

Chapter 4

Network Layer

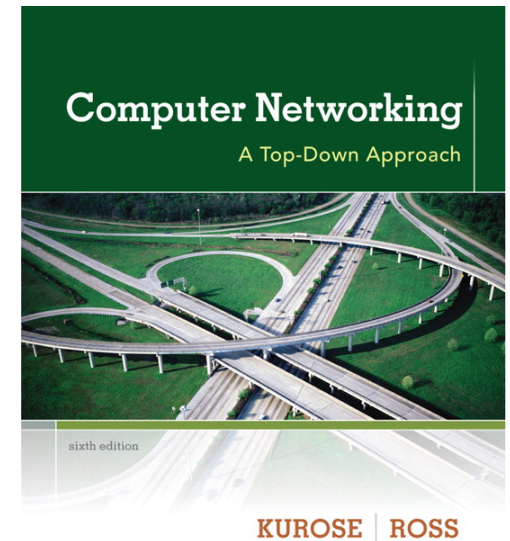
A note on the use of these ppt slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you see the animations; and 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) that you mention their source (after all, we'd like people to use our book!)
- ❖ If you post any slides 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 1996-2013
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)
 - broadcast, multicast
- ❖ instantiation, implementation in the Internet

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

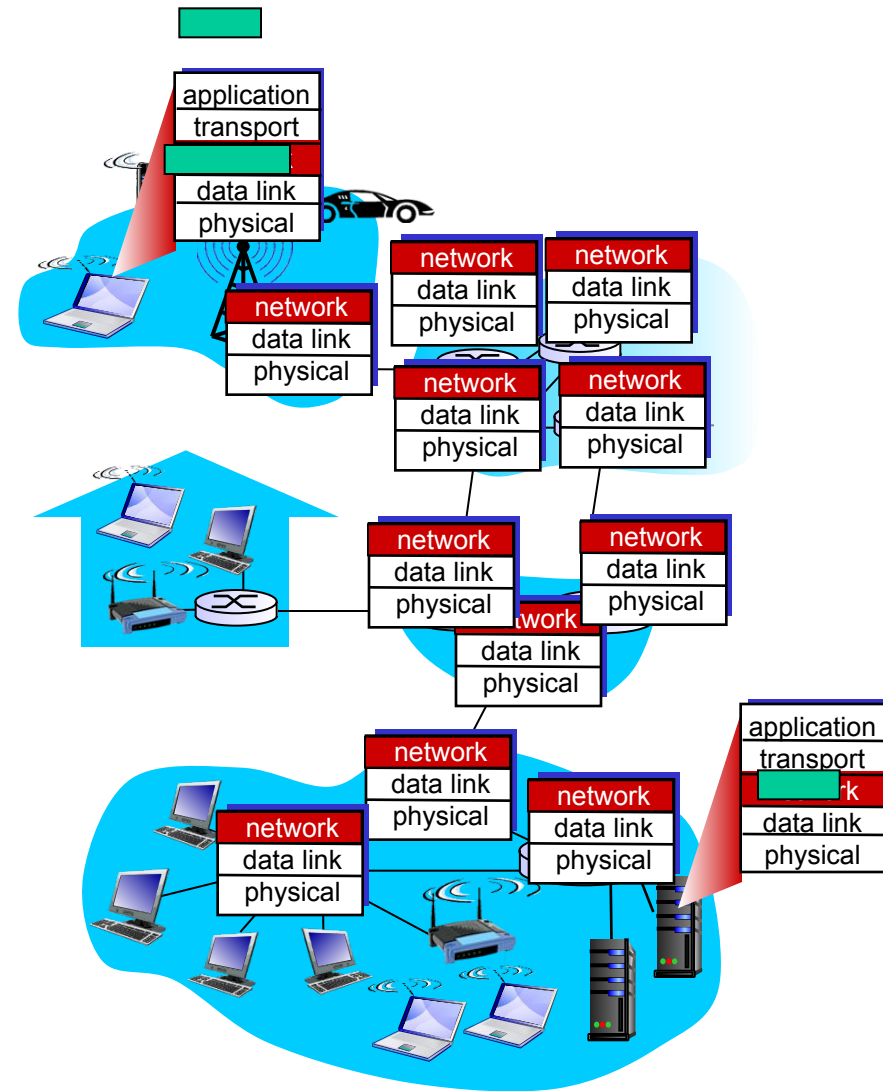
4.6 routing in the Internet

- 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 receiving side, delivers segments to transport layer
- ❖ network layer protocols in *every* host, router
- ❖ router examines header fields in all IP datagrams passing through it



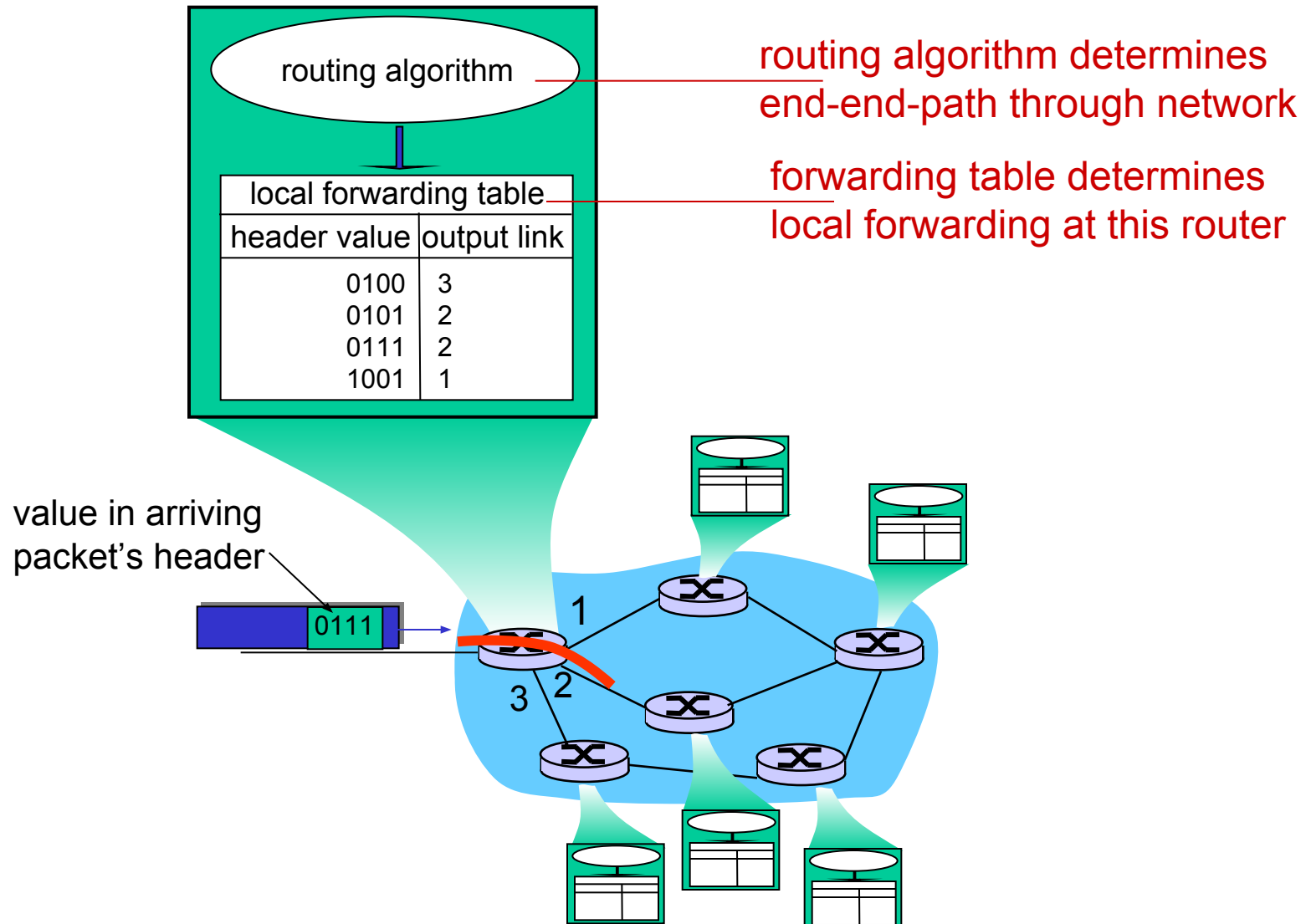
Two key network-layer functions

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

analogy:

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

Interplay between routing and forwarding



Connection setup

- ❖ 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/
guaranteed minimum jitter
- ❖ Security service

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: outline

4.1 introduction

4.2 virtual circuit and datagram networks

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

Connection, connection-less service

- ❖ *datagram* network provides network-layer *connectionless* service
- ❖ *virtual-circuit* network provides network-layer *connection* service
- ❖ analogous to TCP/UDP connection-oriented / connectionless transport-layer services, but:
 - *service*: host-to-host
 - *no choice*: network provides one or the other
 - *implementation*: in network core

Virtual circuits

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

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

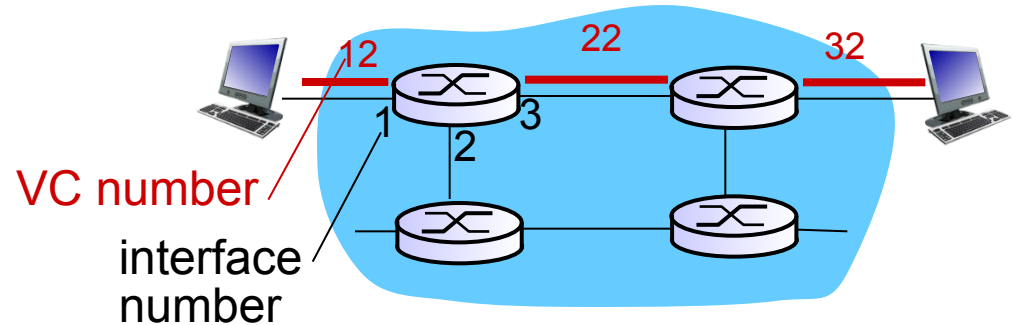
- ❖ call setup, teardown for each call *before* data can flow
- ❖ each packet carries VC identifier (not destination host address)
- ❖ *every* router on source-dest path maintains “state” for each passing connection
- ❖ link, router resources (bandwidth, buffers) may be *allocated* to VC (dedicated resources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

VC implementation

a VC consists of:

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

VC forwarding table



*forwarding table
in*

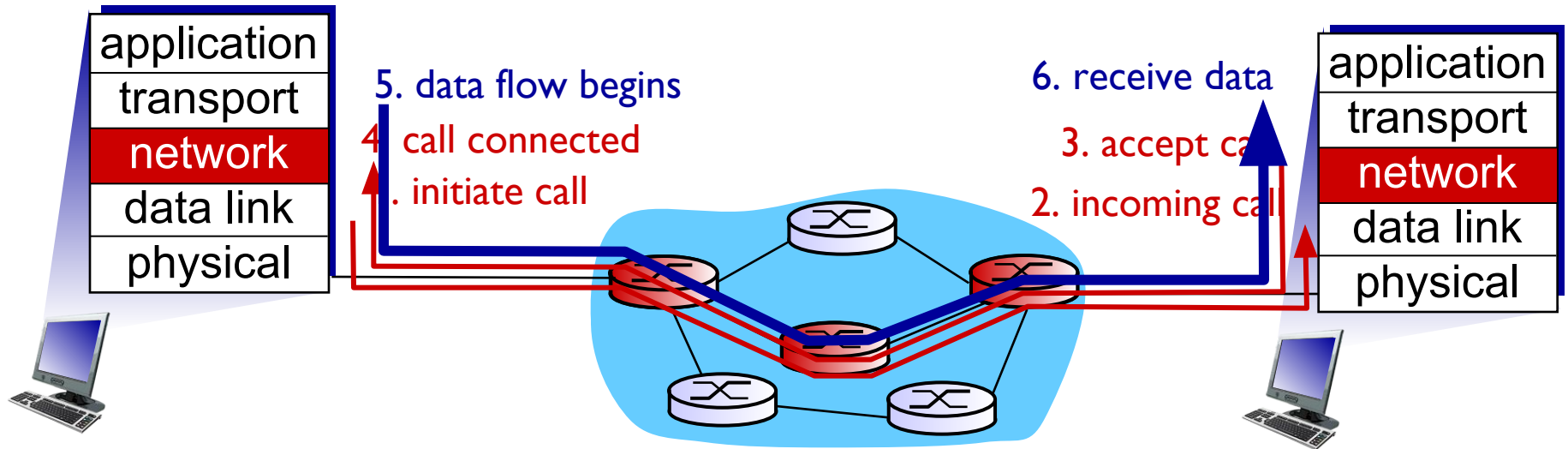
*northwest
router:*

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

*VC routers maintain connection state
information!*

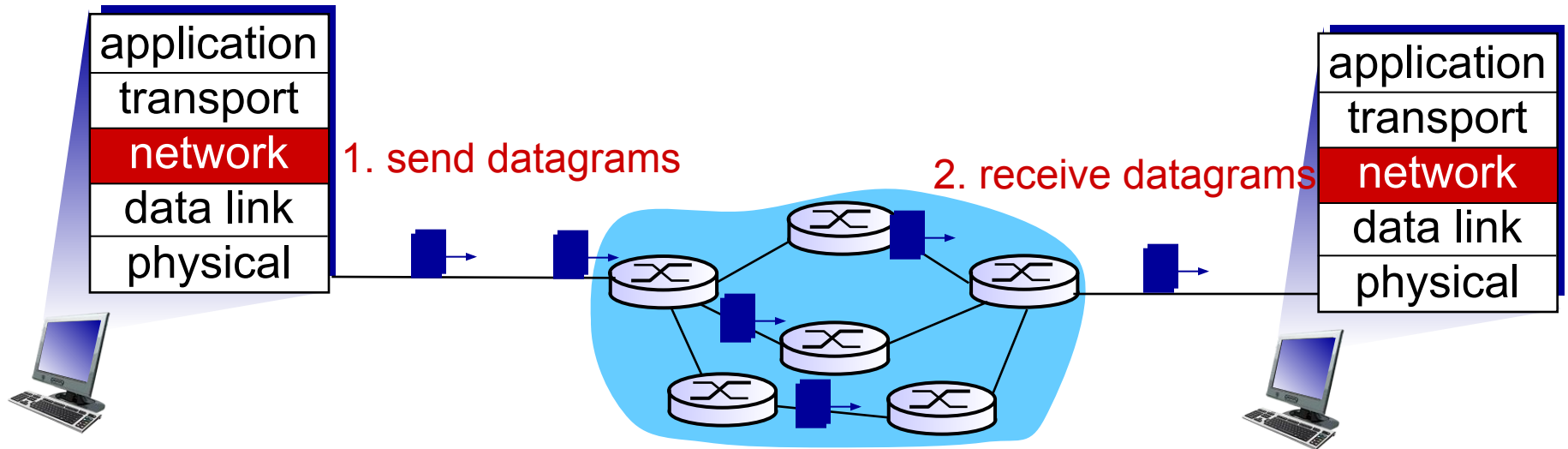
Virtual circuits: signaling protocols

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

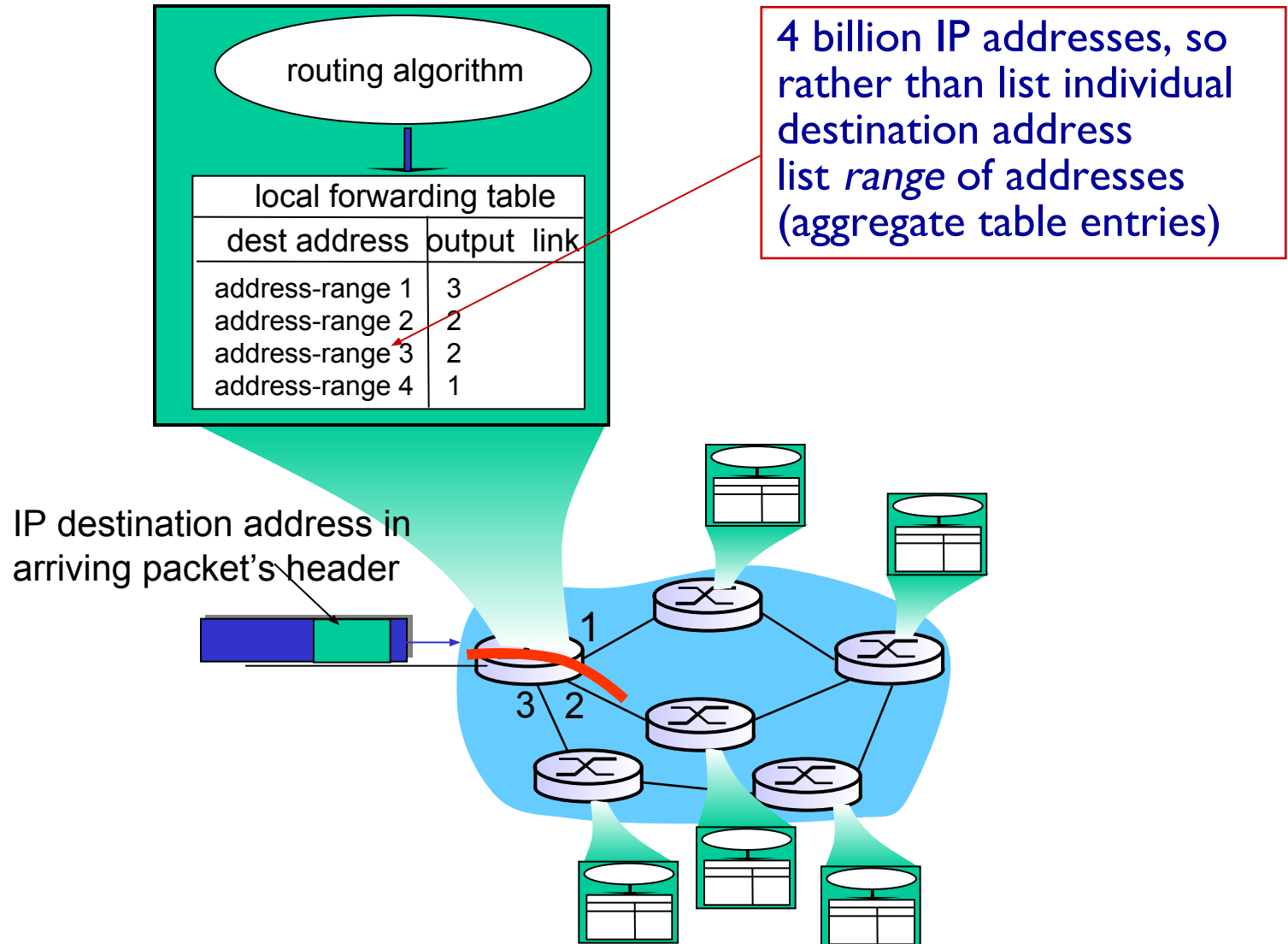


Datagram networks

- ❖ 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



Datagram forwarding table



Datagram forwarding table

Destination Address Range	Link Interface
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

Q: but what happens if ranges don't divide up so nicely?

Longest prefix matching

longest prefix

matching when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

Destination Address Range	Link interface
11001000 00010111 00010*** *****	0
11001000 00010111 00011000 *****	1
11001000 00010111 00011*** *****	2
otherwise	3

examples:

DA: 11001000 00010111 00010110 10100001

which interface?

DA: 11001000 00010111 00011000 10101010

which interface?

Datagram or VC network: why?

Internet (datagram)

- ❖ data exchange among computers
 - “elastic” service, no strict timing req.
- ❖ many link types
 - different characteristics
 - uniform service difficult
- ❖ “smart” end systems (computers)
 - can adapt, perform control, error recovery
 - ***simple inside network, complexity at “edge”***

ATM (VC)

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

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

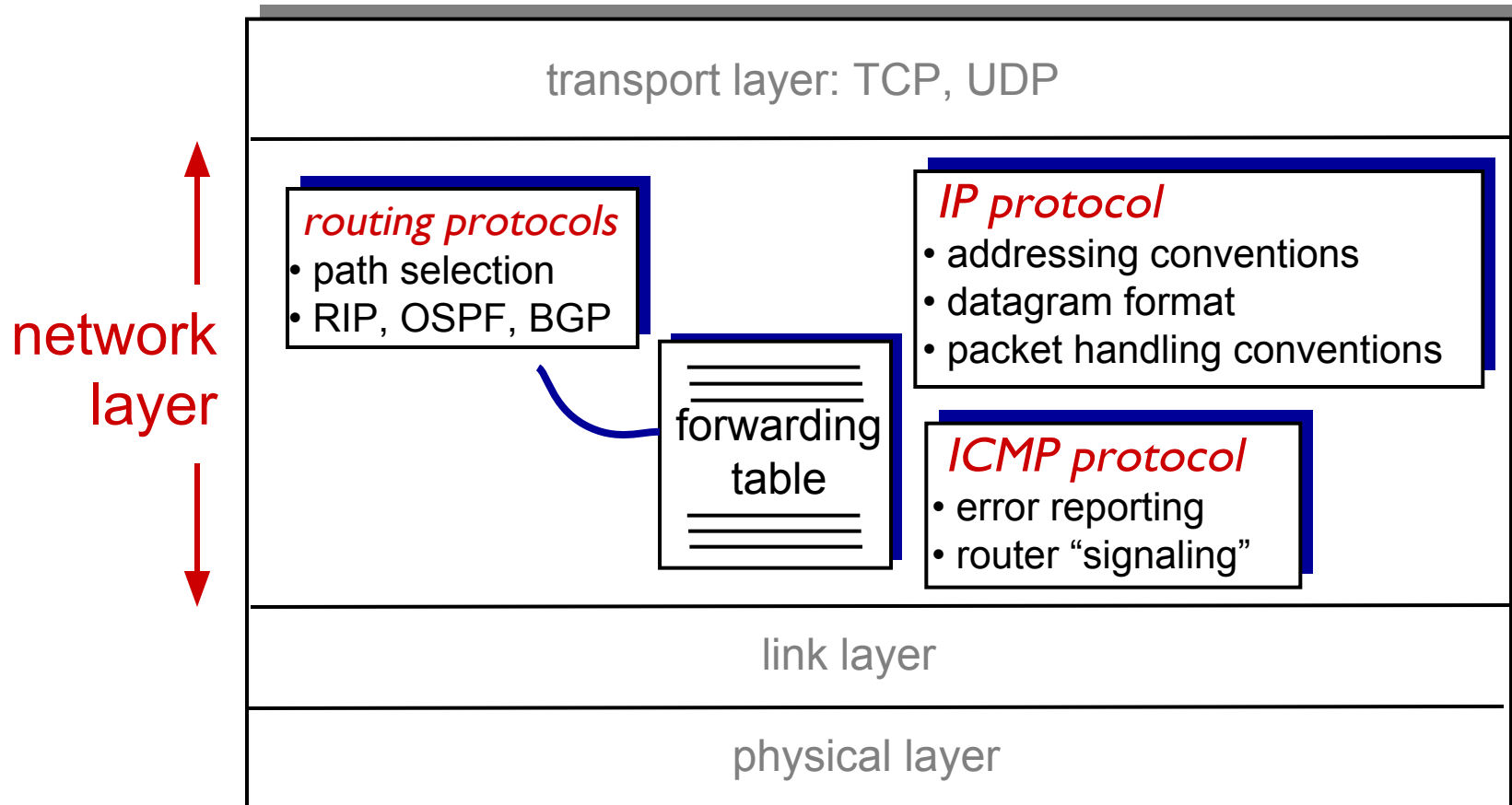
4.6 routing in the Internet

- RIP
- OSPF
- BGP

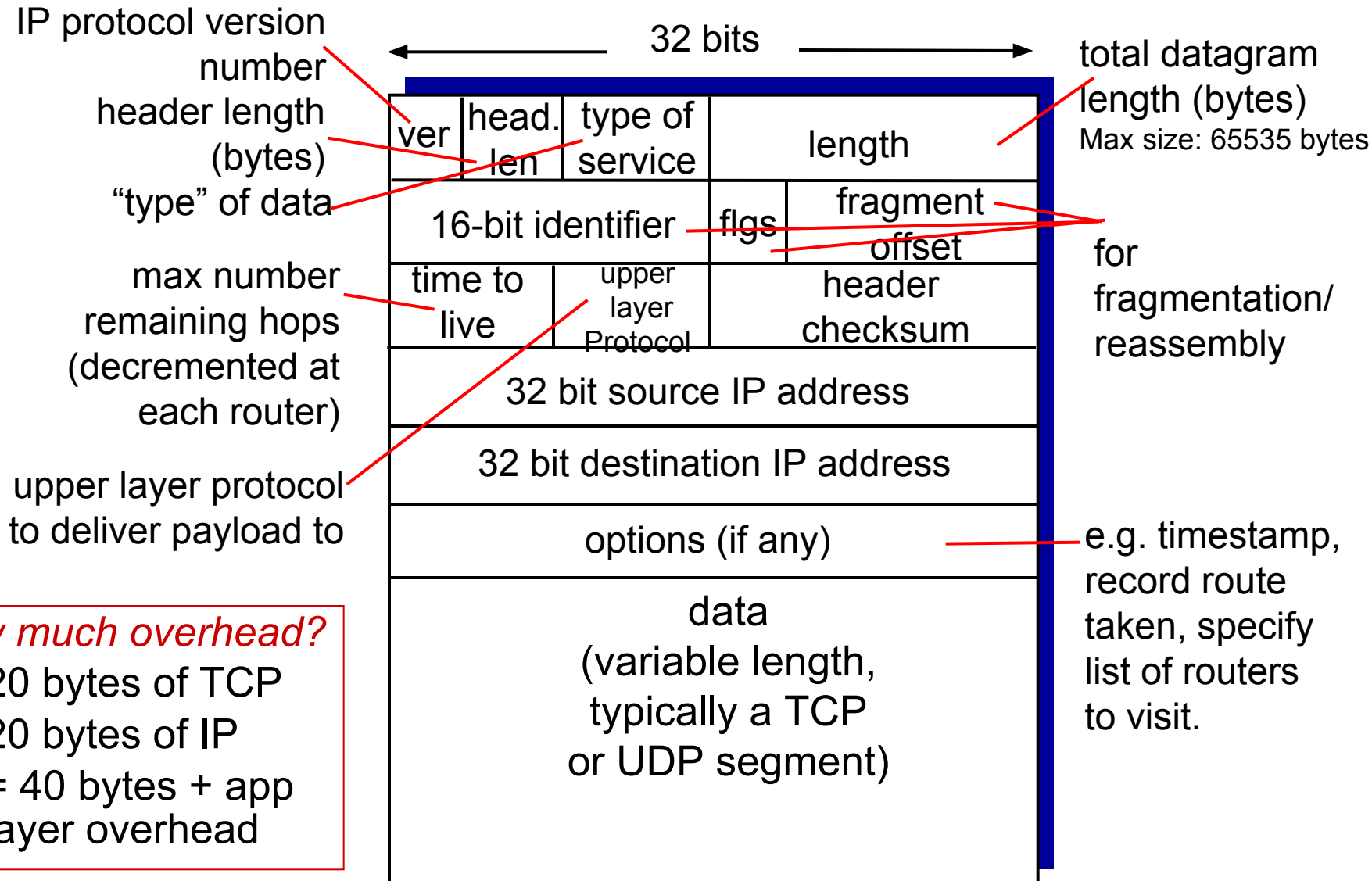
4.7 broadcast and multicast routing

The Internet network layer

host, router network layer functions:

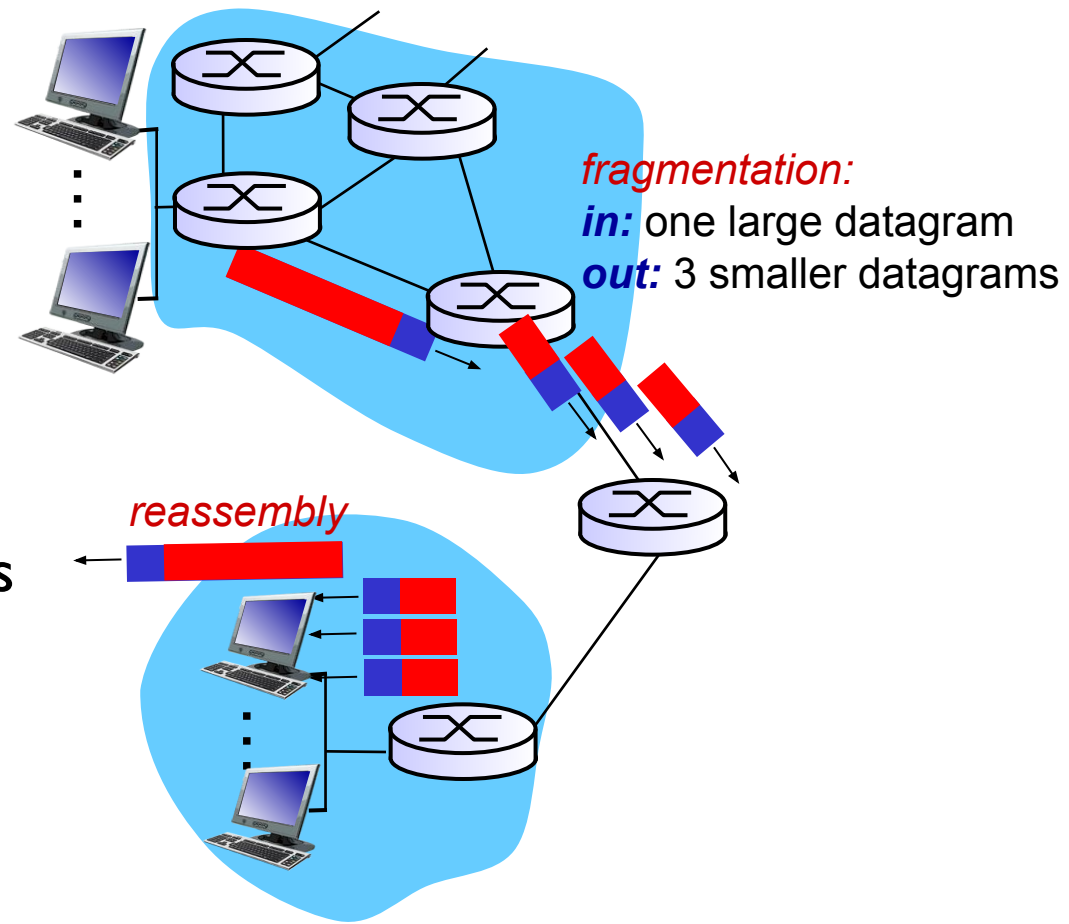


IP datagram format



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, 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: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

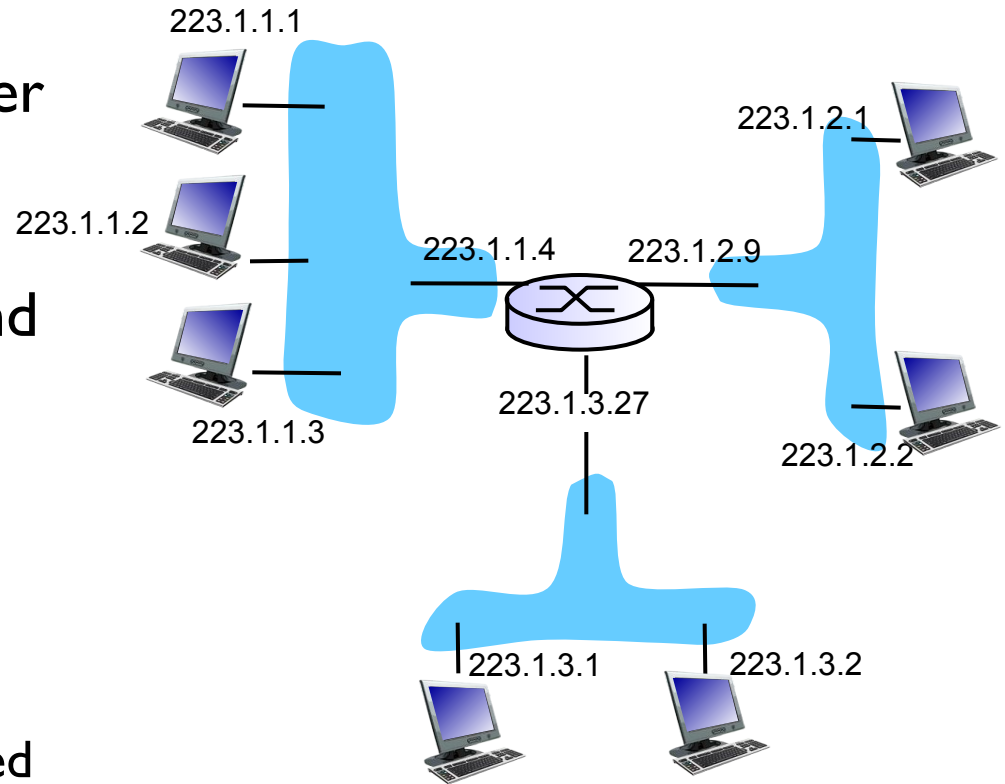
4.6 routing in the Internet

- 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/ Boundary between host and physical link
 - router's typically have multiple interfaces
 - host typically has one or two interfaces (e.g., wired Ethernet, wireless 802.11)
- ❖ **IP addresses associated with each interface**
- ❖ **Written in dotted-decimal notations**



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

IP addressing: introduction

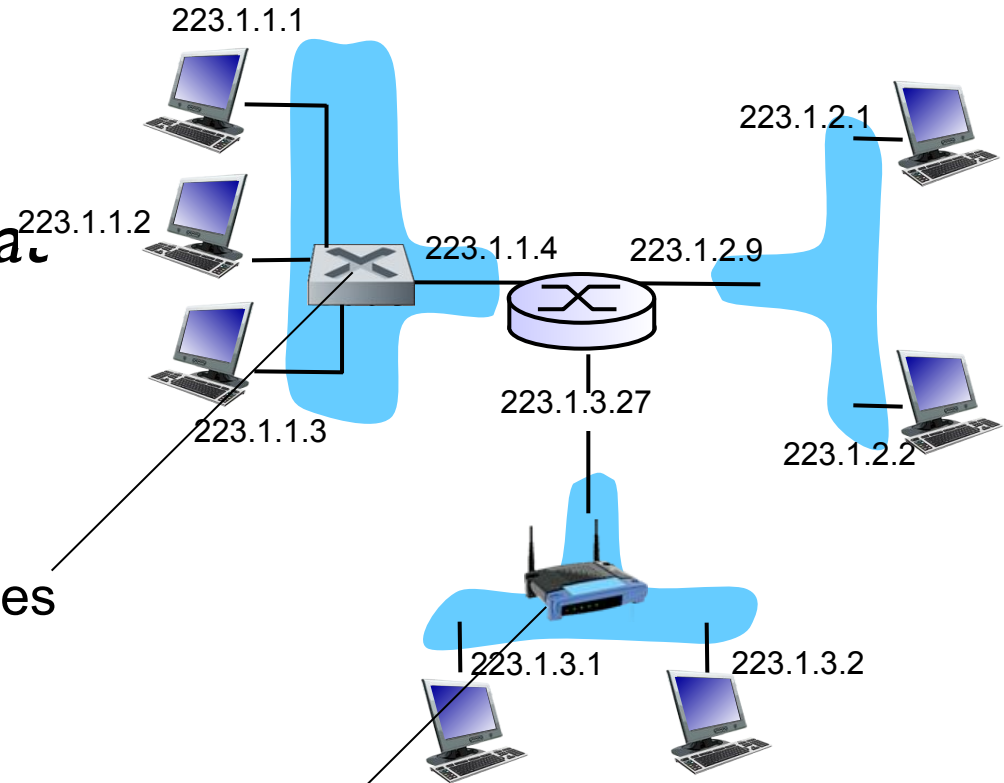
Q: how are interfaces actually connected?

A: we'll learn about that in chapter 5, 6.

A: wired Ethernet interfaces connected by Ethernet switches

For now: don't need to worry about how one interface is connected to another (with no intervening router)

A: wireless WiFi interfaces connected by WiFi base station



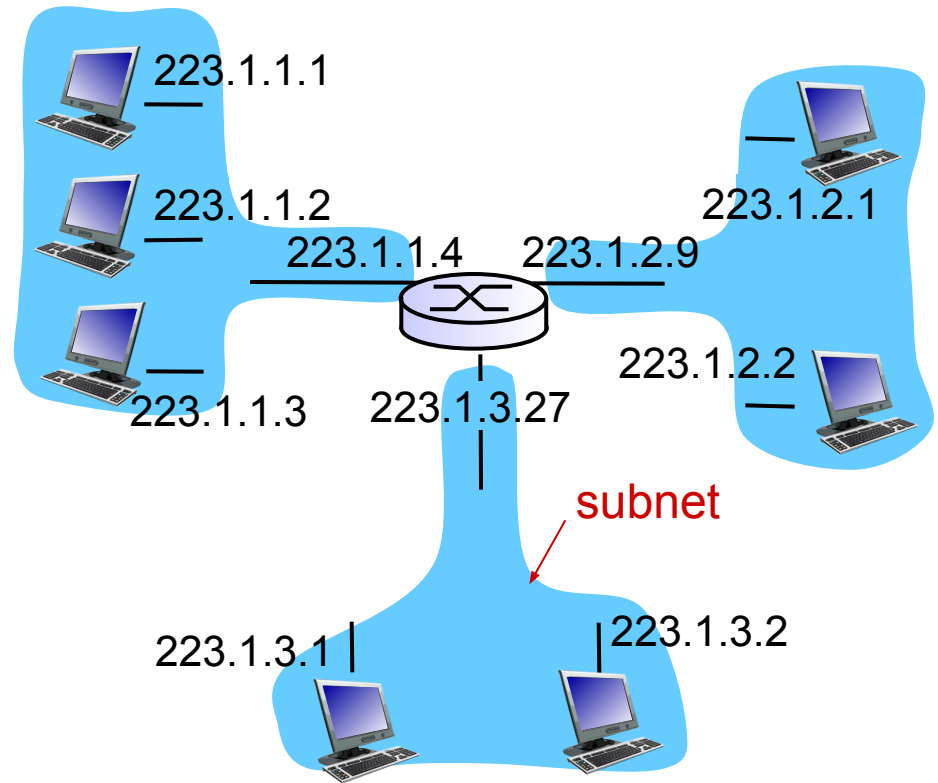
Subnets

❖ 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*

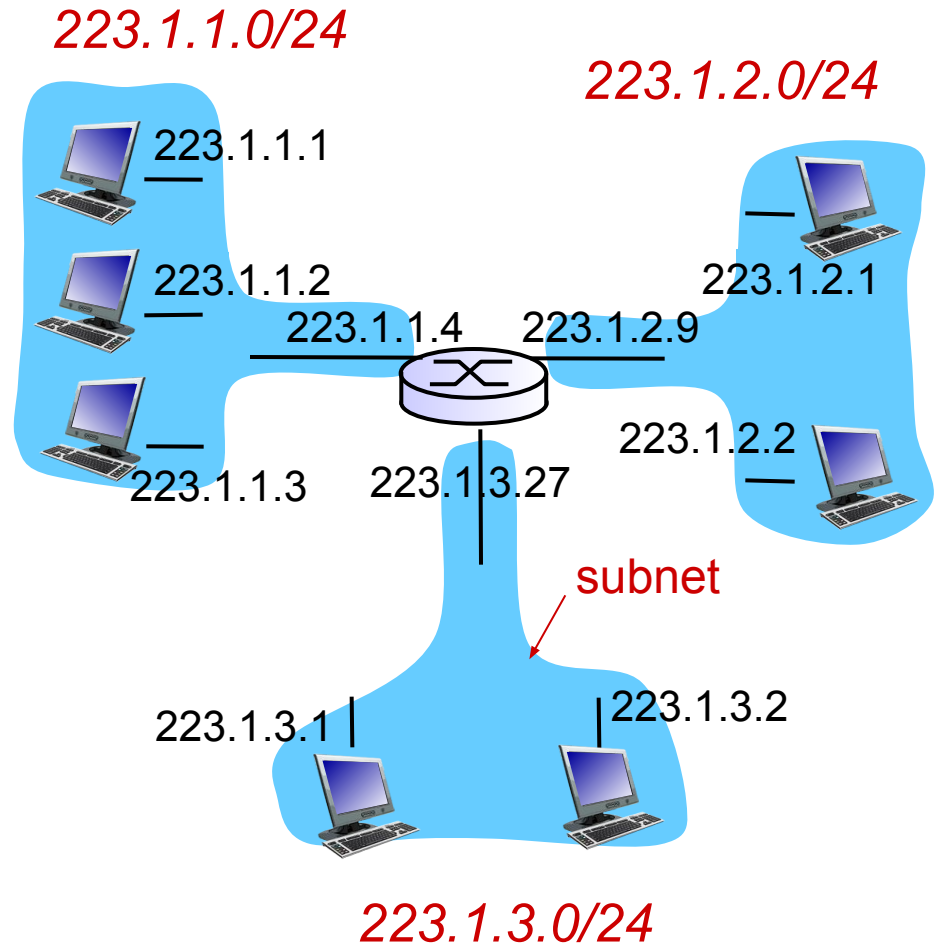


network consisting of 3 subnets

Subnets

recipe

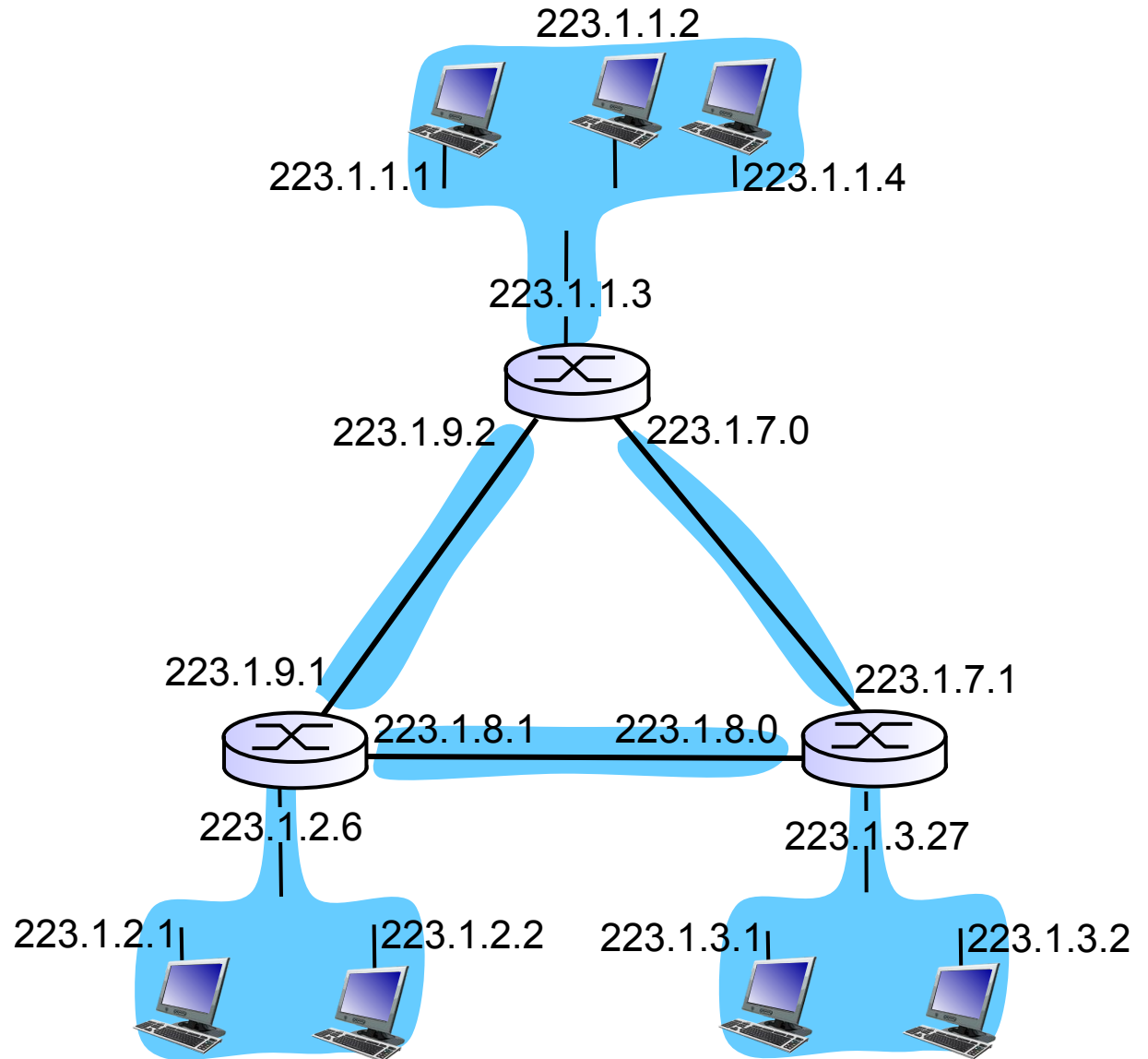
- ❖ to determine the subnets, detach each interface from its host or router, creating islands of isolated networks
- ❖ each isolated network is called a *subnet*



subnet mask: /24

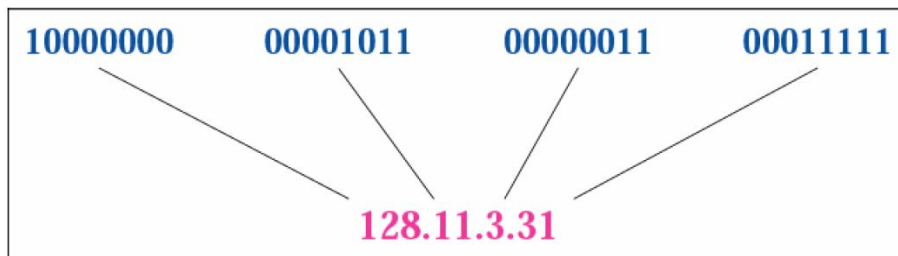
Subnets

how many?



IP addressing

- Notation
 - Binary Notation:
 - IP address is displayed as 32 bits.
 - 01110101 10010101 00011101 11101010
 - Dotted-Decimal Notation
 - IP addresses are written in decimal form with a decimal point (dot) separating the bytes.
 - 128.11.3.31
 - Hexadecimal Notation



Binary Notation

0111 0101 1001 0101 0001 1101 1110 1010
75 95 1D EA
0x75951DEA

Hexadecimal Notation

IP addressing

Find the error, if any, in the following IPv4 addresses.

- a) 111.56.045.78
- b) 221.34.7.8.20
- c) 75.45.301.14
- d) 11100010.23.14.67

Solution

- a) There must be no leading zero (045).
- b) There can be no more than four numbers.
- c) Each number needs to be less than or equal to 255.
- d) A mixture of binary notation and dotted-decimal notation is not allowed.

IP addressing (Classful Addressing)

- IP addresses were divided into 5 classes:
A,B,C,D and E
 - This is the original scheme known as classful addressing
 - From mid-90's, classless addressing is introduced
 - However, classful addressing is still used

	First byte	Second byte	Third byte	Fourth byte
Class A	0			
Class B	10			
Class C	110			
Class D	1110			
Class E	1111			

IP addressing (Classful Addressing)

Exercise:

1. Find the class of these IP addresses:
a) 11000001 10000011 00011011 11111111
b) 10000001 10000011 00011011 11111111
2. How many class B addresses are there altogether?
3. What is the range of class B addresses? Answer this by giving the first and last class B addresses in dotted decimal notation.

IP addressing (Classful Addressing)

Solution:

1. a) First 3 bits are 110 → Class C.
b) First 2 bits are 10 → Class B.
2. Class B addresses: the first two bits are 10 then followed by 30 bits of 1/0
→ 2^{30} addresses
3. The first and last class B addresses in binary are:
10000000 00000000 00000000 00000000
10111111 11111111 11111111 11111111
In dotted decimal notation, they are:
128.0.0.0 and 191.255.255.255

IP addressing (Classful Addressing)

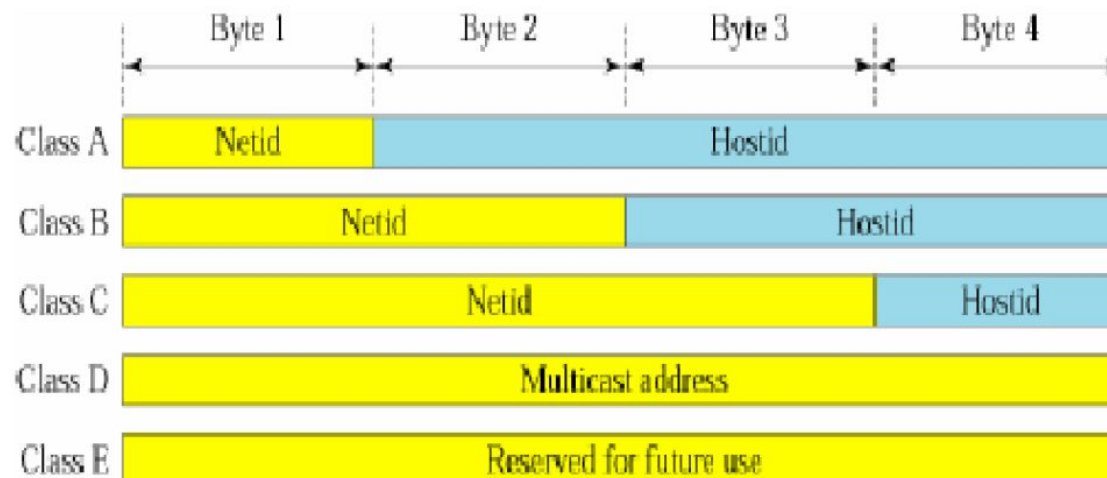
Netid and Hostid

- IP addresses in classes A,B and C are divided into netid and hostid
 - Netid: Identifying the network
 - Hostid: Identifying a host within the network
- Hosts within a network
 - Have the same netid
 - But different hostid

IP addressing (Classful Addressing)

Finding the class in decimal notation

	First byte	Second byte	Third byte	Fourth byte
Class A	0 to 127			
Class B	128 to 191			
Class C	192 to 223			
Class D	224 to 239			
Class E	240 to 255			



IP addressing (Classful Addressing)

Classes and Blocks

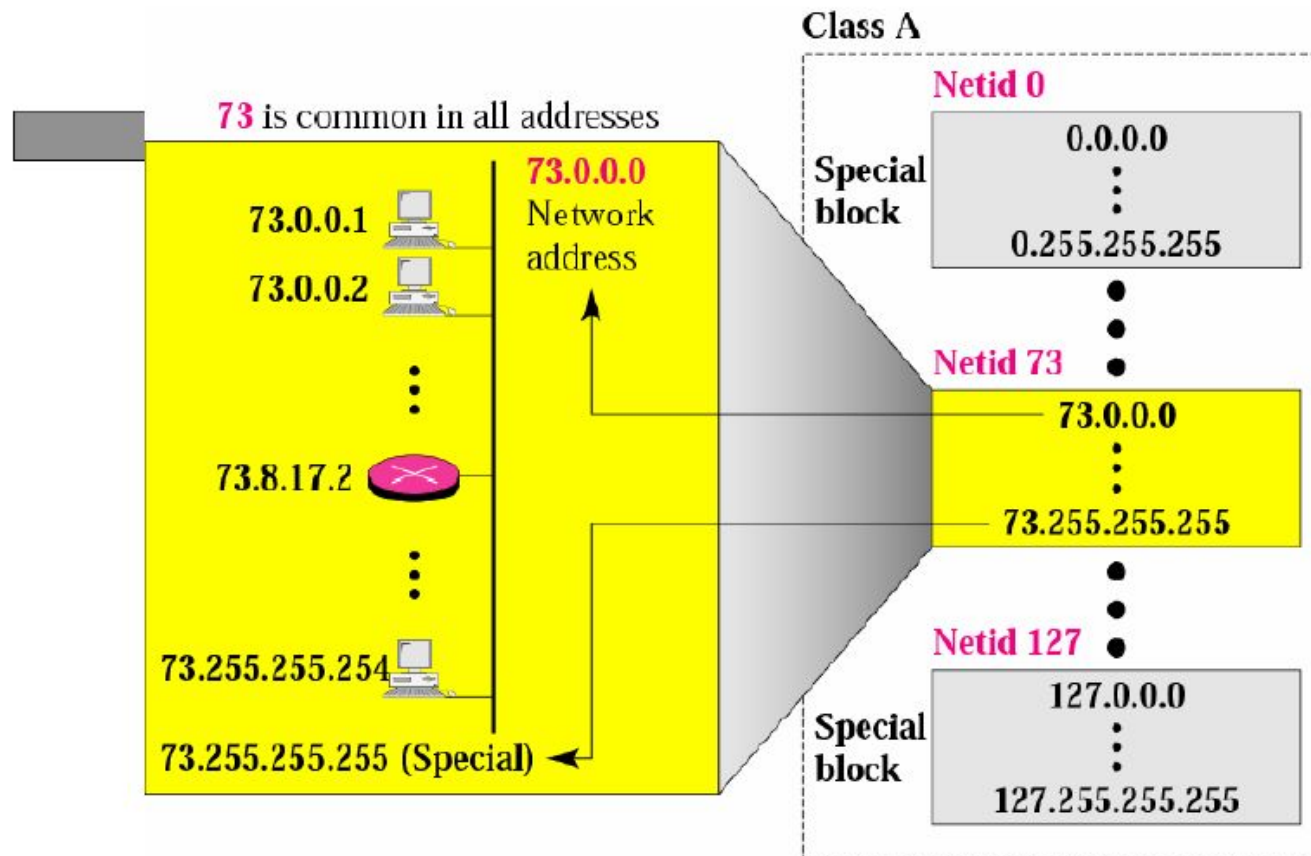
All addresses within a block have the same netid.
For the class A block with netid = 0, the addresses in the block are:

netid		hostid		
00000000	00000000	00000000	00000000	0.0.0.0
00000000	00000000	00000000	00000001	0.0.0.1
00000000	00000000	00000000	00000010	0.0.0.2
...				
00000000	11111111	11111111	11111111	0.255.255.255

- Class A is divided into 128 blocks
 - Each block has a different netid
 - 1st block: 0.0.0.0 to 0.255.255.255 (netid = 0)
 - 2nd block: 1.0.0.0 to 1.255.255.255 (netid = 1)
 - Last block: 127.0.0.0 to 127.255.255.255 (netid = 127)
- *Network address*: the first address of the block

IP addressing (Classful Addressing)

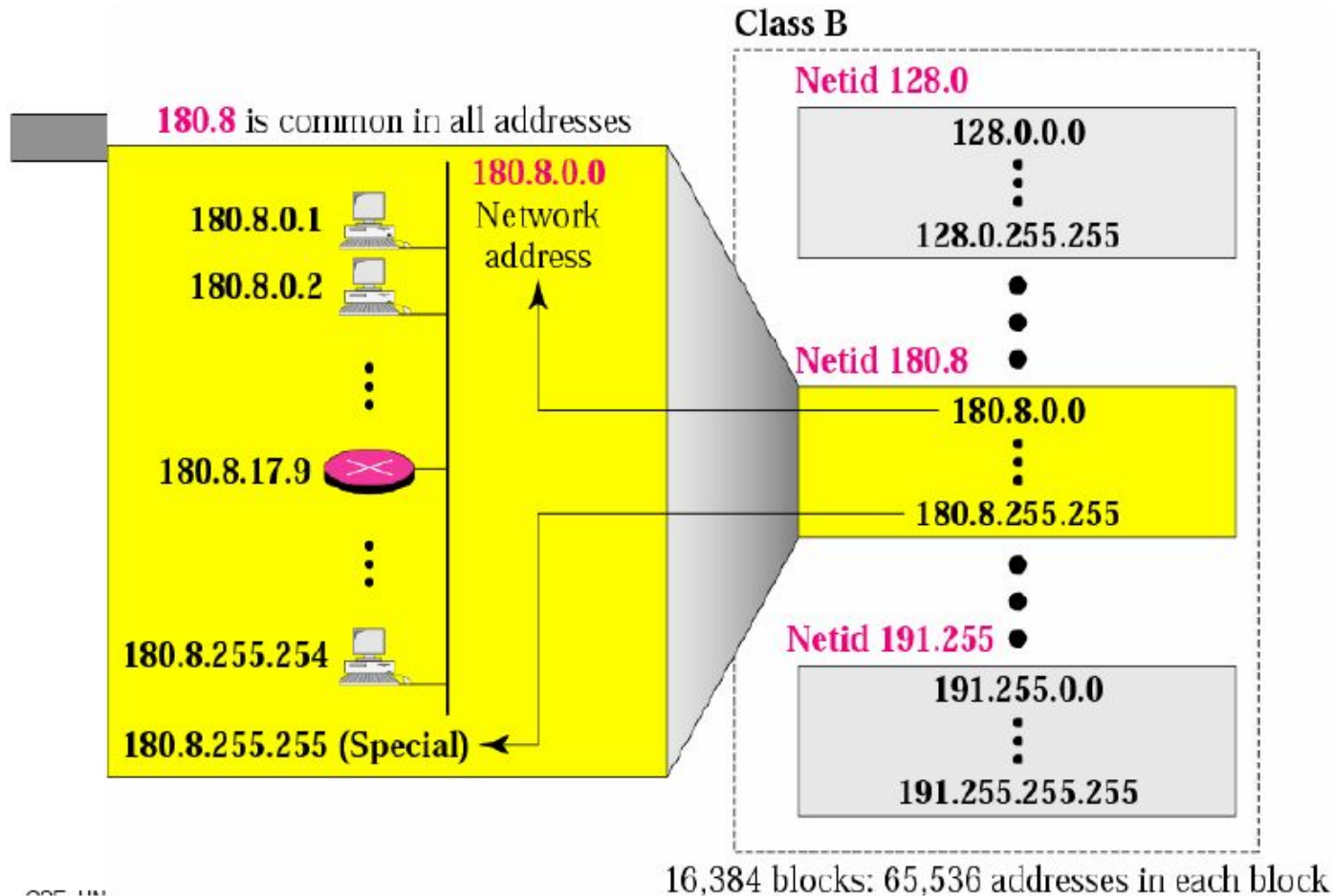
Blocks in Class A



128 blocks: 16,777,216 addresses in each block

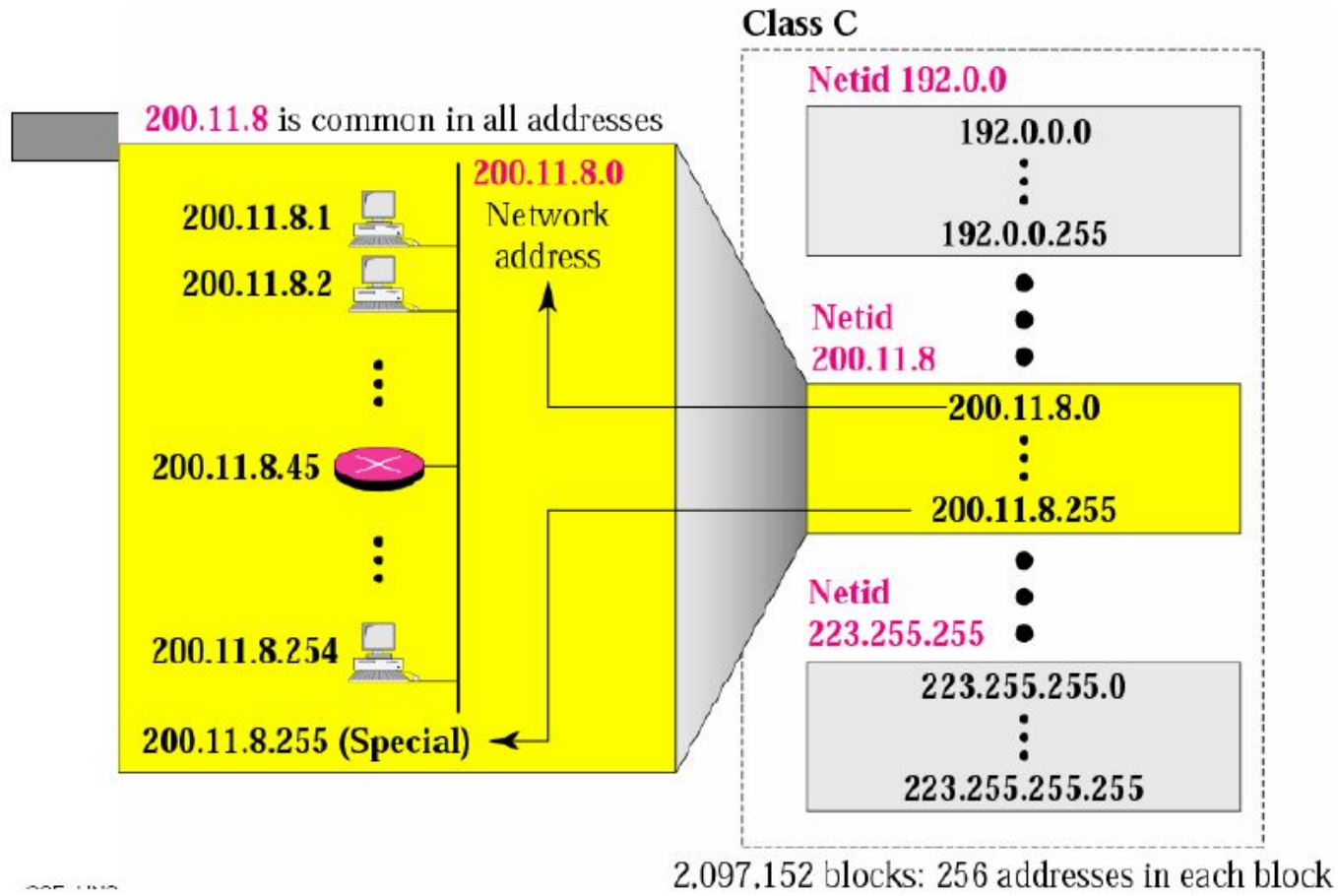
IP addressing (Classful Addressing)

Blocks in Class B



IP addressing (Classful Addressing)

Blocks in Class C



IP addressing (Classful Addressing)

Uses of Addresses

- Classes A, B and C addresses can be assigned to hosts, router ports etc
 - They are also known as unicast addresses
- Class D addresses are for multicast
 - Multicast: One sender, multiple recipients
- Class E addresses are reserved for special purposes

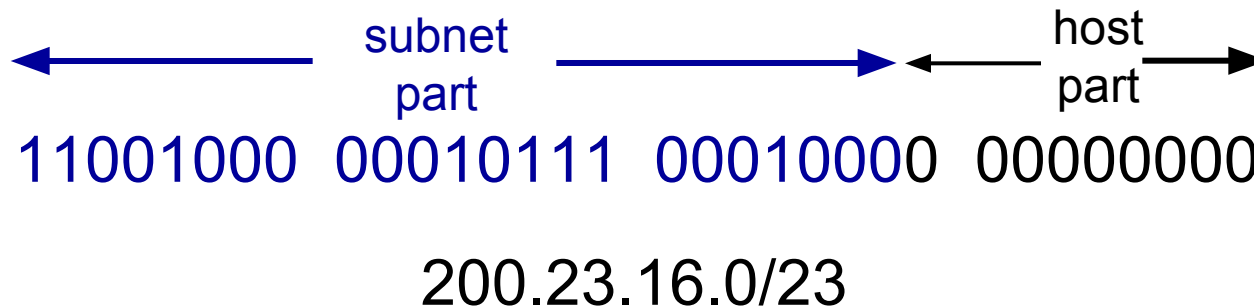
Network Addresses

- The network address is the first address in the block
- The network address defines the network to the rest of the Internet
 - Routers route packets based on network address
- Given the network address, we can find the class of the address and the range of the address in the block

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



IP addressing: Classless Addressing

Note

Limitation of Classful Addressing

In classful addressing, a large part of the available addresses were wasted.

Classful addressing, which is almost obsolete, is replaced with classless addressing.

IP addressing: CIDR

MASK

- ❖ How do we know how many bits represent the network portion and how many bits represent the host portion?
- ❖ The **prefix length** is the number of bits in the address that gives us the network portion. For example, in 172.16.4.0 /24, the /24 is the prefix length - it tells us that the first 24 bits are the network address.
- ❖ another entity that is used to specify the network portion is called **subnet mask**. The subnet mask consists of 32 bits, just as the address does, and uses 1s and 0s to indicate which bits of the address are network bits and which bits are host bits.

IP addressing (CIDR)

Example

A small organization has the network 205.16.37.24/29.
What is the range of the block?

Solution

The binary representation of the given address is

11001111 00010000 00100101 00011000

For first addresses, we set $(32-29) = 3$ rightmost bits to 0, we get

11001111 00010000 00100101 00011000

For last addresses, we set $(32 - 29) = 3$ rightmost bits to 1, we get

11001111 00010000 00100101 00011111

The value of n is 29, which means that the range of the block is 2^{32-29} or 8.

So, There are only 8 addresses in this block

IP addressing (CIDR)

Example

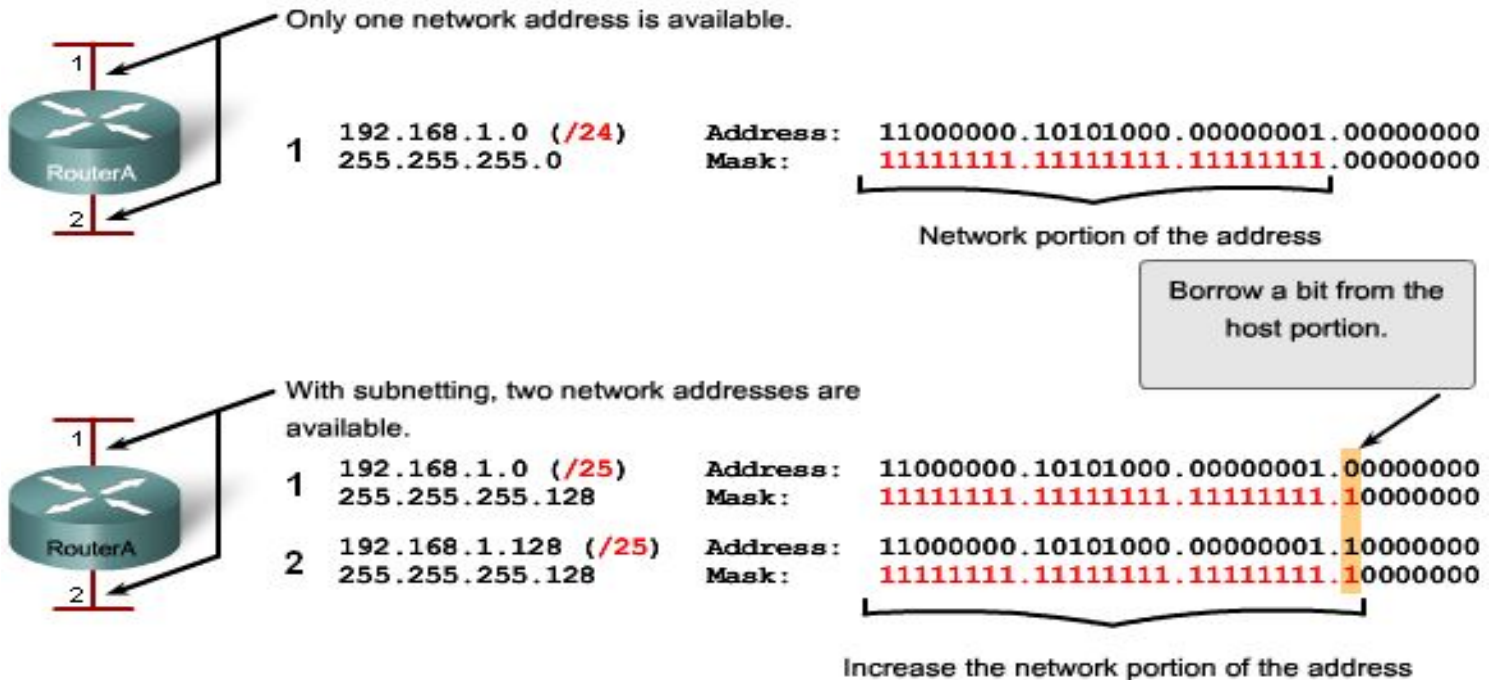
What is the network address if one of the addresses is 167.199.170.82/27?

Solution

The prefix length is 27, which means that we must keep the first 27 bits as is and change the remaining bits (5) to 0s. The 5 bits affect only the last byte. The last byte is 01010010. Changing the last 5 bits to 0s, we get 01000000 or 64. The network address is 167.199.170.64/27.

IP addressing (CIDR)

- ❖ Subnetting allows for creating multiple logical networks from a single address block
- ❖ Subnets are created using one or more of the host bits as network bits
 - done by extending the mask to borrow some of the bits from the host portion to create additional network bits



IP addressing (CIDR)

Calculating Subnets and Hosts

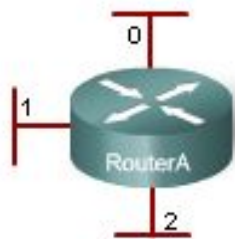
- ❖ The number of subnets is calculated using 2^n , where n is the number of bits borrowed
 - $2^1 = 2$ subnets
 - the more bits borrowed, the more subnets can be defined
- ❖ The number of useable hosts per subnet is calculated using $2^h - 2$ where h is the number of host bits left
 - $2^7 - 2 = 126$ useable hosts per subnet
 - with each bit borrowed, fewer host addresses are available per subnet

Subnet	Network address	Host range	Broadcast address
0	192.168.1.0/25	192.168.1.1 - 192.168.1.126	192.168.1.127
1	192.168.1.128/25	192.168.1.129 - 192.168.1.254	192.168.1.255

IP addressing (CIDR)

Subnetting Example 1

- ❖ Need to borrow a minimum of 2 host bits.
 - $2^2 = 4$ subnets



Borrowing Bits for Subnets

-	192.168.1.0 (/24)	Address:	11000000.10101000.00000001.00000000
	255.255.255.0	Mask:	11111111.11111111.11111111.00000000
0	192.168.1.0 (/26)	Address:	11000000.10101000.00000001.00000000
	255.255.255.192	Mask:	11111111.11111111.11111111.11000000
1	192.168.1.64 (/26)	Address:	11000000.10101000.00000001.01000000
	255.255.255.192	Mask:	11111111.11111111.11111111.11000000
2	192.168.1.128 (/26)	Address:	11000000.10101000.00000001.10000000
	255.255.255.192	Mask:	11111111.11111111.11111111.11000000
3	192.168.1.192 (/26)	Address:	11000000.10101000.00000001.11000000
	255.255.255.192	Mask:	11111111.11111111.11111111.11000000

Two bits are borrowed to provide four subnets.

Unused address in this example.

A 1 in these positions in the mask means that these values are part of the network address.

More subnets are available, but fewer addresses are available per subnet.

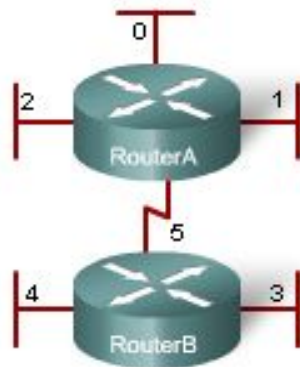
IP addressing (CIDR)

Subnetting Example 2

- ❖ Need to borrow a minimum of 3 host bits.
 - $2^3 = 8$ subnets

Borrowing Bits for Subnets

Start with this address	-	192.168.1.0 (/24) 255.255.255.0	Address: 11000000.10101000.00000001.00000000 Mask: 11111111.11111111.11111111.00000000
Make 8 subnets	0	192.168.1.0 (/27) 255.255.255.224	Address: 11000000.10101000.00000001.00000000 Mask: 11111111.11111111.11111111.11100000
	1	192.168.1.32 (/27) 255.255.255.224	Address: 11000000.10101000.00000001.00100000 Mask: 11111111.11111111.11111111.11100000
	2	192.168.1.64 (/27) 255.255.255.224	Address: 11000000.10101000.00000001.01000000 Mask: 11111111.11111111.11111111.11100000
	3	192.168.1.96 (/27) 255.255.255.224	Address: 11000000.10101000.00000001.01100000 Mask: 11111111.11111111.11111111.11100000
	4	192.168.1.128 (/27) 255.255.255.224	Address: 11000000.10101000.00000001.10000000 Mask: 11111111.11111111.11111111.11100000
	5	192.168.1.160 (/27) 255.255.255.224	Address: 11000000.10101000.00000001.10100000 Mask: 11111111.11111111.11111111.11100000
	6	192.168.1.192 (/27) 255.255.255.224	Address: 11000000.10101000.00000001.11000000 Mask: 11111111.11111111.11111111.11100000
	7	192.168.1.224 (/27) 255.255.255.224	Address: 11000000.10101000.00000001.11100000 Mask: 11111111.11111111.11111111.11100000



Three bits are borrowed to provide eight subnets.

IP addressing (CIDR)

Subnetting Classless Network

- A classless network also be subnetted
- Example:

An organization is granted the block 130.34.12.64/26. The organization needs to have four subnets. What are the subnet addresses and the range of addresses for each subnet?

IP addressing (CIDR)

Solution



Since the prefix length is 26. This means the last 6 bits are available as hostid, the total number of addresses in the block is 64 (2^6). If we create four subnets, each subnet will have 16 addresses.

Solution (Continued)

Let us first find the subnet prefix (subnet mask). We need four subnets, which means we need to add two more 1s to the site prefix. The subnet prefix is then /28.

Subnet 1: 130.34.12.64/28 to 130.34.12.79/28.

Subnet 2 : 130.34.12.80/28 to 130.34.12.95/28.

Subnet 3: 130.34.12.96/28 to 130.34.12.111/28.

Subnet 4: 130.34.12.112/28 to 130.34.12.127/28.

IP addressing (CIDR)

Solution:

The binary representation of the given address 130.34.12.64/26 is

10001000 00100010 00001100 01000000

The organization need to have 4 subnets so, need to borrow a minimum of 2 host bits.

$2^2 = 4$ subnets

Subnet 1:

First address: 10001000 00100010 00001100 01000000 or 130.34.12.64/28

Last address: 10001000 00100010 00001100 01001111 or 130.34.12.79/28

Range: 130.34.12.64/28 - 130.34.12.79/28

Subnet 2:

First address: 10001000 00100010 00001100 01010000 or 130.34.12.80/28

Last address: 10001000 00100010 00001100 01011111 or 130.34.12.95/28

Range: 130.34.12.80/28 - 130.34.12.95/28

Subnet 3:

First address: 10001000 00100010 00001100 01100000 or 130.34.12.96/28

Last address: 10001000 00100010 00001100 01101111 or 130.34.12.111/28

Range: 130.34.12.96/28 - 130.34.12.111/28

Subnet 4:

First address: 10001000 00100010 00001100 01110000 or 130.34.12.112/28

Last address: 10001000 00100010 00001100 01111111 or 130.34.12.127/28

Range: 130.34.12.112/28 - 130.34.12.127/28

So, the range of the addresses for each subnet is $2^{32-28} = 2^4 = 16$

IP addressing (CIDR)

An ISP is granted a block of addresses starting with 190.100.0.0/16 (65,536 addresses). The ISP needs to distribute these addresses to three groups of customers as follows:

- a. The first group has 64 customers; each needs 256 addresses.
- b. The second group has 128 customers; each needs 128 addresses.
- c. The third group has 128 customers; each needs 64 addresses.

Design the subblocks and find out how many addresses are still available after these allocations.

Group 1

For this group, each customer needs 256 addresses. This means that 8 bits are needed to define each host. The prefix length is then $32 - 8 = 24$. The addresses are

<i>1st Customer:</i>	<i>190.100.0.0/24</i>	<i>190.100.0.255/24</i>
<i>2nd Customer:</i>	<i>190.100.1.0/24</i>	<i>190.100.1.255/24</i>
<i>...</i>		
<i>64th Customer:</i>	<i>190.100.63.0/24</i>	<i>190.100.63.255/24</i>
<i>Total = $64 \times 256 = 16,384$</i>		

IP addressing (CIDR)

Group 2

For this group, each customer needs 128 addresses. This means that 7 bits are needed to define each host. The prefix length is then $32 - 7 = 25$. The addresses are

<i>1st Customer:</i>	<i>190.100.64.0/25</i>	<i>190.100.64.127/25</i>
<i>2nd Customer:</i>	<i>190.100.64.128/25</i>	<i>190.100.64.255/25</i>
<i>...</i>		
<i>128th Customer:</i>	<i>190.100.127.128/25</i>	<i>190.100.127.255/25</i>
<i>Total = $128 \times 128 = 16,384$</i>		

Group 3

For this group, each customer needs 64 addresses. This means that 6 bits are needed to each host. The prefix length is then $32 - 6 = 26$. The addresses are

<i>1st Customer:</i>	<i>190.100.128.0/26</i>	<i>190.100.128.63/26</i>
<i>2nd Customer:</i>	<i>190.100.128.64/26</i>	<i>190.100.128.127/26</i>
<i>...</i>		
<i>128th Customer:</i>	<i>190.100.159.192/26</i>	<i>190.100.159.255/26</i>
<i>Total = $128 \times 64 = 8192$</i>		

Number of granted addresses to the ISP: 65,536
Number of allocated addresses by the ISP: 40,960
Number of available addresses: 24,576

IP addresses: how to get one?

Host address can be configured manually.

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

- ❖ hard-coded by system admin in a file
 - Windows:
control-panel->network->configuration->tcp/ip->properties
 - UNIX: /etc/rc.config
- ❖ **DHCP: Dynamic Host Configuration Protocol:**
dynamically get address from a server
 - “plug-and-play”
 - DHCP allows a host to obtain an IP address automatically.

DHCP: Dynamic Host Configuration Protocol

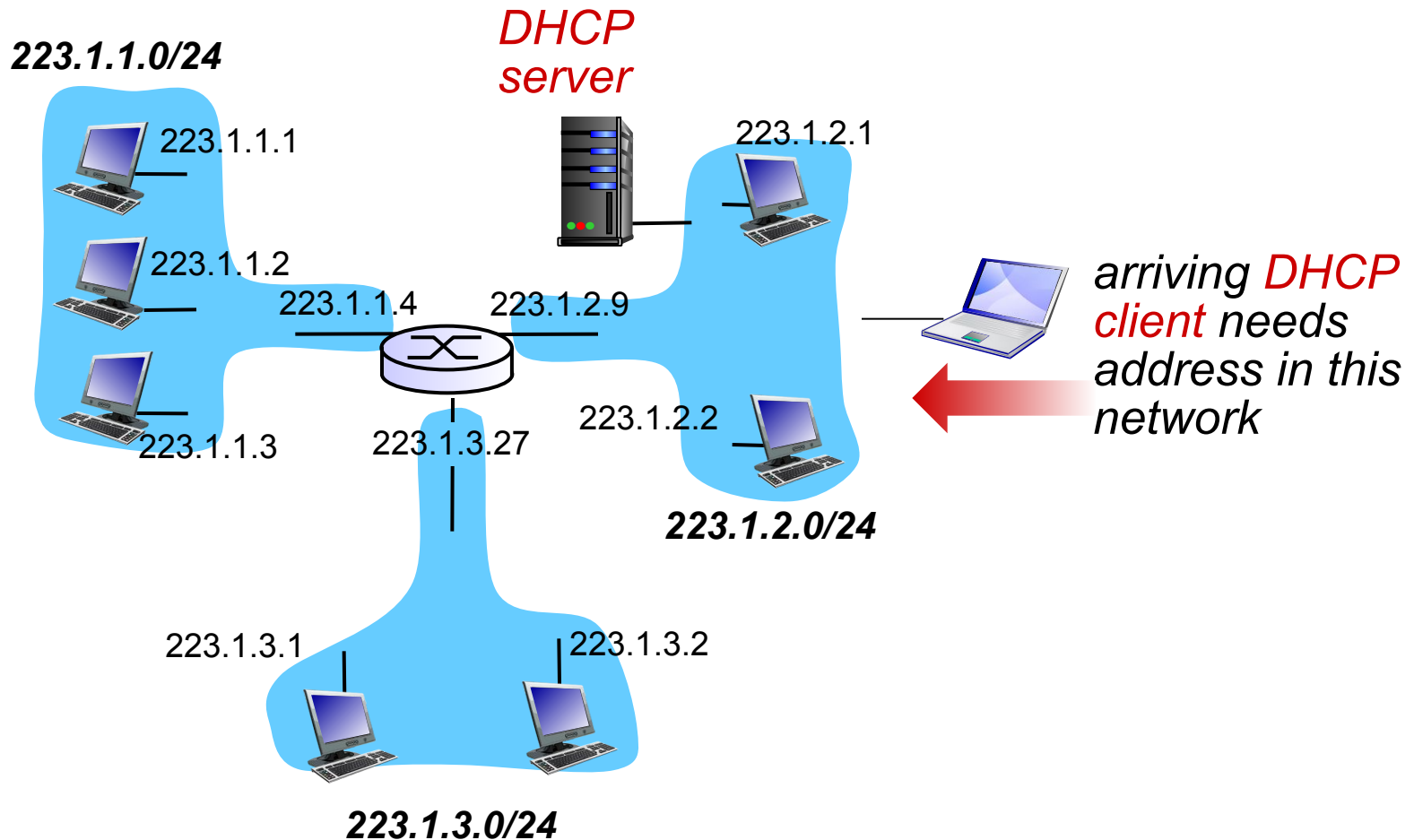
goal: allow host to *dynamically* obtain its IP address from network server when it joins network

- can renew its lease on address in use
- allows reuse of addresses (only hold address while connected/“on”)
- support for mobile users who want to join network (more shortly)
- Is also useful where hosts join and leave the network frequently.

DHCP overview:

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

DHCP client-server scenario



DHCP client-server scenario

DHCP server: 223.1.2.5



DHCP discover

Broadcast: is there a
DHCP server out there?

arriving
client



DHCP offer

Broadcast: I'm a DHCP
server! Here's an IP
address you can use

DHCP request

Broadcast: OK. I'll take
that IP address!

DHCP ACK

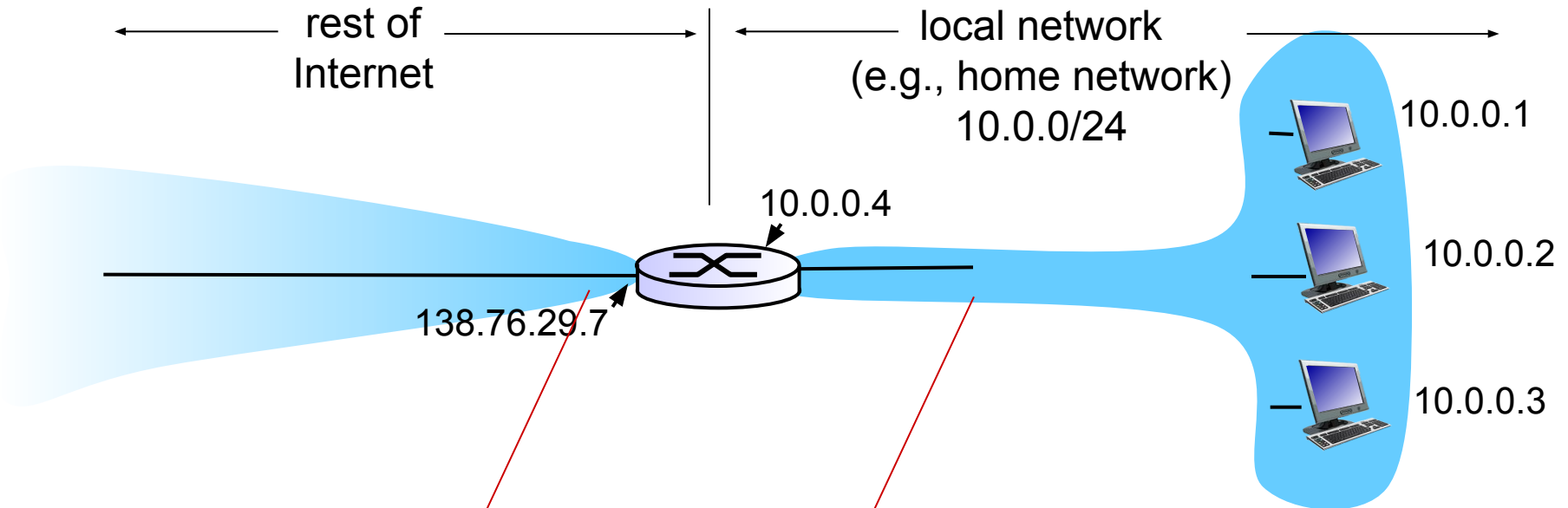
Broadcast: OK. You've
got that IP address!

DHCP: more than IP addresses

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

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

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:

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

NAT: network address translation

Private Addresses

Some IPv4 addresses are designated as private addresses, they are

Prefix	Range	#Addresses
10/8	10.0.0.0- 10.255.255.255	2^{24}
172.16/12	172.16.0.0- 172.31.255.255	2^{20}
192.168/16	192.168.0.0- 192.168.255.255	2^{16}

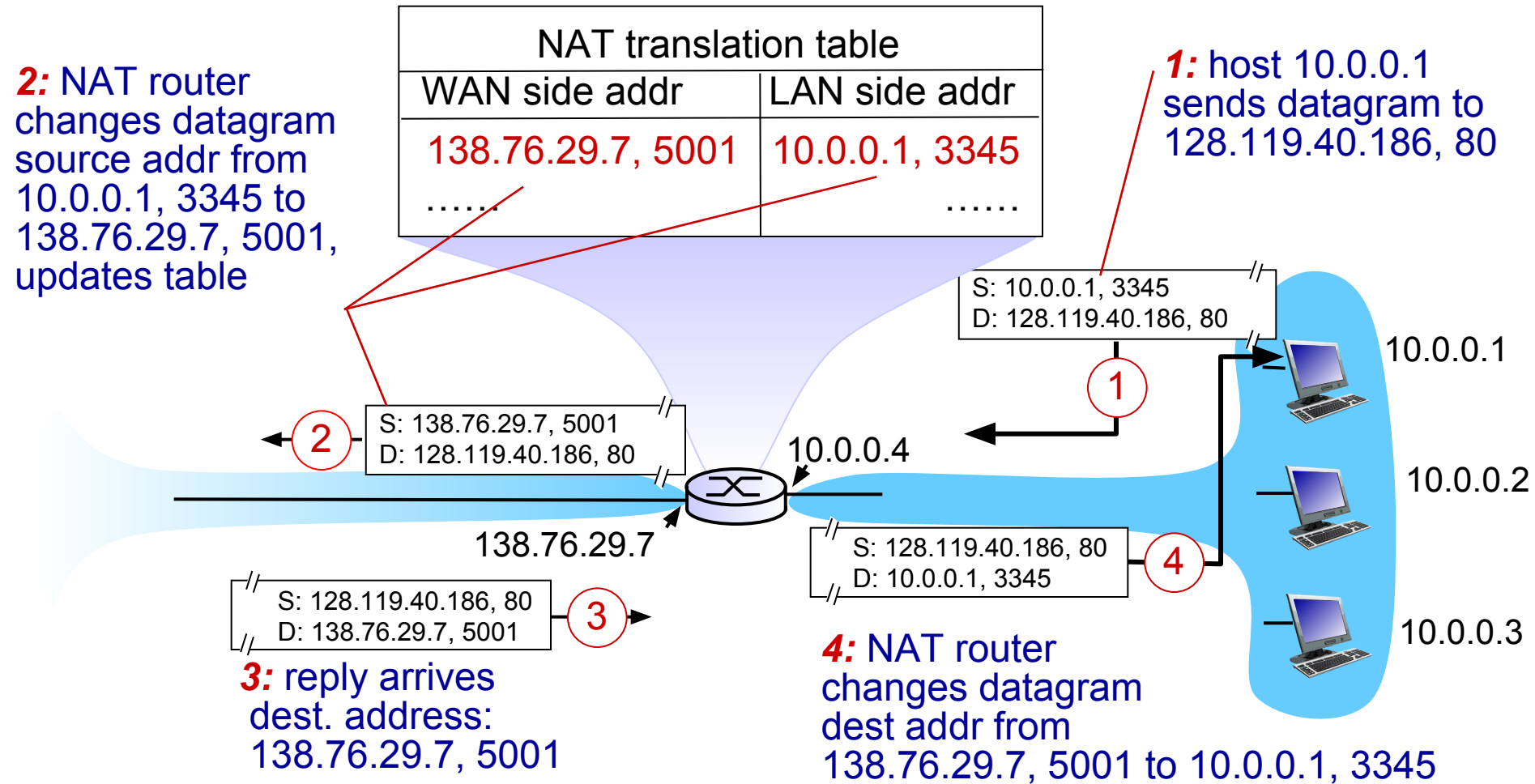
NSW

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

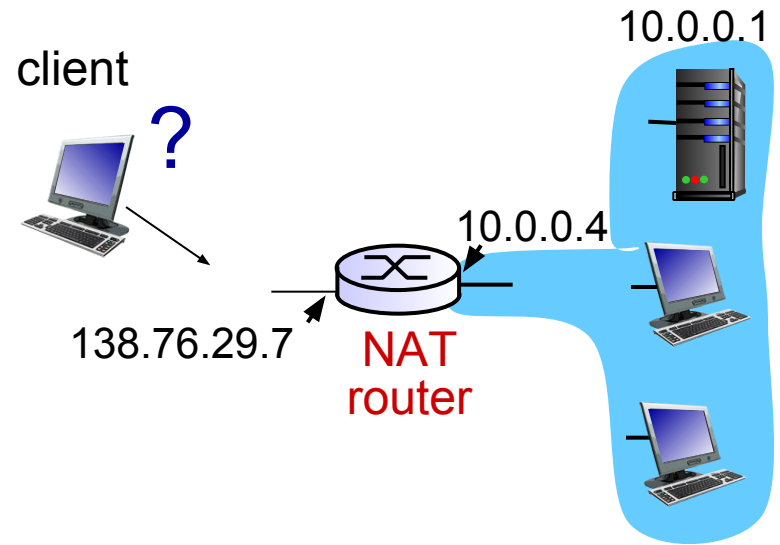


NAT: network address translation

- ❖ 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, e.g., P2P applications
 - address shortage should instead be solved by IPv6

NAT traversal problem

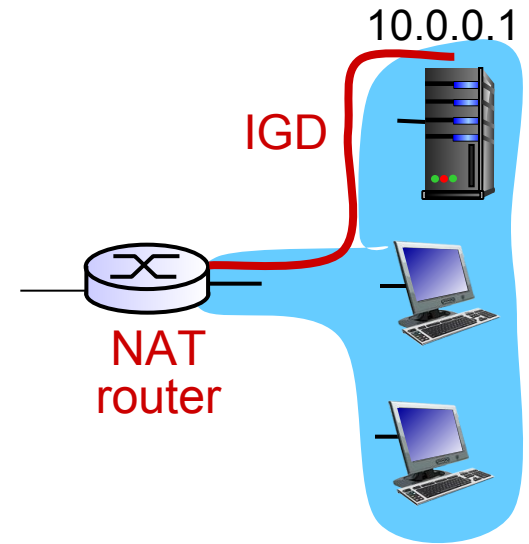
- ❖ client wants to connect to server with address 10.0.0.1
 - server address 10.0.0.1 local to LAN (client can't use it as destination addr)
 - only one externally visible NATed address: 138.76.29.7
- ❖ **solution 1:** statically configure NAT to forward incoming connection requests at given port to 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 NATed host to:
 - ❖ learn public IP address (138.76.29.7)
 - ❖ add/remove port mappings (with lease times)

i.e., automate static NAT port map configuration



Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

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:

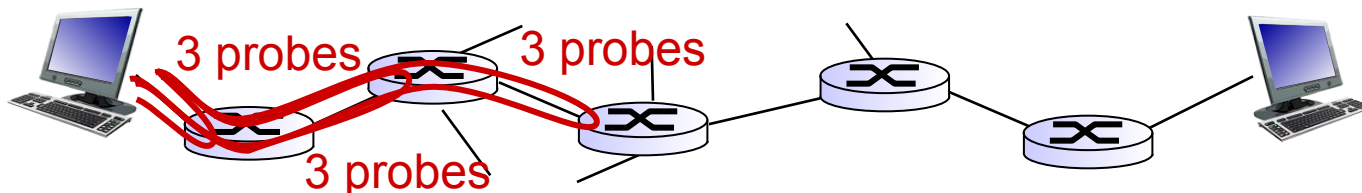
- ICMP msgs carried in IP datagrams

- ❖ **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 set has TTL = 1
 - second set has TTL=2, etc.
 - unlikely port number
 - ❖ when n th set of datagrams arrives to n th router:
 - router discards datagrams
 - and sends source ICMP messages (type 11, code 0)
 - ICMP messages includes name of router & IP address
 - ❖ when ICMP messages arrives, source records RTTs
- stopping criteria:*
- ❖ UDP segment eventually arrives at destination host
 - ❖ destination returns ICMP “port unreachable” message (type 3, code 3)
 - ❖ source stops



IPv6: motivation

- ❖ *initial motivation*: 32-bit address space soon to be completely allocated.
- ❖ additional motivation:
 - header format helps speed processing/forwarding
 - header changes to facilitate QoS
 - Introduces anycast address.

IPv6 datagram format:

- fixed-length 40 byte header
- no fragmentation allowed

IPv6: motivation

- 128 bits means you can have 2^{128} addresses, which is 340,282,366,920,938,463,463,374,607,431,768,211,456
 - This is approximately 3.4×10^{38}
 - Compare with 4×10^9 IPv4 addresses, IPv6 has 10^{29} times more addresses
- Earth's surface area (land + water) is 500×10^{14} sq. metres $\Rightarrow 7 \times 10^{21}$ addresses per sq. metre
- One reason why the address space is so large
 - Address assignment can never achieve 100% efficiency
 - but there are other reasons (to be discussed later)

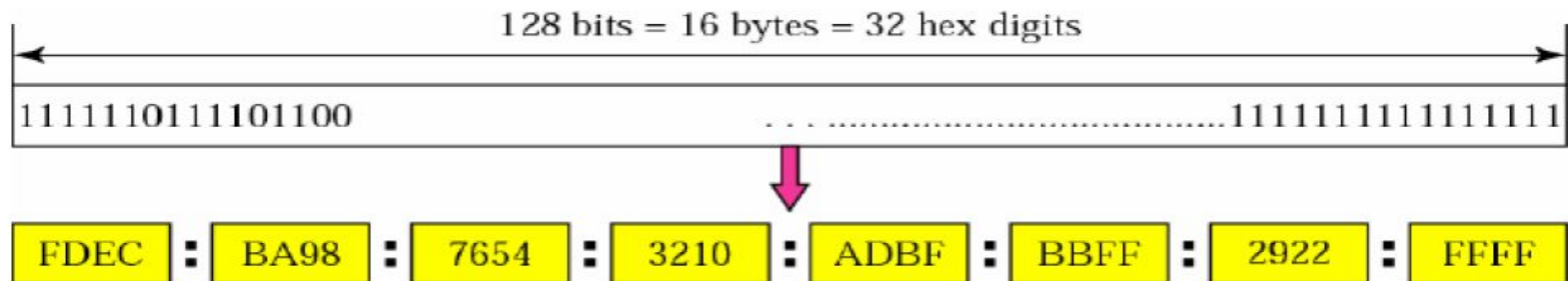


Fig: IPv6 Address Format

IPv6: Abbreviated Address

Unabbreviated

FDEC : BA98 : 0074 : 3210 : 000F : BBFF : 0000 : FFFF



FDEC : BA98 : 74 : 3210 : F : BBFF : 0 : FFFF

Abbreviated

Abbreviated address with consecutive zeros

Abbreviated

FDEC : 0 : 0 : 0 : 0 : BBFF : 0 : FFFF



FDEC :: BBFF : 0 : FFFF

More Abbreviated

The address can also be abbreviated as
FDEC:0:0:0:0:BBFF::FFFF

IPv6: Abbreviated Address

Exercises

- What is the complete IPv6 address for
 - ABBA:CAB:1234::FEED:3:BEEF

Solution

The complete address for

ABBA:CAB:1234::FEED:3:BEEF is

ABBA:0CAB:1234:0000:0000:FEED:0003:BEEF

IPv6: Abbreviated Address

more exercise

Exercise: What is the complete address for

- ABBA::7::FEED?

Solution

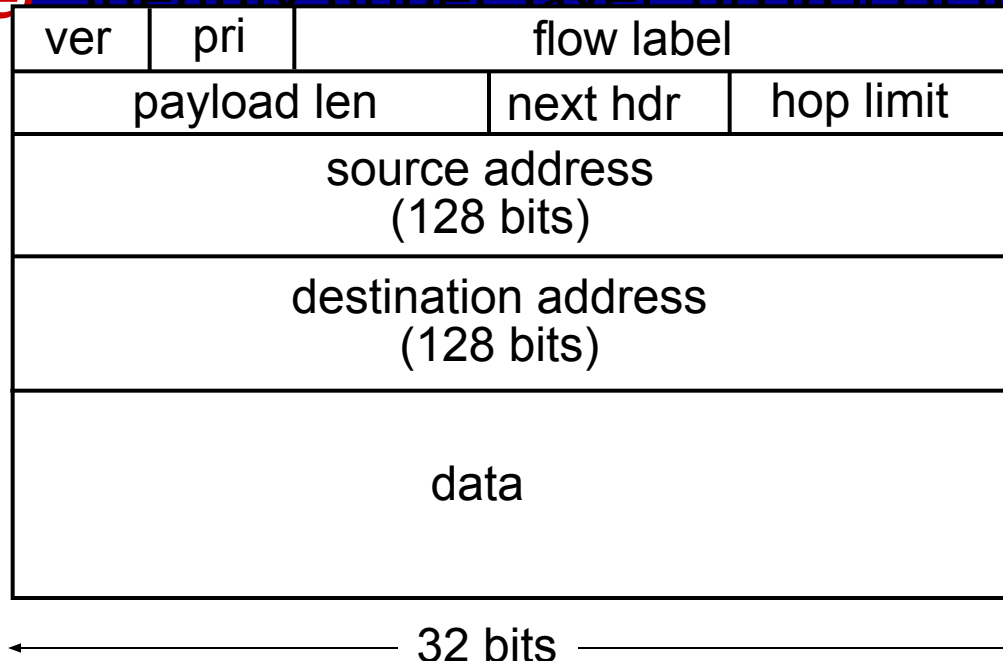
Solution: The address is invalid. Only one group of zeros can be suppressed.

IPv6 datagram format

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

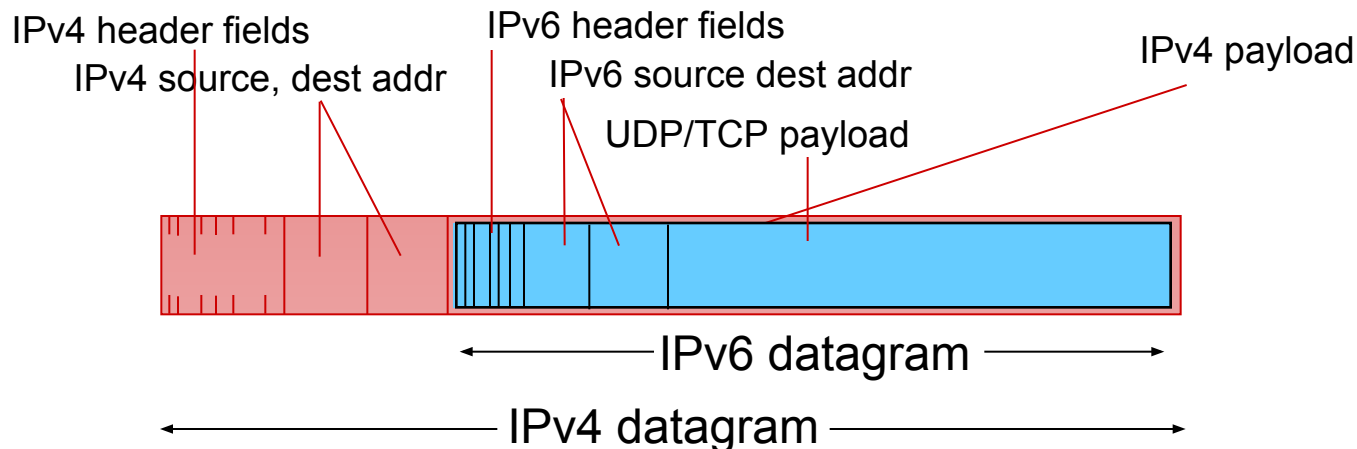


Other changes from IPv4

- ❖ *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

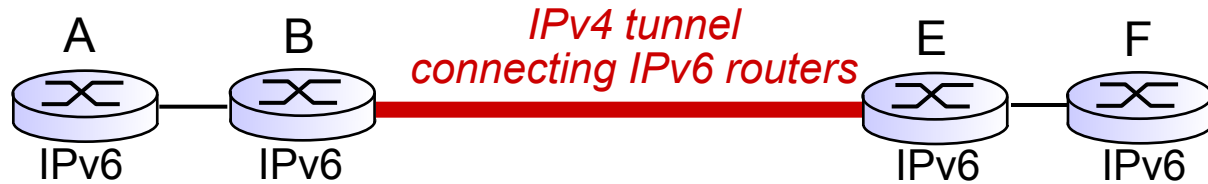
Transition from IPv4 to IPv6

- ❖ not all routers can be upgraded simultaneously
 - no “flag days”
 - how will network operate with mixed IPv4 and IPv6 routers?
- ❖ **tunneling**: IPv6 datagram carried as *payload* in IPv4 datagram among IPv4 routers

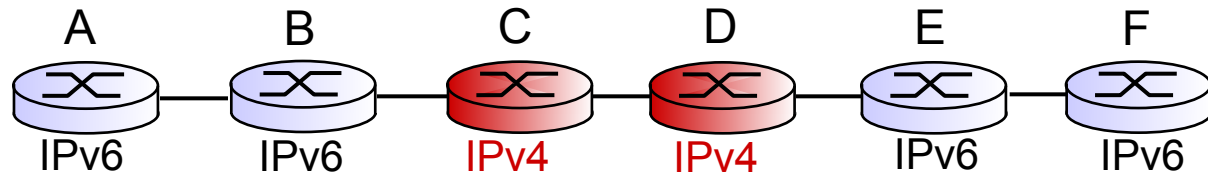


Tunneling

logical view:

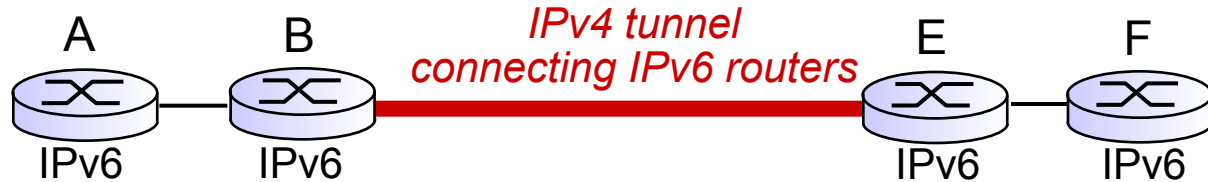


physical view:

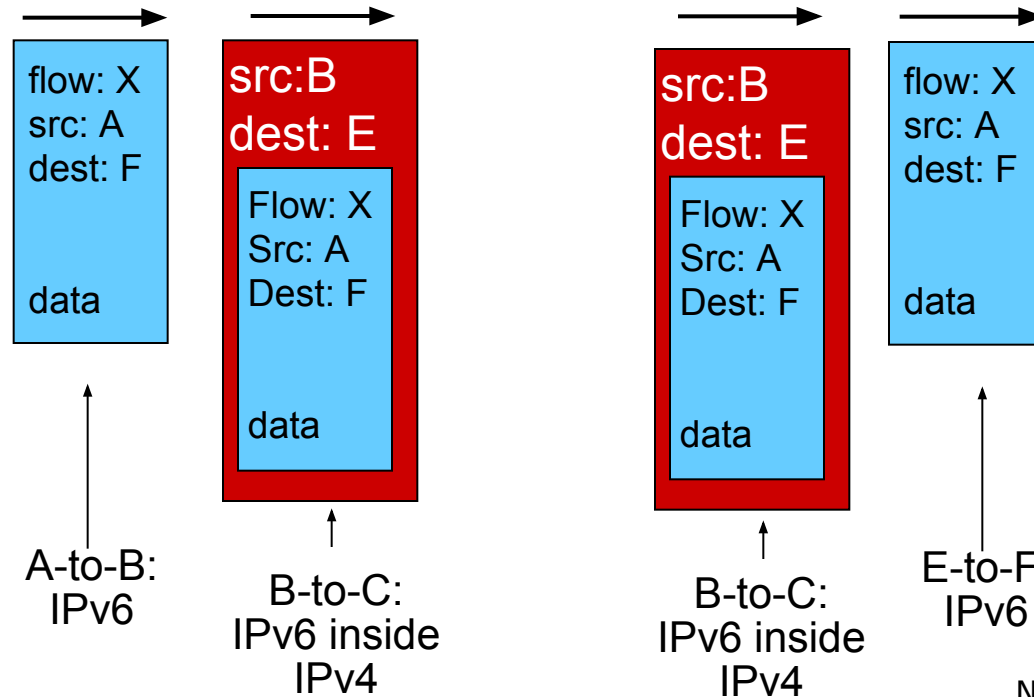
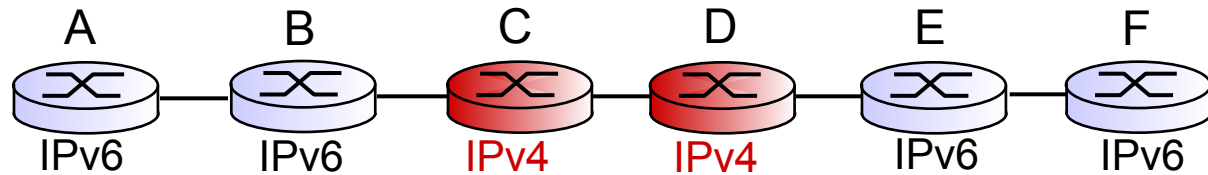


Tunneling

logical view:



physical view:



IPv6: adoption

- ❖ US National Institutes of Standards estimate [2013]:
 - ~3% of industry IP routers
 - ~11% of US gov't routers
- ❖ *Long (long!) time for deployment, use*
 - 20 years and counting!
 - think of application-level changes in last 20 years: WWW, Facebook, ...
 - *Why?*

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

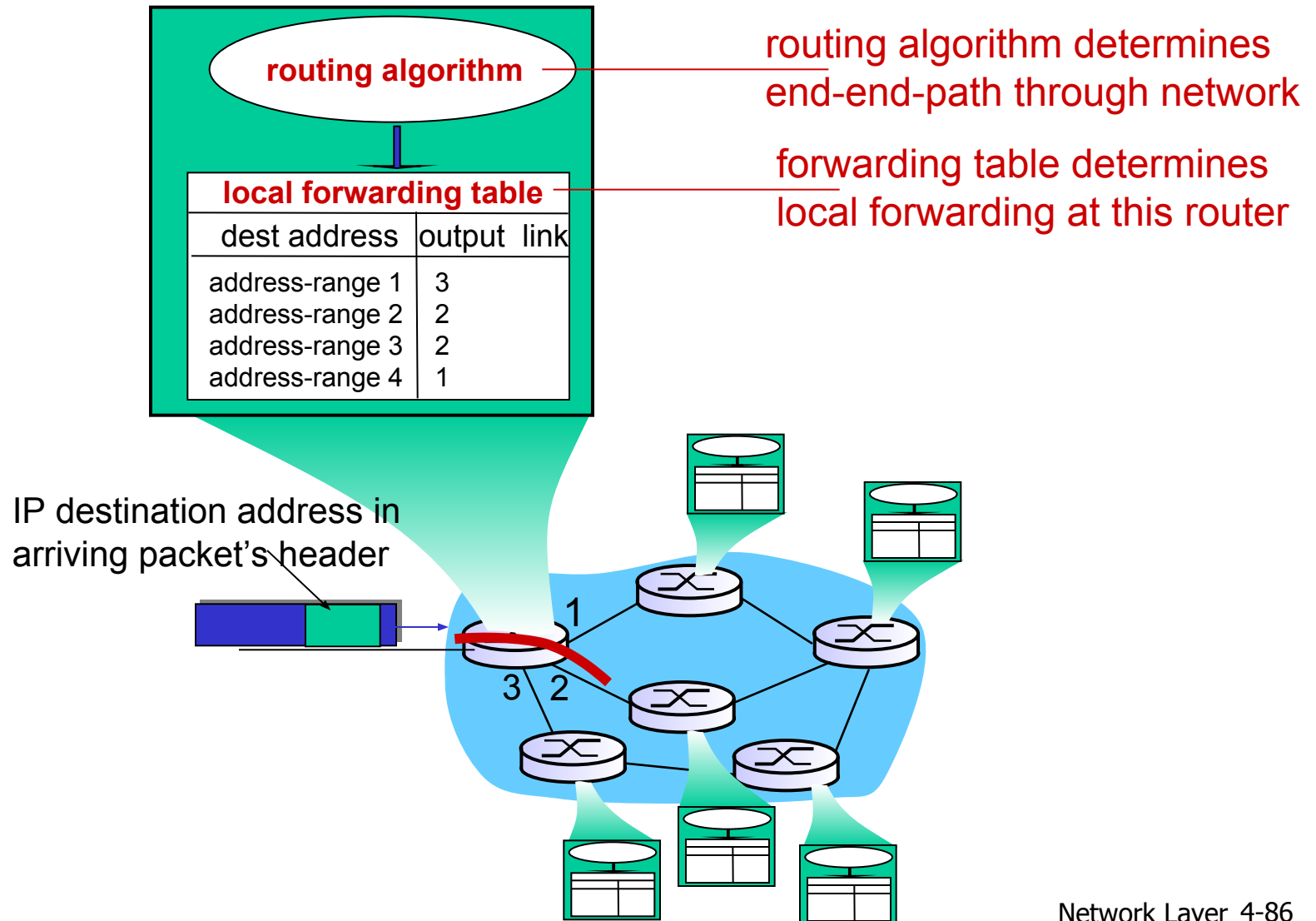
- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet

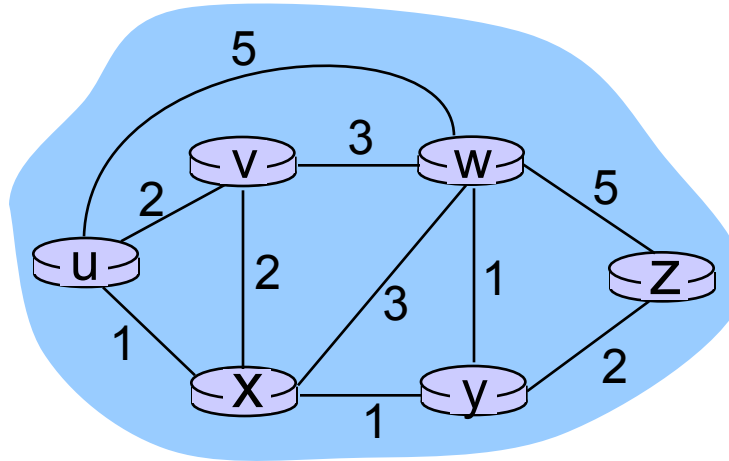
- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Interplay between routing, forwarding



Graph abstraction



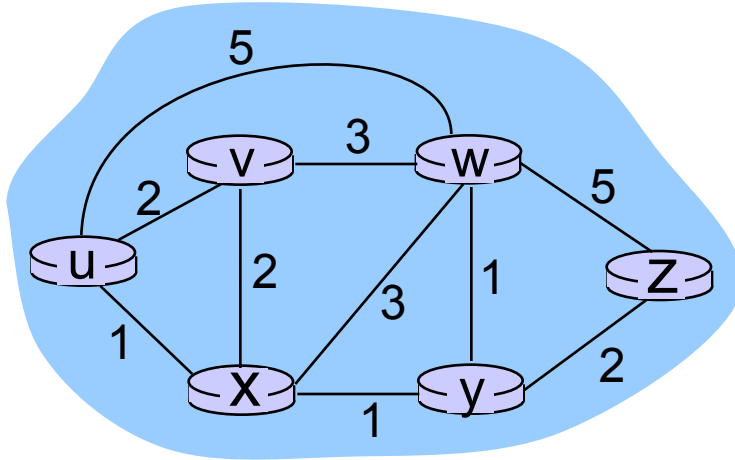
graph: $G = (N, E)$

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

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

aside: graph abstraction is useful in other network contexts, e.g., 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 $(x_1, x_2, x_3, \dots, x_p) = c(x_1, x_2) + c(x_2, x_3) + \dots + c(x_{p-1}, x_p)$

key question: what is the least-cost path between u and z?

routing algorithm: algorithm that finds that least cost path

Routing algorithm classification

Q: 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

Q: static or dynamic?

static:

- ❖ routes change slowly over time

dynamic:

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

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector

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 */

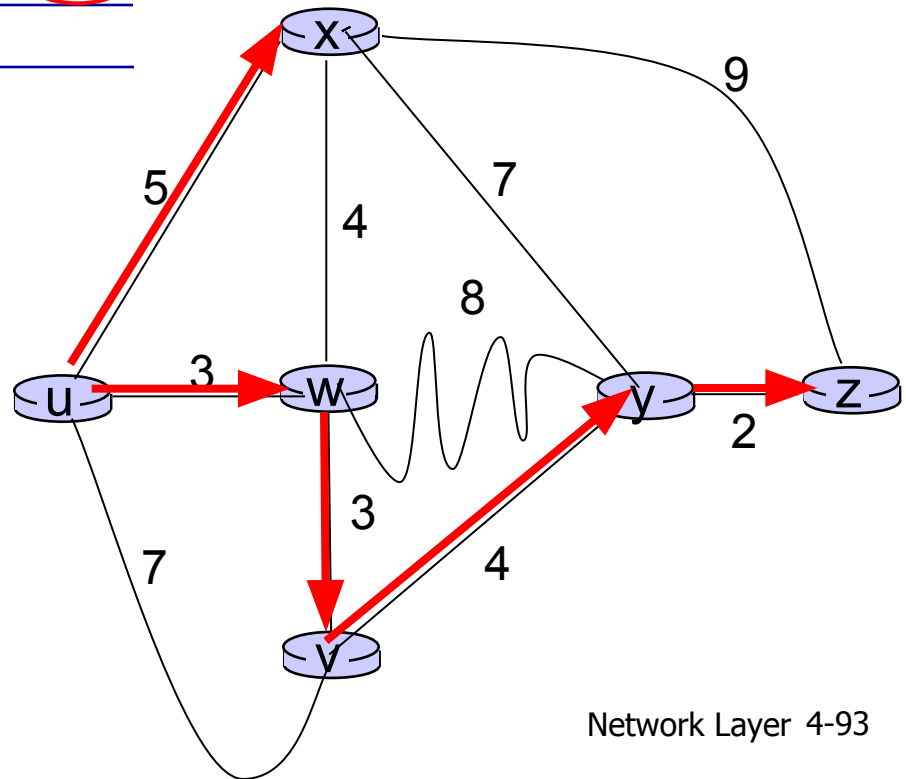
15 **until all nodes in N'**

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	7,u	3,u	5,u	∞	∞
1	uw	6,w		5,u	11,w	∞
2	uwx	6,w			11,w	14,x
3	uwxv				10,v	14,x
4	uwxvy					12,y
5	uwxvyz					

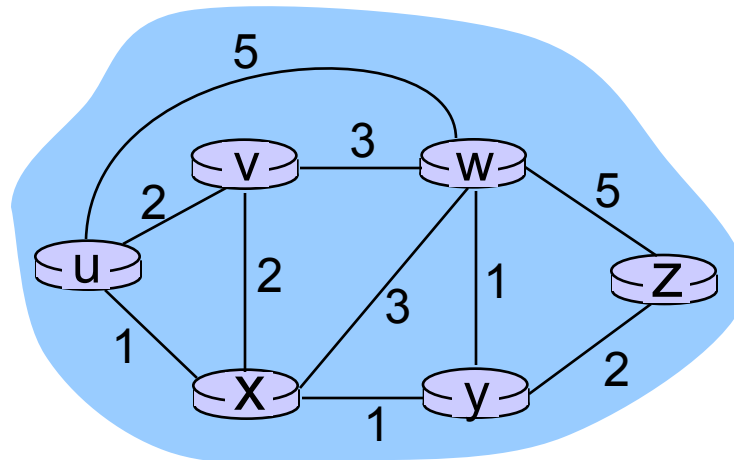
notes:

- ❖ construct shortest path tree by tracing predecessor nodes
- ❖ ties can exist (can be broken arbitrarily)



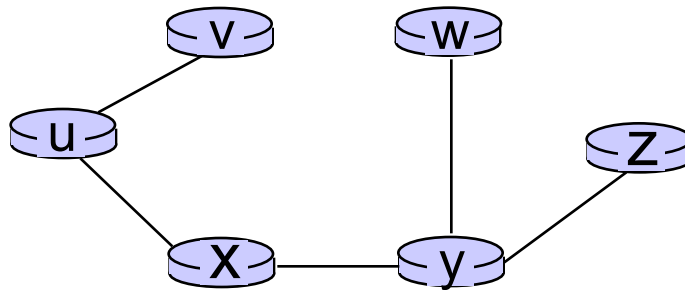
Dijkstra's algorithm: another example

Step	N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0	u	2,u	5,u	1,u	∞	∞
1	ux	2,u	4,x		2,x	∞
2	uxy	2,u	3,y			4,y
3	uxyv		3,y			4,y
4	uxyvw					4,y
5	uxyvwz					



Dijkstra's algorithm: example (2)

resulting shortest-path tree from u:



resulting forwarding table in u:

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

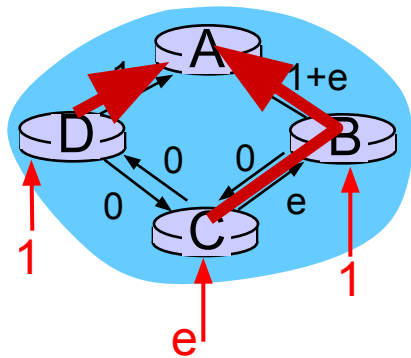
Dijkstra's algorithm, discussion

algorithm complexity: n nodes

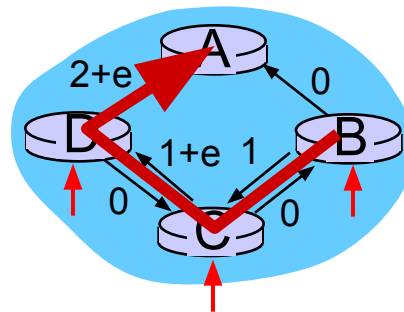
- ❖ each iteration: need to check all nodes, w, not in N
- ❖ $n(n+1)/2$ comparisons: $O(n^2)$
- ❖ more efficient implementations possible: $O(n \log n)$

oscillations possible:

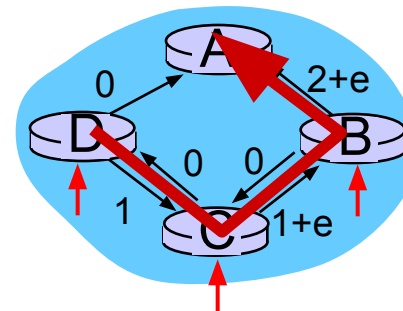
- ❖ e.g., support link cost equals amount of carried traffic:



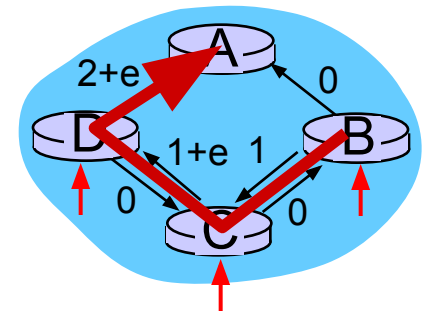
initially



given these costs,
find new routing....
resulting in new costs



given these costs,
find new routing....
resulting in new costs



given these costs,
find new routing....
resulting in new costs

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector

4.6 routing in the Internet

- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Distance vector algorithm

Bellman-Ford equation (dynamic programming)

let

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

then

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

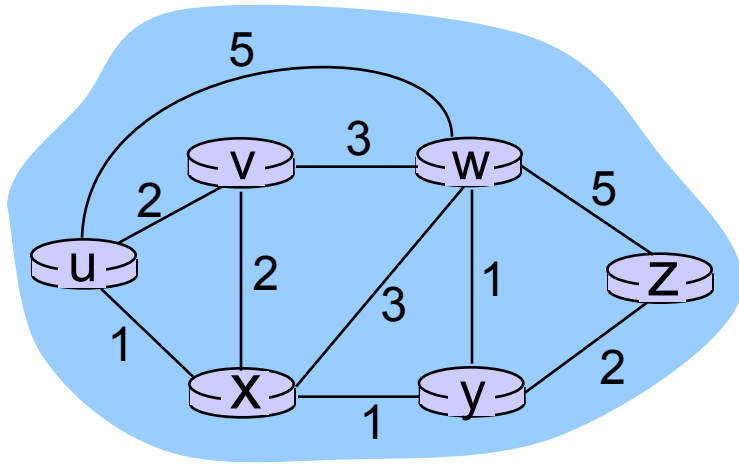
cost from neighbor v to destination y

cost to neighbor v

\min 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 achieving minimum is next
hop in shortest path, used in forwarding table

Distance vector algorithm

- ❖ $D_x(y)$ = estimate of least cost from x to y
 - x maintains distance vector $\mathbf{D}_x = [D_x(y): y \in N]$
- ❖ node x :
 - knows cost to each neighbor v : $c(x,v)$
 - maintains its neighbors' distance vectors. For each neighbor v , x maintains $\mathbf{D}_v = [D_v(y): y \in N]$

Distance vector algorithm

key idea:

- ❖ from time-to-time, each node sends its own distance vector estimate to neighbors
- ❖ when x receives new DV estimate from neighbor, it updates its own DV using B-F equation:

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

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

Distance vector algorithm

iterative, asynchronous:

each local iteration
caused by:

- ❖ local link cost change
- ❖ DV update message from neighbor

distributed:

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

each

node:

wait for (change in local link cost or msg from neighbor)

recompute estimates

if DV to any dest has changed, *notify* neighbors

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

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

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

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

**node x
table**

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

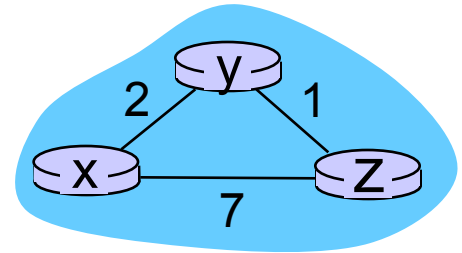
		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	7	1	0

**node y
table**

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

**node z
table**

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



time

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

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

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

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

**node x
table**

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

**node y
table**

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

**node z
table**

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

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

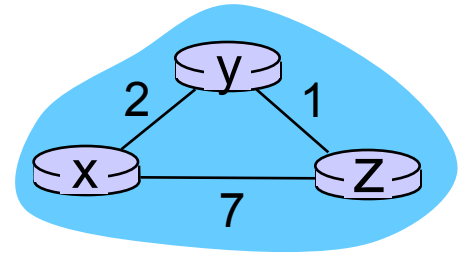
	cost to		
	x	y	z
from x	0	2	7
from y	2	0	1
from z	7	1	0

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

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

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

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

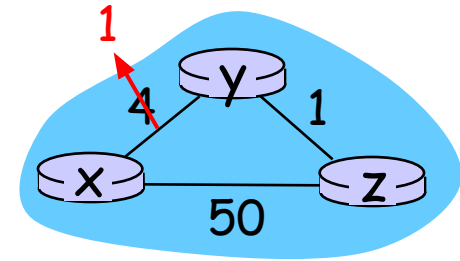


time →

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”

t_0 : y detects link-cost change, updates its DV, informs its neighbors.

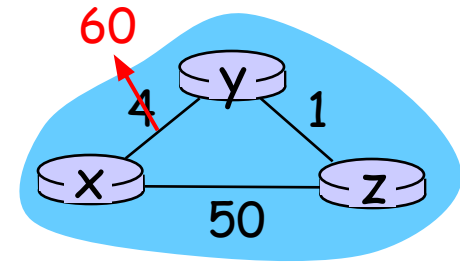
t_1 : z receives update from y, updates its table, computes new least cost to x, sends its neighbors its DV.

t_2 : y receives z's update, updates its distance table. y's least costs do *not* change, so y does *not* send a message to z.

Distance vector: link cost changes

link cost changes:

- ❖ node detects local link cost change
- ❖ *bad news travels slow* - “count to infinity” problem!
- ❖ 44 iterations before algorithm stabilizes: see text



poisoned reverse:

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

Comparison of LS and DV algorithms

message complexity

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

speed of convergence

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

robustness: what happens if router malfunctions?

LS:

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

DV:

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

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector

4.6 routing in the Internet

- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Intra-AS Routing

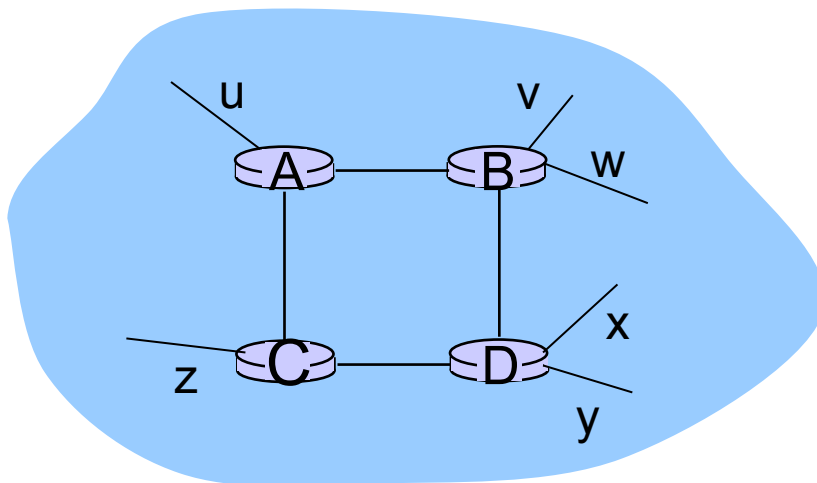
- ❖ also known as *interior gateway protocols (IGP)*
- ❖ most common intra-AS routing protocols:
 - RIP: Routing Information Protocol
 - OSPF: Open Shortest Path First
 - IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

Three standard Internet routing protocols:

Routing protocol	Based on routing algorithm	IGP/EGP
Routing Information Protocol (RIP)	Distance vector	IGP
Open shortest path first (OSPF)	Link state	IGP
Border Gateway Protocol (BGP)	Path vector	EGP

RIP (Routing Information Protocol)

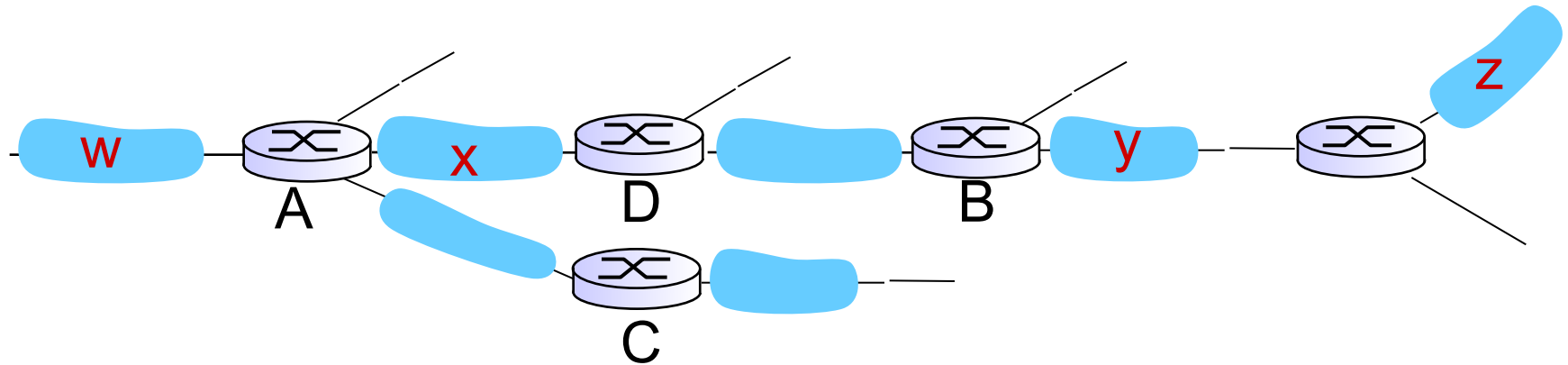
- ❖ included in BSD-UNIX distribution in 1982
- ❖ distance vector algorithm
 - distance metric: # hops (max = 15 hops), each link has cost 1
 - DVs exchanged with neighbors every 30 sec in response message (aka **advertisement**)
 - each advertisement: list of up to 25 destination **subnets** (in IP addressing sense)



from router A to destination **subnets**:

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

RIP: example



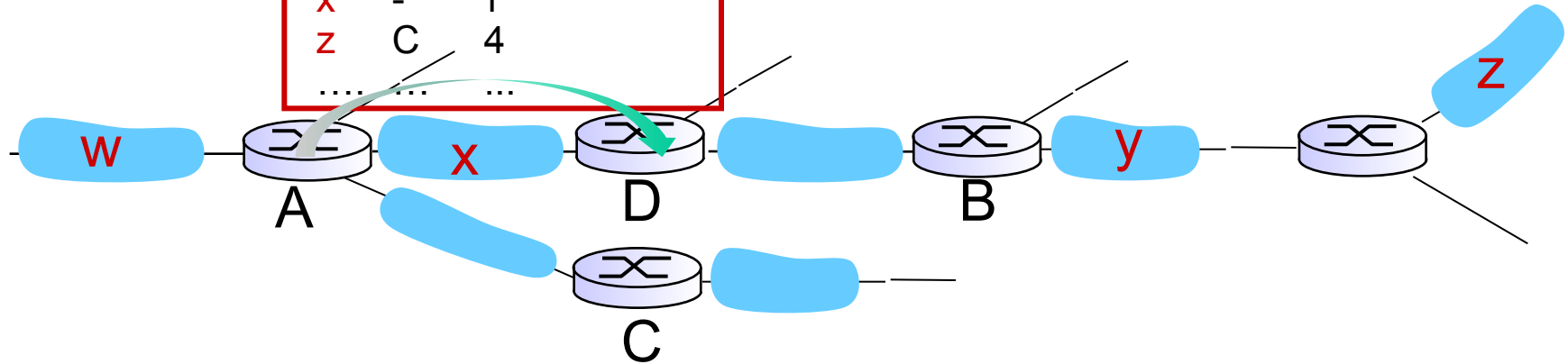
routing table in router D

destination	subnet	next router	# hops to dest
w	A	2	
y	B	2	
z	B	7	
x	--	1	
....	

RIP: example

A-to-D advertisement

dest	next	hops
W	-	1
X	-	1
Z	C	4
...



routing table in router D

destination	subnet	next router	# hops to dest
w	A	2	
y	B	2	
z	B	7	
x	--	1	
....	

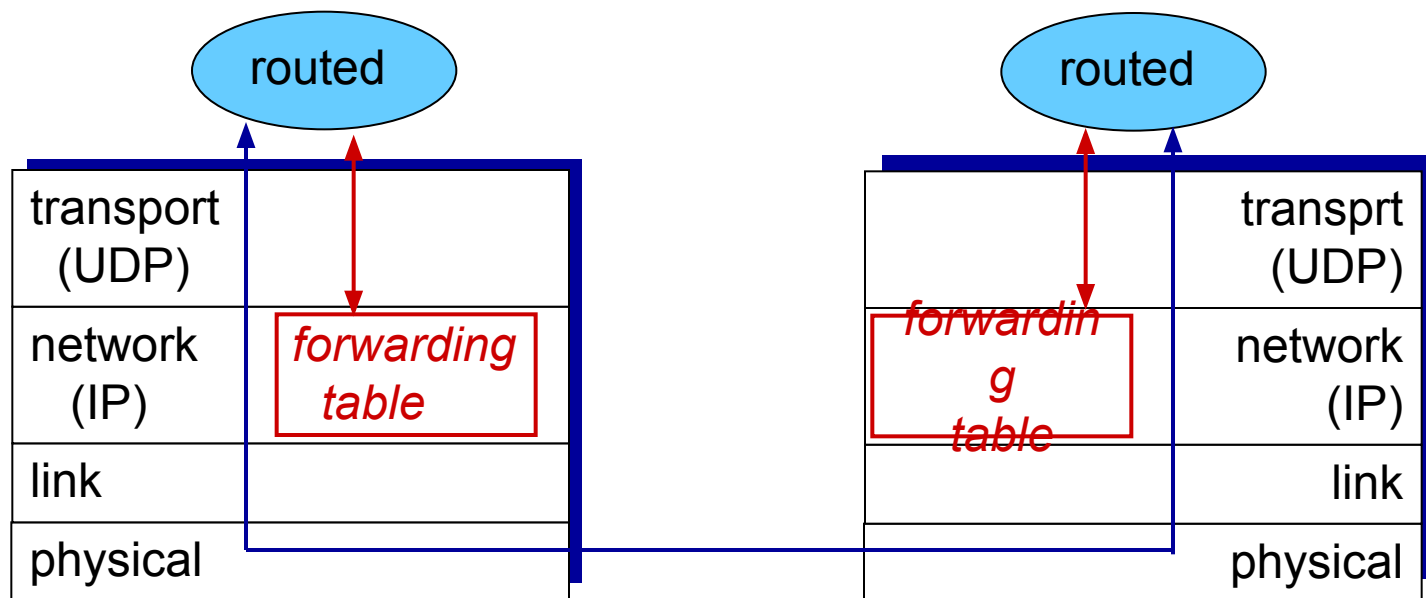
RIP: link failure, recovery

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



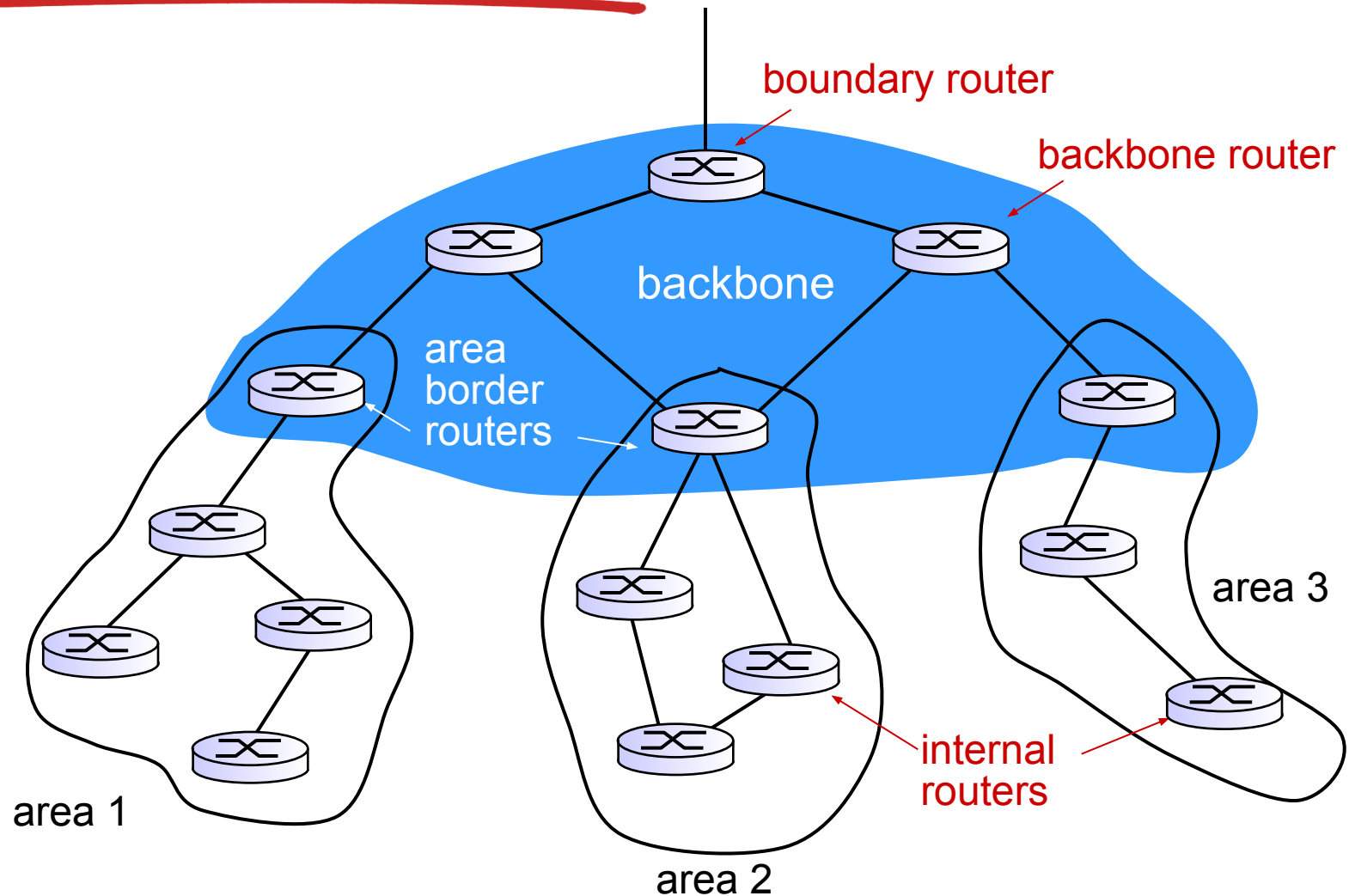
OSPF (Open Shortest Path First)

- ❖ “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
- ❖ advertisements flooded to *entire* AS
 - carried in OSPF messages directly over IP (rather than TCP or UDP)
- ❖ *IS-IS routing* protocol: nearly identical to OSPF

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 ToS; high for real time ToS)
- ❖ integrated uni- and **multicast** support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- ❖ **hierarchical** OSPF in large domains.

Hierarchical OSPF



Hierarchical OSPF

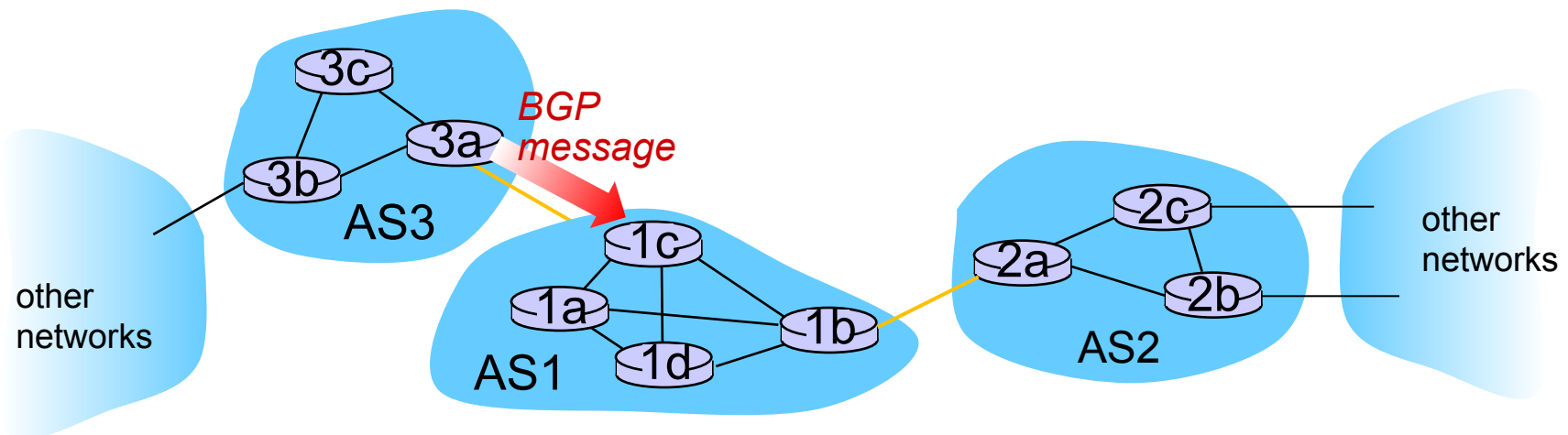
- ❖ *two-level hierarchy*: local area, backbone.
 - link-state advertisements only in area
 - each nodes has detailed area topology; only know direction (shortest path) to nets in other areas.
- ❖ *area border routers*: “summarize” distances to nets in own area, advertise to other Area Border routers.
- ❖ *backbone routers*: run OSPF routing limited to backbone.
- ❖ *boundary routers*: connect to other AS's.

Internet inter-AS routing: BGP

- ❖ **BGP (Border Gateway Protocol):** *the de facto* inter-domain routing protocol
 - “glue that holds the Internet together”
- ❖ BGP provides each AS a means to:
 - **eBGP:** obtain subnet reachability information from neighboring ASs.
 - **iBGP:** propagate reachability information to all AS-internal routers.
 - determine “good” routes to other networks based on reachability information and policy.
- ❖ allows subnet to advertise its existence to rest of Internet: *“I am here”*

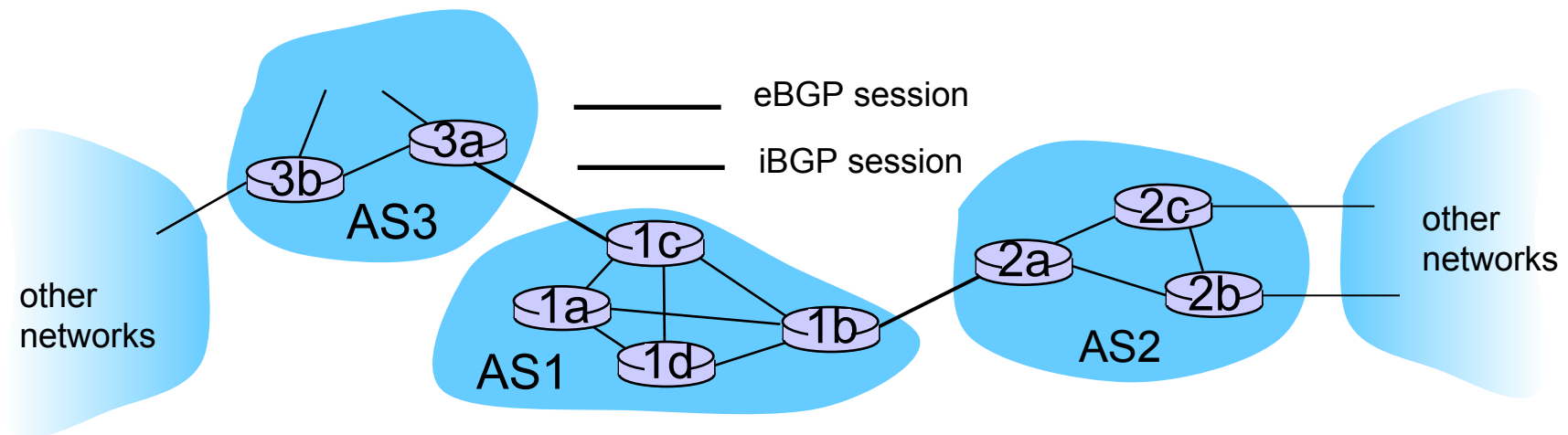
BGP basics

- ❖ **BGP session:** two BGP routers (“peers”) exchange BGP messages:
 - advertising *paths* to different destination network prefixes (“path vector” protocol)
 - exchanged over semi-permanent TCP connections
- ❖ **when AS3 advertises a prefix to AS1:**
 - *AS3 promises it will forward datagrams towards that prefix*
 - *AS3 can aggregate prefixes in its advertisement*



BGP basics: distributing path information

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



Path attributes and 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)
- ❖ gateway router receiving route advertisement uses **import policy** to accept/decline
 - e.g., never route through AS x
 - *policy-based* routing

BGP route selection

- ❖ router may learn about more than 1 route to destination AS, selects route based on:
 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 between peers over TCP connection
- ❖ 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

Chapter 4: *done!*

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

- ❖ understand principles behind network layer services:
 - network layer service models, forwarding versus routing how a router works, routing (path selection), broadcast, multicast
- ❖ instantiation, implementation in the Internet