> =1040	<i>ID</i> =χ	<i>fragflag</i> =0	<i>offset</i> =370	
--	------------------------	-----------------	-----------------------	-----------------------	--

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*
 - *router's typically have multiple interfaces*
 - *host typically has one interface*
 - *IP addresses associated with each interface*



$$223.1.1.1 = \underbrace{11011111}_{223} \underbrace{00000001}_1 \underbrace{00000001}_1 \underbrace{00000001}_1$$

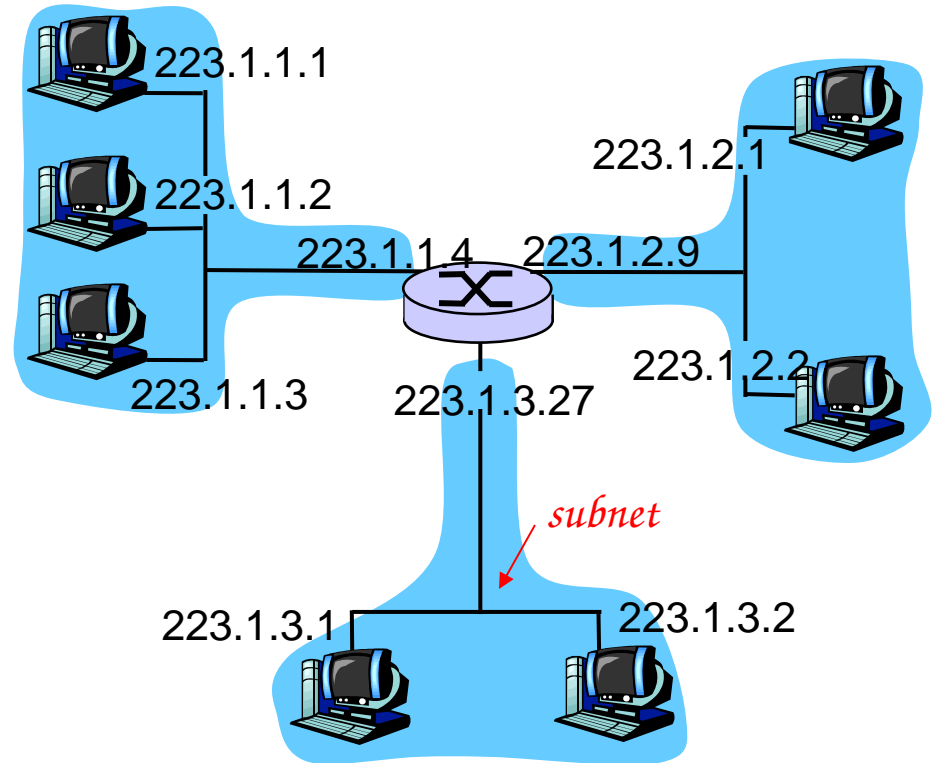
Subnets

□ *IP address:*

- *subnet part (high order bits)*
- *host part (low order bits)*

□ *What's a subnet ?*

- *device interfaces with same subnet part of IP address*
- *can physically reach each other without intervening router*

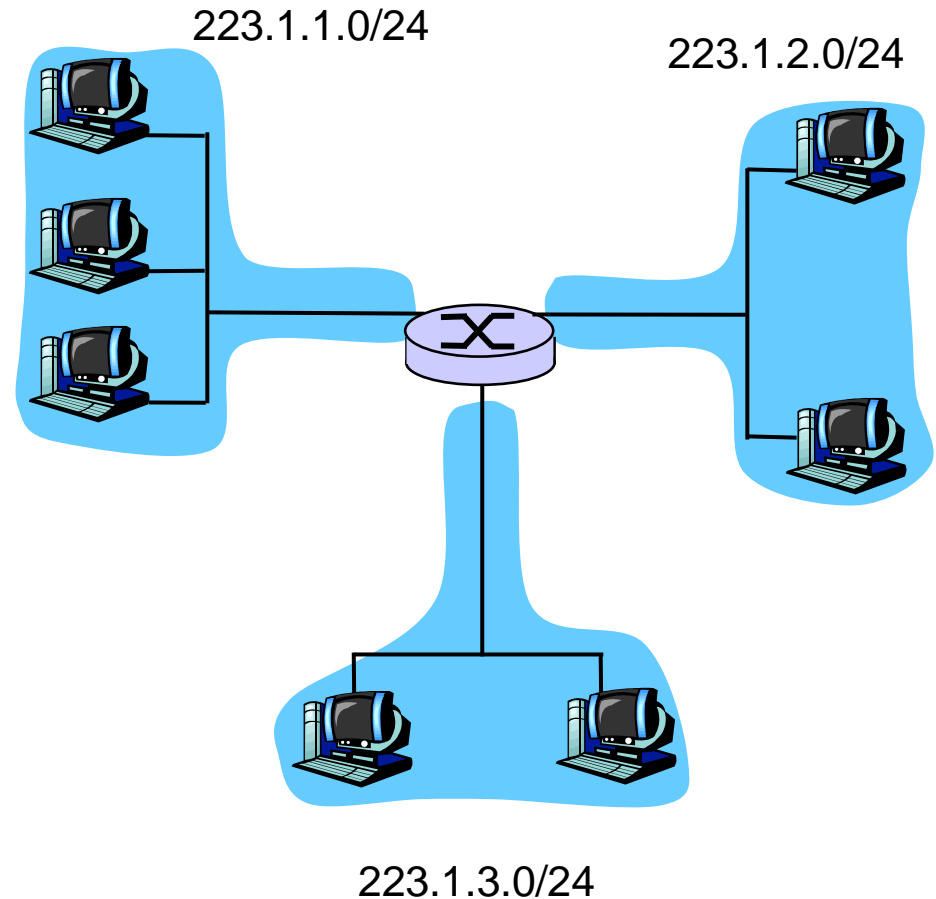


network consisting of 3 subnets

Subnets

Recipe

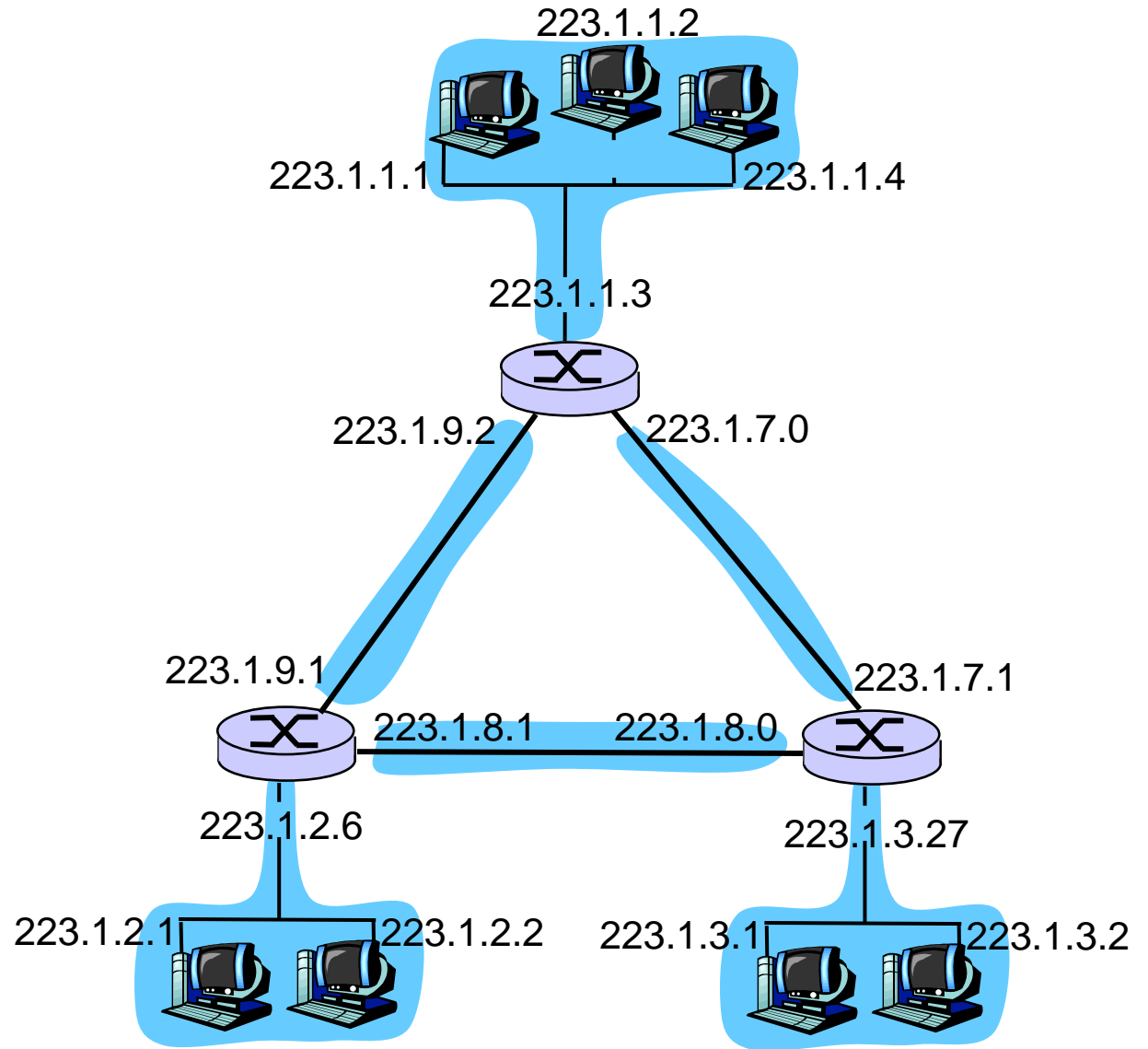
- To determine the subnets, detach each interface from its host or router, creating islands of isolated networks. Each isolated network is called a *subnet*.



Subnet mask: /24

Subnets

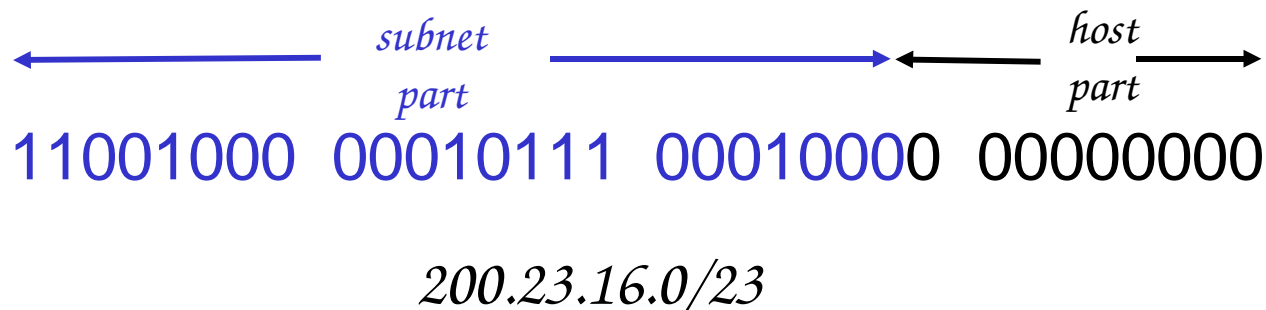
How many?



IP addressing: CIDR

CIDR: Classless InterDomain Routing

- subnet portion of address of arbitrary length
- address format: *a.b.c.d/χ*, where χ is # bits in subnet portion of address



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**: **D**ynamic **H**ost **C**onfiguration **P**rotocol: dynamically get address from as server*
 - *“plug-and-play”*

(more in next chapter)

IP addresses: how to get one?

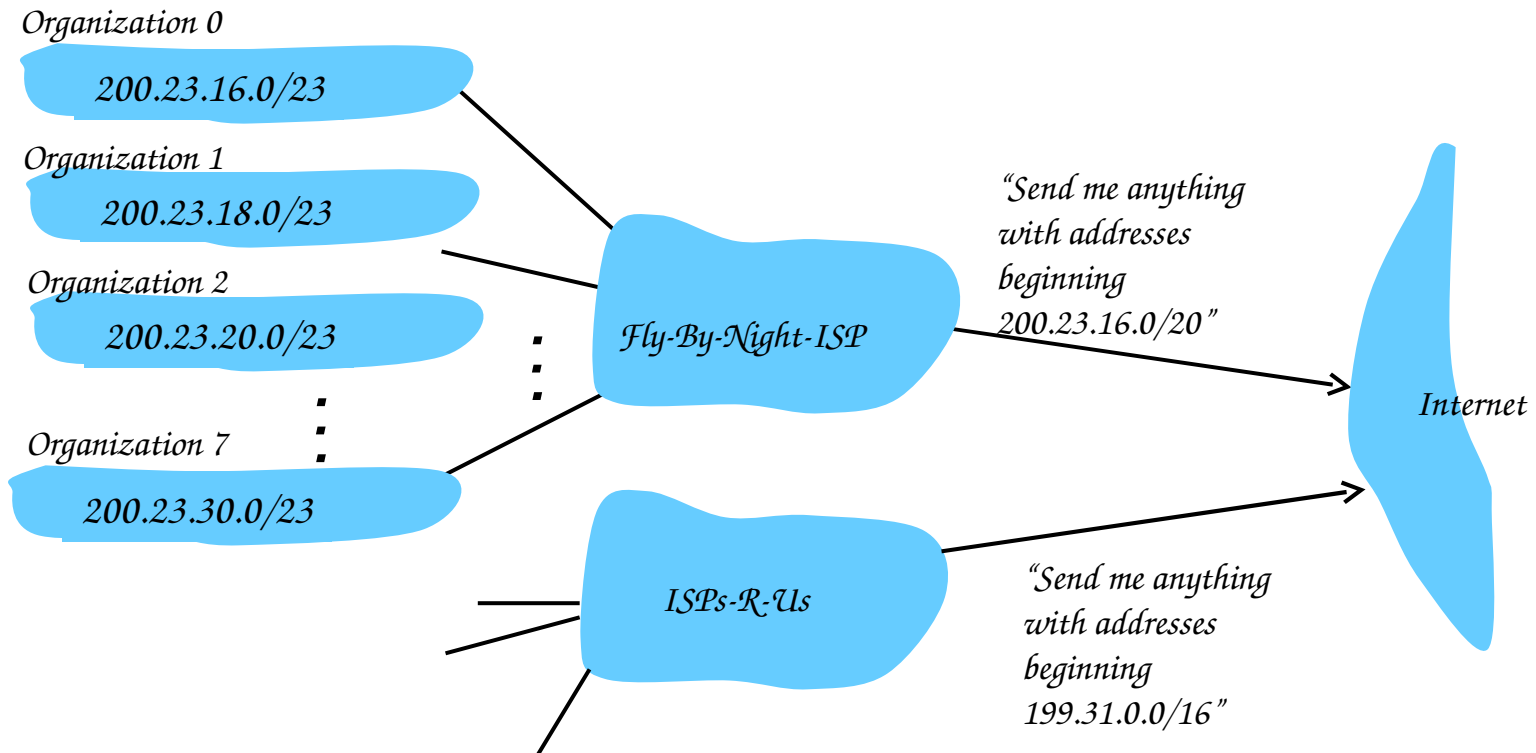
Q: *How does network get subnet part of IP addr?*

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

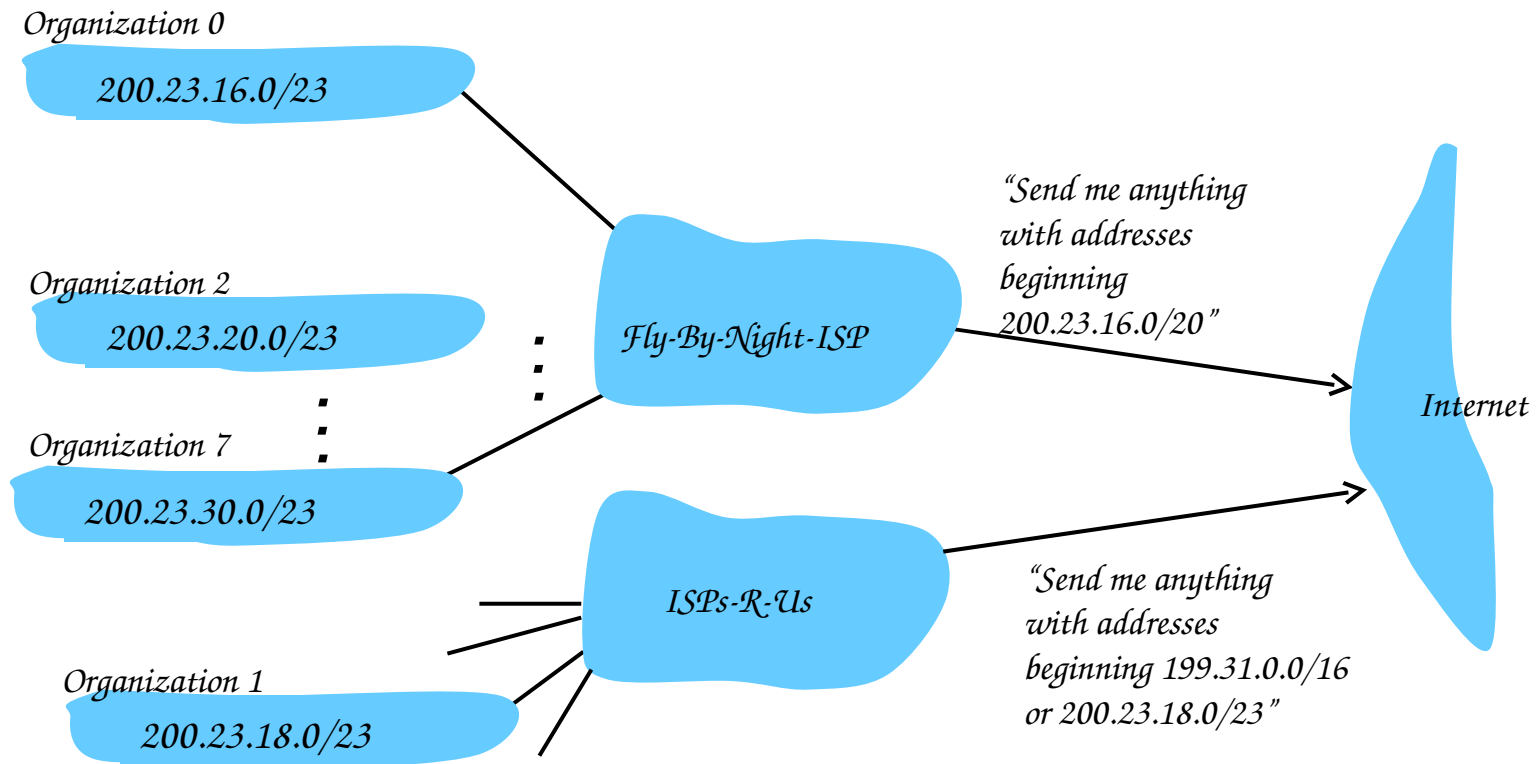
Hierarchical addressing: route aggregation

Hierarchical addressing allows efficient advertisement of routing information:



Hierarchical addressing: more specific routes

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



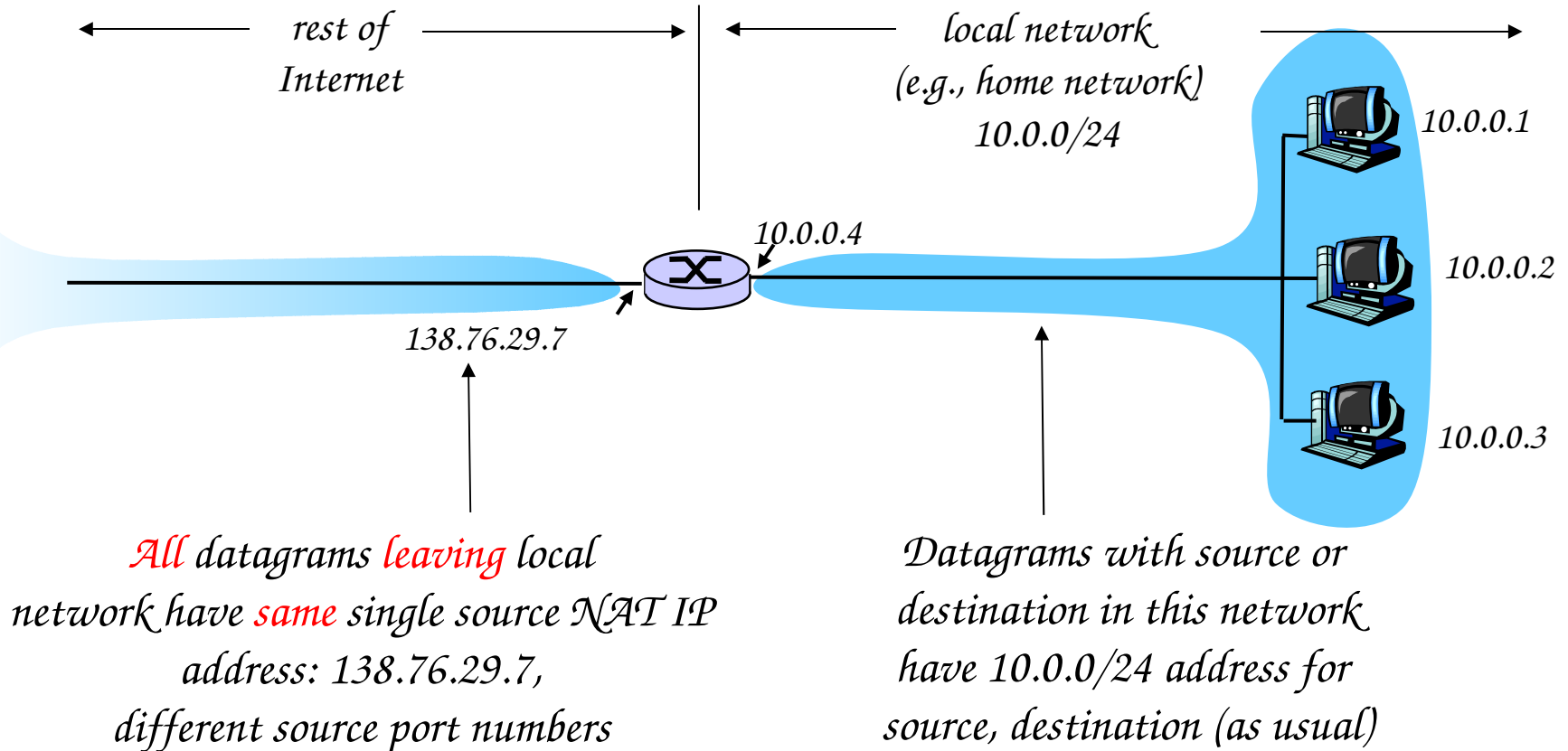
IP addressing: the last word...

Q: *How does an ISP get block of addresses?*

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

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

NAT: Network Address Translation



NAT: Network Address Translation

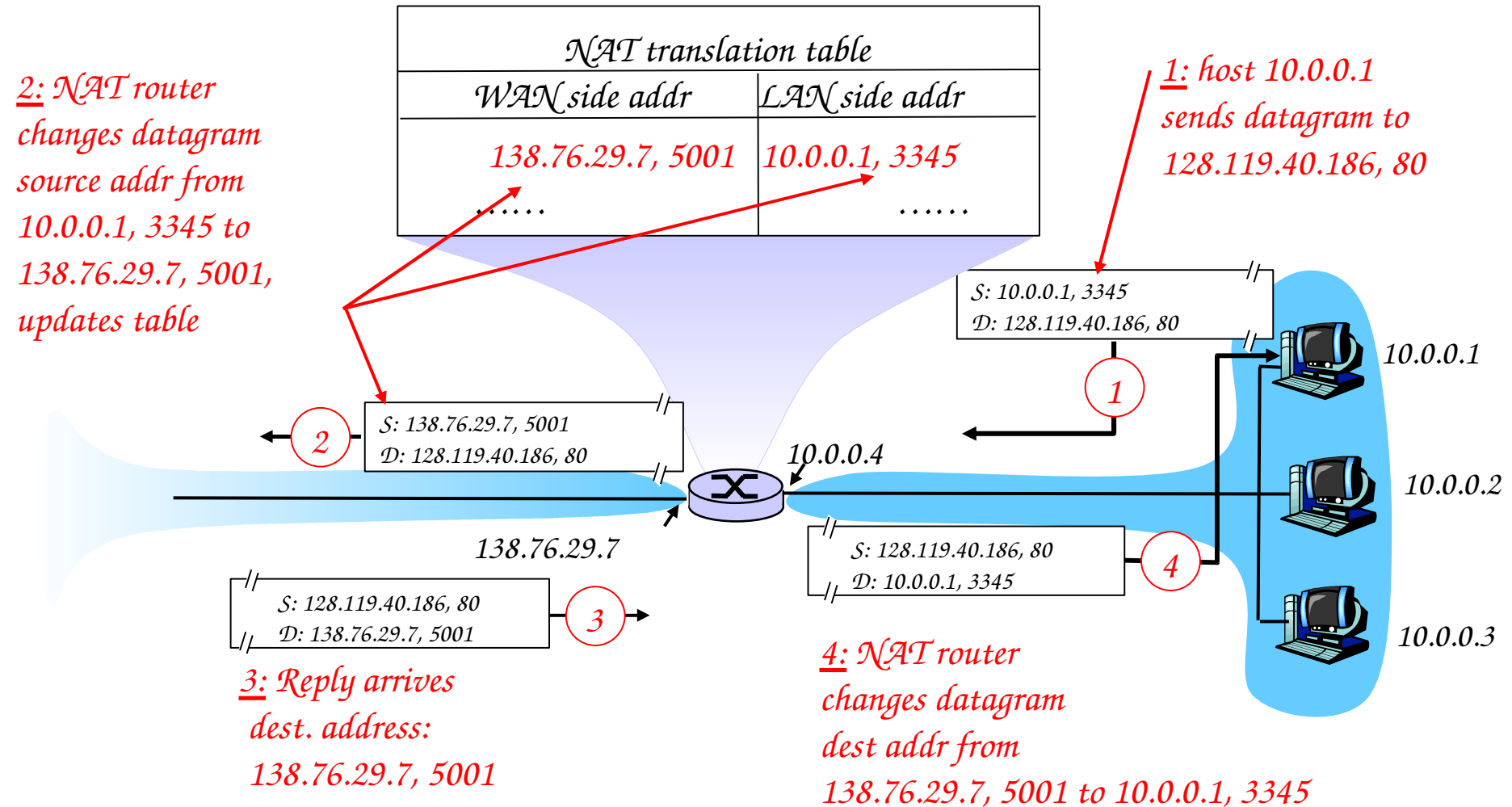
- **Motivation:** *local network uses just one IP address as far as outside world is concerned:*
 - *no need to be allocated range of addresses from ISP: - just one IP address is used for all devices*
 - *can change addresses of devices in local network without notifying outside world*
 - *can change ISP without changing addresses of devices in local network*
 - *devices inside local net not explicitly addressable, visible by outside world (a security plus).*

NAT: Network Address Translation

Implementation: NAT router must:

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

NAT: Network Address Translation



NAT: 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*

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 n th datagram arrives to n th 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*

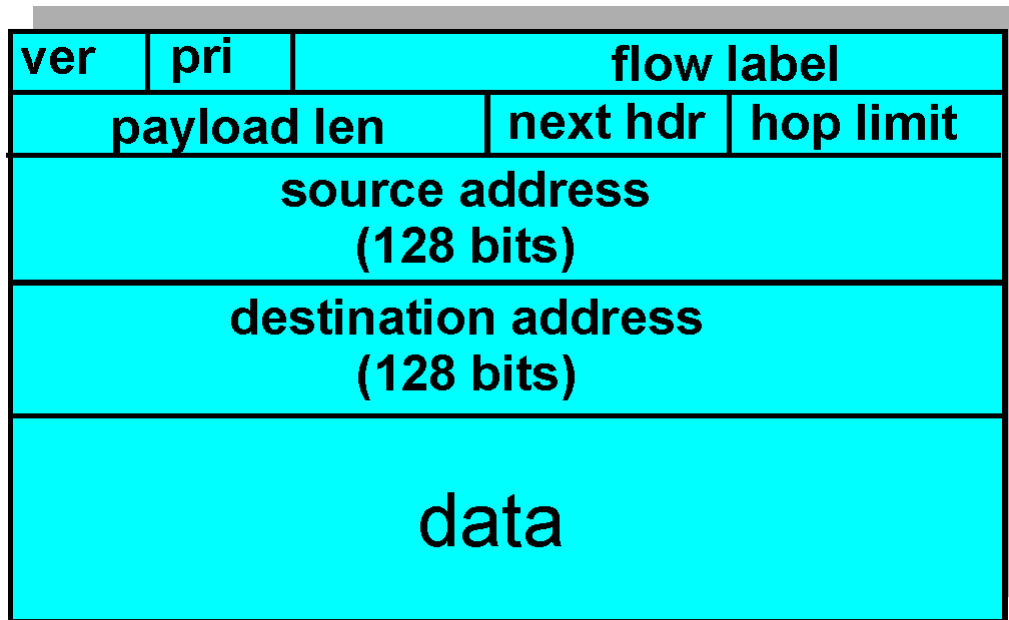
IPv6 Header (Cont)

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



← 32 bits →

Other Changes from IPv4

- ❑ *Checksum*: removed entirely to reduce processing time at each hop
- ❑ *Options*: allowed, but outside of header, indicated by “Next Header” field
- ❑ *ICMPv6*: new version of ICMP
 - additional message types, e.g. “Packet Too Big”
 - multicast group management functions

Transition From IPv4 To IPv6

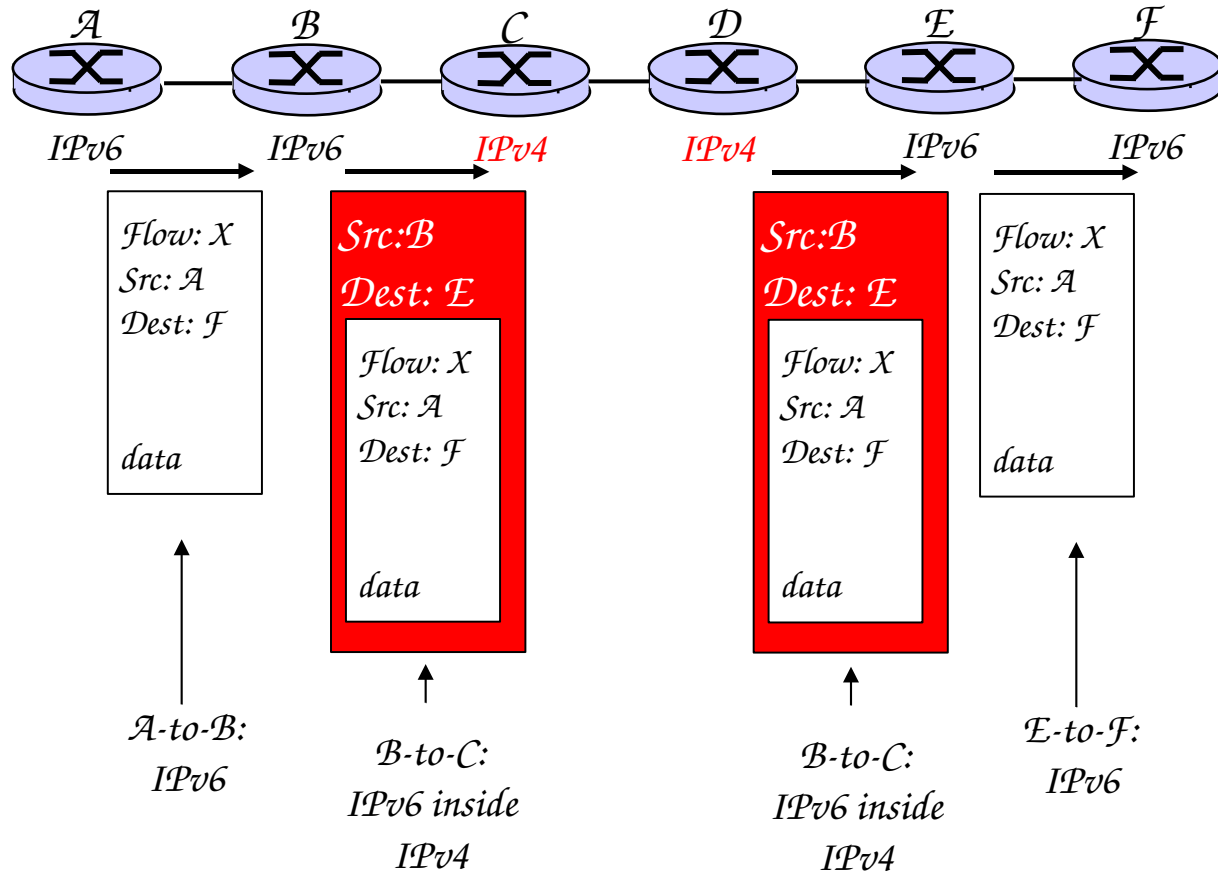
- ❑ *Not all routers can be upgraded simultaneous*
 - *no “flag days”*
 - *How will the network operate with mixed IPv4 and IPv6 routers?*
- ❑ ***Tunneling:** IPv6 carried as payload in IPv4 datagram among IPv4 routers*

Tunneling

Logical view:



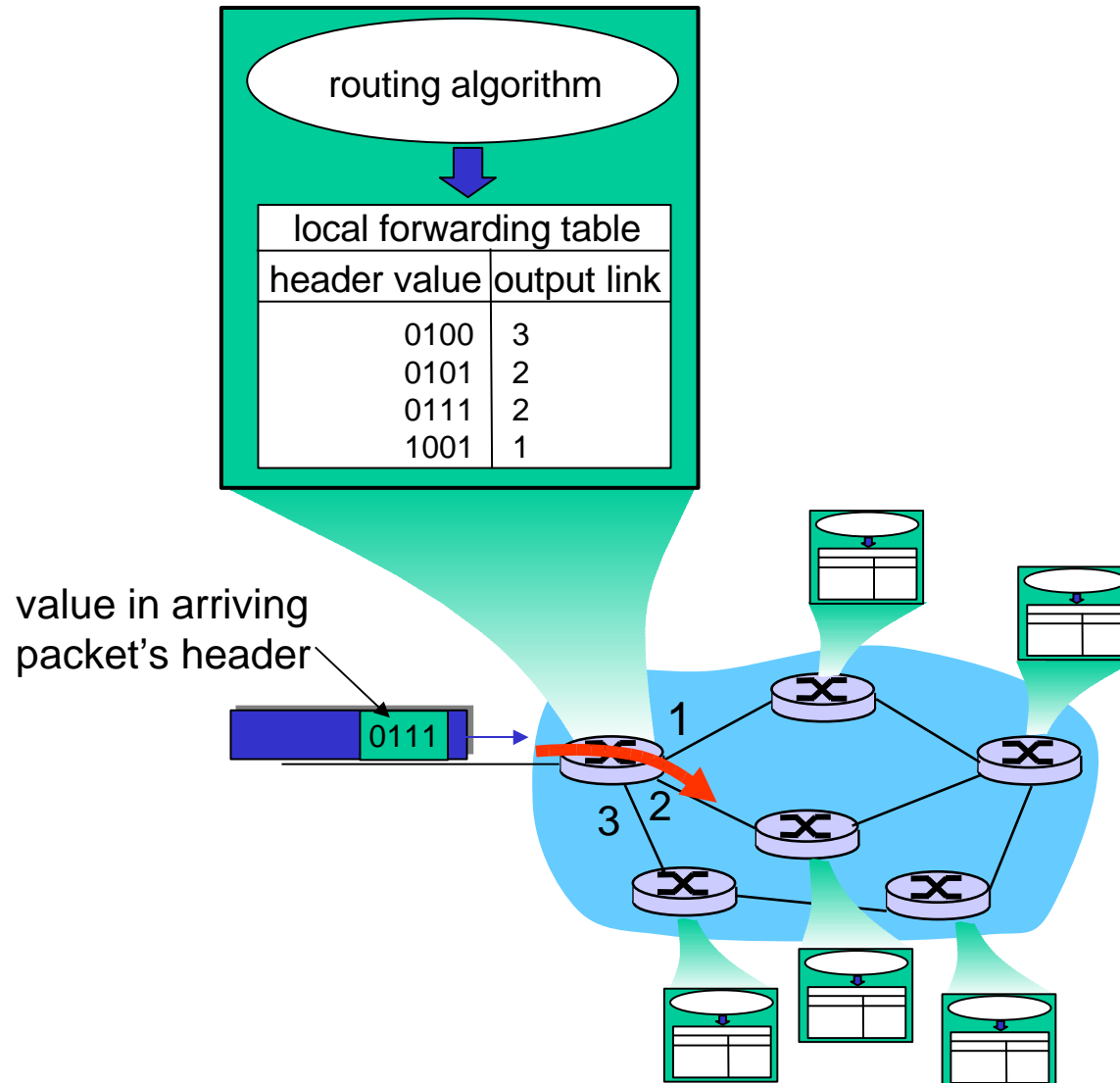
Physical view:



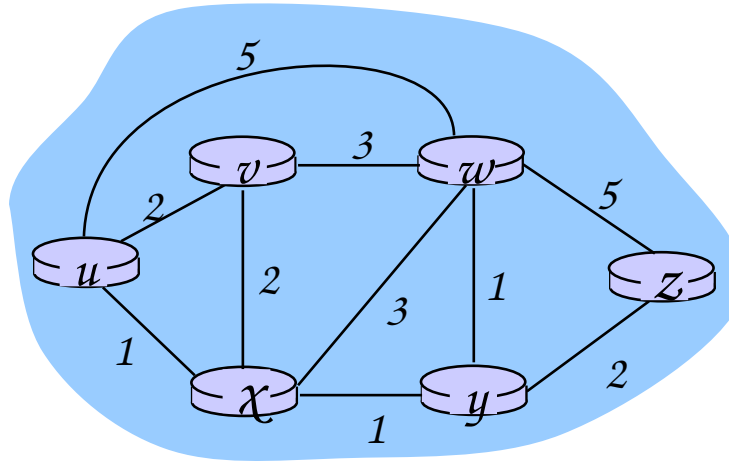
Chapter 4: Network Layer

- 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 and forwarding



Graph abstraction



Graph: $G = (N, E)$

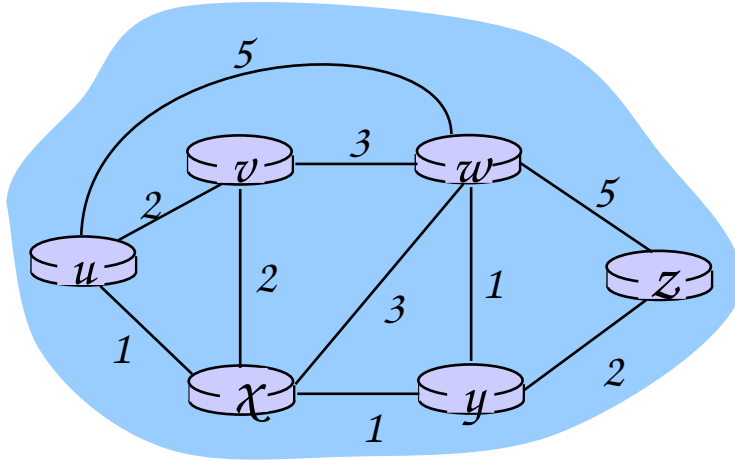
N = set of routers = $\{ u, v, w, x, y, z \}$

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 \mathcal{N} is set of peers and \mathcal{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 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - *Link state*
 - Distance Vector
 - Hierarchical routing
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast and multicast routing

A Link-State Routing Algorithm

Dijkstra's algorithm

- net topology, link costs known to all nodes
 - accomplished via “link state broadcast”
 - all nodes have same info
- computes least cost paths from one node (‘source’) to all other nodes
 - gives *forwarding table* for that node
- iterative: after k iterations, know least cost path to k 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 $D(v) = \min(D(v), D(w) + c(w,v))$

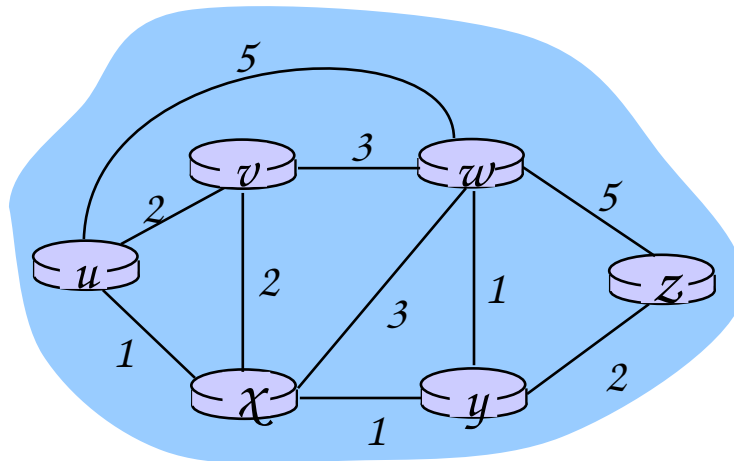
13 /* new cost to v is either old cost to v or known

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

15 **until all nodes in N'**

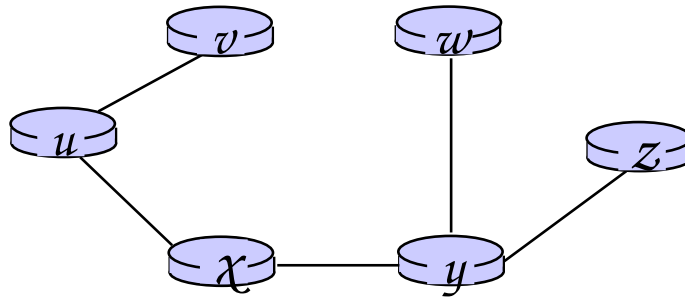
Dijkstra's algorithm: example

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	∞	∞
1	ux	2,u	4,x		2,x	∞
2	uxy	2,u	3,y			4,y
3	uxyv		3,y			4,y
4	uxyvw					4,y
5	uxyvwz					



Dijkstra's algorithm: example (2)

Resulting shortest-path tree from u :



Resulting forwarding table in u :

<i>destination</i>	<i>link</i>
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

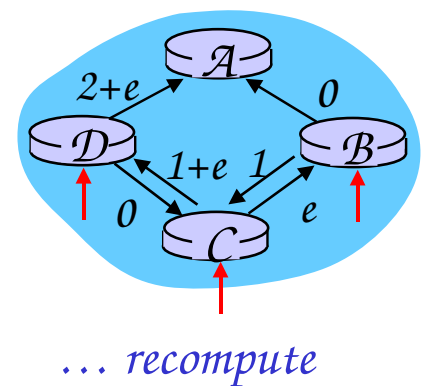
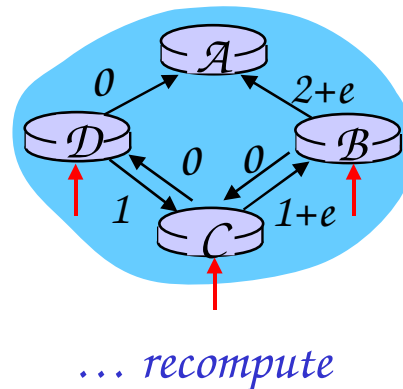
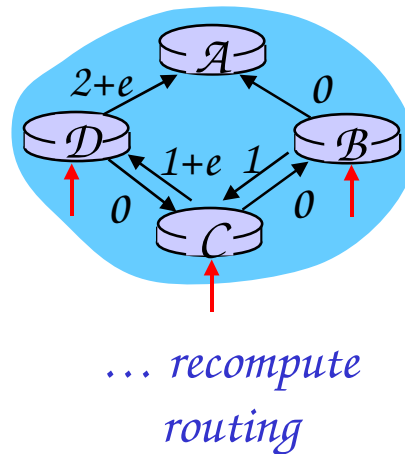
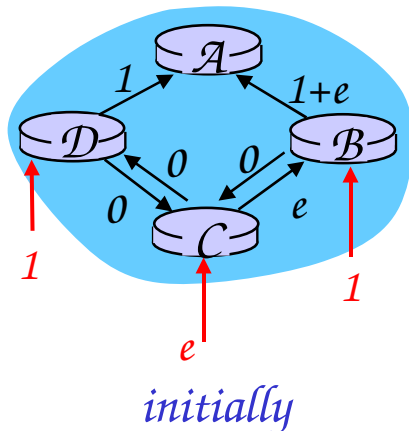
Dijkstra's algorithm, discussion

Algorithm complexity: n nodes

- each iteration: need to check all nodes, w , not in \mathcal{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 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

Distance Vector Algorithm

Bellman-Ford Equation (dynamic programming)

Define

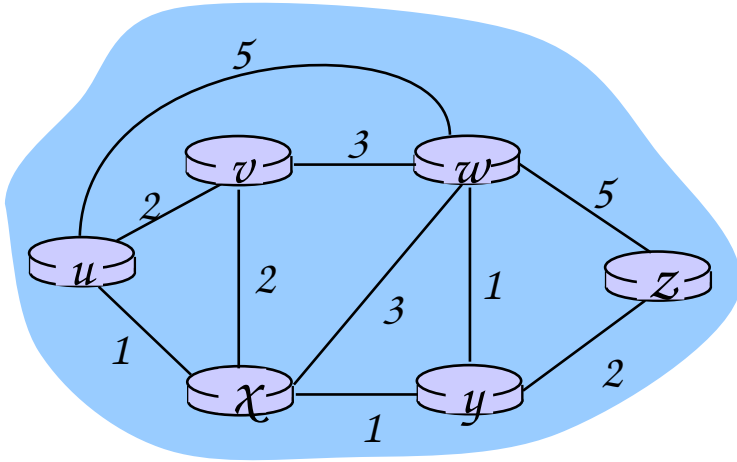
$d_x(y) := \text{cost of least-cost path from } x \text{ to } y$

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

- $\mathcal{D}_x(y)$ = estimate of least cost from x to y
- Distance vector: $\mathbf{D}_x = [\mathcal{D}_x(y): y \in \mathcal{N}]$
- Node x knows cost to each neighbor v : $c(xv)$
- Node x maintains $\mathbf{D}_x = [\mathcal{D}_x(y): y \in \mathcal{N}]$
- Node x also maintains its neighbors' distance vectors
 - For each neighbor v , x maintains $\mathbf{D}_v = [\mathcal{D}_v(y): y \in \mathcal{N}]$

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 $d_x(y)$

Distance Vector Algorithm (5)

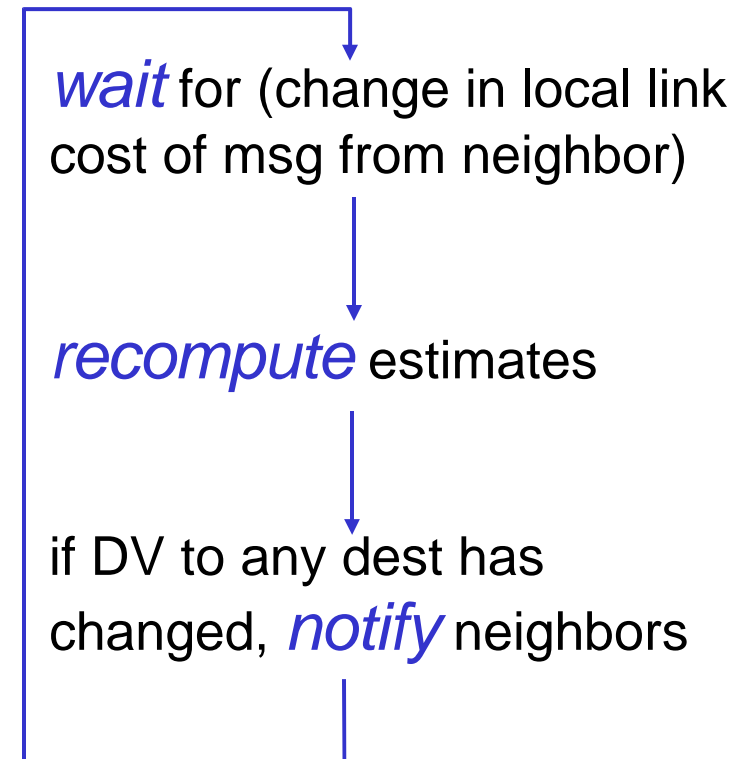
Iterative, asynchronous: each local iteration caused by:

- ❑ 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:



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

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

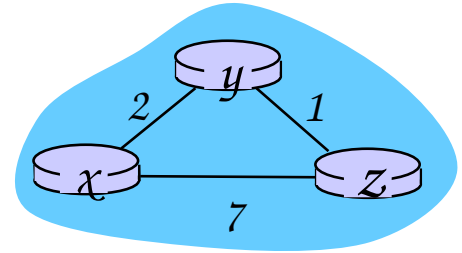
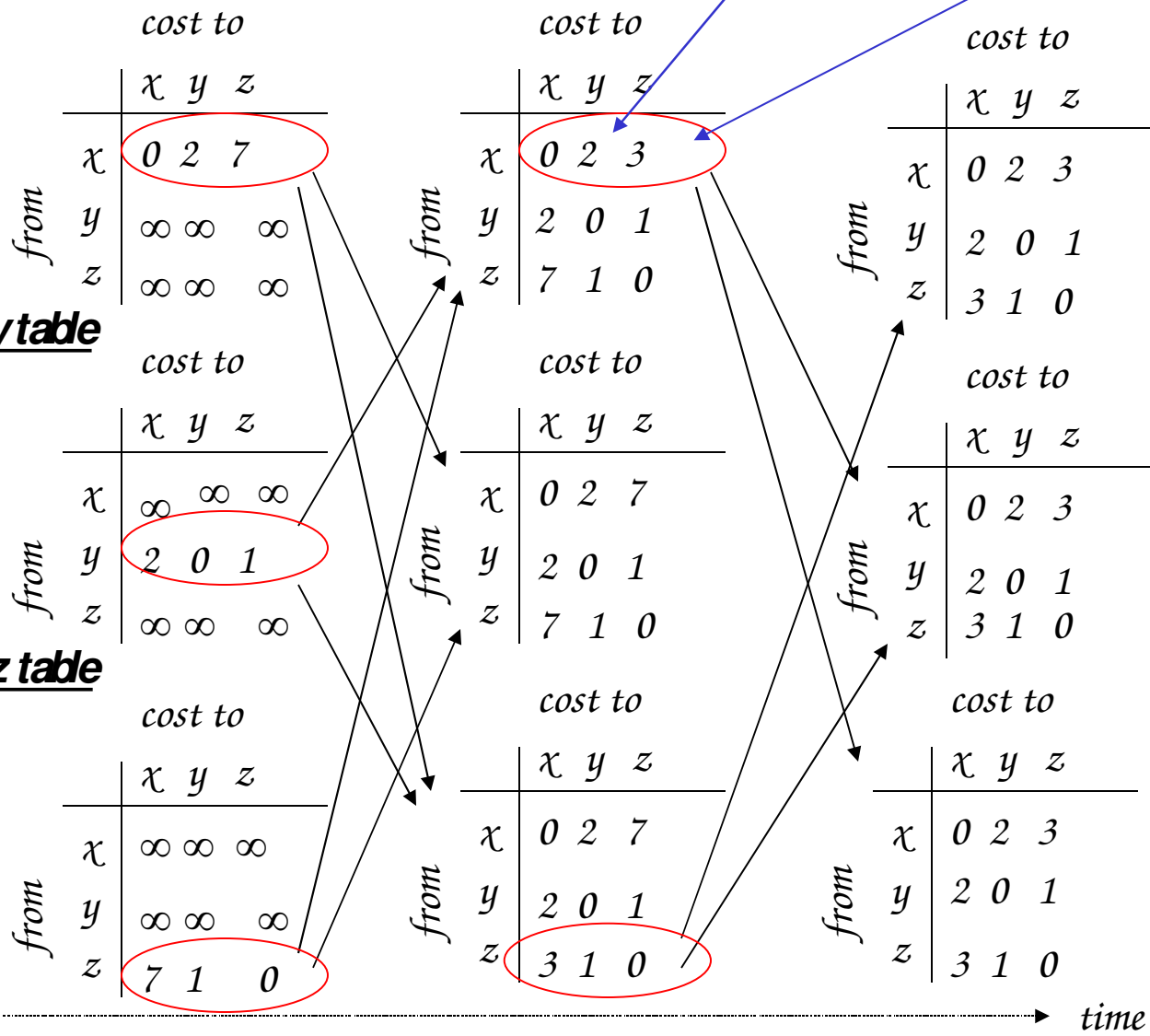
$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\}$$

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

nodex table

nodey table

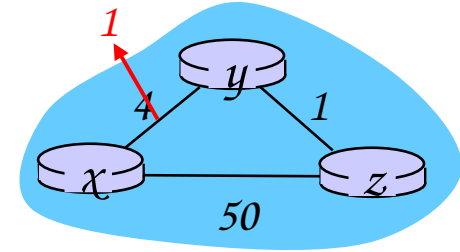
nodez table



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



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 .

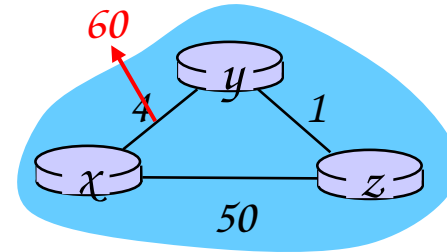
At time t_2 , y receives z 's update and updates its distance table. y 's least costs do not change and hence y does not send any message to z .

“good
news
travels
fast”

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 \mathcal{LS} and \mathcal{DV} algorithms

Message complexity

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

Speed of Convergence

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

Robustness: what happens if router malfunctions?

\mathcal{LS} :

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

\mathcal{DV} :

- \mathcal{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*

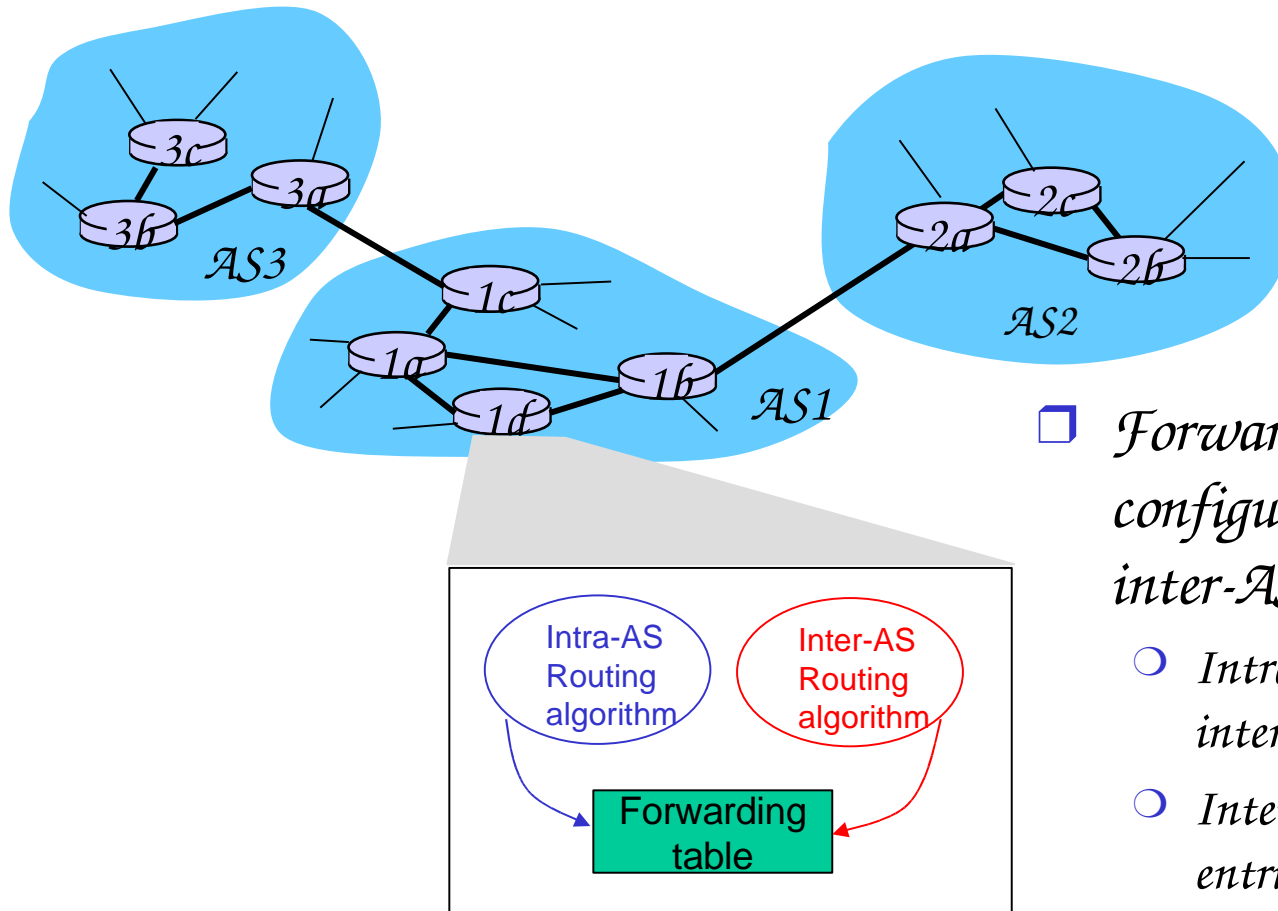
Hierarchical Routing

- 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 router

- Direct link to router in another AS

Interconnected ASes



- Forwarding table is configured by both intra- and inter-AS routing algorithm
 - Intra-AS sets entries for internal dests
 - Inter-AS & Intra-As sets entries for external dests

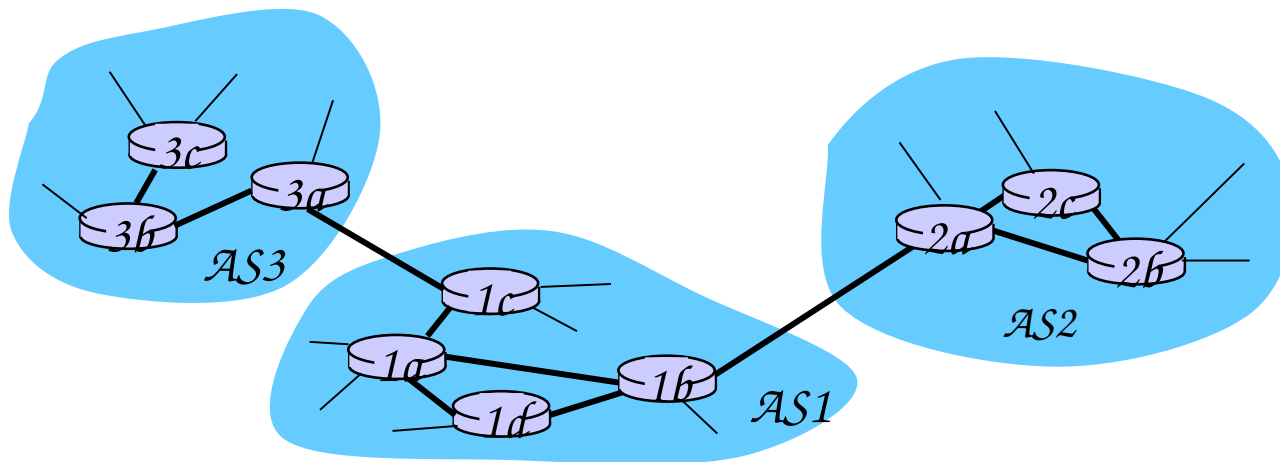
Inter-AS tasks

- Suppose router in AS1 receives datagram for which dest is outside of AS1
 - Router should forward packet towards one of the gateway routers, but which one?

AS1 needs:

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

Job of inter-AS routing!

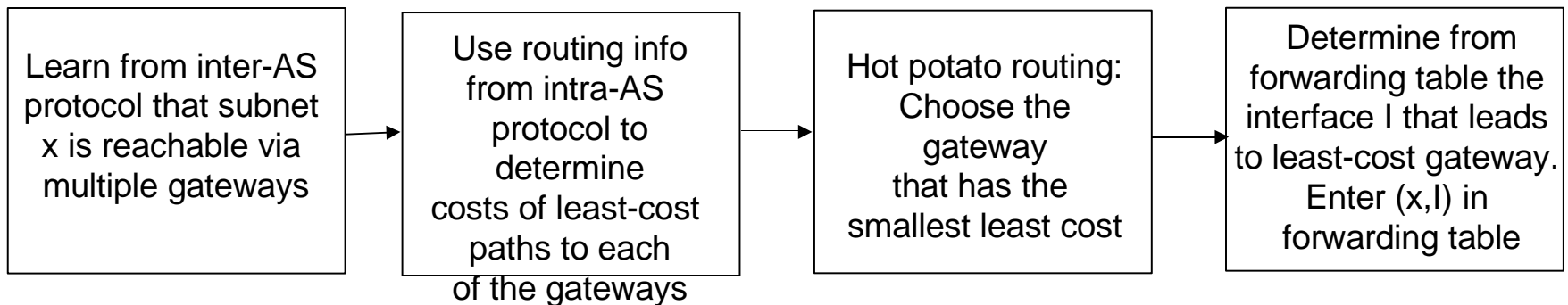


Example: Setting forwarding table in router 1d

- Suppose $\mathcal{AS}1$ learns from the inter- \mathcal{AS} protocol that subnet χ is reachable from $\mathcal{AS}3$ (gateway 1c) but not from $\mathcal{AS}2$.
- Inter- \mathcal{AS} protocol propagates reachability info to all internal routers.
- Router 1d determines from intra- \mathcal{AS} routing info that its interface I is on the least cost path to 1c.
- Puts in forwarding table entry (χ, I) .

Example: Choosing among multiple ASes

- Now suppose AS1 learns from the inter-AS protocol that subnet x is reachable from AS3 and from AS2.
- To configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest x .
- This is also the job on 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

Intra-AS Routing

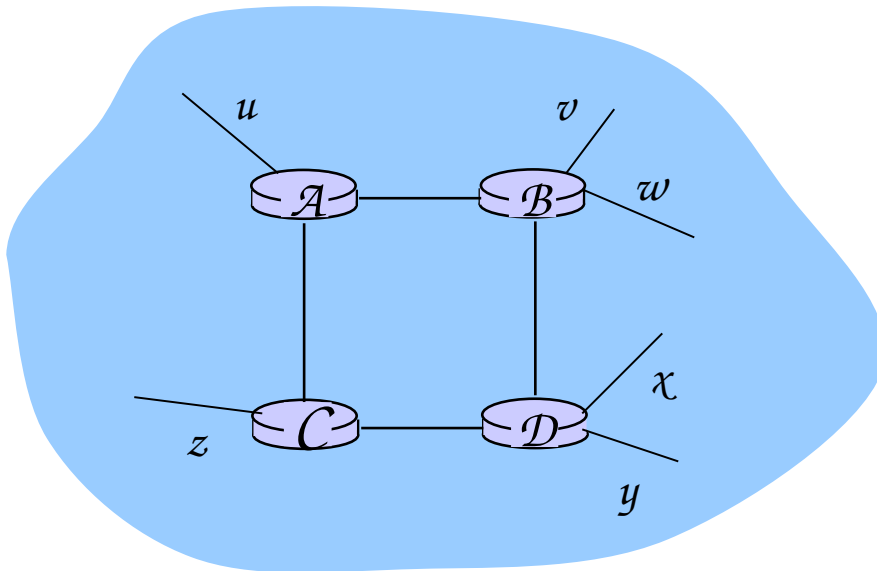
- ❑ Also known as *Interior Gateway Protocols (IGP)*
- ❑ Most common Intra-AS routing protocols:
 - *RIP: Routing Information Protocol*
 - *OSPF: Open Shortest Path First*
 - *IGRP: Interior Gateway Routing Protocol (Cisco proprietary)*

Chapter 4: Network Layer

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

RIP (Routing Information Protocol)

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



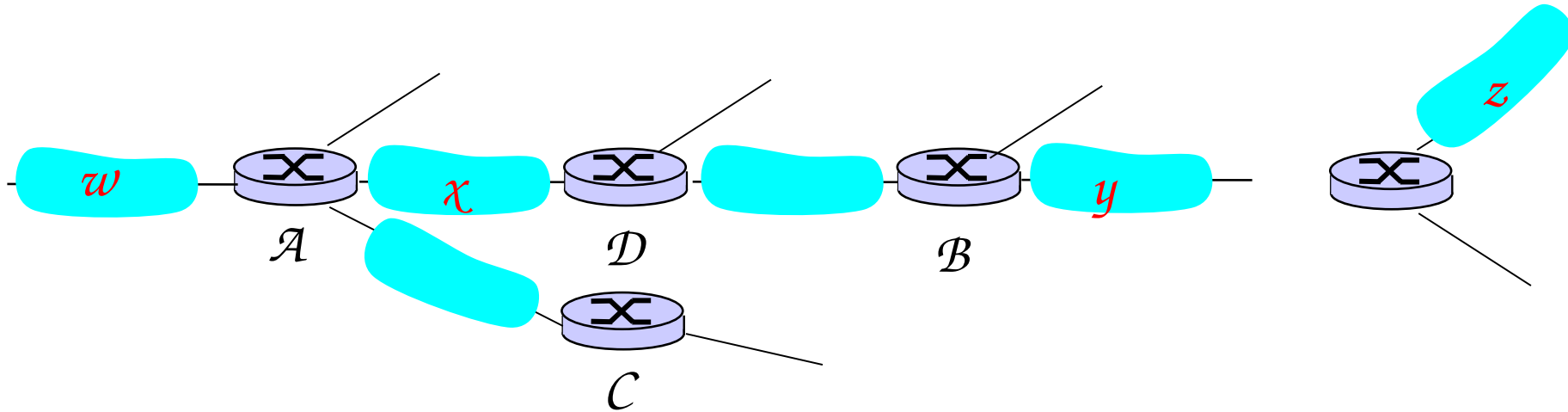
From router A to subsets:

<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



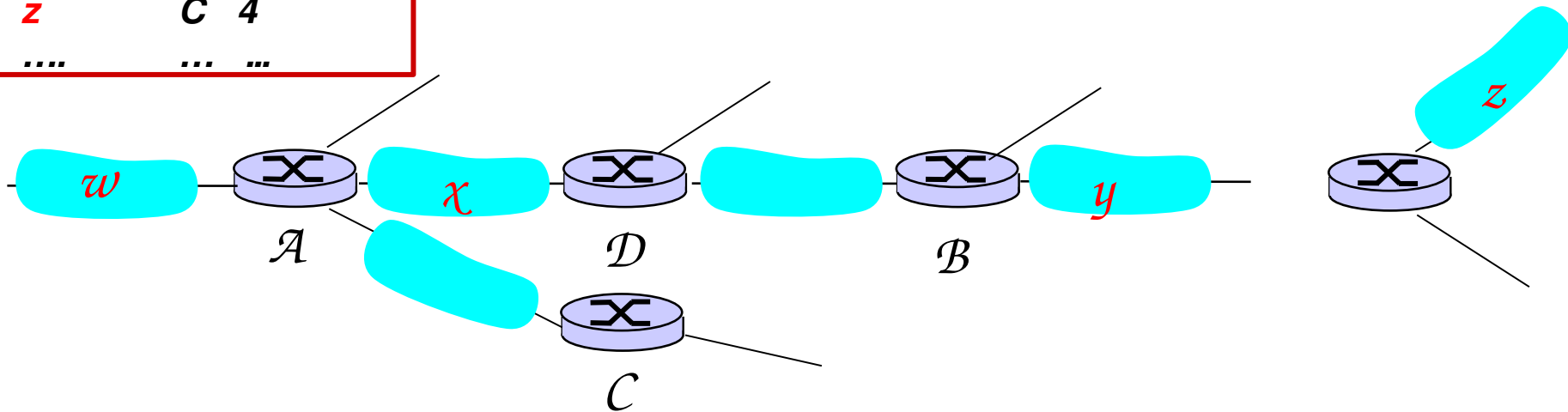
<i>Destination Network</i>	<i>Next Router</i>	<i>Num of hops to dest.</i>
w	A	2
y	B	2
z	B	7
x	—	1
....

Routing table in D

RIP: Example

Dest	Next	hops
<i>w</i>	-	1
<i>x</i>	-	1
<i>z</i>	C	4
....

Advertisement
from A to D



Destination Network	Next Router	Num of hops to dest.
<i>w</i>	A	2
<i>y</i>	B	2
<i>z</i>	BA	75
<i>x</i>	-	1
....

Routing table in D

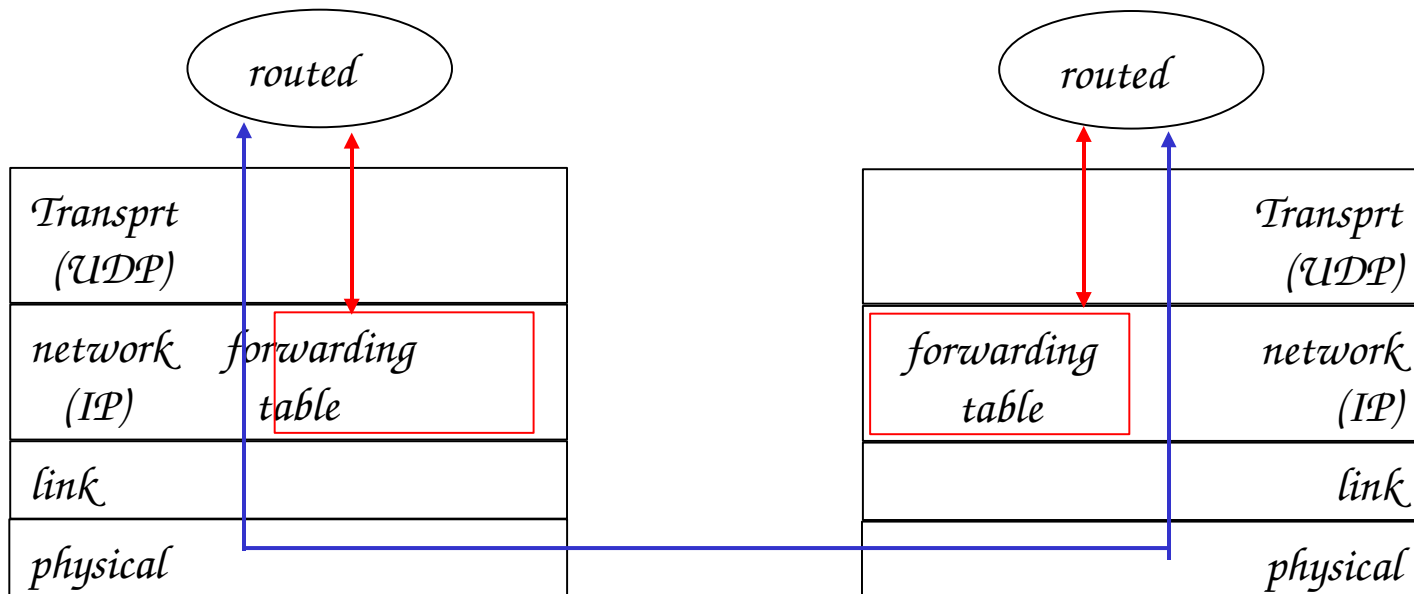
RIP: Link Failure and Recovery

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

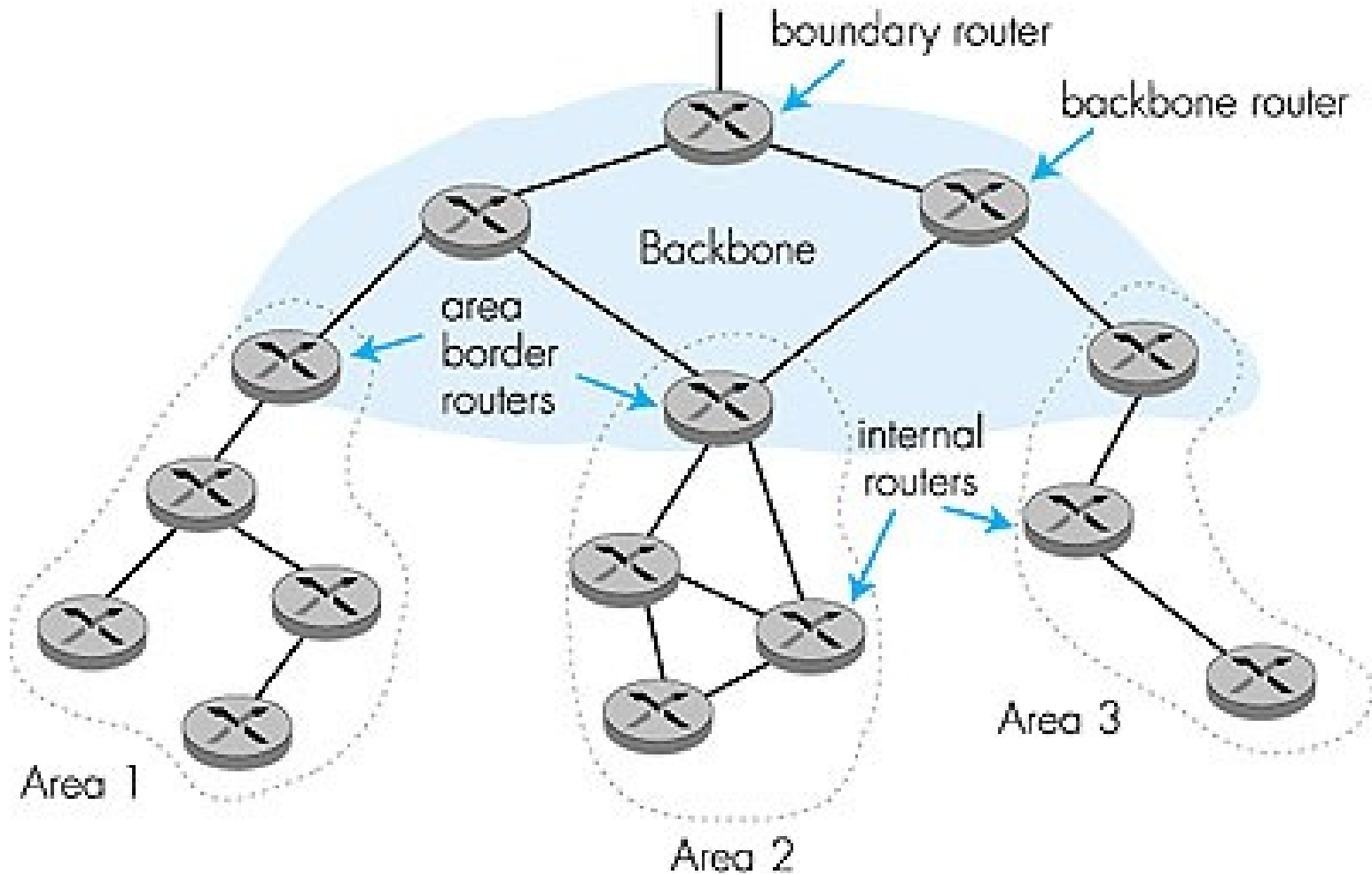
OSPF (Open Shortest Path First)

- ❑ “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

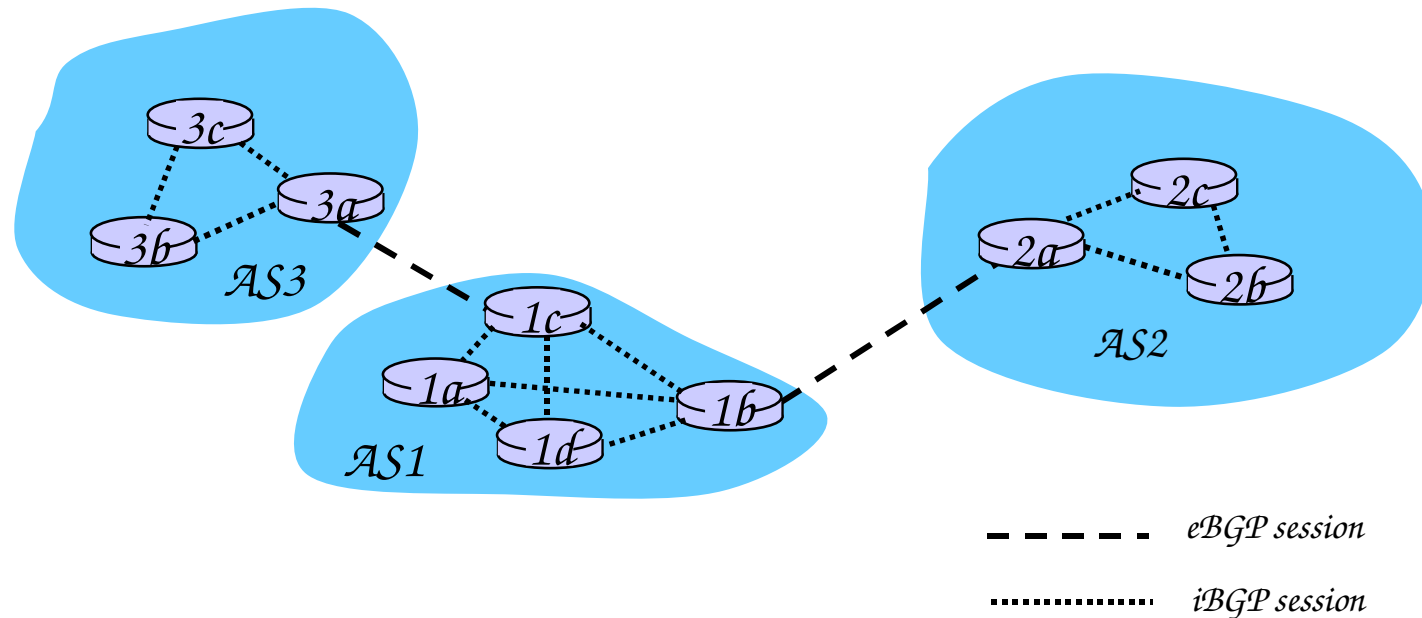
- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
 - Hierarchical routing
- 4.6 Routing in the Internet
 - RIP
 - OSPF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

Internet inter-AS routing: BGP

- ❑ *BGP (Border Gateway Protocol): the de facto standard*
- ❑ *BGP provides each AS a means to:*
 1. *Obtain subnet reachability information from neighboring ASs.*
 2. *Propagate the reachability information to all routers internal to the AS.*
 3. *Determine “good” routes to subnets based on reachability information and policy.*
- ❑ *Allows a subnet to advertise its existence to rest of the Internet: “I am here”*

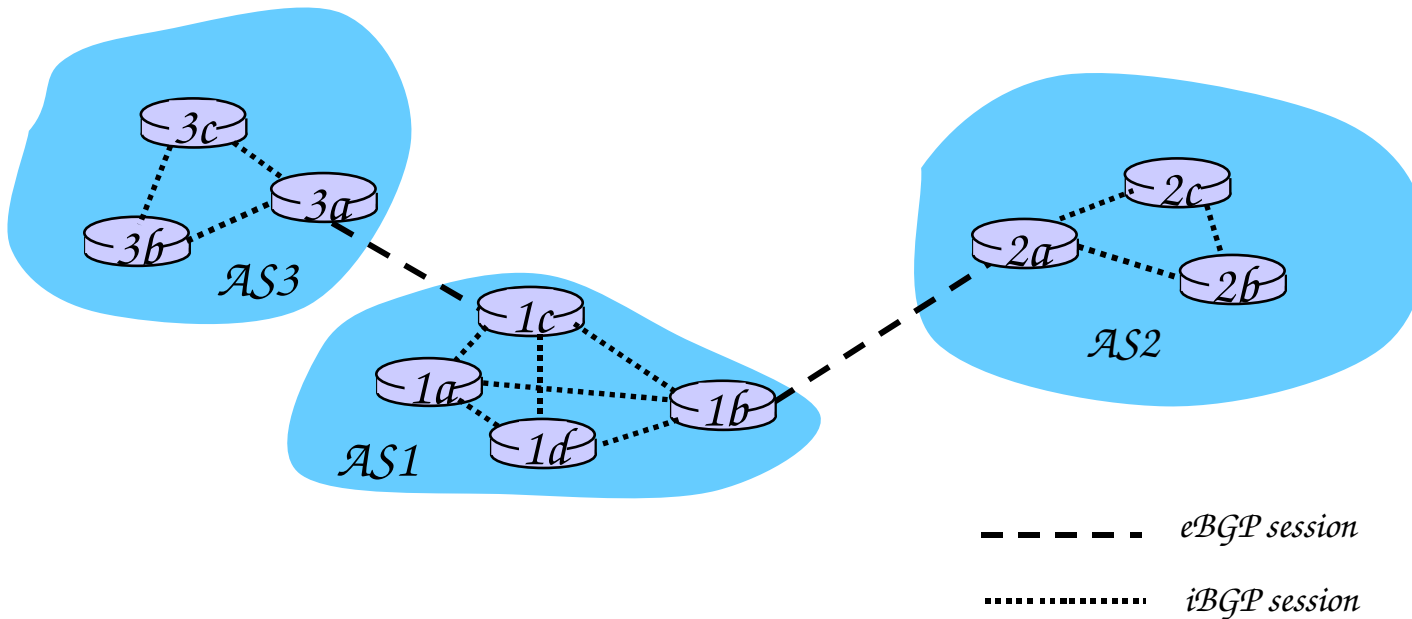
BGP basics

- Pairs of routers (BGP peers) exchange routing info over semi-permanent TCP conctns: **BGP sessions**
- Note that BGP sessions do not correspond to physical links.
- When AS2 advertises a prefix to AS1, AS2 is **promising** it will forward any datagrams destined to that prefix towards the prefix.
 - AS2 can aggregate prefixes in its advertisement



Distributing reachability info

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



Path attributes & BGP routes

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

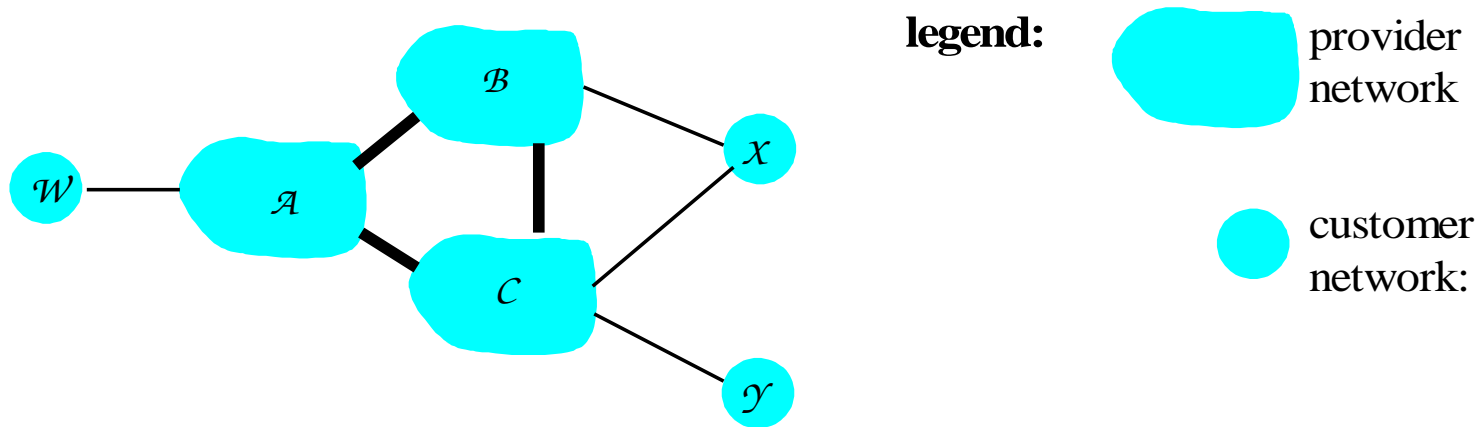
BGP route selection

- ❑ Router may learn about more than 1 route to some prefix. Router must select route.
- ❑ Elimination rules:
 1. Local preference value attribute: policy decision
 2. Shortest AS-PATH
 3. Closest NEXT-HOP router: hot potato routing
 4. Additional criteria

BGP messages

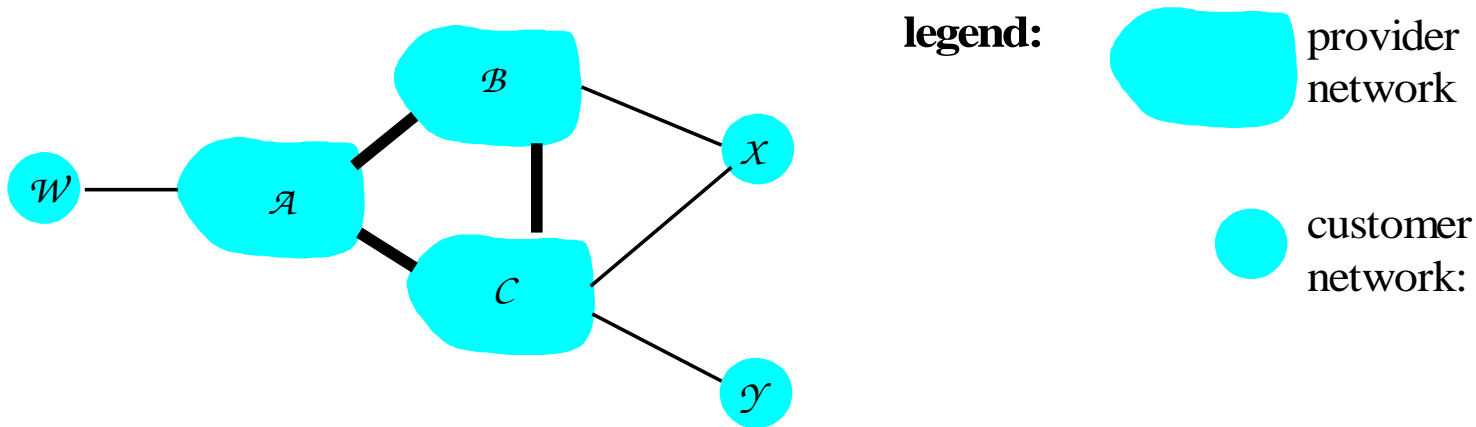
- *BGP messages exchanged using TCP.*
- *BGP messages:*
 - *OPEN*: opens TCP connection to peer and authenticates sender
 - *UPDATE*: advertises new path (or withdraws old)
 - *KEEPALIVE* keeps connection alive in absence of *UPDATES*; also *ACKs* *OPEN* request
 - *NOTIFICATION*: reports errors in previous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

BGP routing policy



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

Why different Intra- and Inter-AS routing ?

Policy:

- ❑ *Inter-AS: 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

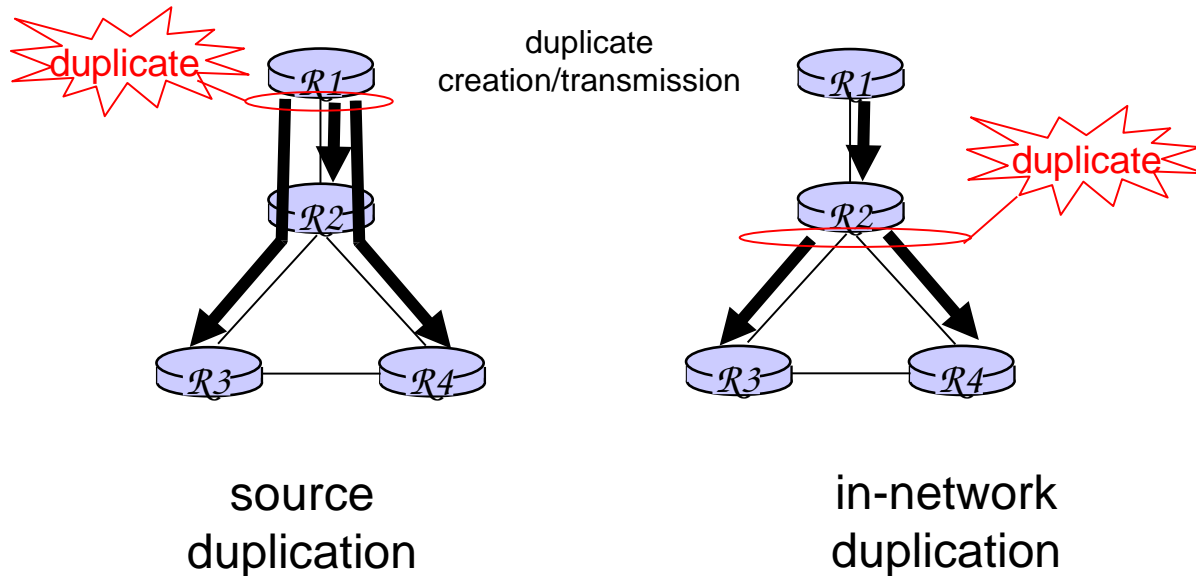
- ❑ *Intra-AS: can focus on performance*
- ❑ *Inter-AS: policy may dominate over performance*

Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 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

Broadcast Routing

- Deliver packets from srce to all other nodes
- Source duplication is inefficient:



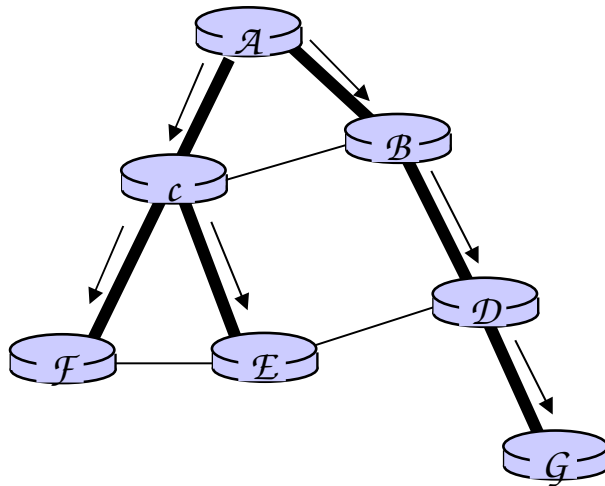
- Source duplication: how does source determine recipient addresses

In-network duplication

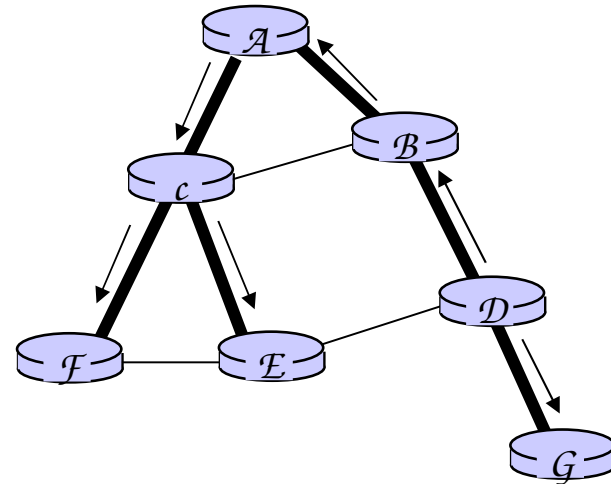
- ❑ *Flooding: when node receives brdcst pkt, sends copy to all neighbors*
 - *Problems: cycles & broadcast storm*
- ❑ *Controlled flooding: node only brdcsts pkt if it hasn't brdcst same packet before*
 - *Node keeps track of pkt ids already brdcsted*
 - *Or reverse path forwarding (RPF): only forward 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*



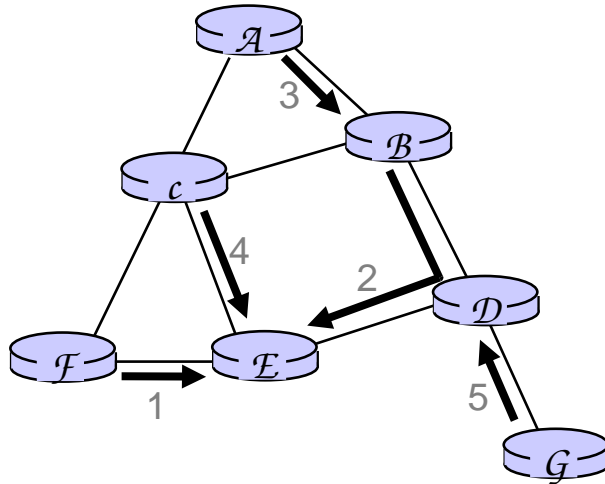
(a) Broadcast initiated at A



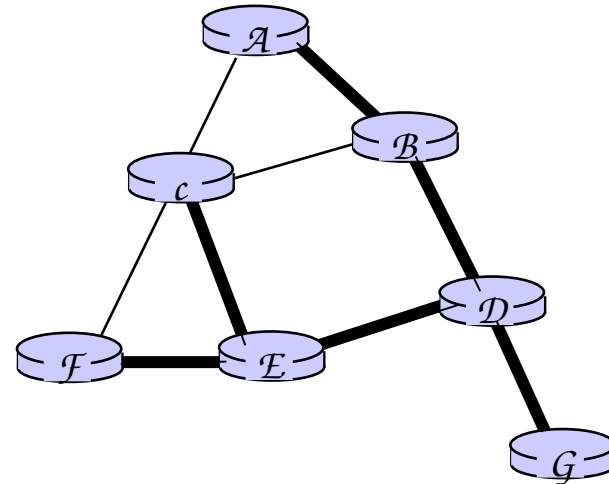
(b) Broadcast initiated at D

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



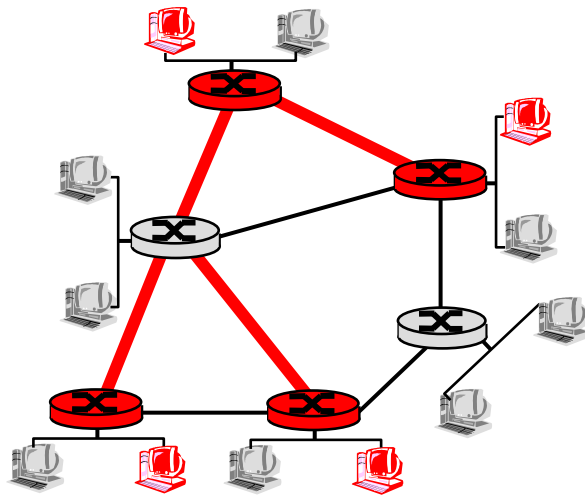
(a) Stepwise construction of spanning tree



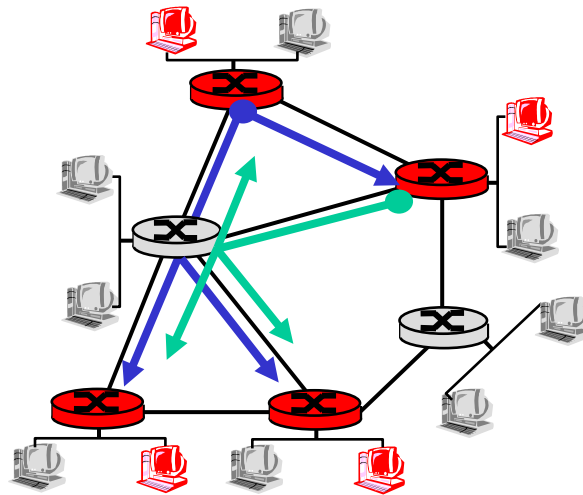
(b) Constructed spanning tree

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*

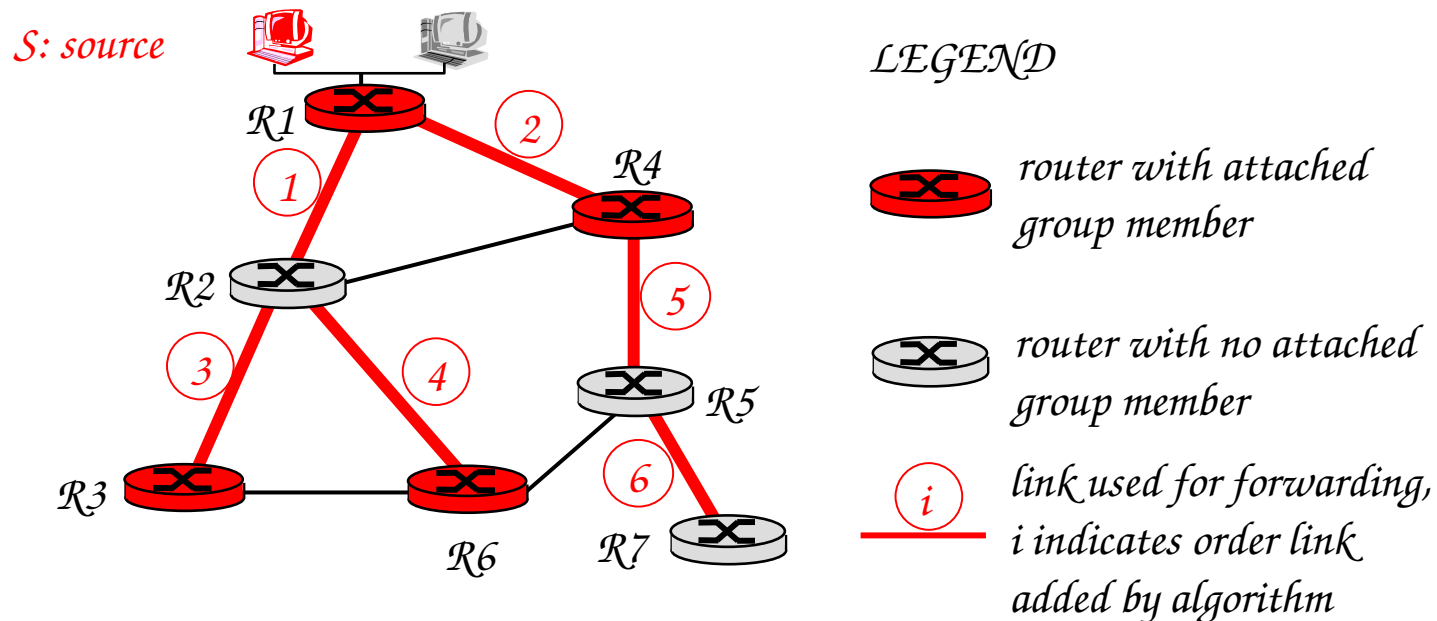
□ *group-shared tree: group uses one tree*

- *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

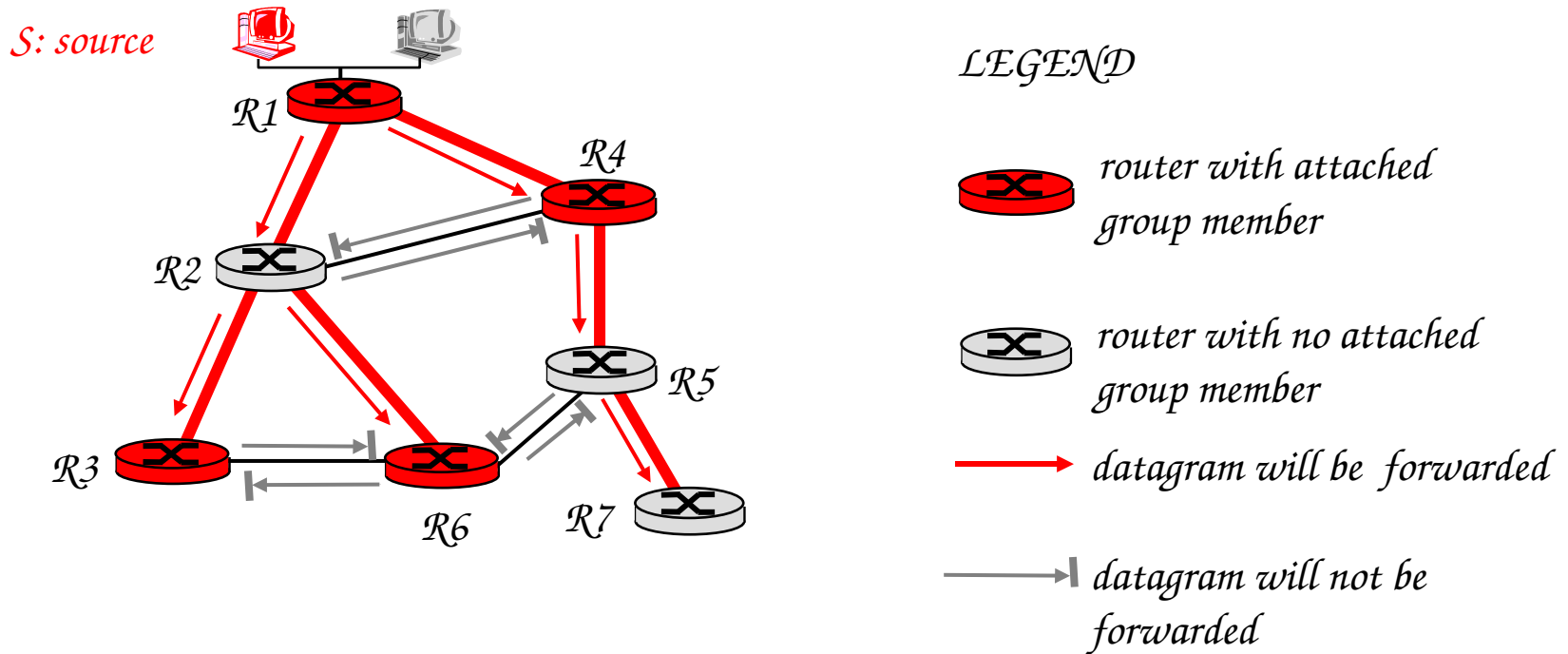
each router has simple forwarding behavior:

if (mcast datagram received on incoming link on shortest path back to center)

then flood datagram onto all outgoing links

else ignore datagram

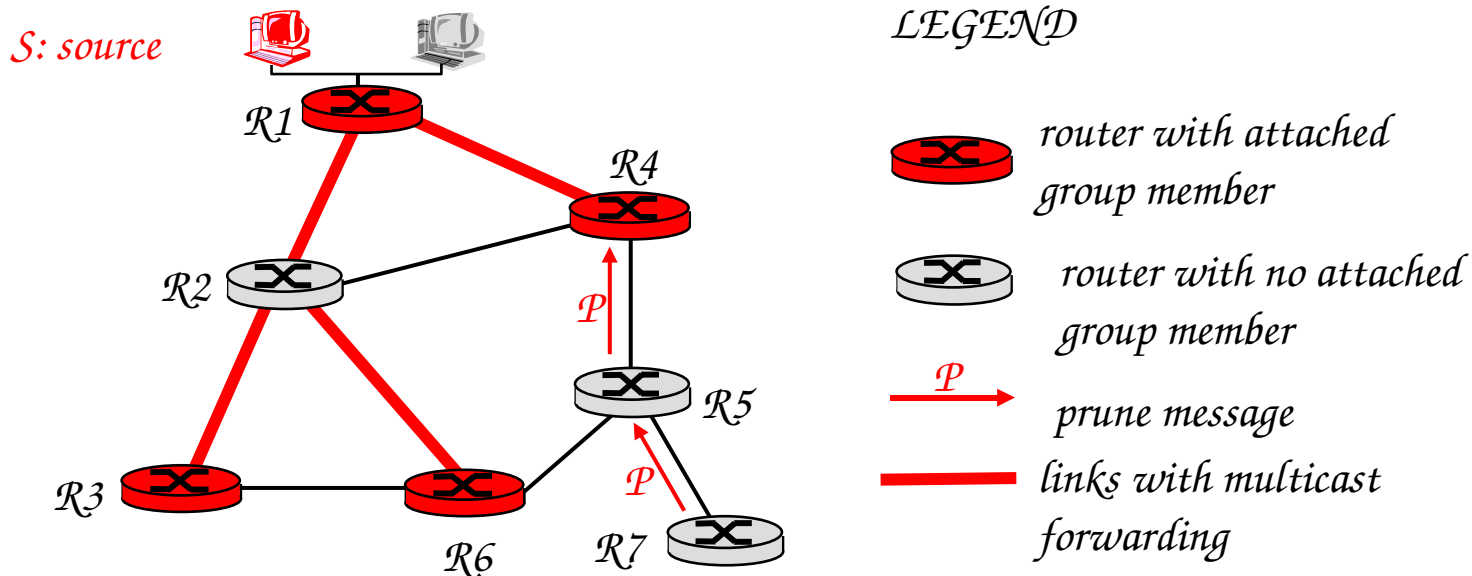
Reverse Path Forwarding: example



- 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



Shared-Tree: Steiner Tree

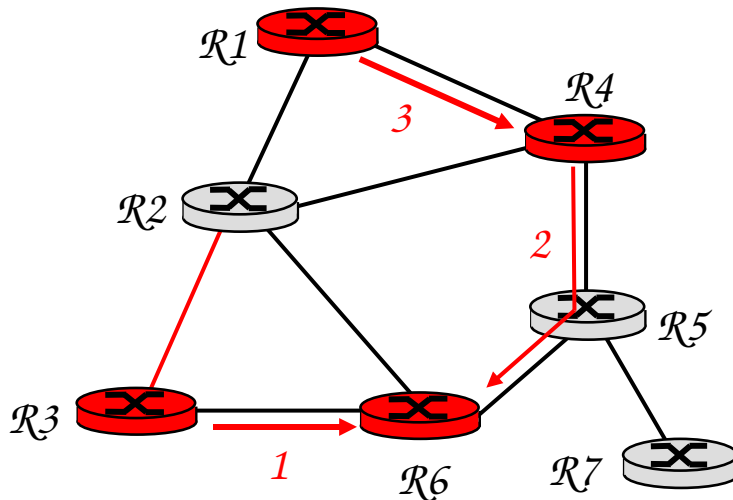
- ❑ *Steiner Tree*: minimum cost tree connecting all routers with attached group members
- ❑ problem is \mathcal{NP} -complete
- ❑ excellent heuristics exists
- ❑ not used in practice:
 - computational complexity
 - information about entire network needed
 - monolithic: rerun whenever a router needs 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 R_6 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

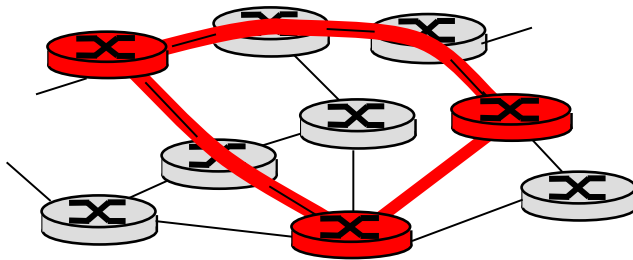
- *DVMRP: distance vector multicast routing protocol, RFC1075*
- *flood and prune: reverse path forwarding, source-based tree*
 - *RPF tree based on DVMRP's own routing tables constructed by communicating DVMRP routers*
 - *no assumptions about underlying unicast*
 - *initial datagram to mcast group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

DVMRP: continued...

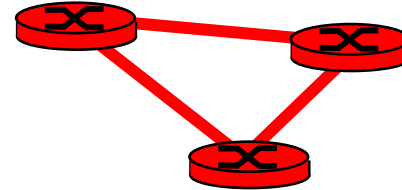
- soft state: DVMRP router periodically (1 min.) “forgets”
branches are pruned:
 - mcast data again flows down unpruned branch
 - downstream router: reprune or else continue to receive data
- routers can quickly regraft to tree
 - following IGMF join at leaf
- odds and ends
 - commonly implemented in 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:

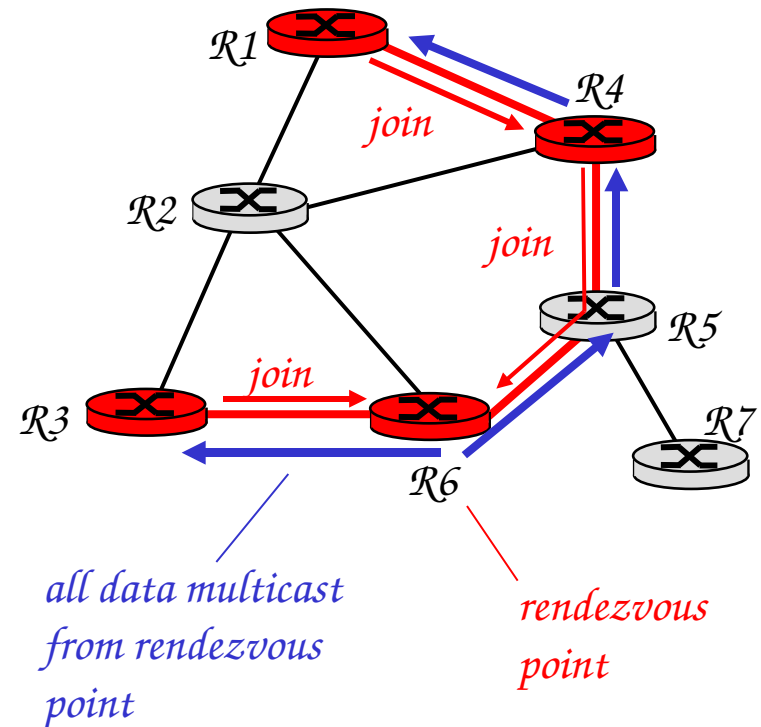
- *no membership until routers explicitly join*
- *receiver-driven construction of mcast tree (e.g., center-based)*
- *bandwidth and non-group-router processing conservative*

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 \mathcal{RP} , which distributes down \mathcal{RP} -rooted tree
- \mathcal{RP} can extend mcast tree upstream to source
- \mathcal{RP} can send stop msg if no attached receivers
 - “no one is listening!”

