

# Network Layer

Chapter 4

# 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

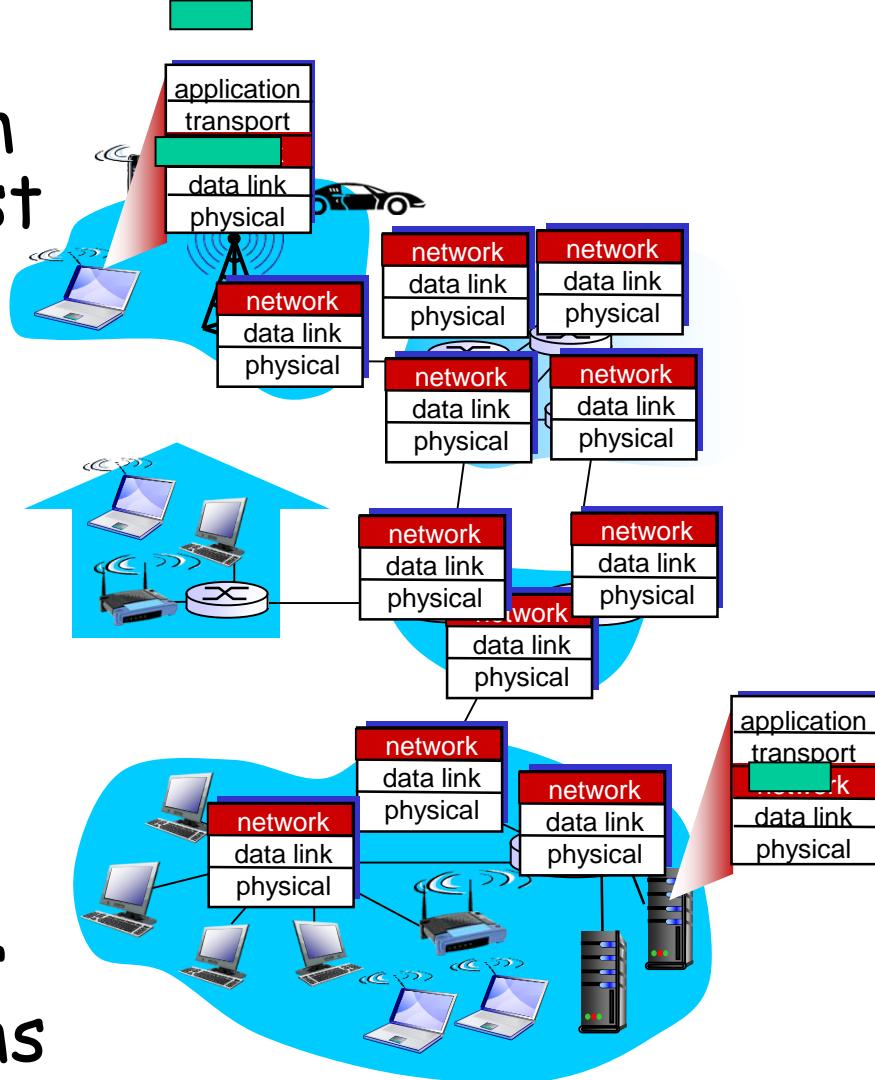
- link state
- distance vector

## 4.6 routing in the Internet

- RIP
- OSPF
- BGP

# Network layer

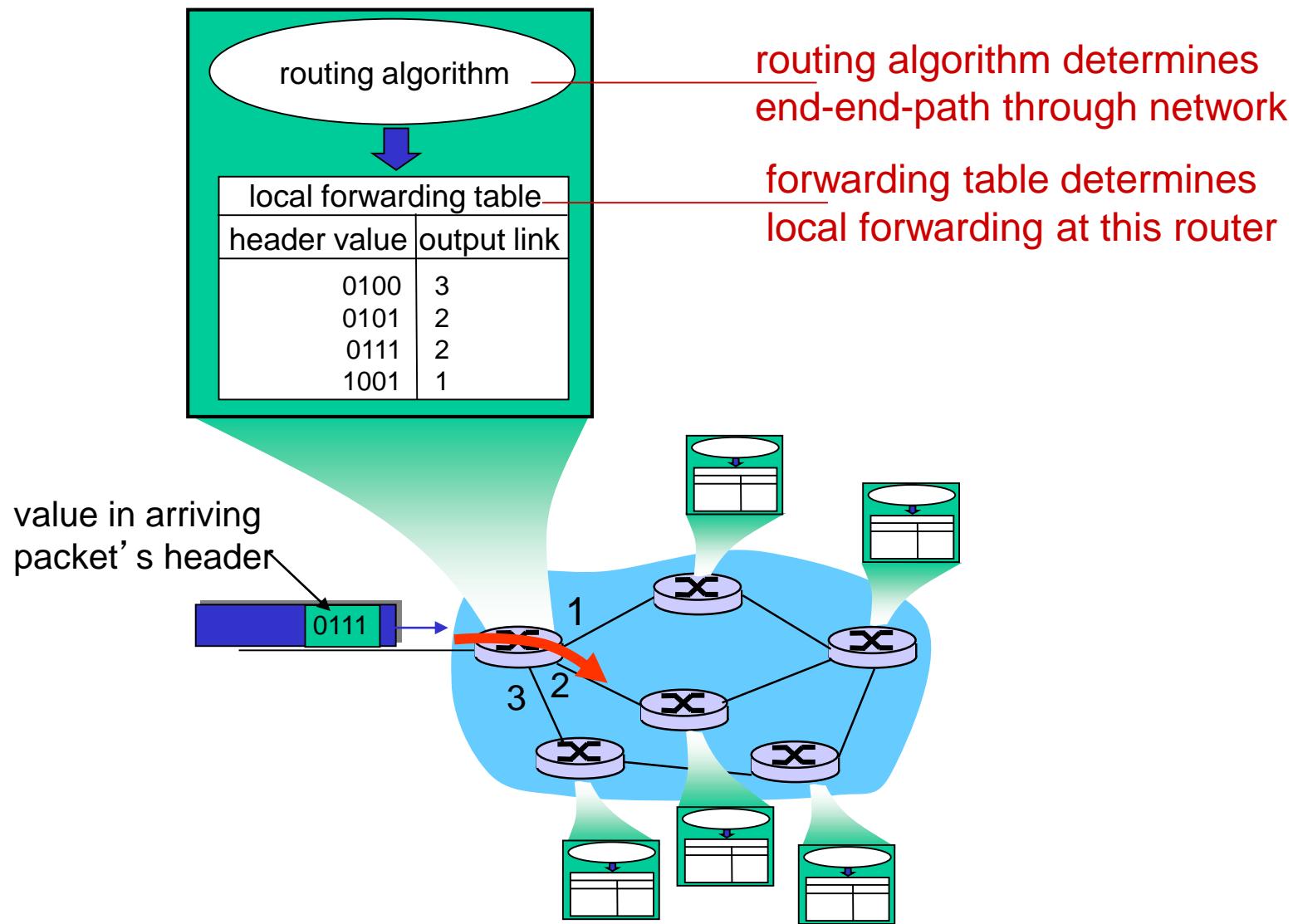
- ❖ transport segment from sending to receiving host
- ❖ on sending side encapsulates segments into datagrams
- ❖ on receiving side, delivers segments to transport layer
- ❖ network layer protocols in **every** host, router
- ❖ router examines header fields in all IP datagrams passing through it



# 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*

# Interplay between routing and forwarding



# Connection setup

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

# Network service model

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

*example services for individual datagrams:*

- ❖ guaranteed delivery
- ❖ guaranteed delivery with less than 40 msec (bounded)delay

*example services for a flow of datagrams:*

- ❖ in-order datagram delivery
- ❖ guaranteed minimum bandwidth to flow.
- ❖ Security services

# Network layer service models:

Network Architecture	Service Model	Bandwidth Guarantee	No-Loss Guarantee	Ordering	Timing	Congestion Indication
Internet	Best Effort	None	None	Any order possible	Not maintained	None
ATM	CBR	Guaranteed constant rate	Yes	In order	Maintained	Congestion will not occur
ATM	ABR	Guaranteed minimum	None	In order	Not maintained	Congestion indication provided

# Chapter 4: outline

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

# Chapter 4: outline

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

# 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)

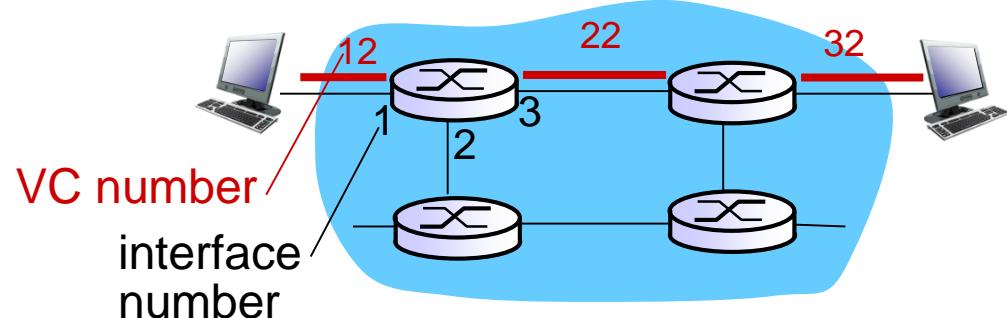
# VC implementation

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:*

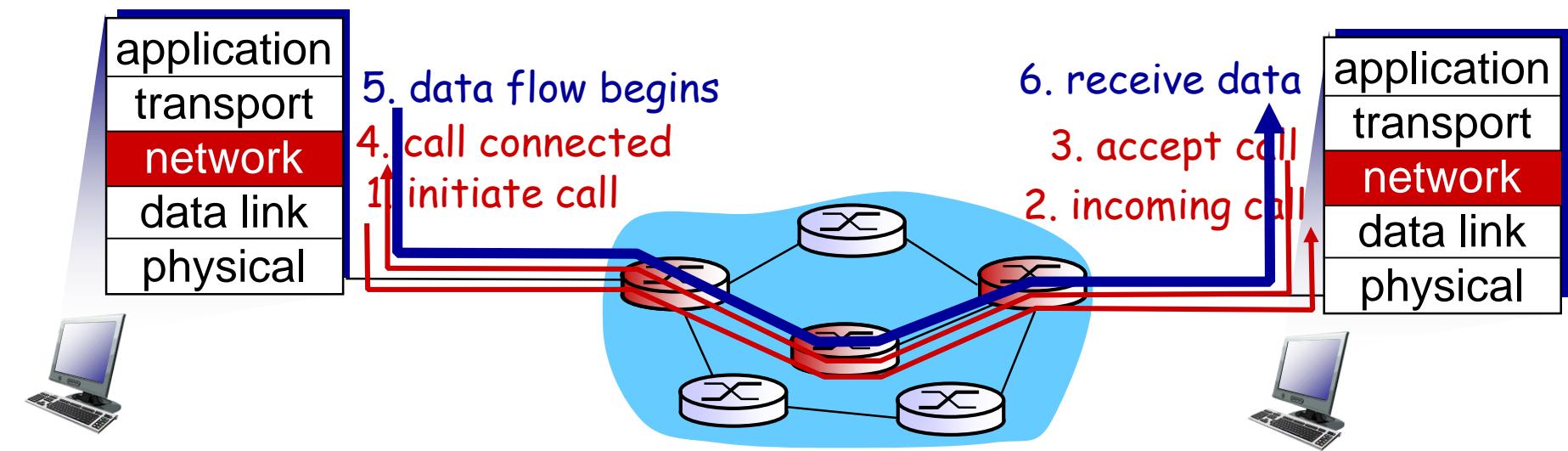


Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	3	22
2	63	1	18
3	7	2	17
1	97	3	87
...	...	...	...

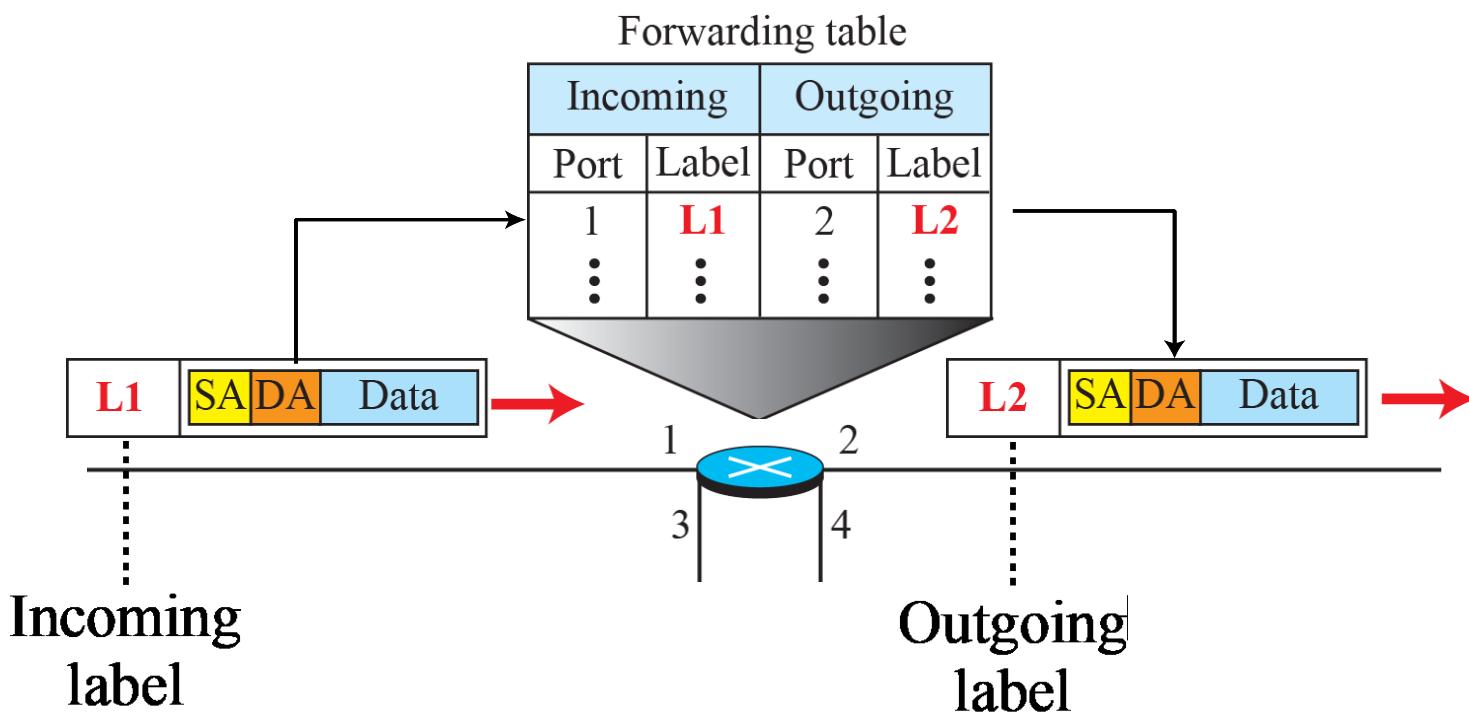
*VC routers maintain connection state information!*

# Virtual circuits: signaling protocols

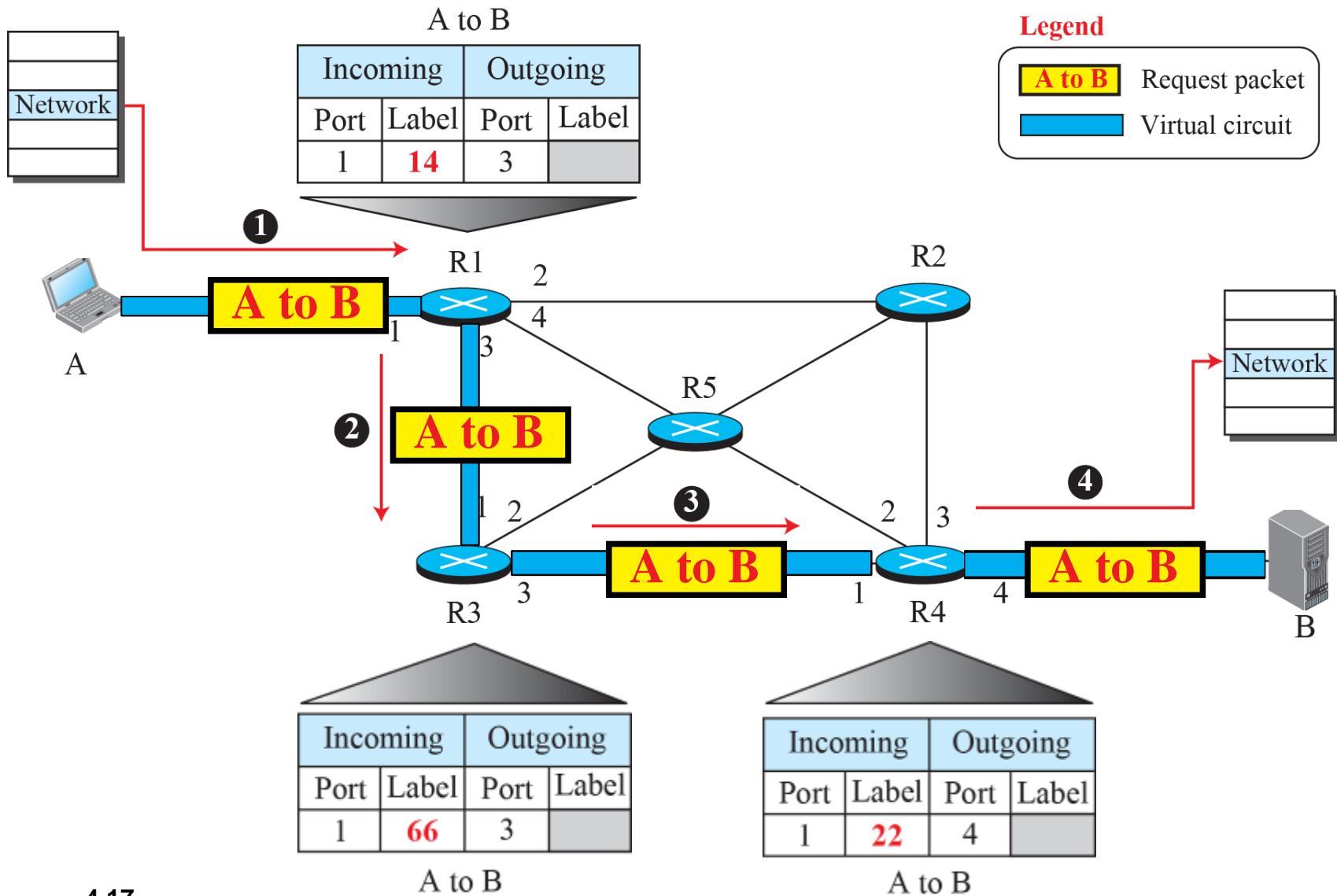
- ❖ used to setup, maintain teardown VC
- ❖ used in ATM, frame-relay, X.25
- ❖ not used in today's Internet



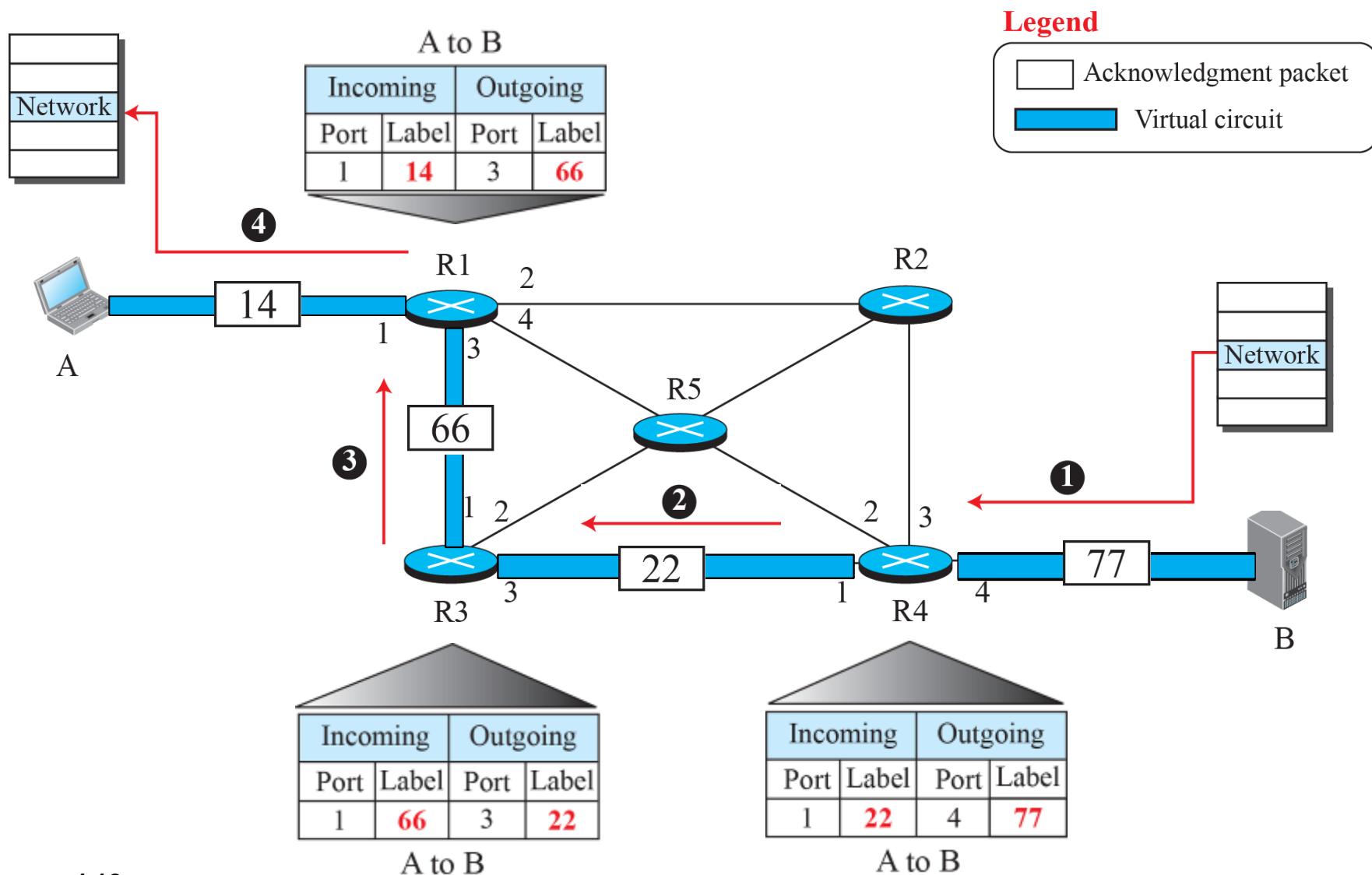
**Figure 4.6: Forwarding process in a router when used in a virtual circuit network**



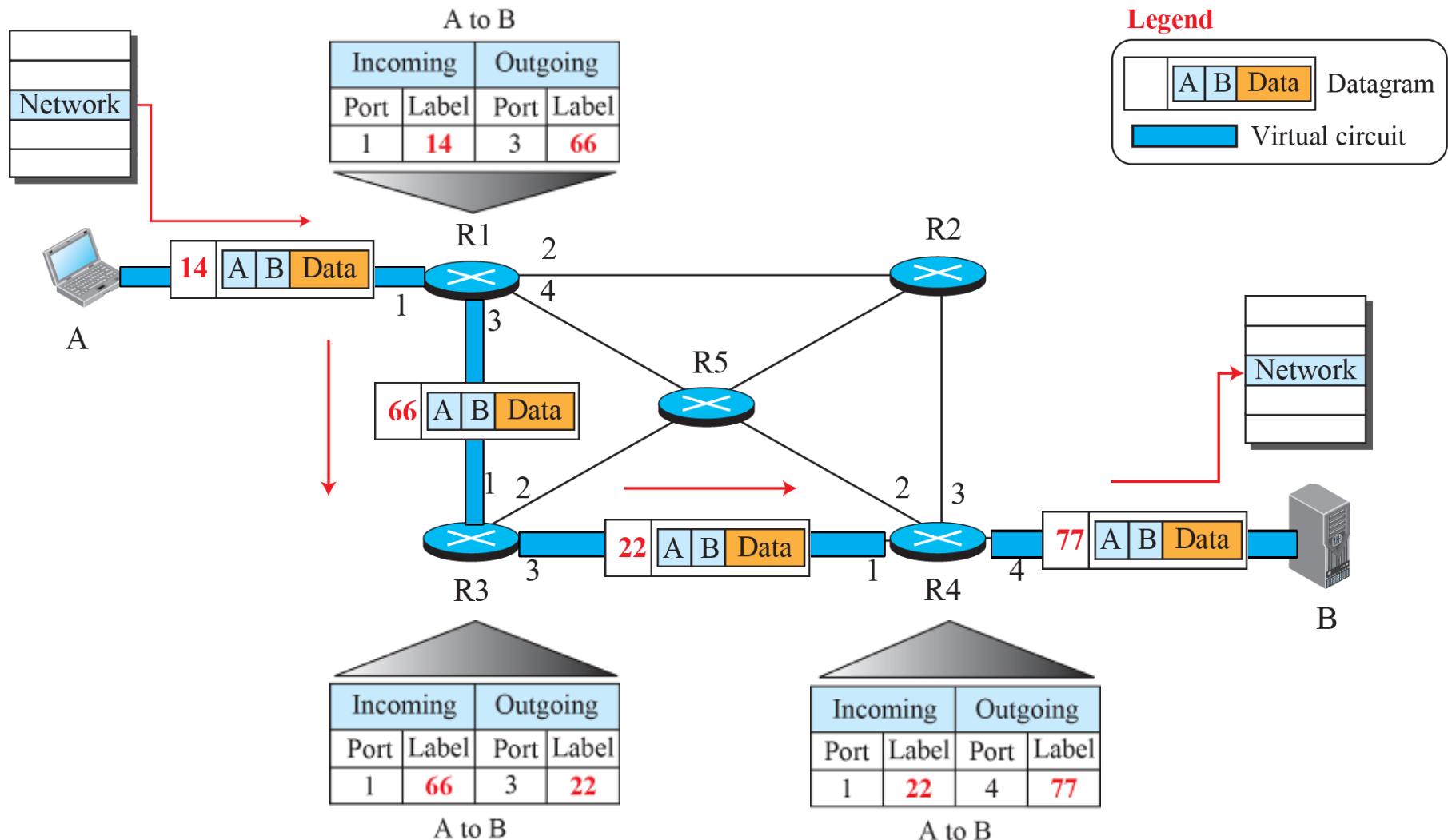
**Figure 4.7: Sending request packet in a virtual-circuit network**



**Figure 4.8: Sending acknowledgments in a virtual-circuit network**

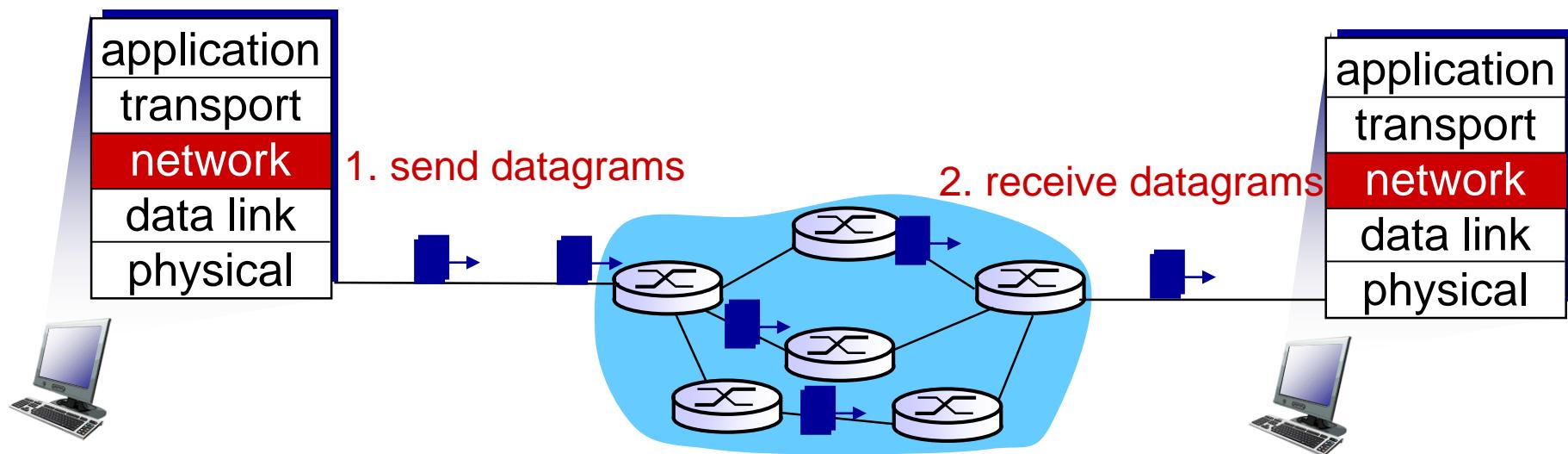


**Figure 4.8: Sending acknowledgments in a virtual-circuit network**

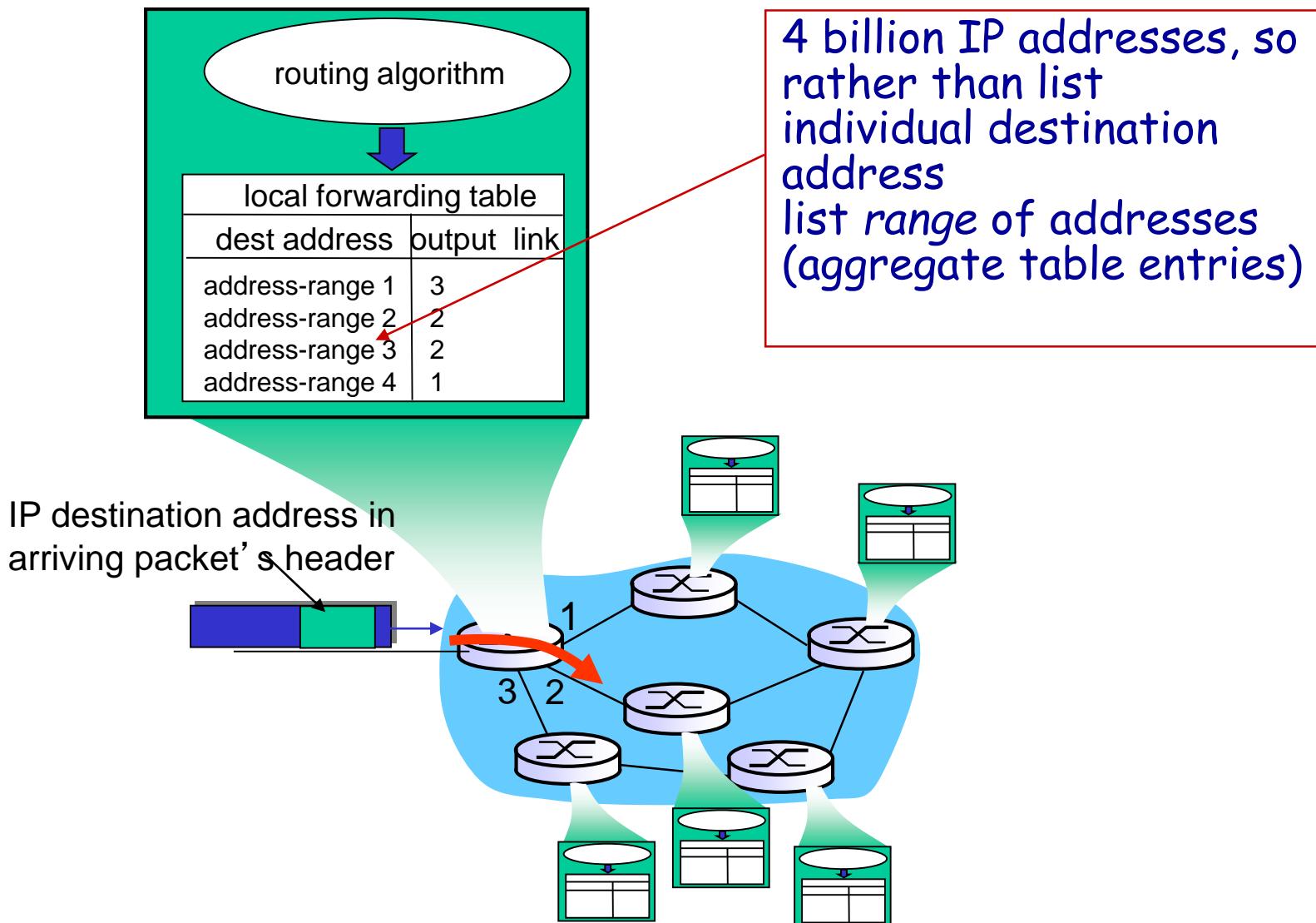


# Datagram networks

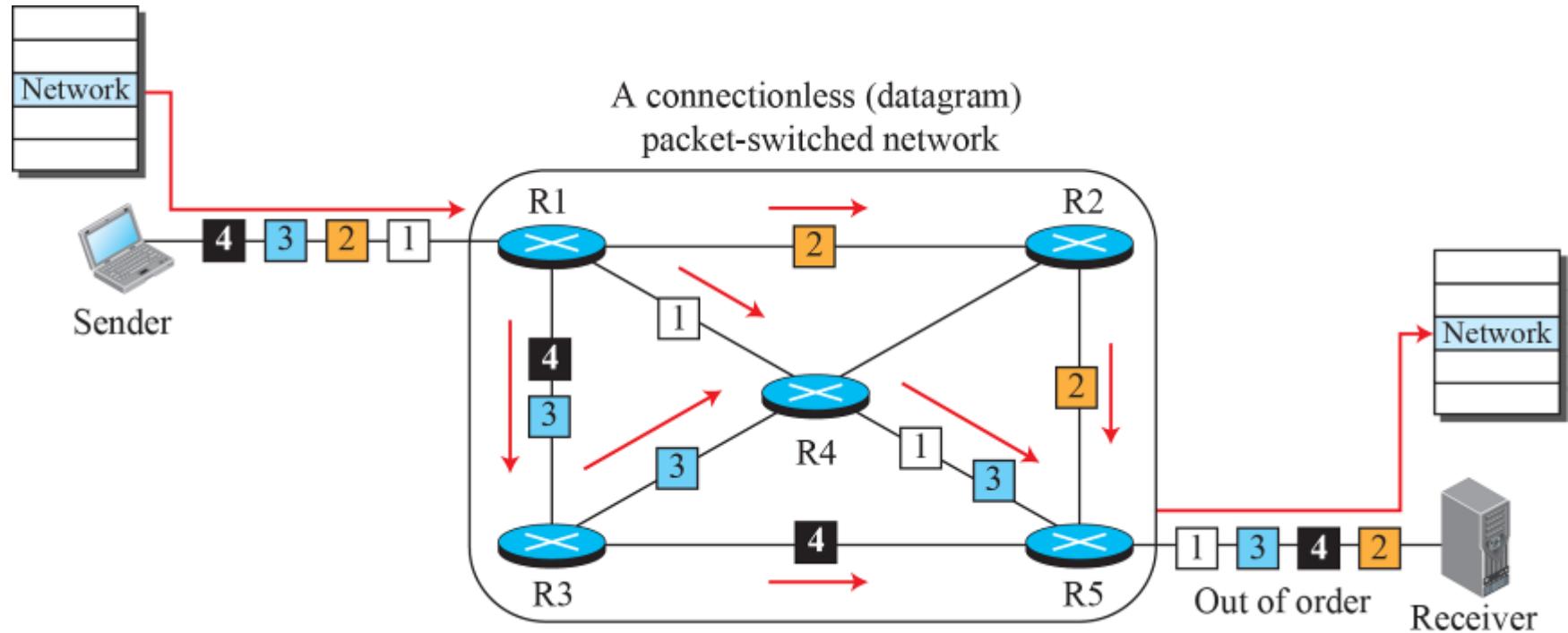
- ❖ no call setup at network layer
- ❖ routers: no state about end-to-end connections
  - no network-level concept of “connection”
- ❖ packets forwarded using destination host address



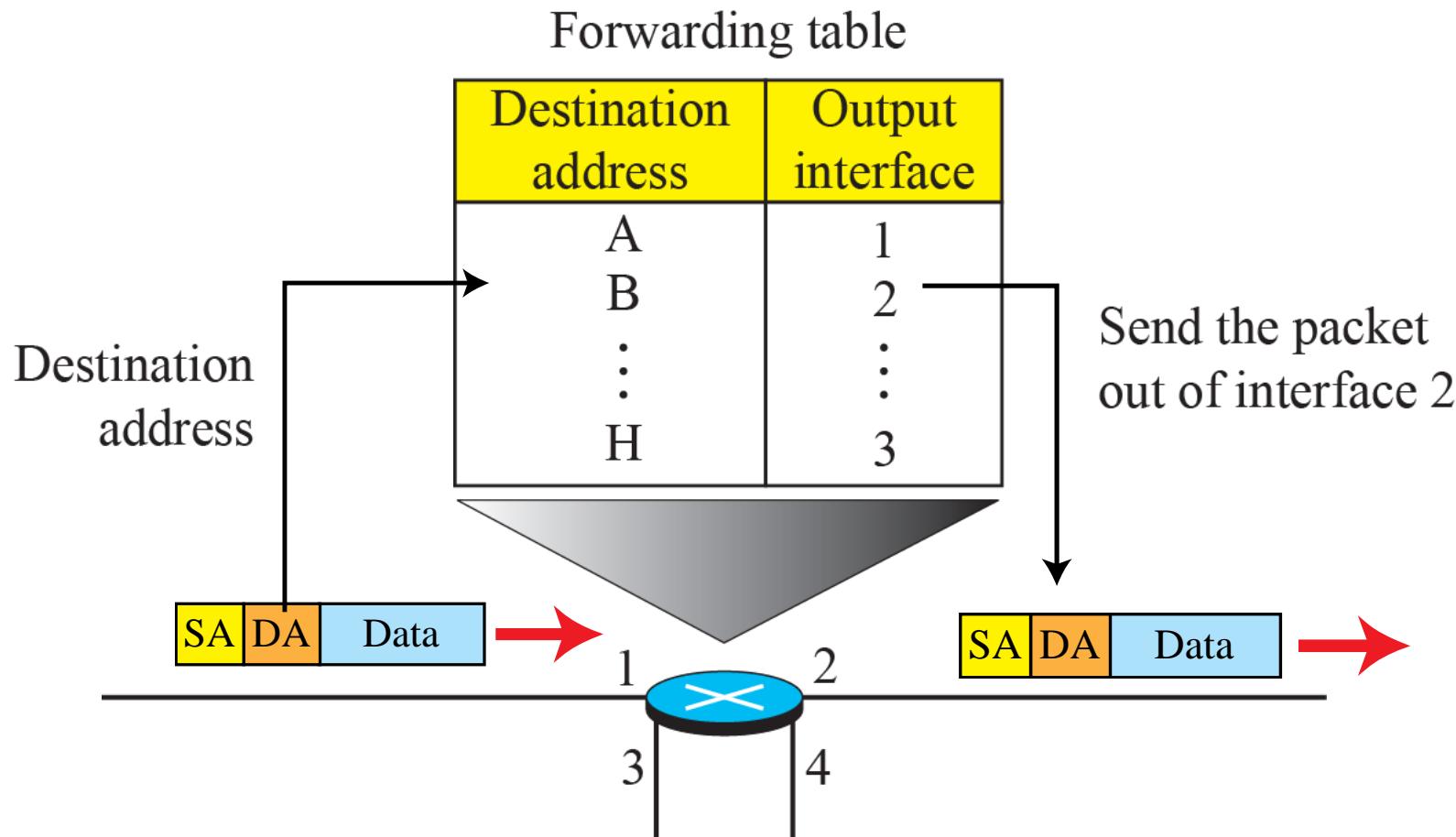
# Datagram forwarding table



**Figure 4.3: A connectionless packet-switched network**



**Figure 4.4: Forwarding process in a router when used in a connectionless network**



# 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

4.3 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

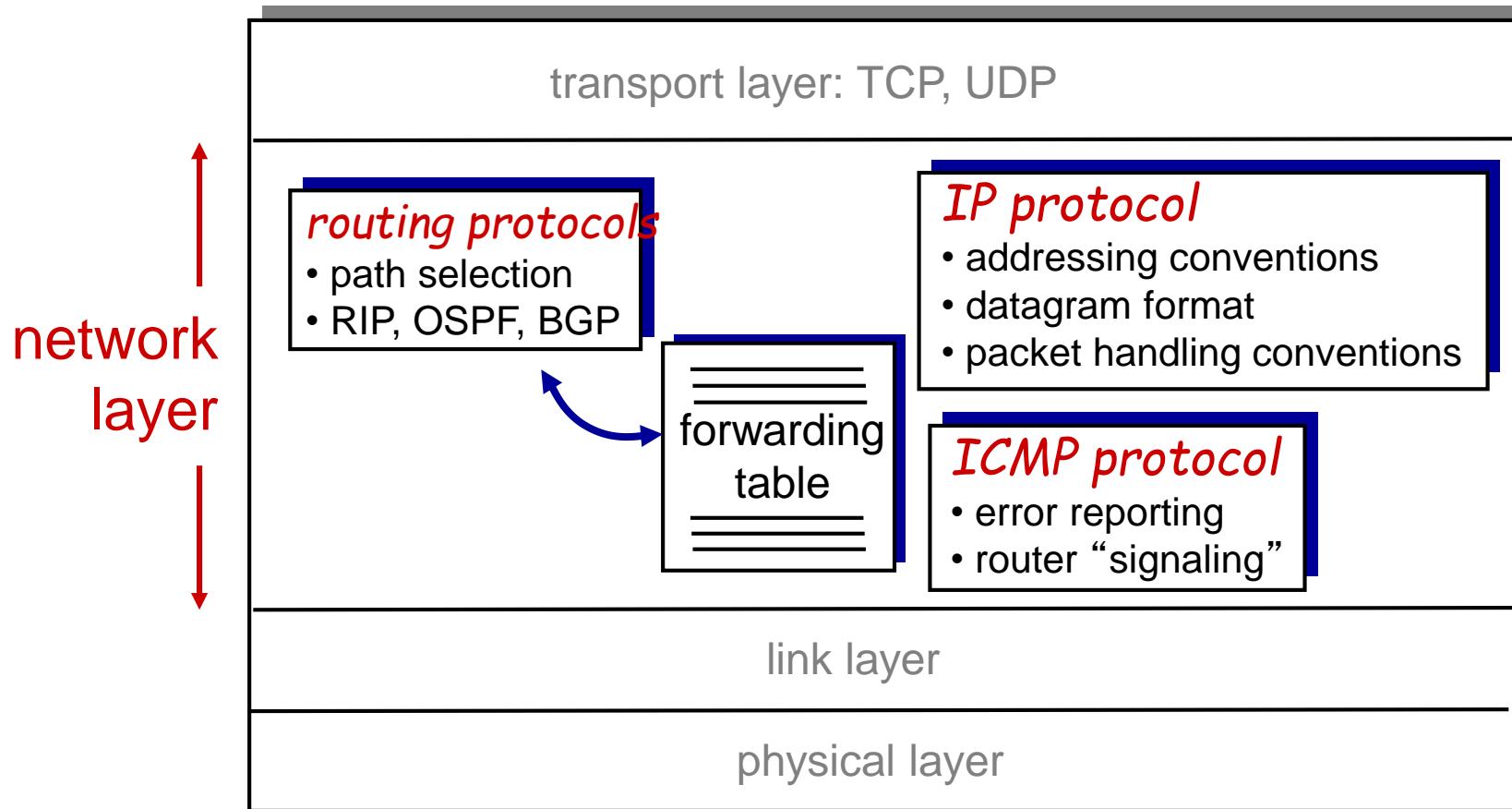
4.6 routing in the  
Internet

- RIP
- OSPF
- BGP

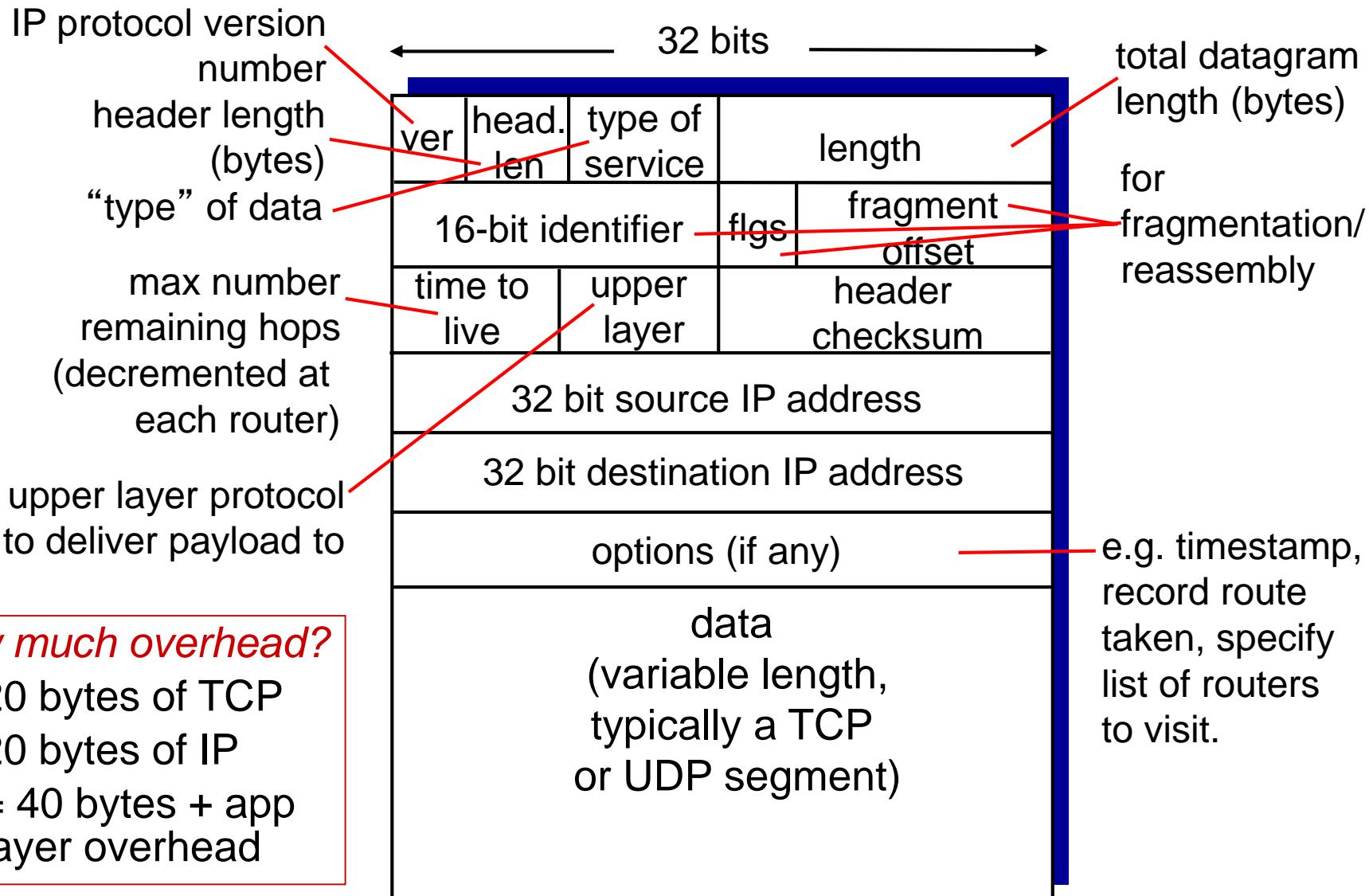
4.7 broadcast and  
multicast routing

# The Internet network layer

host, router network layer functions:

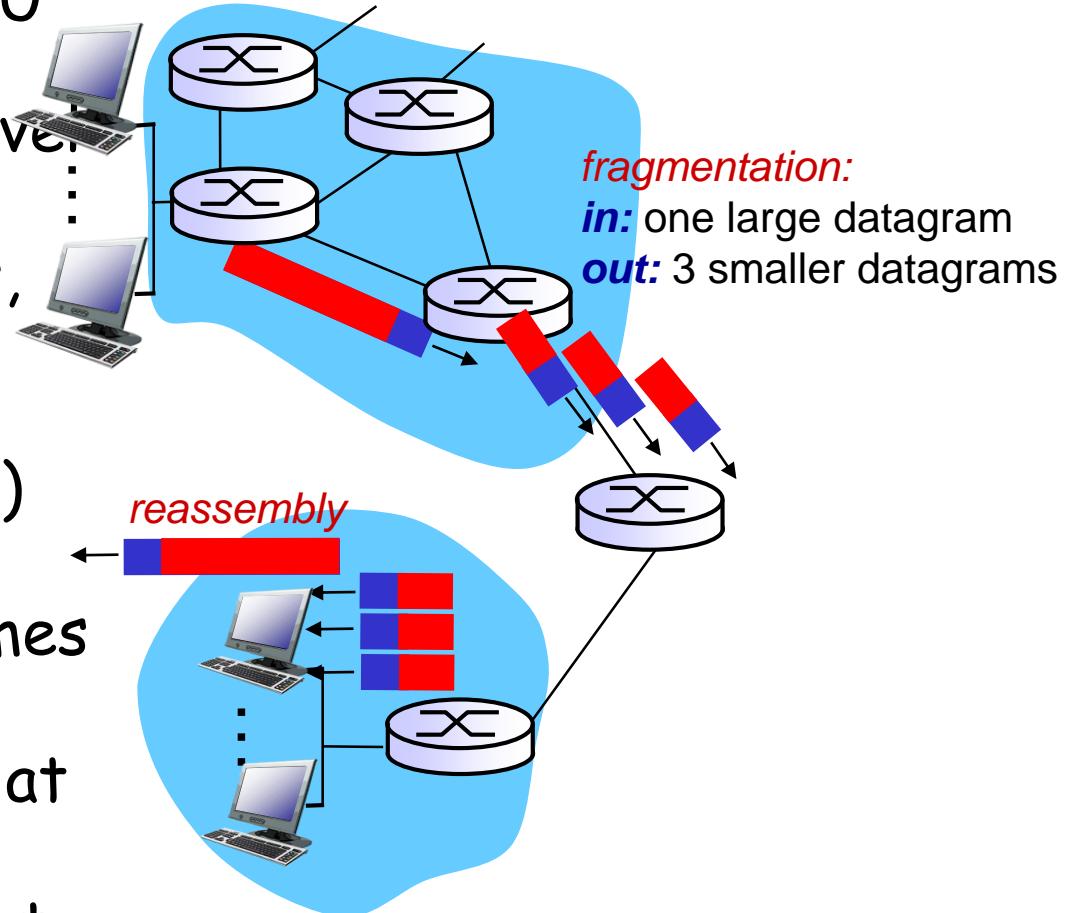


# IP datagram format



# 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

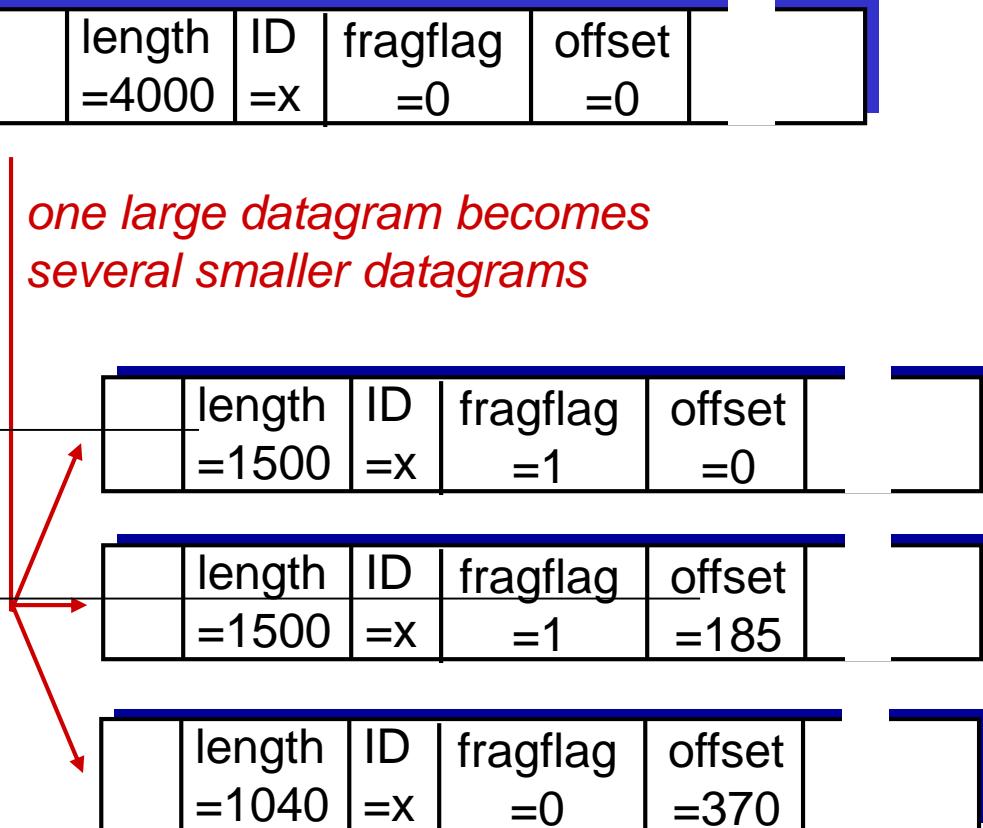
# IP fragmentation, reassembly

*example:*

- ❖ 4000 byte datagram
- ❖ MTU = 1500 bytes

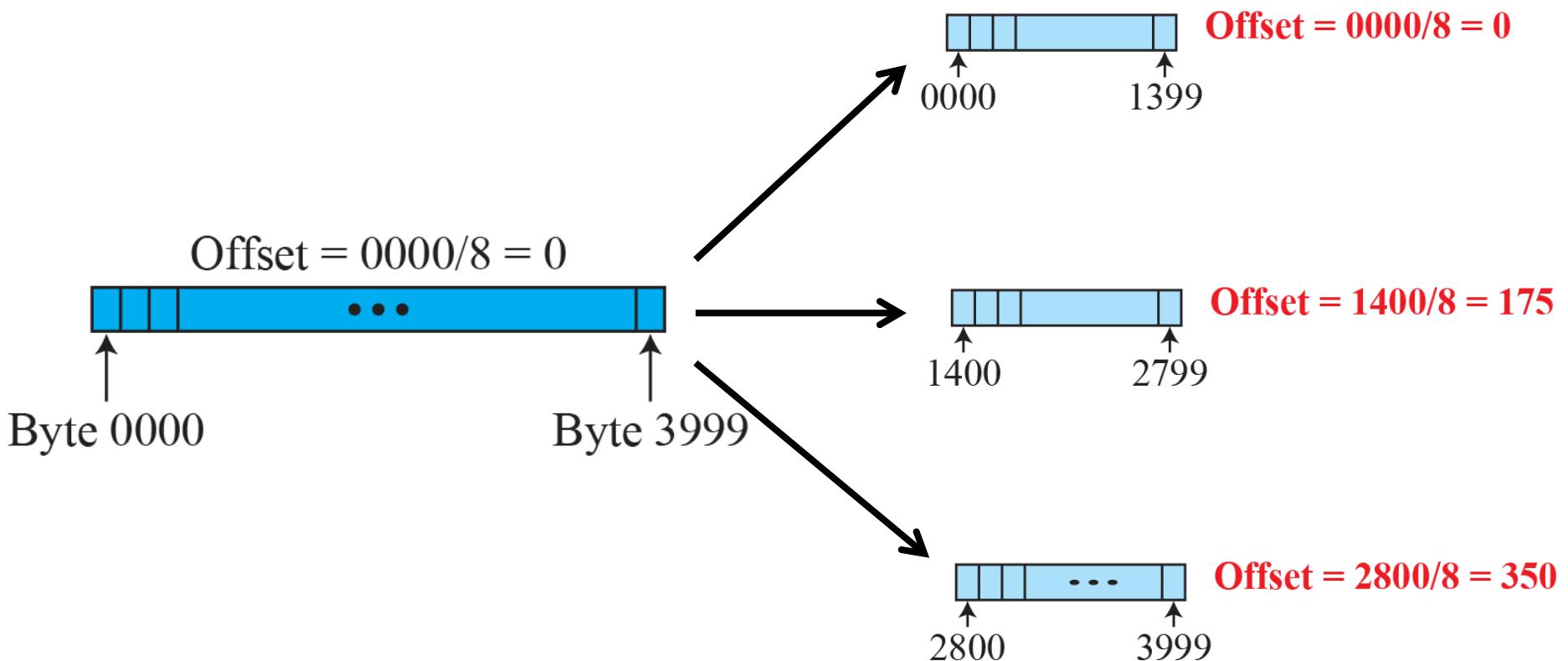
1480 bytes in data field

offset =  
 $1480/8$

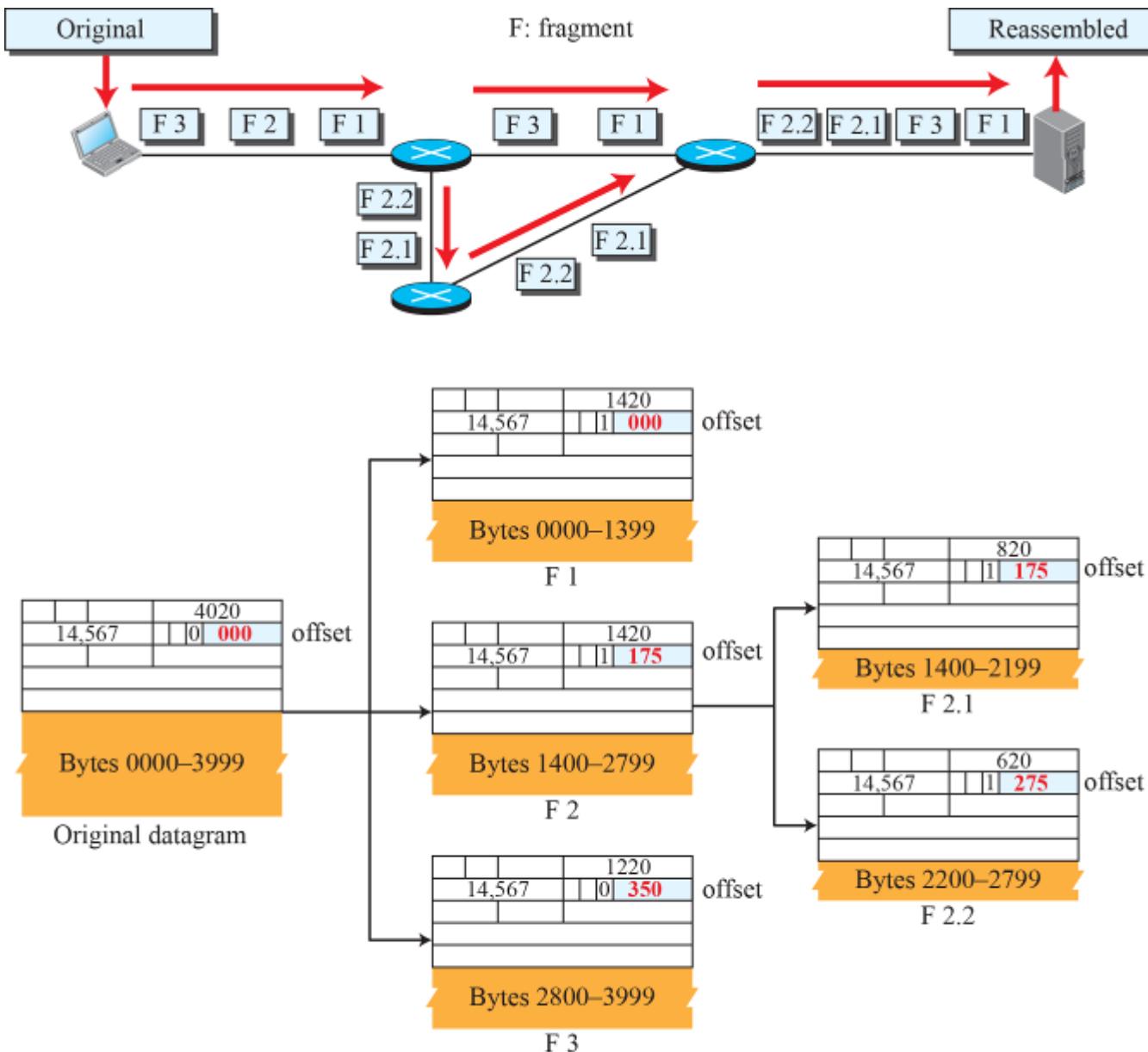


Total data length=?

**Figure 4.27:** Fragmentation example



**Figure 4.28: Detailed fragmentation example**



# Reassemble the original datagram

- ❖ Identify first fragment :  
offset=0
- ❖ Divide the length of first fragment by 8 :  
second fragment offset
- ❖ Divide the length of first and second :  
fragment by 8 : Third fragment offset
- ❖ Continue the process until last fragment :  
M bit=0

# Chapter 4: outline

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

## 4.2.2 IPv4 Addresses

The identifier used in the IP layer.

Identifies the connection of each device to the Internet  
is called the Internet address or IP address.

An IPv4 address is a 32-bit address

The IP address is the address of the connection, not  
the host or the router, because if the device is moved  
to another network, the IP address may be changed.

## 4.2.2 (*continued*)

- ❑ Address Space
- ❑ Notation
- ❑ Hierarchy in Addressing
- ❑ Classful Addressing
  - ❖ Address Depletion
  - ❖ Subnetting and Supernetting
  - ❖ Advantage of Classful Addressing

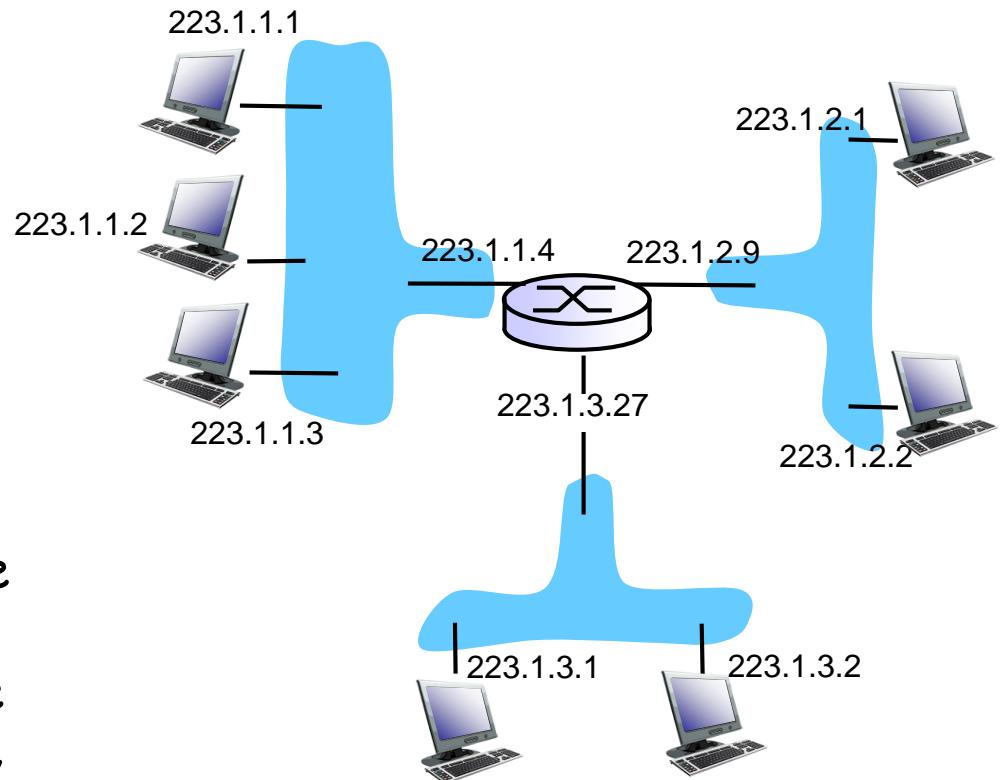
## 4.2.2 (*continued*)

### Classless Addressing

- ❖ Prefix Length: Slash Notation
- ❖ Extracting information from an address
- ❖ Address Mask
- ❖ Network Address
- ❖ Block Allocation
- ❖ Subnetting
- ❖ Address Aggregation
- ❖ Special Addresses

# IP addressing: introduction

- ❖ **IP address:** 32-bit identifier for host, router *interface*
- ❖ **interface:** connection between host/router and physical link
  - routers typically have multiple interfaces
  - host typically has one active interface (e.g., wired Ethernet, wireless 802.11)
- ❖ **one IP address associated with each interface**



$223.1.1.1 = \underbrace{11011111}_{223} \underbrace{00000001}_{1} \underbrace{00000001}_{1} \underbrace{00000001}_{1}$

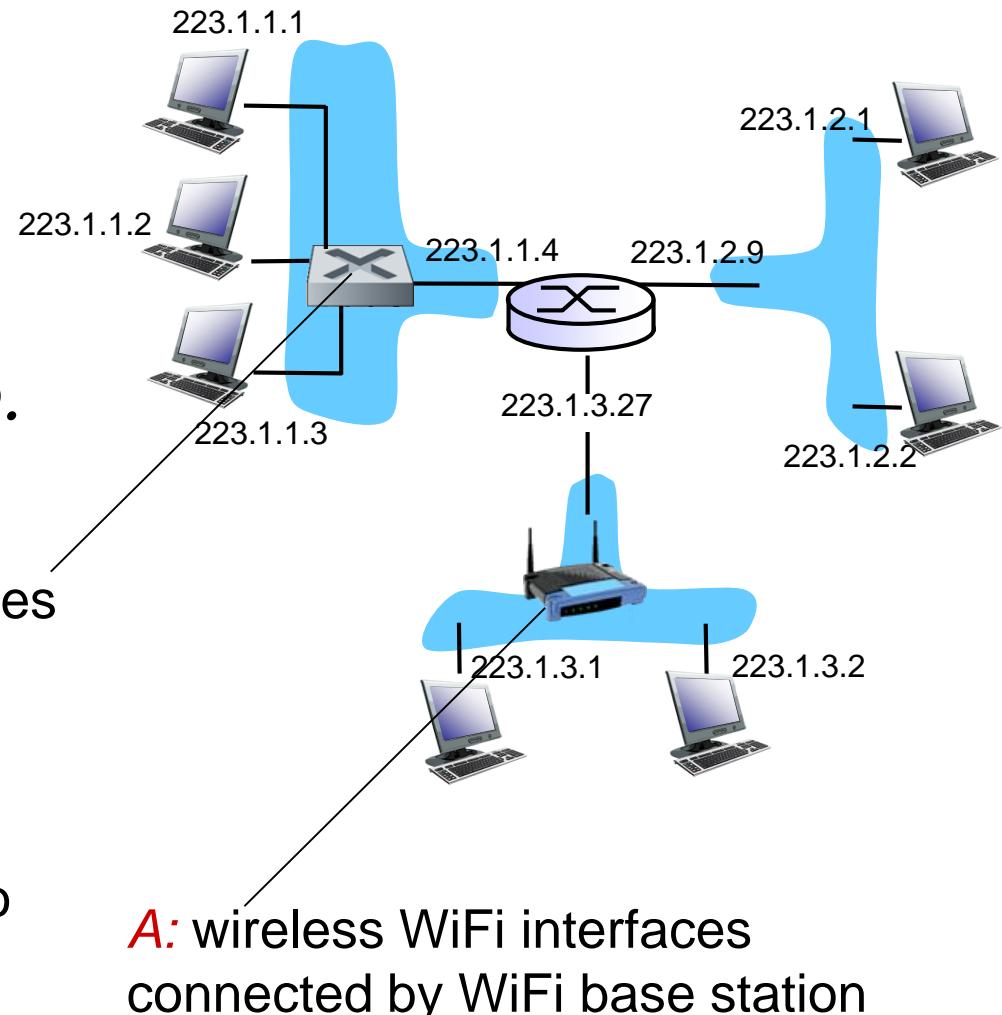
# IP addressing: introduction

**Q:** how are interfaces actually connected?

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

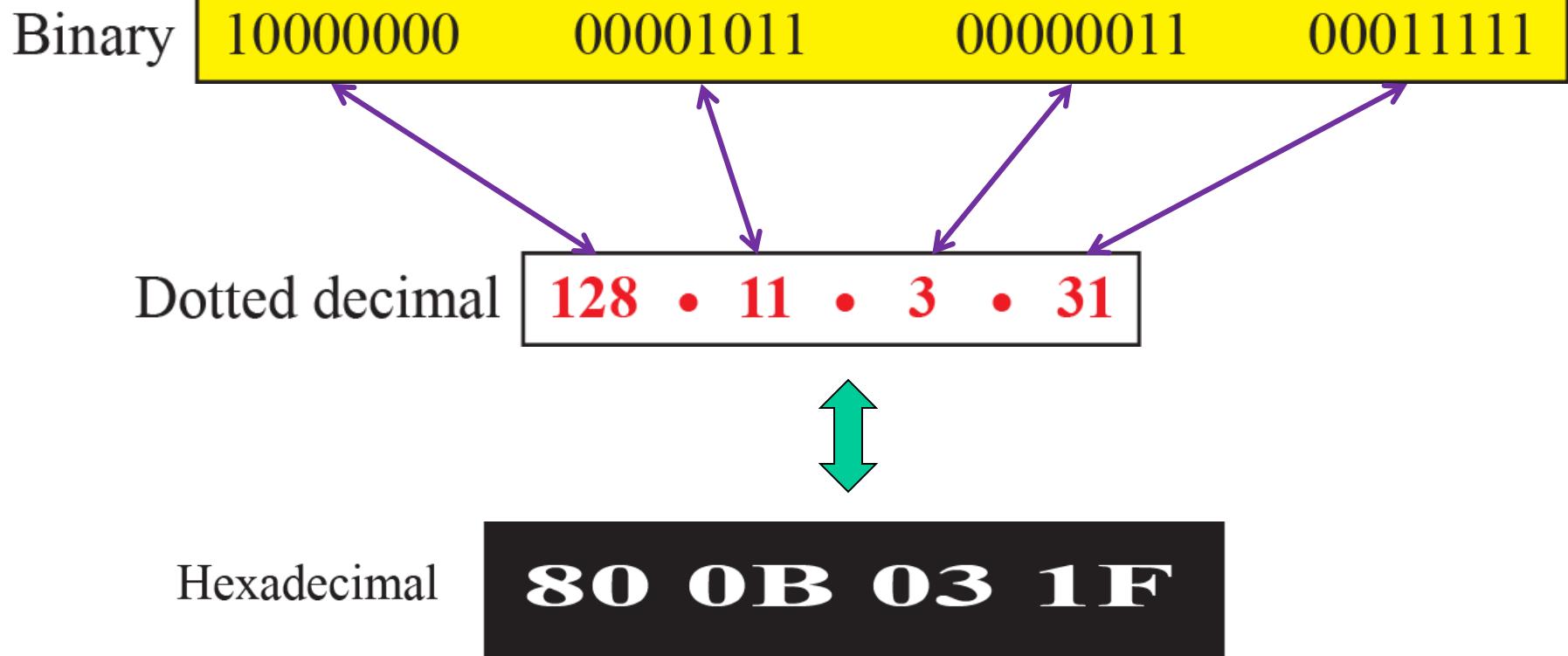
**A:** wired Ethernet interfaces connected by Ethernet switches

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

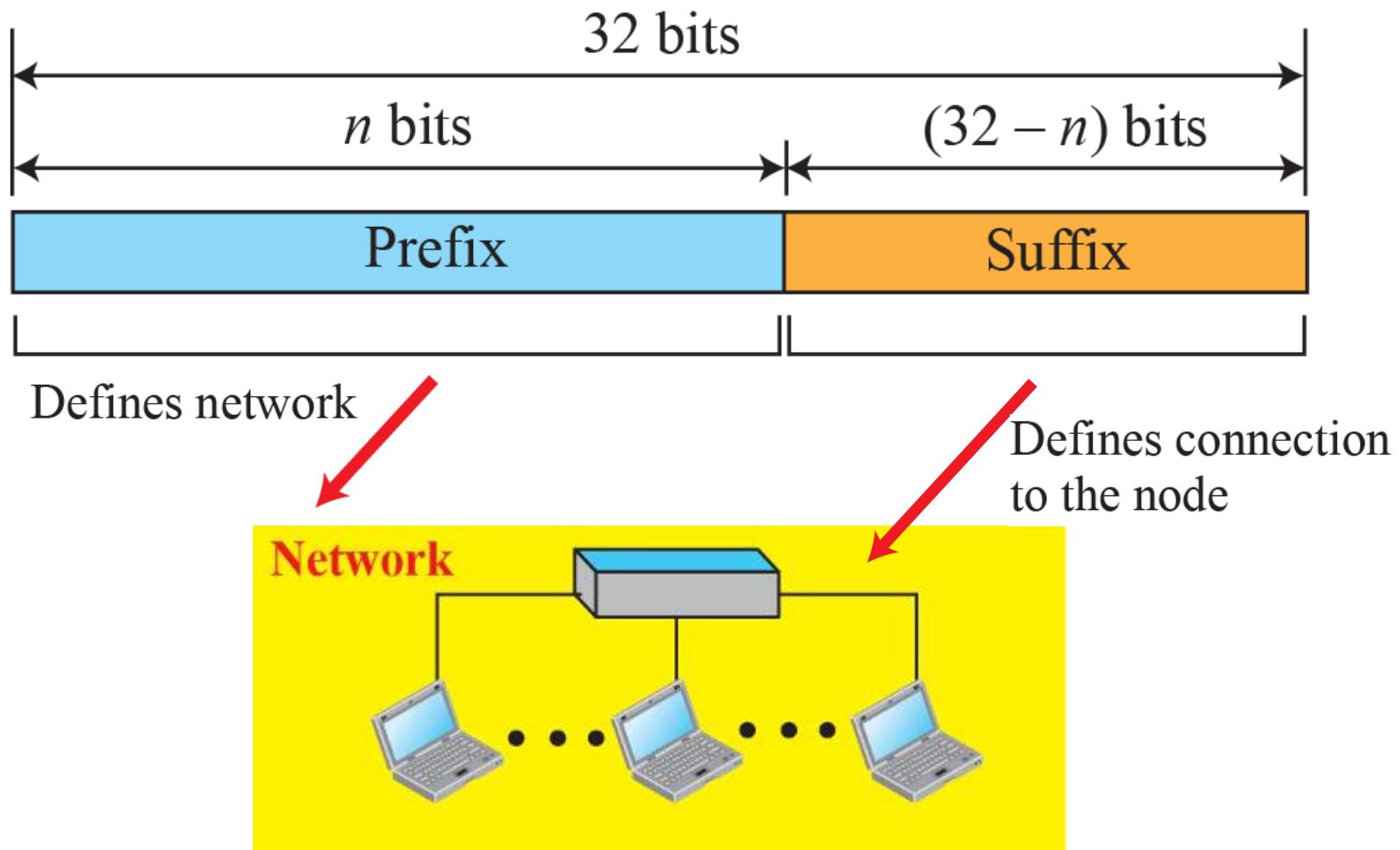


**A:** wireless WiFi interfaces connected by WiFi base station

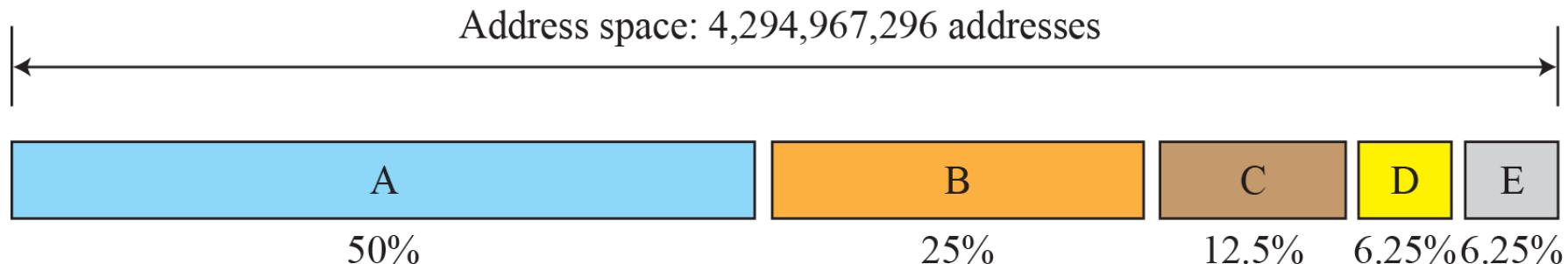
**Figure 4.29:** Three different notations in IPv4 addressing



**Figure 4.30: Hierarchy in addressing**



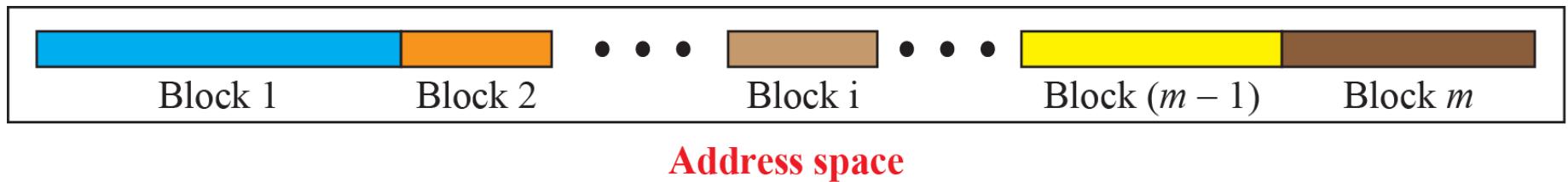
**Figure 4.31: Occupation of the address space in classful addressing**



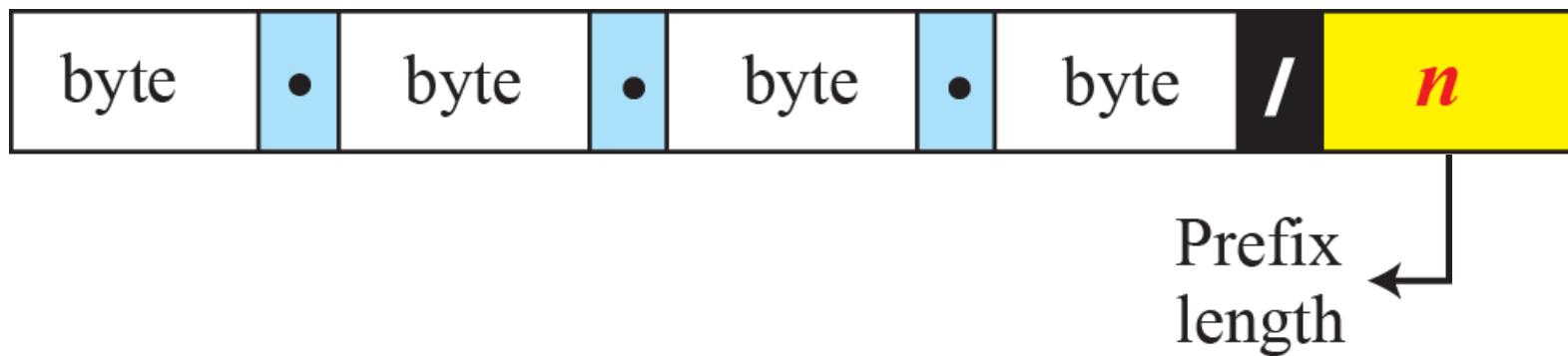
Class A	0 Prefix	Suffix
Class B	10 Prefix	Suffix
Class C	110 Prefix	Suffix
Class D	1110 Multicast addresses	
Class E	1111 Reserved for future use	

Class	Prefixes	First byte
A	$n = 8$ bits	0 to 127
B	$n = 16$ bits	128 to 191
C	$n = 24$ bits	192 to 223
D	Not applicable	224 to 239
E	Not applicable	240 to 255

**Figure 4.32:** Variable-length blocks in classless addressing



## **Slash notation (CIDR)**



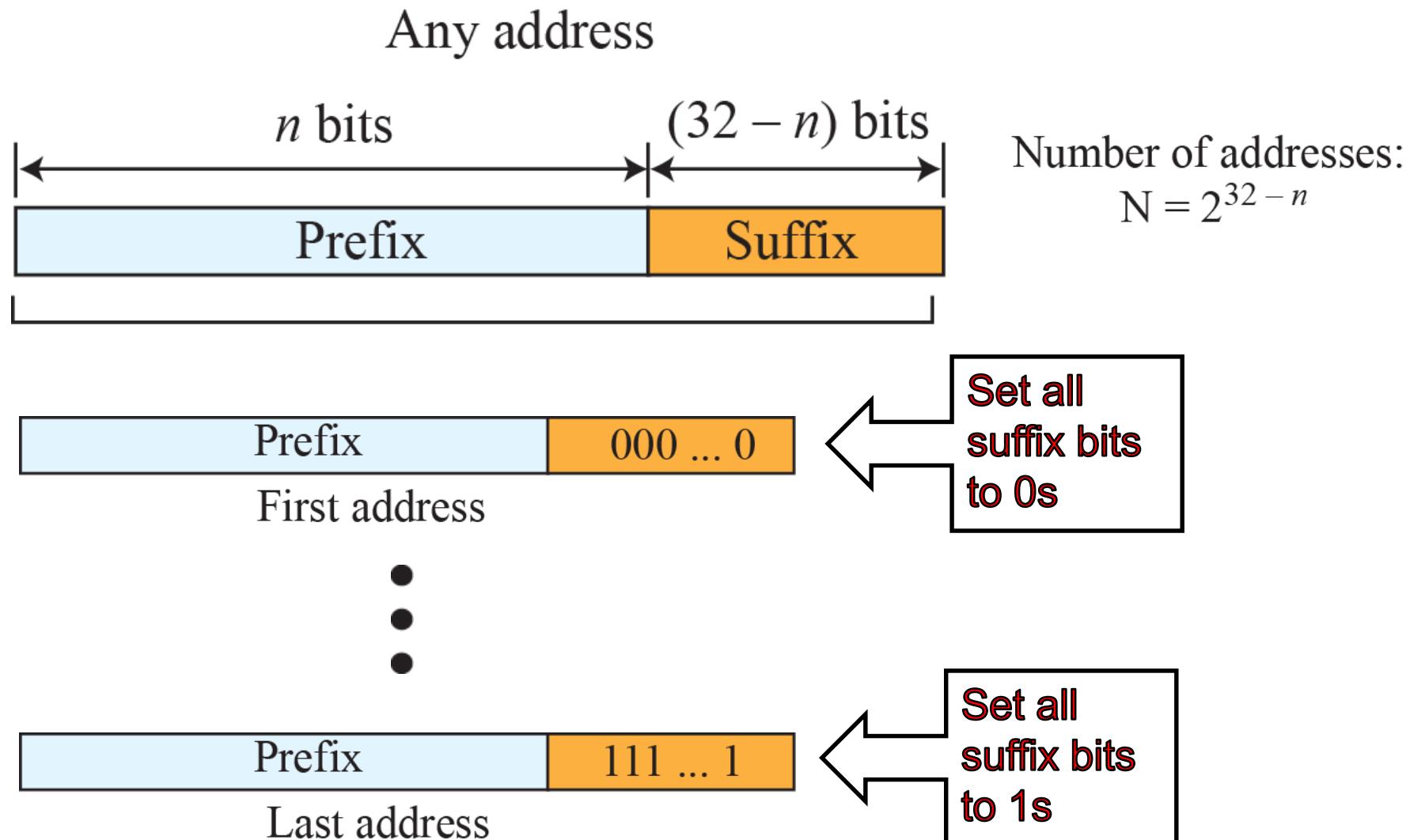
### **Examples:**

**12.24.76.8/8**

**23.14.67.92/12**

**220.8.24.255/25**

**Figure 4.34: Information extraction in classless addressing**



## Example 4.1

A classless address is given as 167.199.170.82/27. We can find the above three pieces of information as follows. The number of addresses in the network is  $2^{32-n} = 2^5 = 32$  addresses. The first address can be found by keeping the first 27 bits and changing the rest of the bits to 0s.

Address: 167.199.170.82/ <b>27</b>	10100111	11000111	10101010	01010010
First address: 167.199.170.64/ <b>27</b>	10100111	11000111	10101010	010 <b>00000</b>

The last address can be found by keeping the first 27 bits and changing the rest of the bits to 1s.

Address: 167.199.170.82/ <b>27</b>	10100111	11000111	10101010	01011111
Last address: 167.199.170.95/ <b>27</b>	10100111	11000111	10101010	010 <b>11111</b>

We repeat Example 4.1 using the mask. The mask in dotted-decimal notation is 255.255.255.224. The AND, OR, and NOT operations can be applied to individual bytes using calculators and applets at the book website.

Number of addresses in the block:  $N = \text{NOT}(\text{mask}) + 1 = 0.0.0.31 + 1 = 32 \text{ addresses}$

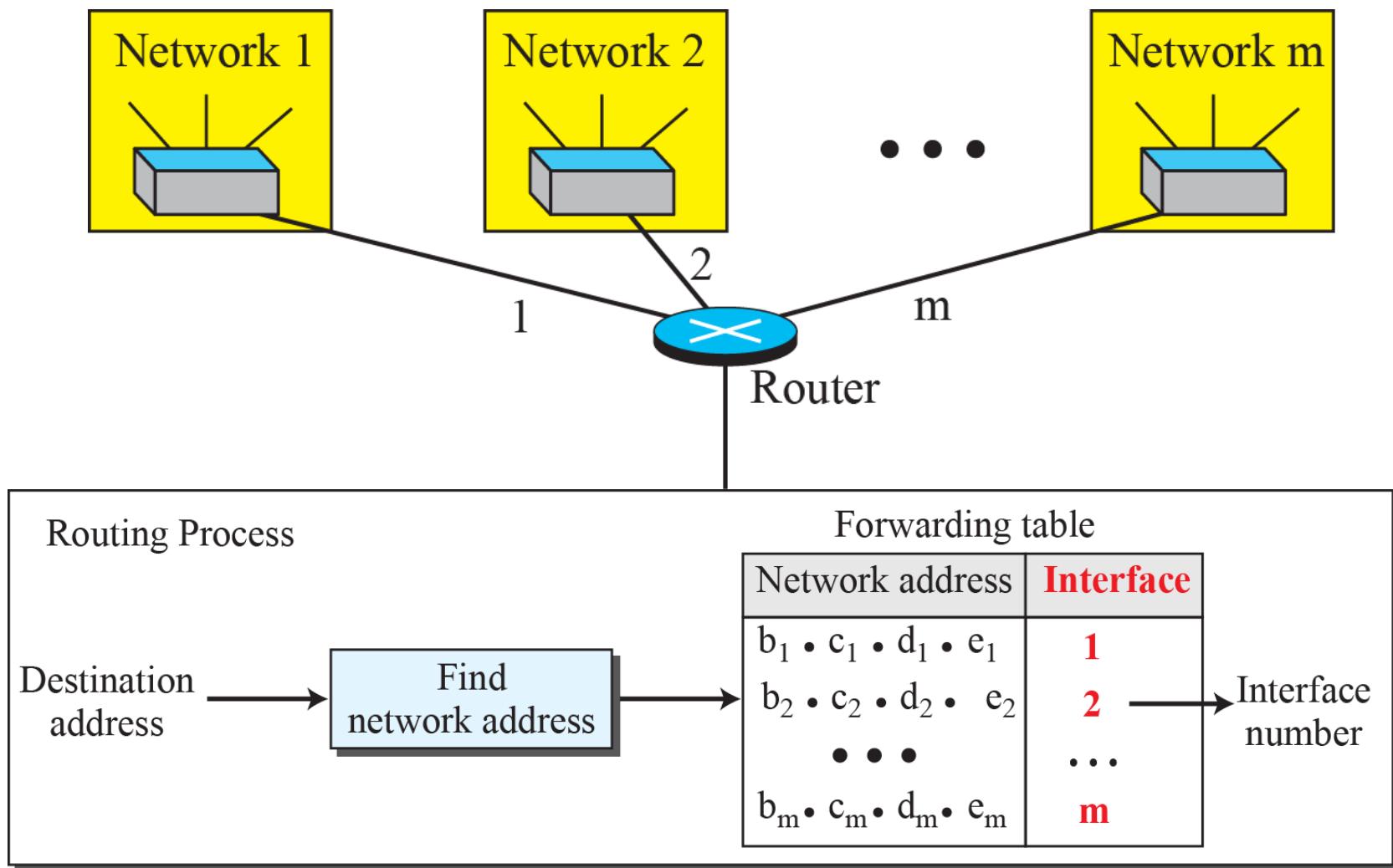
First address:  $\text{First} = (\text{address}) \text{ AND } (\text{mask})$

Last address:  $\text{Last} = (\text{address}) \text{ OR } (\text{NOT mask})$

In classless addressing, an address cannot per se define the block the address belongs to. For example, the address 230.8.24.56 can belong to many blocks. Some of them are shown below with the value of the prefix associated with that block.

Prefix length:16	→	Block:	230.8.0.0	to	230.8.255.255
Prefix length:20	→	Block:	230.8.16.0	to	230.8.31.255
Prefix length:26	→	Block:	230.8.24.0	to	230.8.24.63
Prefix length:27	→	Block:	230.8.24.32	to	230.8.24.63
Prefix length:29	→	Block:	230.8.24.56	to	230.8.24.63
Prefix length:31	→	Block:	230.8.24.56	to	230.8.24.57

**Figure 4.35: Network address**



By ICANN ( Internet Corporation for Assigned Names & Numbers )

1. No of addresses N: power of 2, n to be an int.
2. Contiguous block . But selection such that first address needs to be divisible by N. ( I address : suffix zero, i.e 32-n zeros. The decimal value then equal to the decimal equivalent of prefix with suffix all zeros.

i.e  $FA = \text{prefix} * 2^{32-n} = \text{prefix in decimal} \times N$

An ISP has requested a block of 1000 addresses. Since 1000 is not a power of 2, 1024 addresses are granted. The prefix length is calculated as  $n = 32 - \log_2 1024 = 22$ . An available block, 18.14.12.0/**22**, is granted to the ISP. It can be seen that the first address in decimal is 302,910,464, which is divisible by 1024.

# Subnetting

- ❖ Divide range into several subranges & assign each one to different subnetwork.
- ❖ Subnetwork -- into several sub-subnetwork.
- ❖ Design: to enable routing.
- ❖ Let Total no of addresses is  $N$  with prefix  $n$ .
- ❖ Total no of addresses is  $N_{\text{sub}}$  with prefix  $n_{\text{sub}}$ .

Rule: 1.  $N_{\text{sub}}$  is power of 2

$$2. n_{\text{sub}} = 32 - \log_2 N_{\text{sub}}$$

3. FA in each subnet should be divisible by  $N_{\text{sub}}$ .

An organization is granted a block of addresses with the beginning address 14.24.74.0/**24**. The organization needs to have 3 subblocks of addresses to use in its three subnets: one subblock of 10 addresses, one subblock of 60 addresses, and one subblock of 120 addresses. Design the subblocks.

## Solution

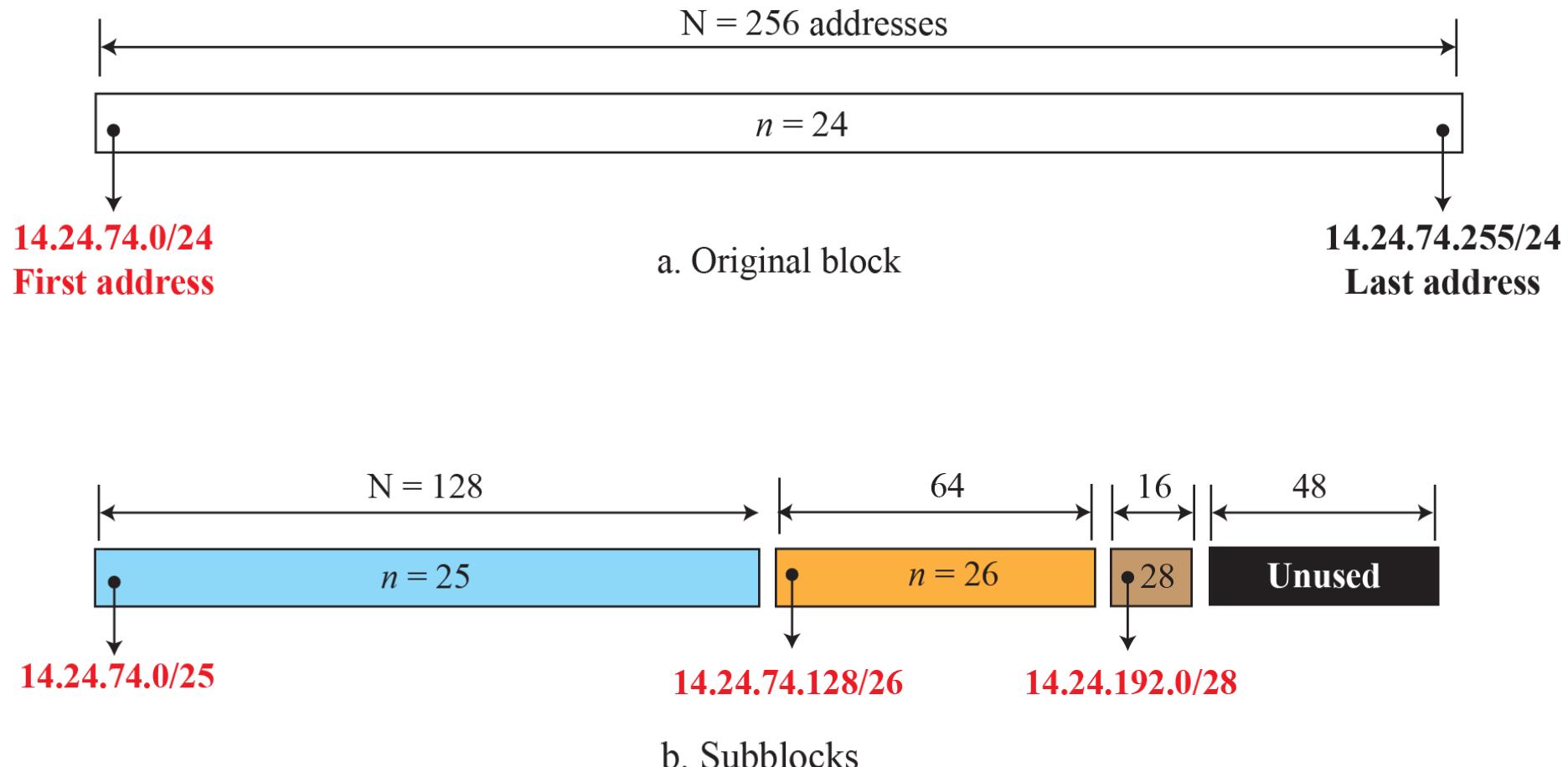
There are  $2^{32-24} = 256$  addresses in this block. The first address is 14.24.74.0/**24**; the last address is 14.24.74.255/**24**. To satisfy the third requirement, we assign addresses to subblocks, starting with the largest and ending with the smallest one.

- a. The number of addresses in the largest subblock, which requires 120 addresses, is not a power of 2. We allocate 128 addresses. The subnet mask for this subnet can be found as  $n_1 = 32 - \log_2 128 = 25$ . The first address in this block is 14.24.74.0/**25**; the last address is 14.24.74.127/**25**.
- b. The number of addresses in the second largest subblock, which requires 60 addresses, is not a power of 2 either. We allocate 64 addresses. The subnet mask for this subnet can be found as  $n_2 = 32 - \log_2 64 = 26$ . The first address in this block is 14.24.74.128/**26**; the last address is 14.24.74.191/**26**.

c. The number of addresses in the smallest subblock, which requires 10 addresses, is not a power of 2. We allocate 16 addresses. The subnet mask for this subnet can be found as  $n_1 = 32 - \log_2 16 = 28$ . The first address in this block is 14.24.74.192/28; the last address is 14.24.74.207/28.

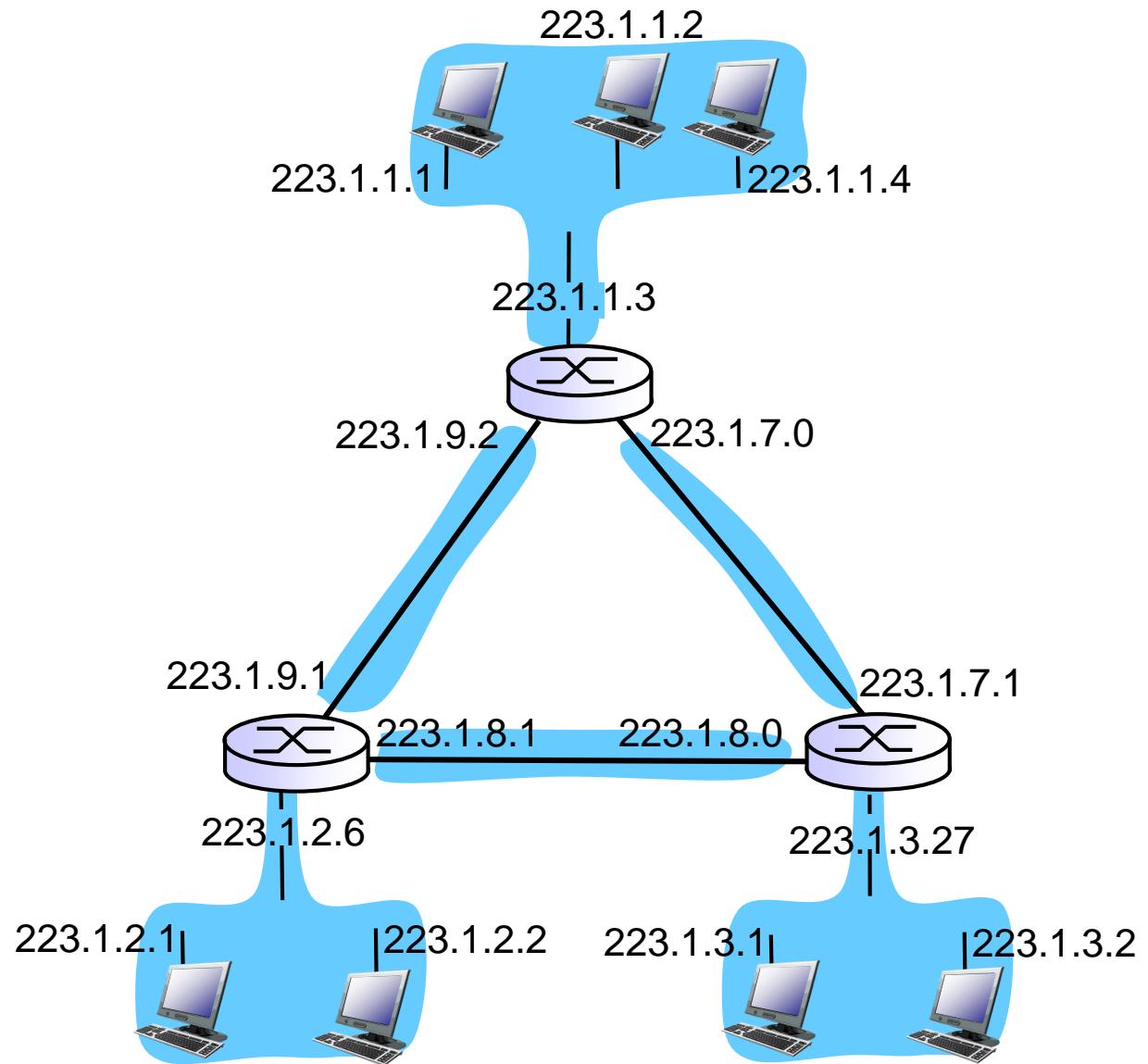
If we add all addresses in the previous subblocks, the result is 208 addresses, which means 48 addresses are left in reserve. The first address in this range is 14.24.74.208. The last address is 14.24.74.255. We don't know about the prefix length yet. Figure 4.36 shows the configuration of blocks. We have shown the first address in each block.

**Figure 4.36: Solution to Example 4.5**



# Subnets

how many?



# IP addressing: CIDR

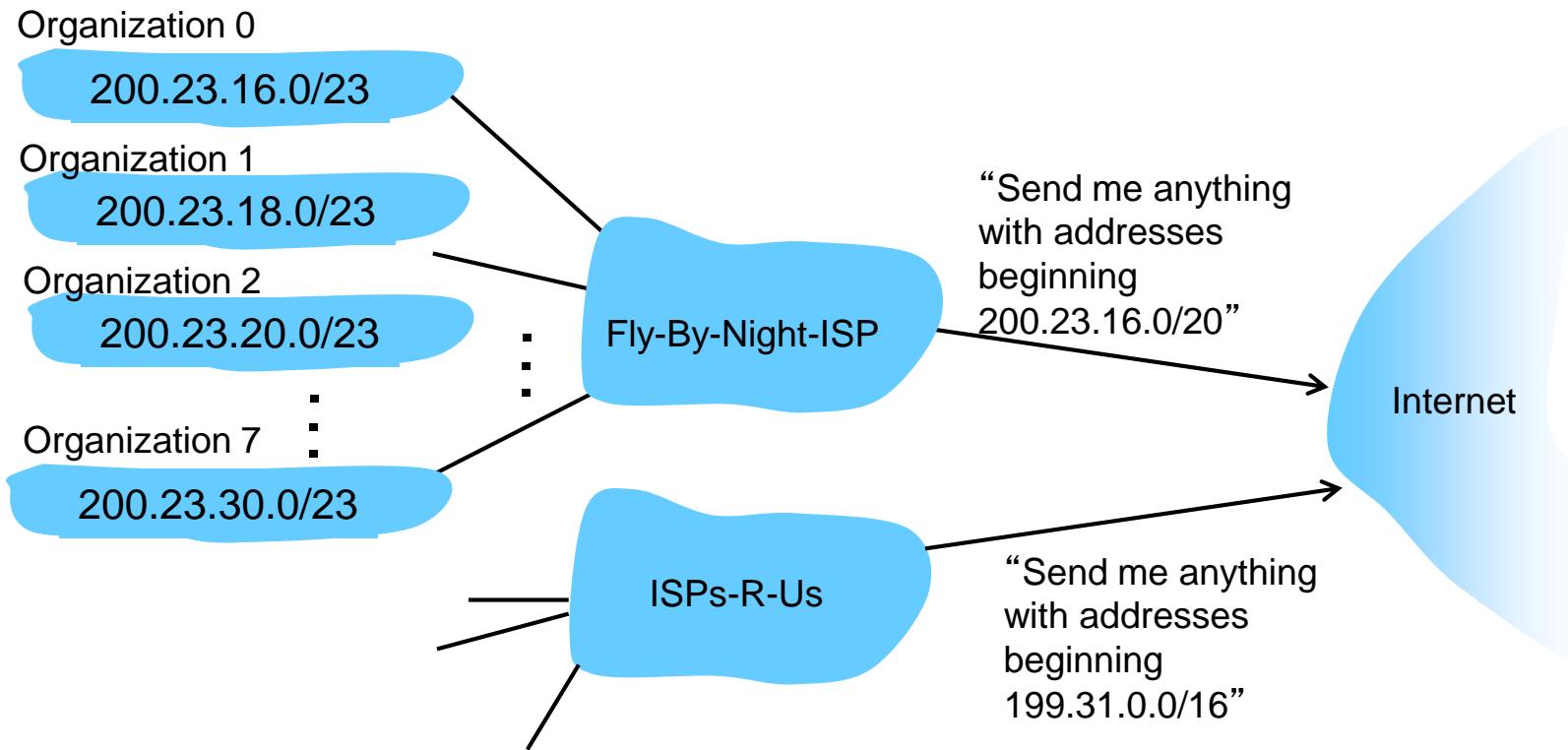
## CIDR: Classless InterDomain Routing

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



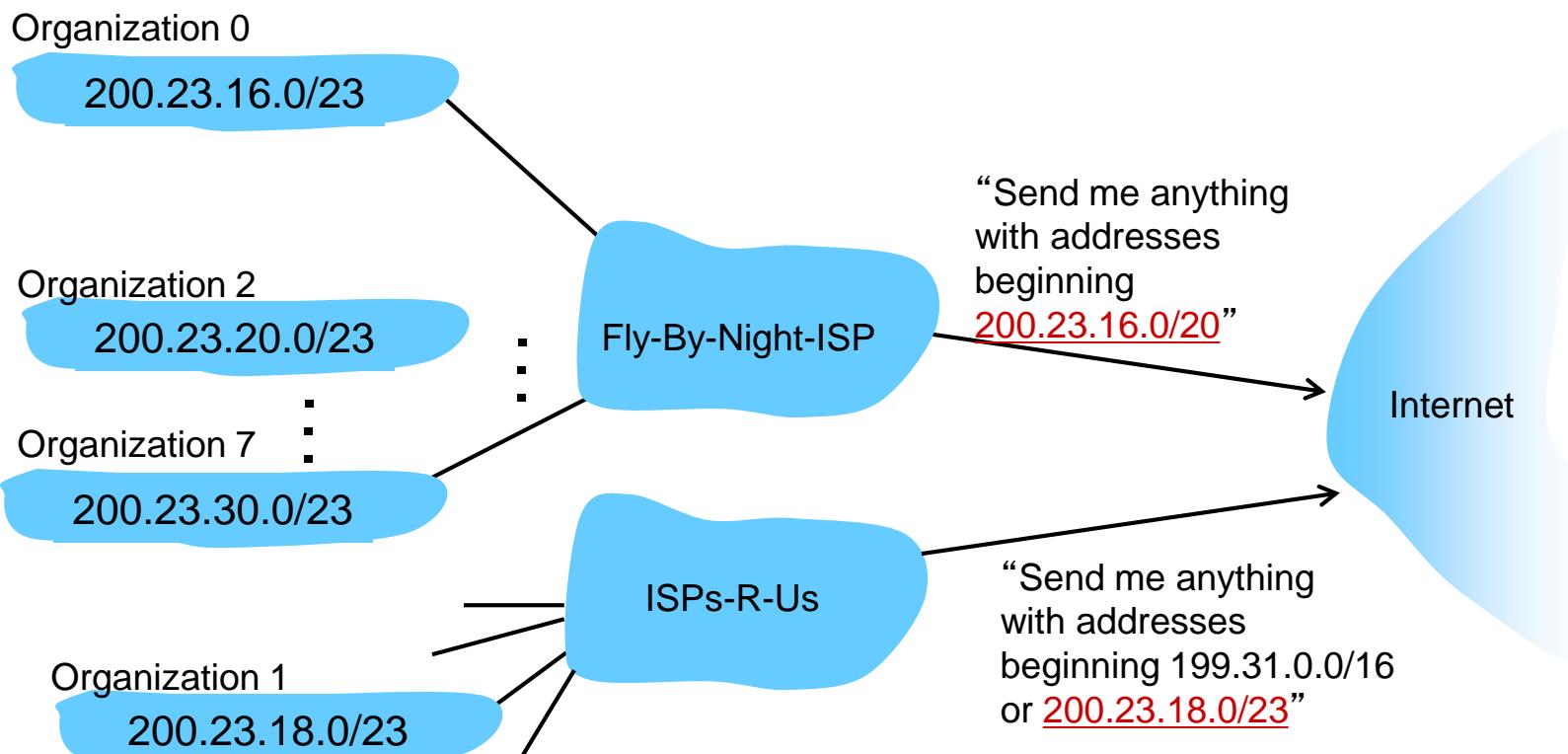
# Hierarchical addressing: route aggregation

hierarchical addressing allows efficient advertisement of routing information:



# Hierarchical addressing: more specific routes

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



# Special Addresses

- ❖ 1. This host Address: 0.0.0.0
- ❖ 2. Limited broadcast address:  
255.255.255.255, all hosts in a network
- ❖ 3. Loopback address: 127.0.0.0/8.
- ❖ Packet in the same host. Software testing
- ❖ 4. Private address: 4 Blocks
- ❖ 17.16.0.0/12, 10.0.0.0/8, 169.254.0.0/16,  
192.168.0.0/16

# IP addressing: how to get a block?

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

**A:** ICANN: Internet Corporation for Assigned Names and Numbers <http://www.icann.org/>

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

# IP addresses: how to get one?

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

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

# DHCP: Dynamic Host Configuration Protocol

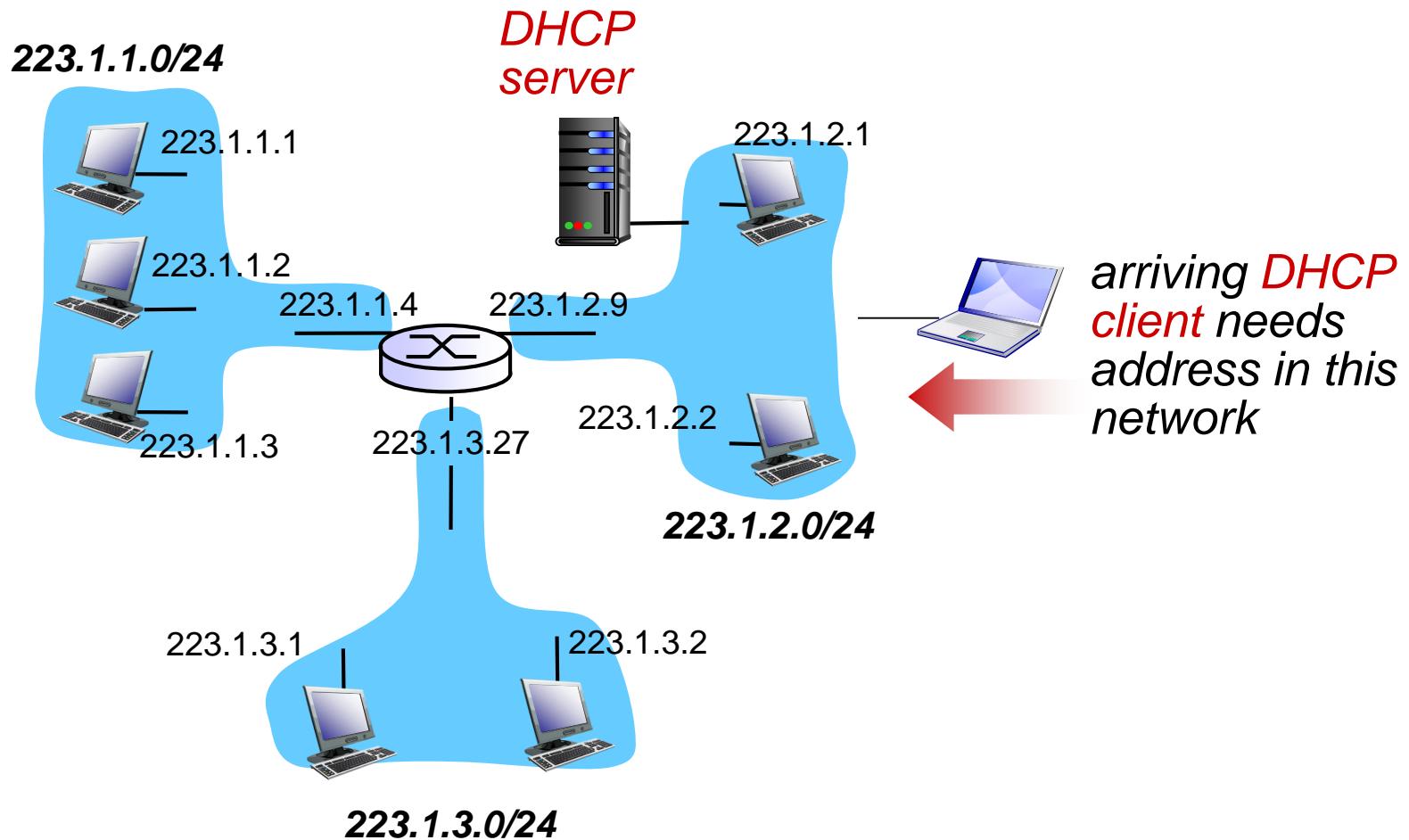
*goal:* allow host to dynamically obtain its IP address from network server when it joins network

- can renew its lease on address in use
- allows reuse of addresses (only hold address while connected/“on”)
- support for mobile users who want to join network

*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

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

arriving  
client



DHCP offer

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

DHCP request

```
src: 0.0.0.0, 68  
dest:: 255.255.255.255, 67  
yiaddr: 223.1.2.4  
transaction ID: 655  
lifetime: 3600 secs
```

DHCP ACK

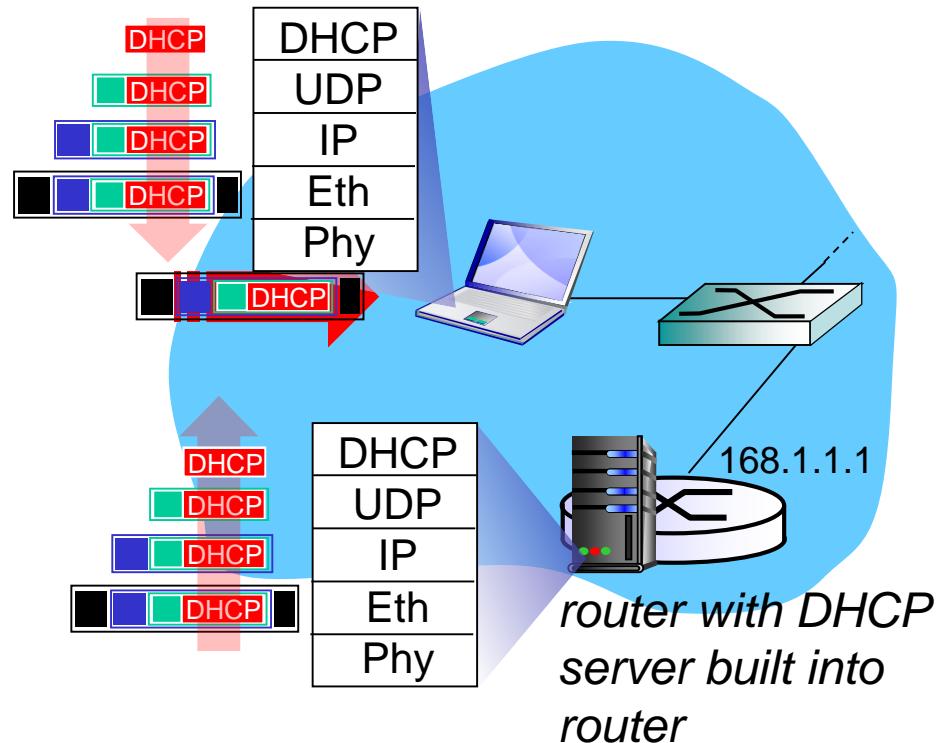
```
src: 223.1.2.5, 67  
dest: 255.255.255.255, 68  
yiaddr: 223.1.2.4  
transaction ID: 655  
lifetime: 3600 secs
```

# DHCP: more than IP addresses

## DHCP returns:

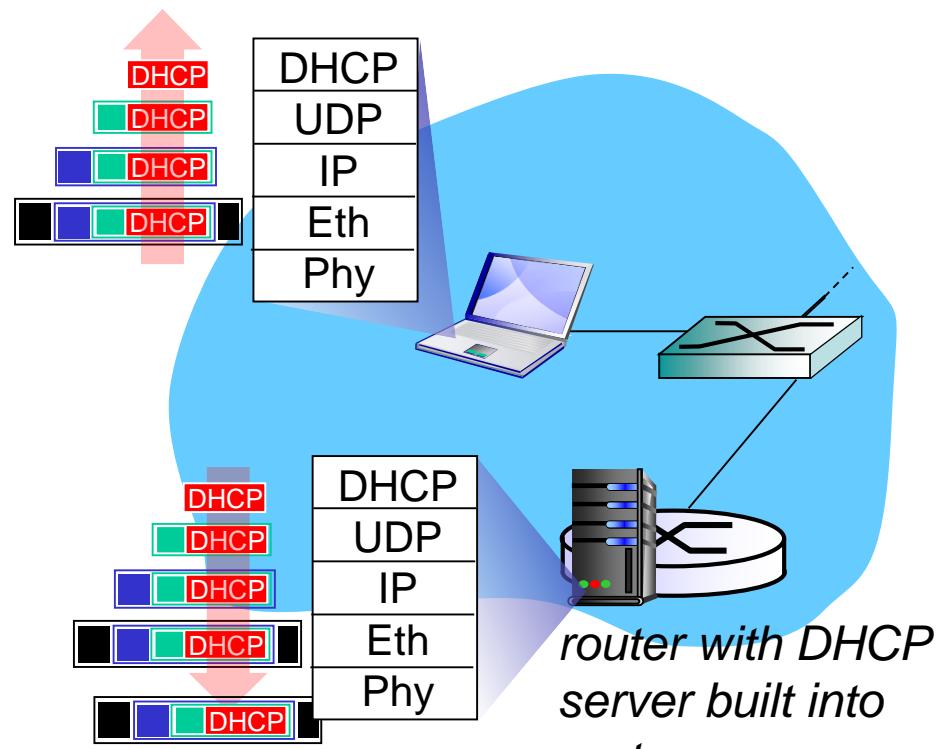
- IP address
- address of first-hop router for client
- name and IP address of DNS sever
- network mask (indicating network versus host portion of address)

# DHCP: example



- ❖ connecting laptop needs its IP address, addr of first-hop router, addr of DNS server: use DHCP
- ❖ DHCP request encapsulated in UDP, encapsulated in IP, encapsulated in 802.3 Ethernet
- ❖ Ethernet frame broadcast (dest: FFFFFFFFFFFF) on LAN, received at router running DHCP server
- ❖ Ethernet demuxed to IP demuxed, UDP demuxed to DHCP

# DHCP: example



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

# DHCP: Wireshark output (home LAN)

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

Hardware type: Ethernet

Hardware address length: 6

Hops: 0

**Transaction ID: 0x6b3a11b7**

Seconds elapsed: 0

Bootp flags: 0x0000 (Unicast)

Client IP address: 0.0.0.0 (0.0.0.0)

Your (client) IP address: 0.0.0.0 (0.0.0.0)

Next server IP address: 0.0.0.0 (0.0.0.0)

Relay agent IP address: 0.0.0.0 (0.0.0.0)

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

Server host name not given

Boot file name not given

Magic cookie: (OK)

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

Option: (61) Client identifier

Length: 7; Value: 010016D323688A;

Hardware type: Ethernet

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

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

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

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

Length: 11; Value: 010F03062C2E2F1F21F92B

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

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

44 = NetBIOS over TCP/IP Name Server

request

reply

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

Hardware type: Ethernet

Hardware address length: 6

Hops: 0

**Transaction ID: 0x6b3a11b7**

Seconds elapsed: 0

Bootp flags: 0x0000 (Unicast)

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

Your (client) IP address: 0.0.0.0 (0.0.0.0)

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

Relay agent IP address: 0.0.0.0 (0.0.0.0)

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

Server host name not given

Boot file name not given

Magic cookie: (OK)

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

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

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

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

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

Length: 12; Value: 445747E2445749F244574092;

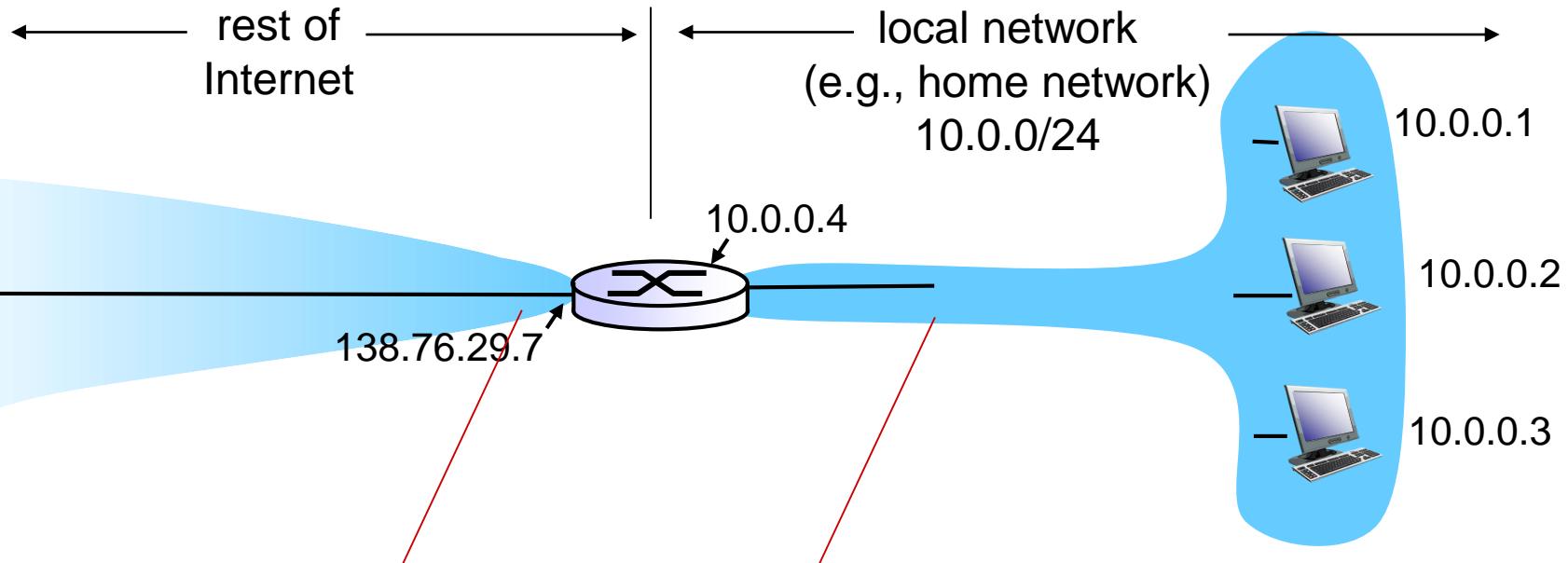
IP Address: 68.87.71.226;

IP Address: 68.87.73.242;

IP Address: 68.87.64.146

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

# 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

**implementation:** NAT router must:

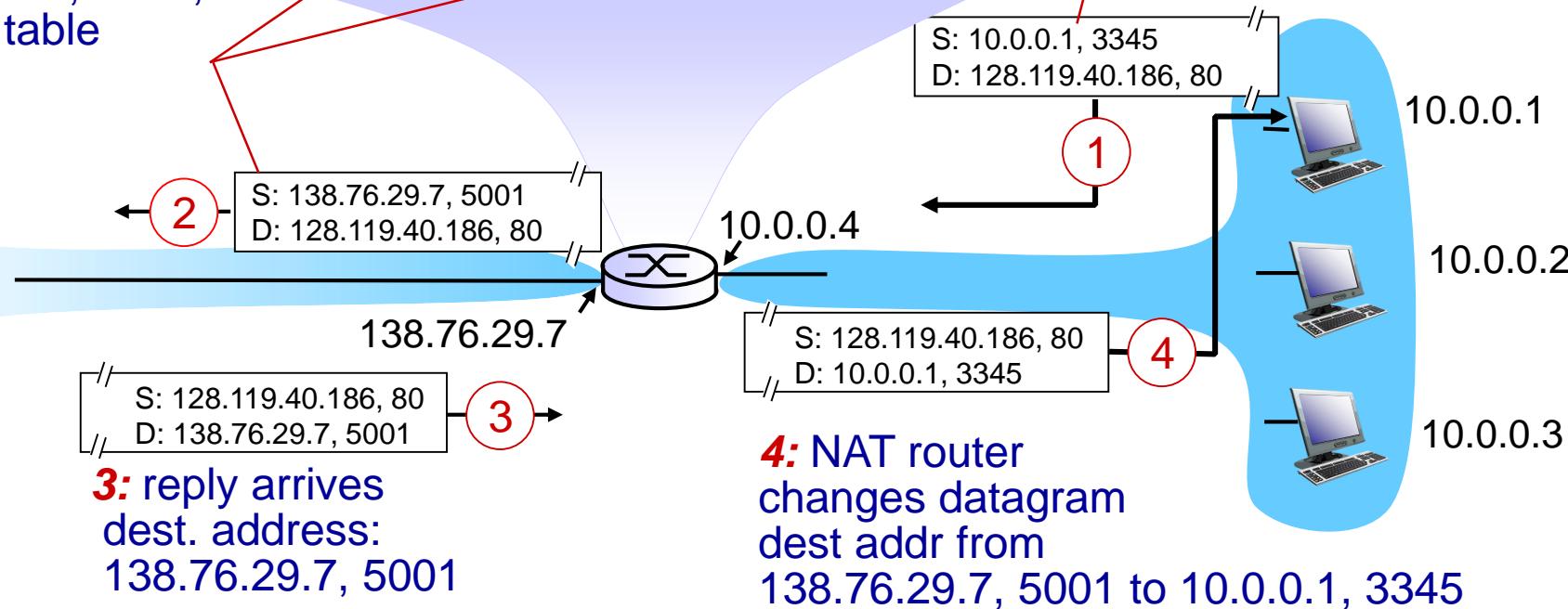
- *outgoing datagrams: replace* (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #)  
. . . remote clients/servers will respond using (NAT IP address, new port #) as destination addr
- *remember (in NAT translation table) every (source IP address, port #) to (NAT IP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

# NAT: network address translation

2: NAT router changes datagram source addr from 10.0.0.1, 3345 to 138.76.29.7, 5001, updates table

NAT translation table	
WAN side addr	LAN side addr
138.76.29.7, 5001	10.0.0.1, 3345
.....	.....

1: host 10.0.0.1 sends datagram to 128.119.40.186, 80



3: reply arrives  
dest. address:  
138.76.29.7, 5001

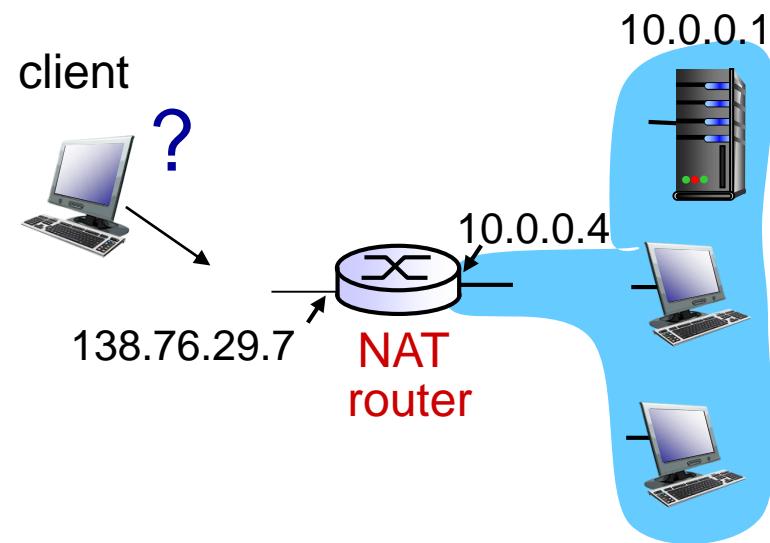
4: NAT router  
changes datagram  
dest addr from  
138.76.29.7, 5001 to 10.0.0.1, 3345

# NAT: network address translation

- ❖ 16-bit port-number field:
  - 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
    - ⑩ NAT possibility must be taken into account by app designers, e.g., P2P applications
  - address shortage should instead be solved by IPv6

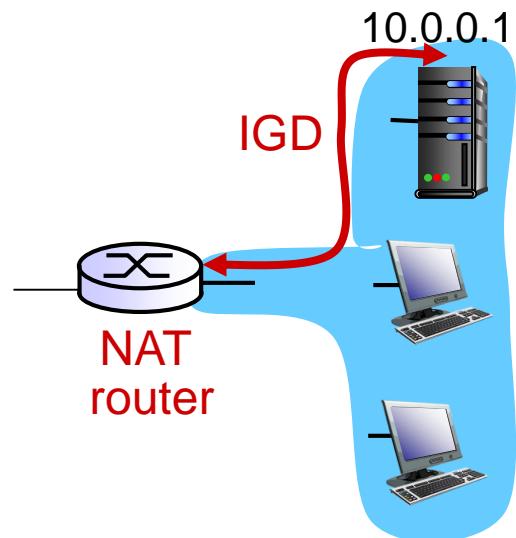
# NAT traversal problem

- ❖ 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
- ❖ **solution1:** statically configure NAT to forward incoming connection requests at given port to server
  - e.g., (123.76.29.7, port 25000) always forwarded to 10.0.0.1 port 25000



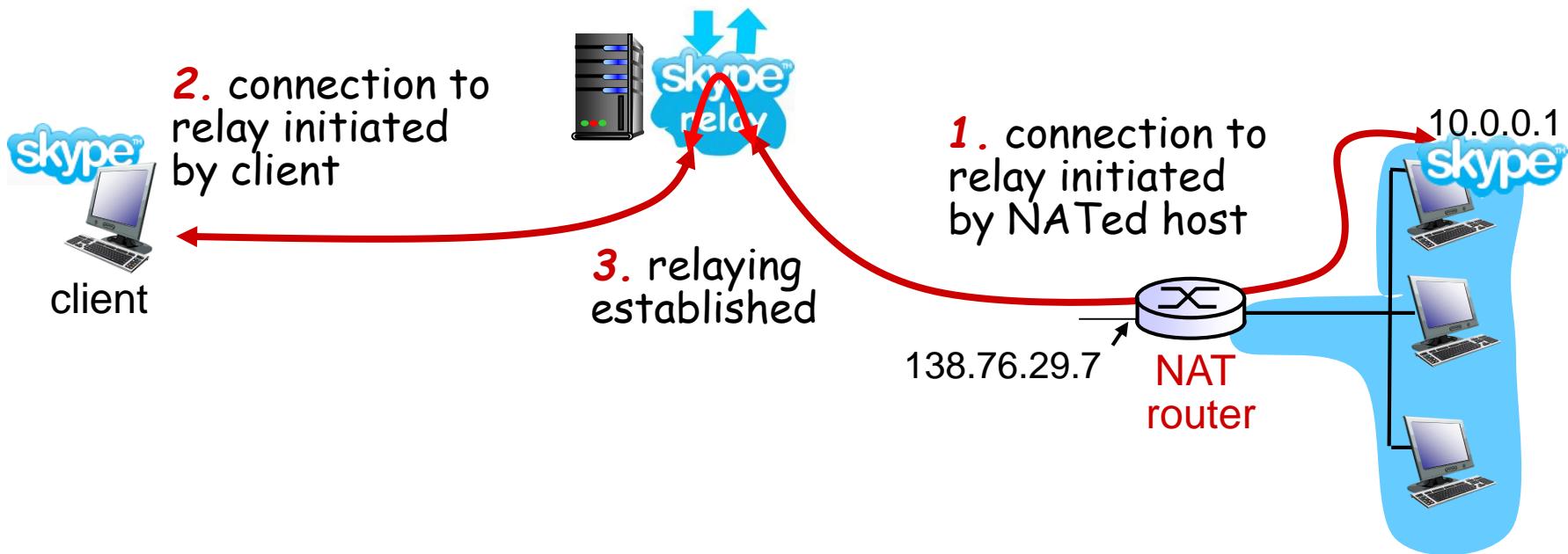
# NAT traversal problem

- ❖ **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



# NAT traversal problem

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



# Chapter 4: outline

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:

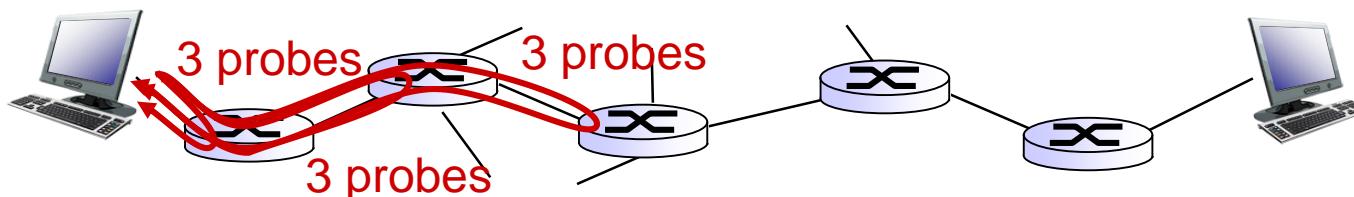
- ICMP msgs carried in IP datagrams

- ❖ **ICMP message:** type, code plus first 8 bytes of IP datagram causing error

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

# Traceroute and ICMP

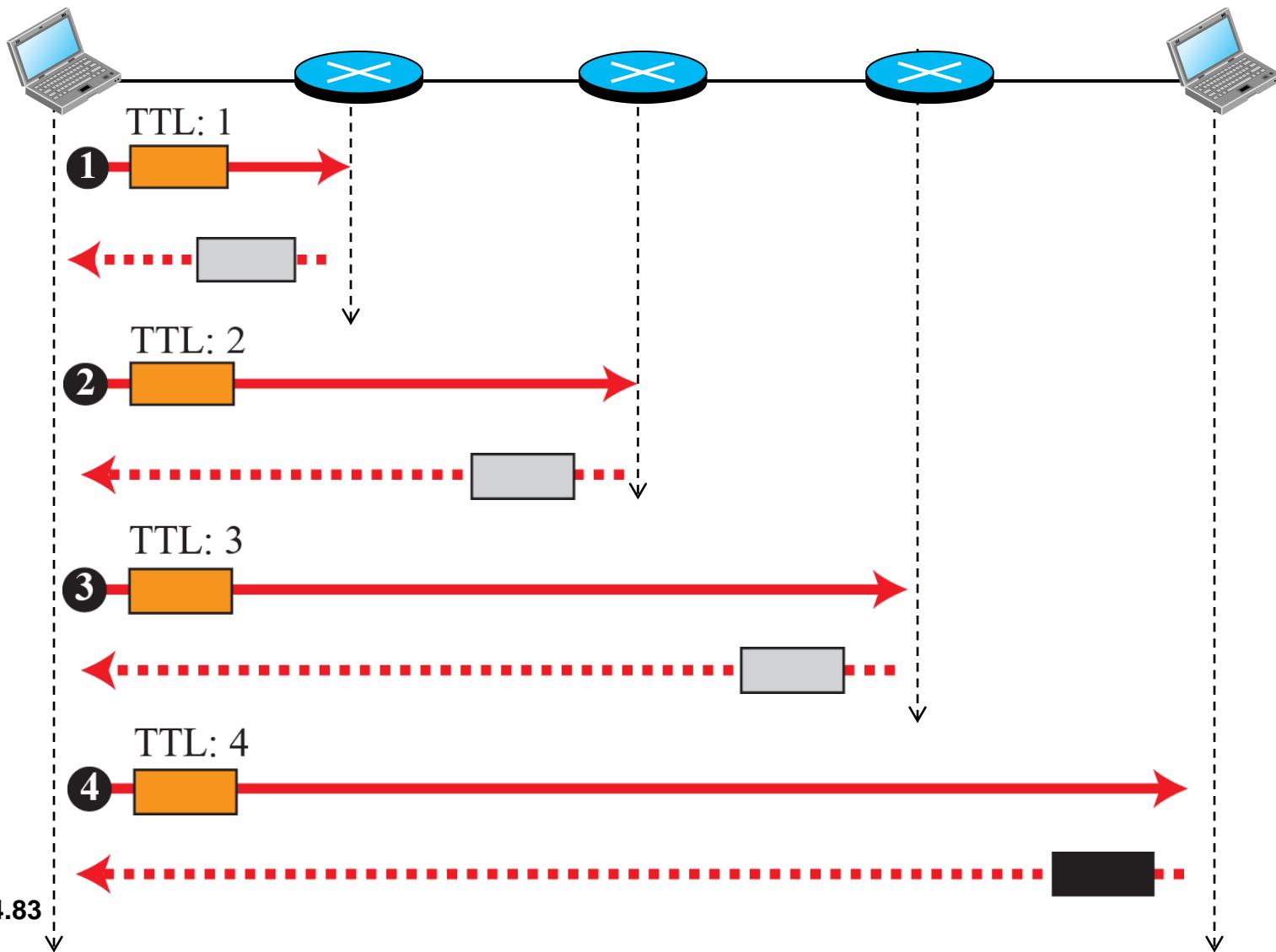
- ❖ source sends series of UDP segments to dest
    - first set has TTL =1
    - second set has TTL=2, etc.
    - unlikely port number
  - ❖ when  $n$ th set of datagrams arrives to  $n$ th router:
    - router discards datagrams
    - and sends source ICMP messages (type 11, code 0)
    - ICMP messages 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



**Figure 4.55: Example of traceroute program**

**Message types**

- Traceroute
- Time-exceeded
- Destination-unreachable



# 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

## *IPv6 datagram format:*

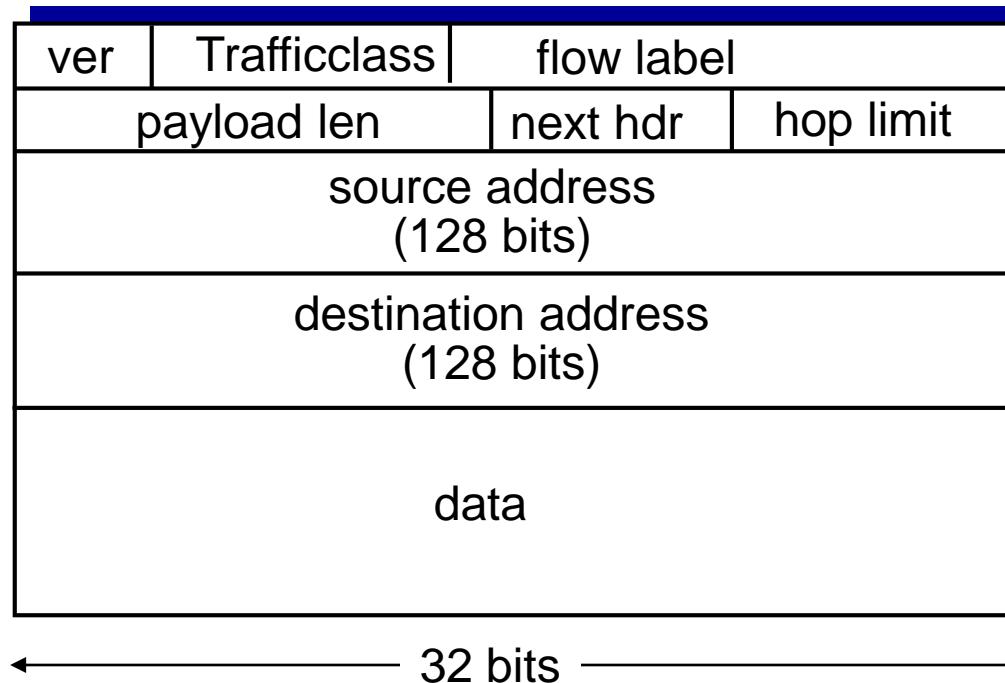
- fixed-length 40 byte header
- no fragmentation allowed

# IPv6 datagram format

**priority:** identify priority among datagrams in flow

**flow Label:** identify datagrams in same “flow.”

**next header:** identify upper layer protocol for data

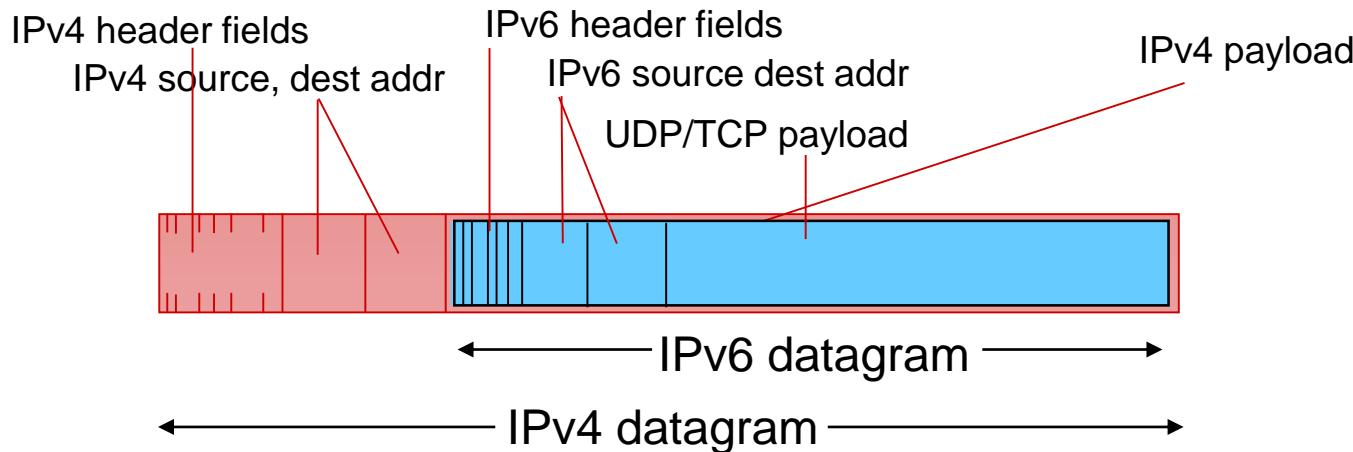


# Other changes from IPv4

- ❖ **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

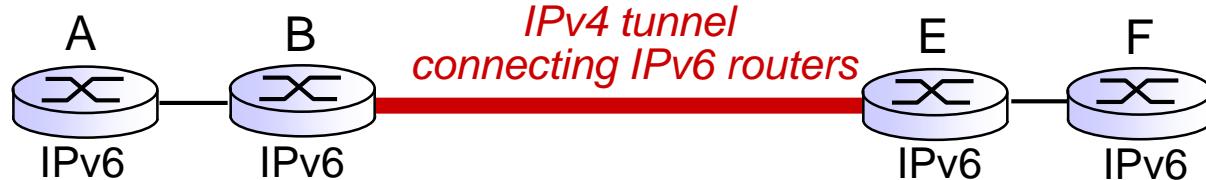
# Transition from IPv4 to IPv6

- ❖ 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

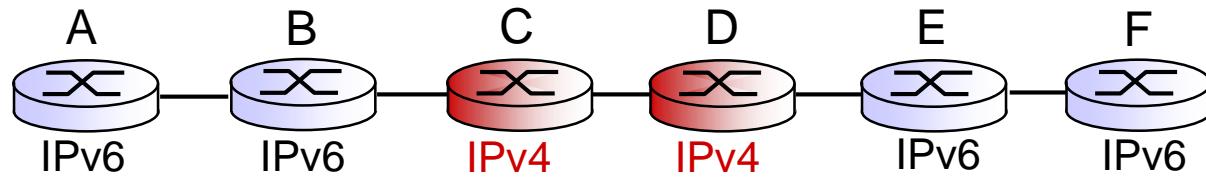


# Tunneling

logical view:

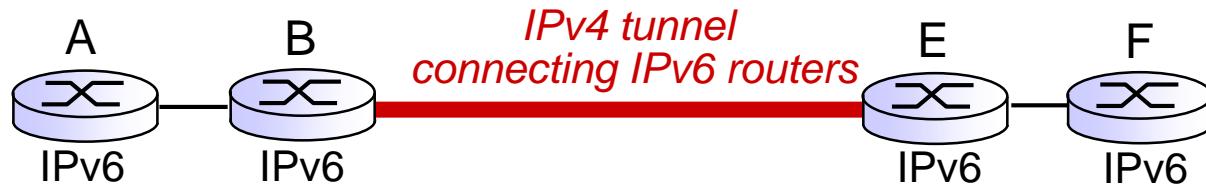


physical view:

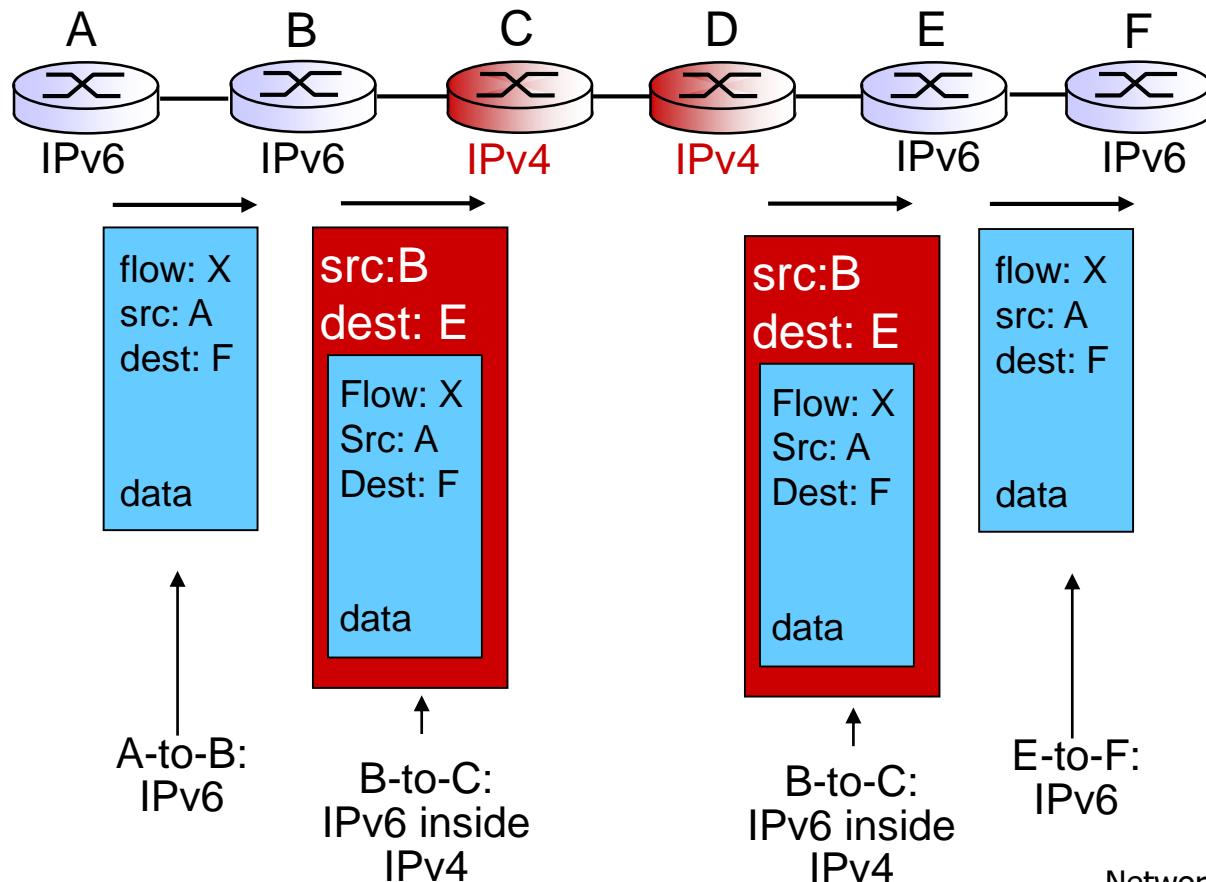


# Tunneling

logical view:



physical view:



# Chapter 4: outline

4.1 introduction

4.2 virtual circuit and  
datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

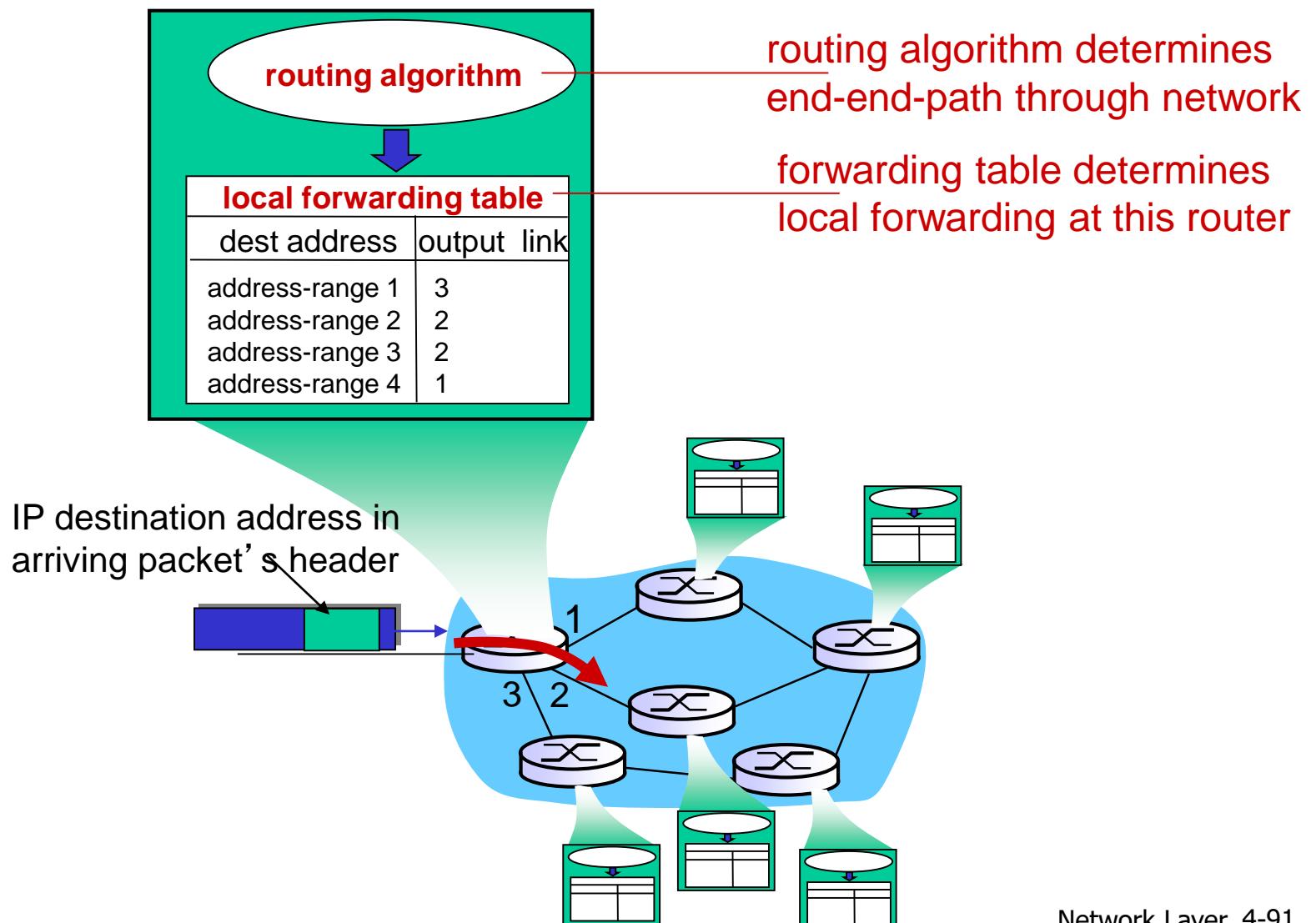
- link state
- distance vector
- hierarchical routing

4.6 routing in the  
Internet

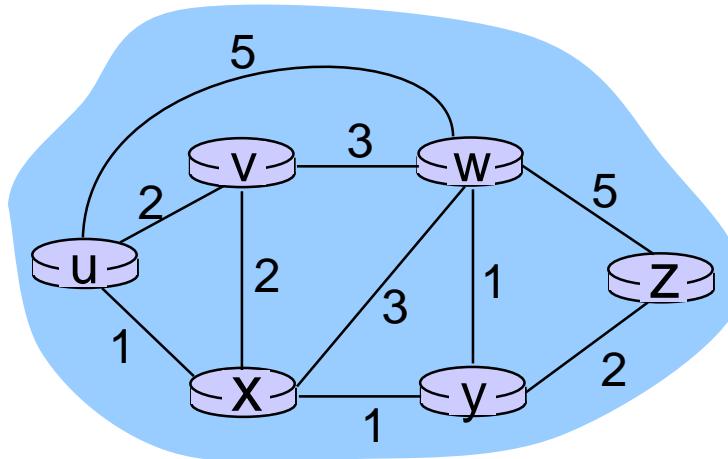
- RIP
- OSPF
- BGP

4.7 broadcast and  
multicast routing

# Interplay between routing, forwarding



# Graph abstraction



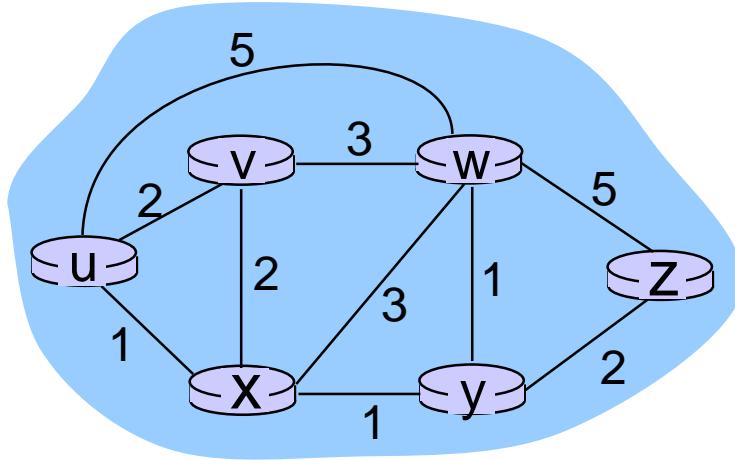
graph:  $G = (N, E)$

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

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

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

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

**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

4.3 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  $N'$**

# Chapter 4: outline

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)

let

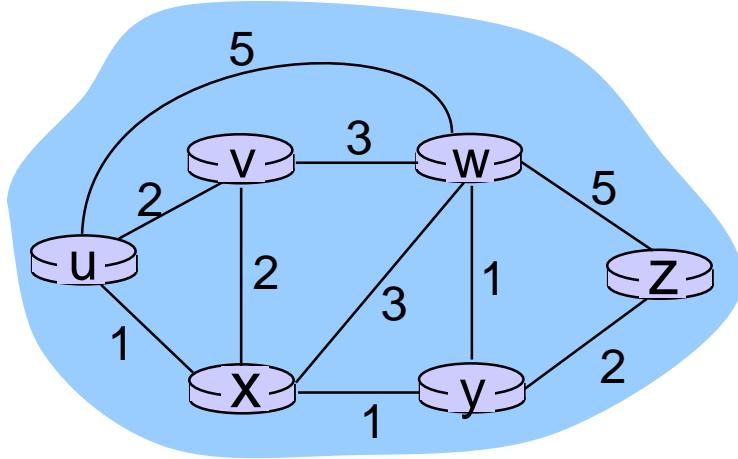
$d_x(y) :=$  cost of least-cost path from  $x$  to  $y$

then

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

cost from neighbor  $v$  to destination  
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  $D_x = [D_x(y): y \in N]$
- ❖ node  $x$ :
  - knows cost to each neighbor  $v$ :  $c(x,v)$
  - maintains its neighbors' distance vectors. For each neighbor  $v$ ,  $x$  maintains  $D_v = [D_v(y): y \in N]$

# 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

## Distance-Vector (DV) Algorithm

At each node,  $x$ :

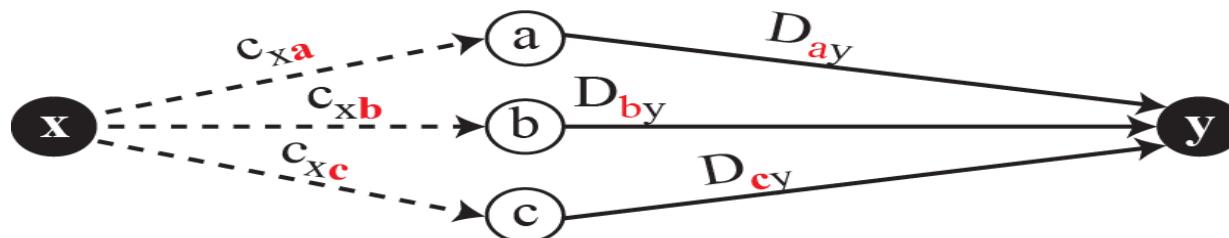
- 1 Initialization:
- 2     for all destinations  $y$  in  $N$ :  
3          $D_x(y) = c(x,y)$  /\* if  $y$  is not a neighbor then  $c(x,y) = \infty$  \*/
- 4     for each neighbor  $w$   
5          $D_w(y) = ?$  for all destinations  $y$  in  $N$
- 6     for each neighbor  $w$   
7         send distance vector  $D_x = [D_x(y): y \text{ in } N]$  to  $w$
- 8

```
9   loop
10    wait (until I see a link cost change to some neighbor w or
11        until I receive a distance vector from some neighbor w)
12
13    for each y in N:
14       $D_x(y) = \min_v \{c(x,v) + D_v(y)\}$ 
15
16    if  $D_x(y)$  changed for any destination y
17      send distance vector  $D_x = [D_x(y); y \in N]$  to all neighbors
18
19 forever
```

# Summary

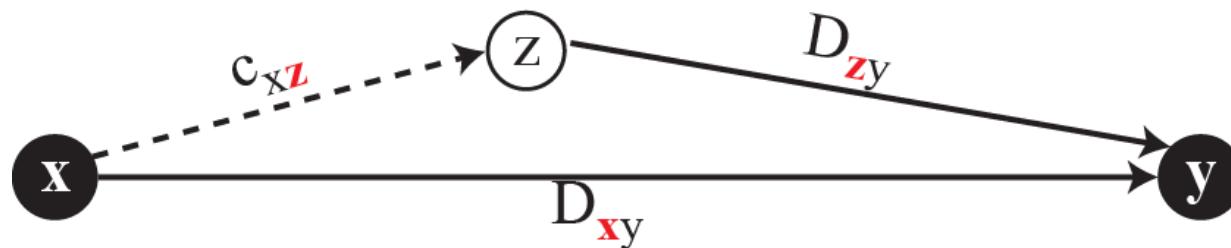
- ❖ Each node creates its own least cost tree with the help of its immediate neighbours.
- ❖ Incomplete tree exchanged with immediate neighbours to make tree more and more complete.
- ❖ i.e router continuously tell about what it knows about internet.
- ❖ Uses Bellman -Ford equation & distance vector.

**Figure 4.58: Graphical idea behind Bellman-Ford equation**



a. General case with three intermediate nodes

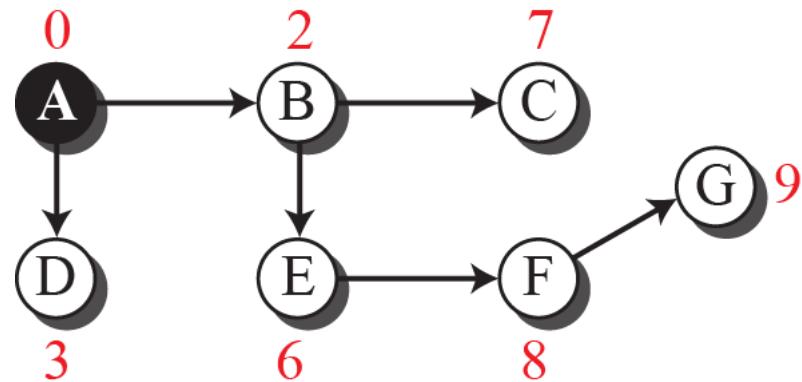
$$D_{xy} = \min\{ (C_{xa} + D_{ay}), (C_{xb} + D_{by}), (C_{xc} + D_{cy}) \dots \}$$



b. Updating a path with a new route

$$D_{xy} = \min\{ D_{xy}, (C_{xz} + D_{zy}) \}$$

**Figure 4.59: The distance vector corresponding to a tree**



a. Tree for node A

A	
A	0
B	2
C	7
D	3
E	6
F	8
G	9

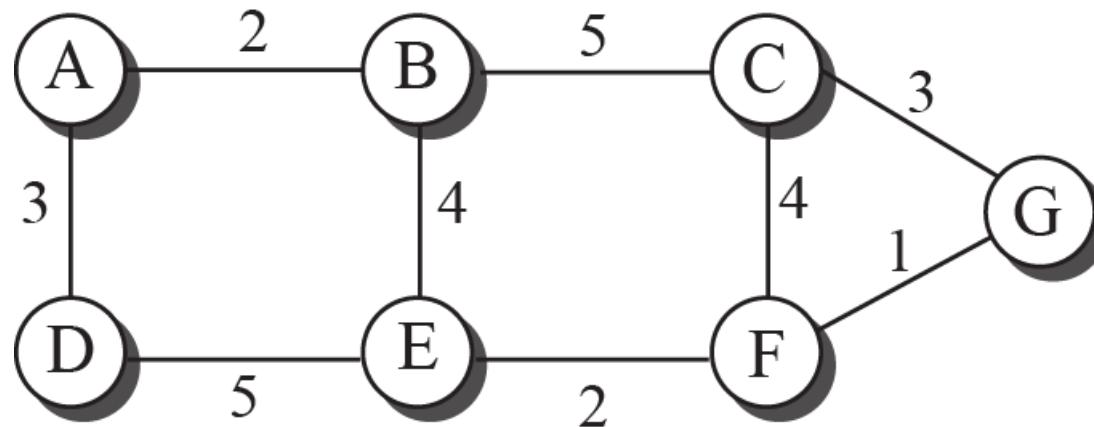
b. Distance vector for node A

**Figure 4.60:** The first distance vector for an internet

A	0
B	2
C	$\infty$
D	3
E	$\infty$
F	$\infty$
G	$\infty$

A	2
B	0
C	5
D	$\infty$
E	4
F	$\infty$
G	$\infty$

A	$\infty$
B	5
C	0
D	$\infty$
E	$\infty$
F	4
G	3



A	$\infty$
B	$\infty$
C	3
D	$\infty$
E	$\infty$
F	1
G	0

A	3
B	$\infty$
C	$\infty$
D	0
E	5
F	$\infty$
G	$\infty$

A	$\infty$
B	4
C	$\infty$
D	5
E	0
F	2
G	$\infty$

A	$\infty$
B	$\infty$
C	4
D	$\infty$
E	2
F	0
G	1

4.109

1.109

## Figure 4.61: Updating distance vectors

New B	Old B	A
A 2	A 2	A 0
B 0	B 0	B 2
C 5	C 5	C $\infty$
D 5	D $\infty$	D 3
E 4	E 4	E $\infty$
F $\infty$	F $\infty$	F $\infty$
G $\infty$	G $\infty$	G $\infty$

$B[ ] = \min(B[ ], 2 + A[ ])$

Note:  
 $X[ ]$ : the whole vector

a. First event: B receives a copy of A's vector.

New B	Old B	E
A 2	A 2	A $\infty$
B 0	B 0	B 4
C 5	C 5	C $\infty$
D 5	D 5	D 5
E 4	E 4	E 0
F 6	F $\infty$	F 2
G $\infty$	G $\infty$	G $\infty$

$B[ ] = \min(B[ ], 4 + E[ ])$

b. Second event: B receives a copy of E's vector.

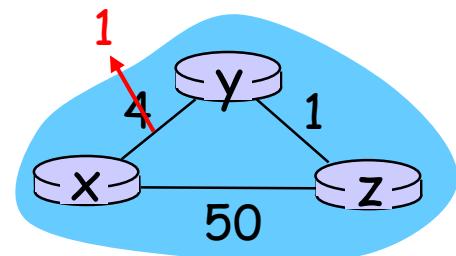
# Distance vector: link cost changes

## 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

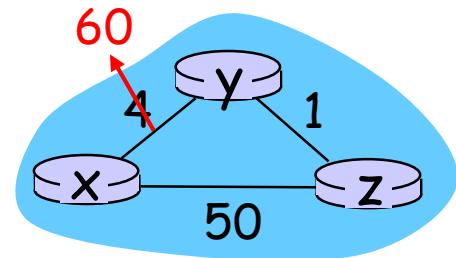
## link cost changes:

- ❖ node detects local link cost change
- ❖ **bad news travels slow** - “count to infinity” problem!

- ❖ 44 iterations before algorithm stabilizes:

## 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
  - ① error propagate thru network

# Chapter 4: outline

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

# Hierarchical routing

our routing study thus far - idealization

- ❖ all routers identical
- ❖ network “flat”
- ... *not true in practice*

*scale:* with 600 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

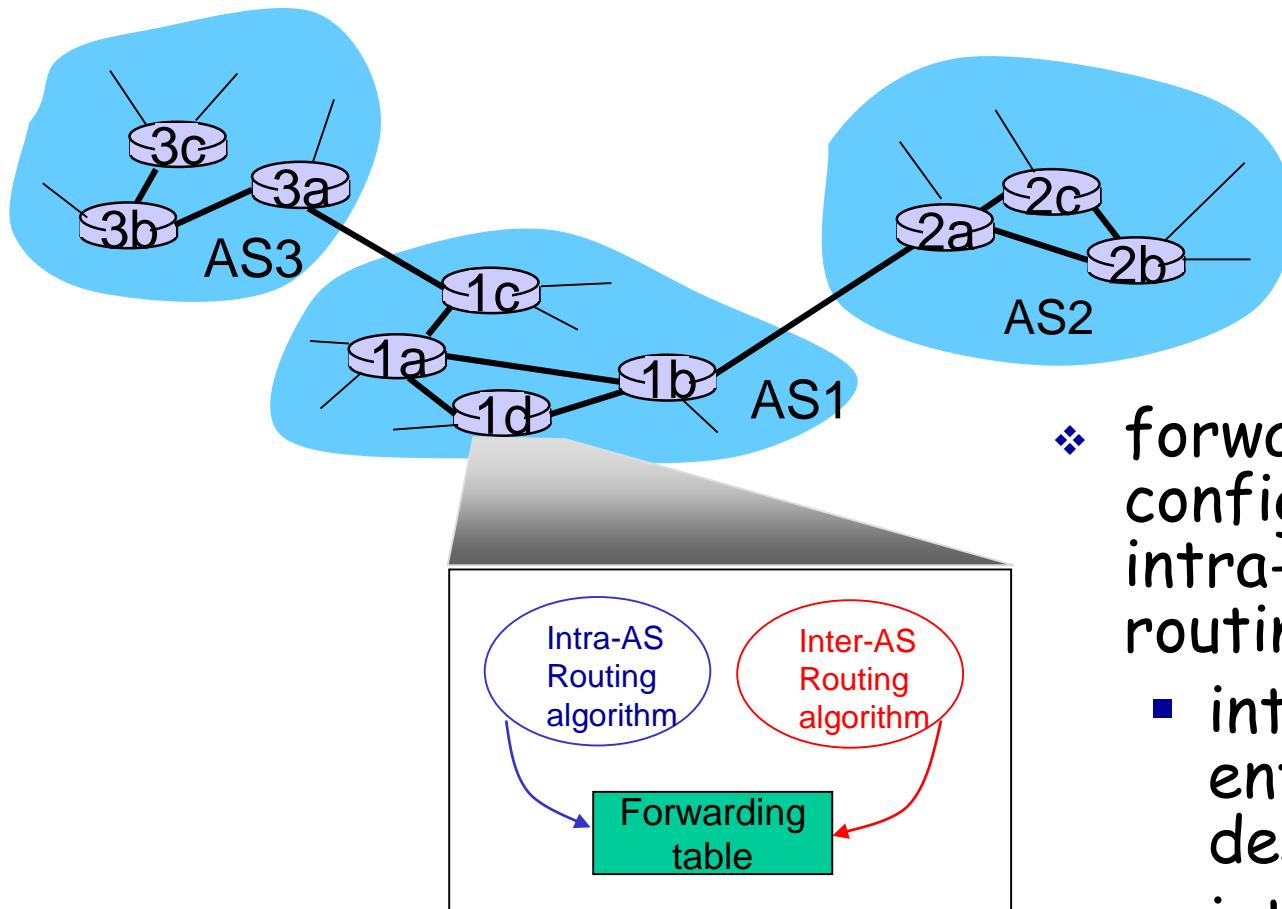
# Hierarchical routing

- ❖ 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:*

- ❖ at “edge” of its own AS
- ❖ has link to router in another AS

# Interconnected ASes



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

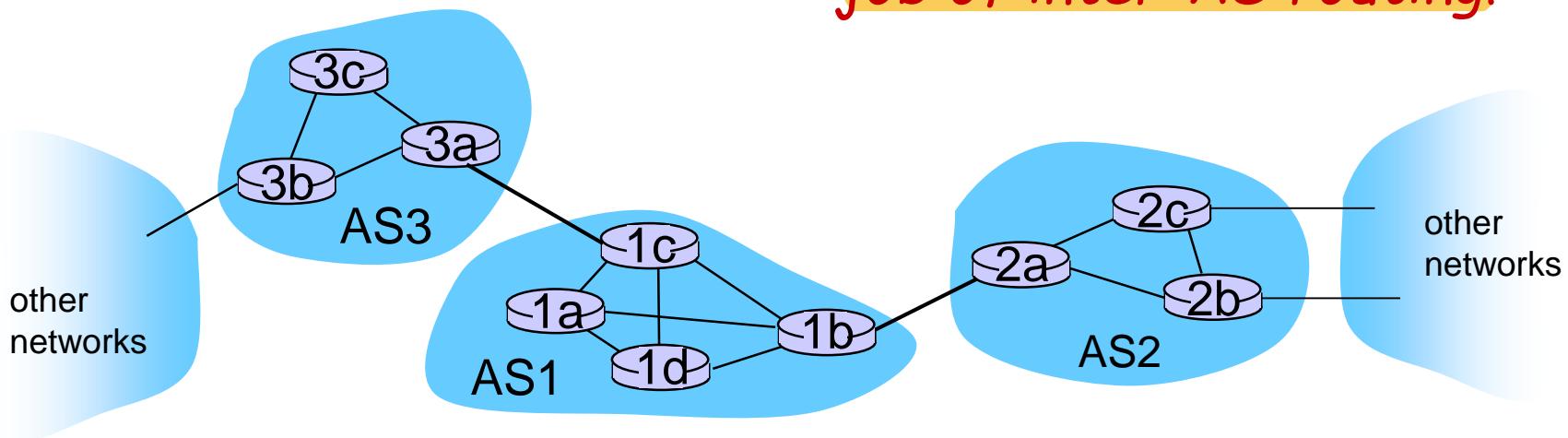
# Inter-AS tasks

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

*AS1 must:*

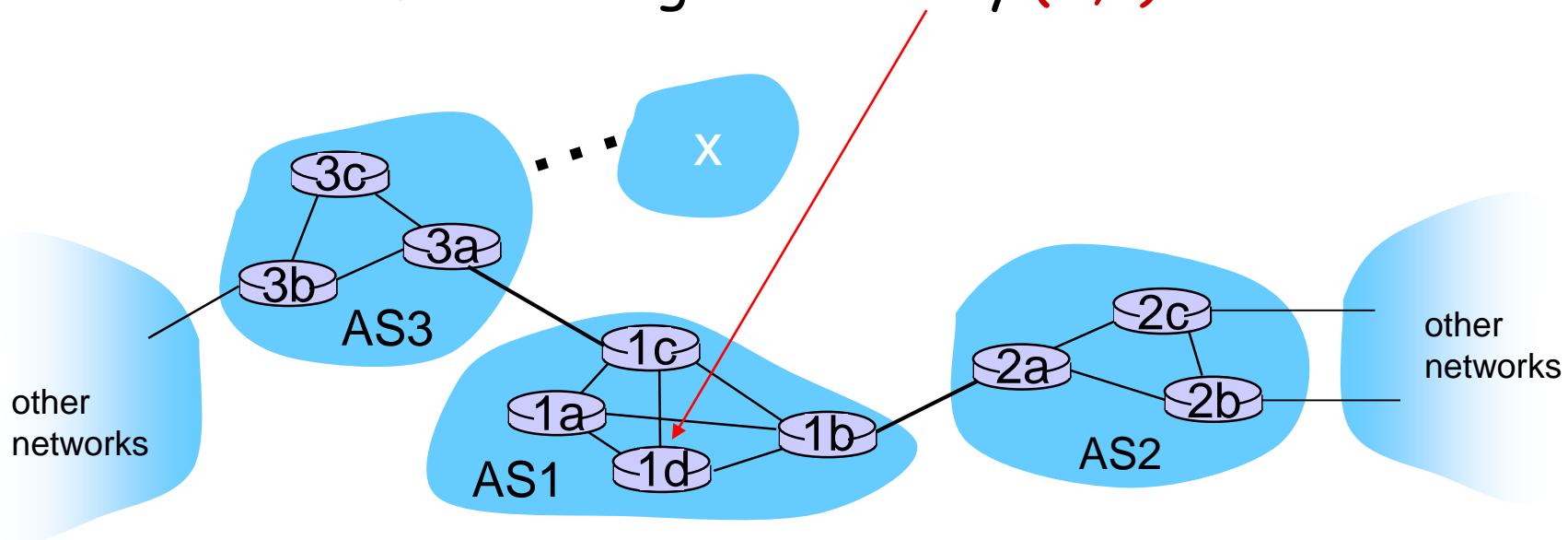
1. learn which dests are reachable through AS2, which through AS3
2. propagate this reachability info to all routers in AS1

*job of inter-AS routing!*



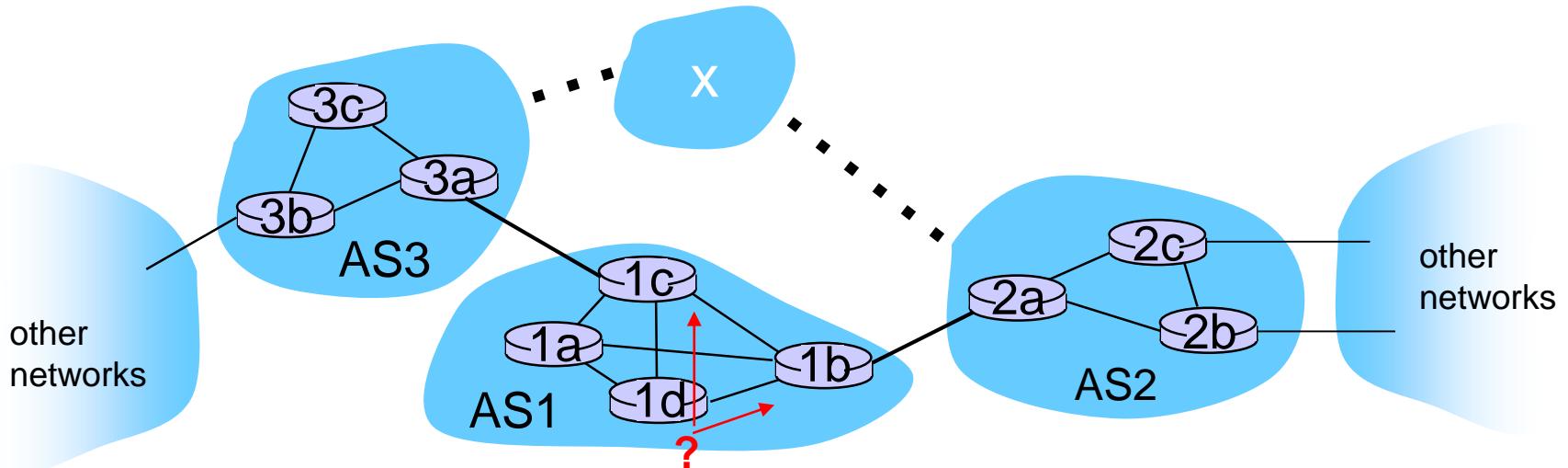
## Example: setting forwarding table in router 1d

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



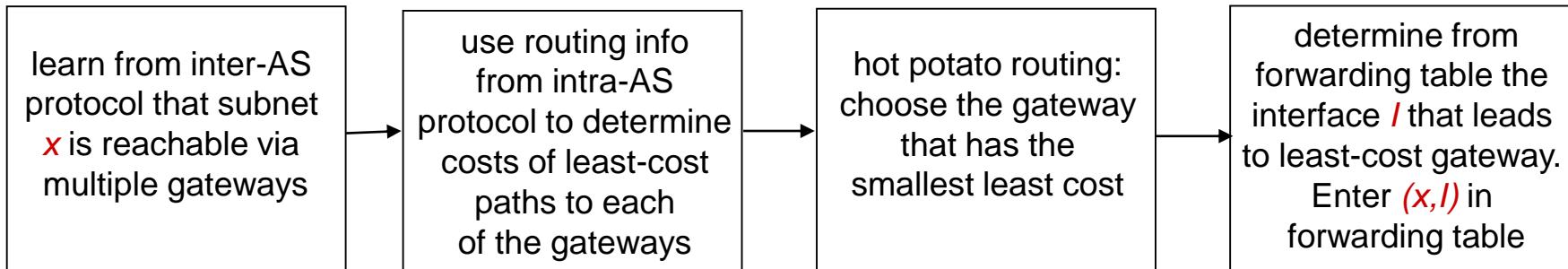
# Example: choosing among multiple ASes

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



# Example: choosing among multiple ASes

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



# Chapter 4: outline

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

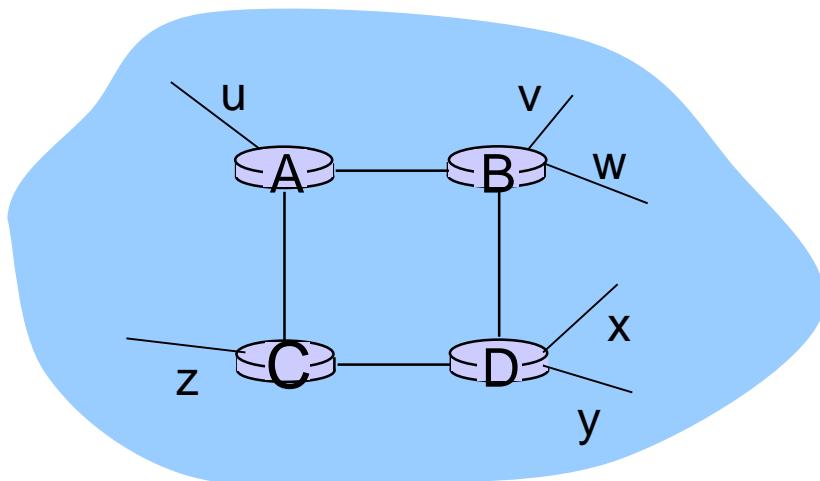
4.7 broadcast and  
multicast routing

# 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)

# RIP ( Routing Information Protocol)

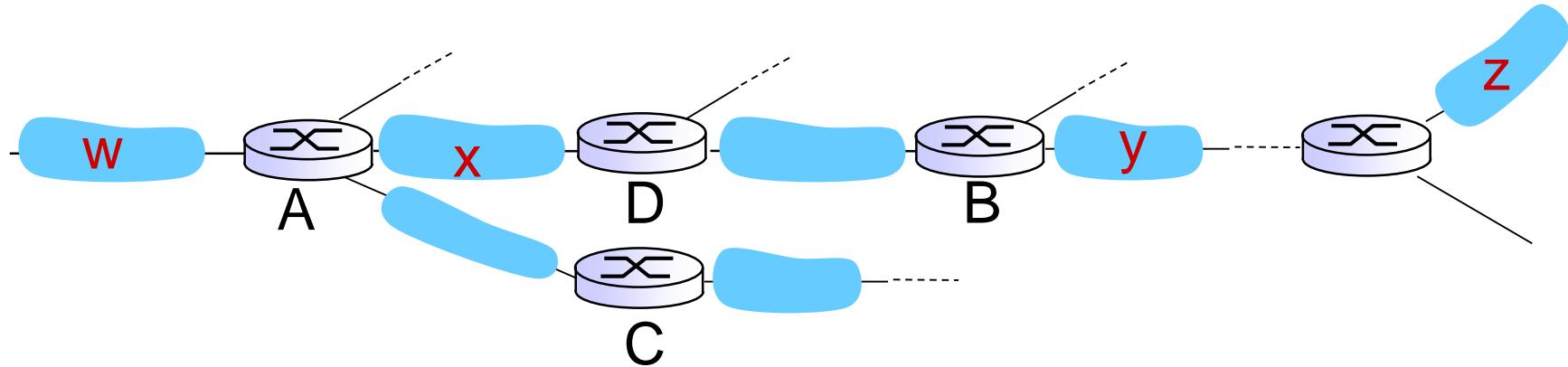
- ❖ included in BSD-UNIX distribution in 1982
- ❖ distance vector algorithm
  - distance metric: # hops (max = 15 hops), each link has cost 1
  - DVs exchanged with neighbors every 30 sec in response message (aka **advertisement**)
  - each advertisement: list of up to 25 destination **subnets** (in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

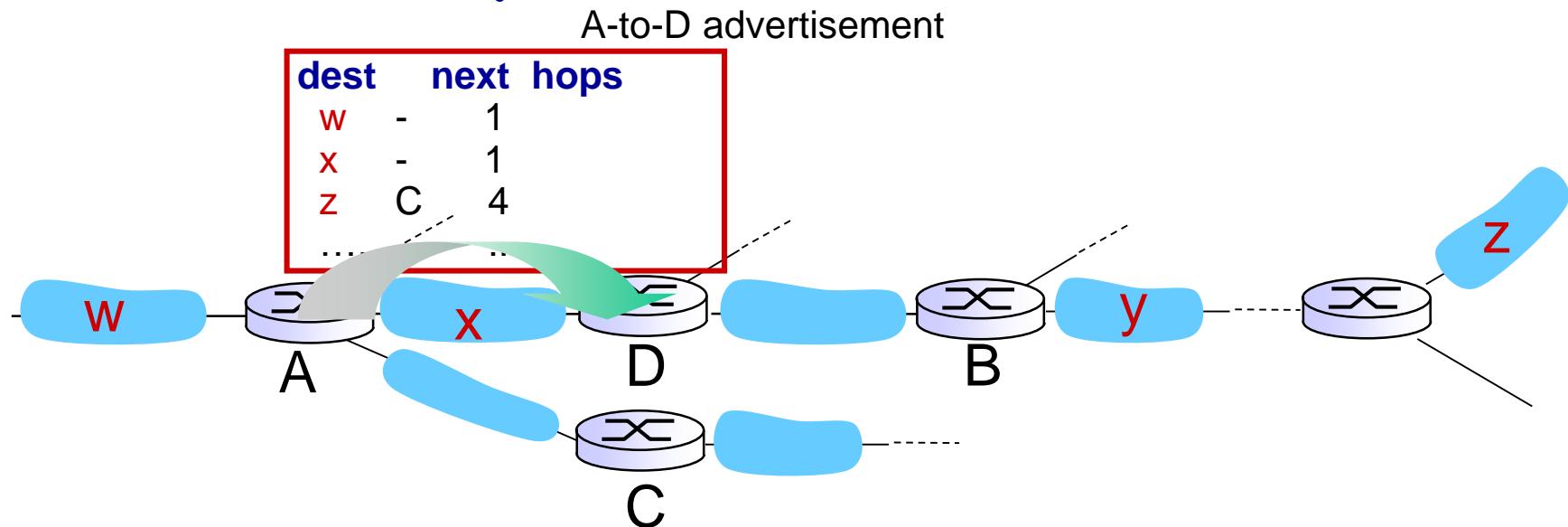
# RIP: example



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



routing table in router D

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

A  
5

# RIP: link failure, recovery

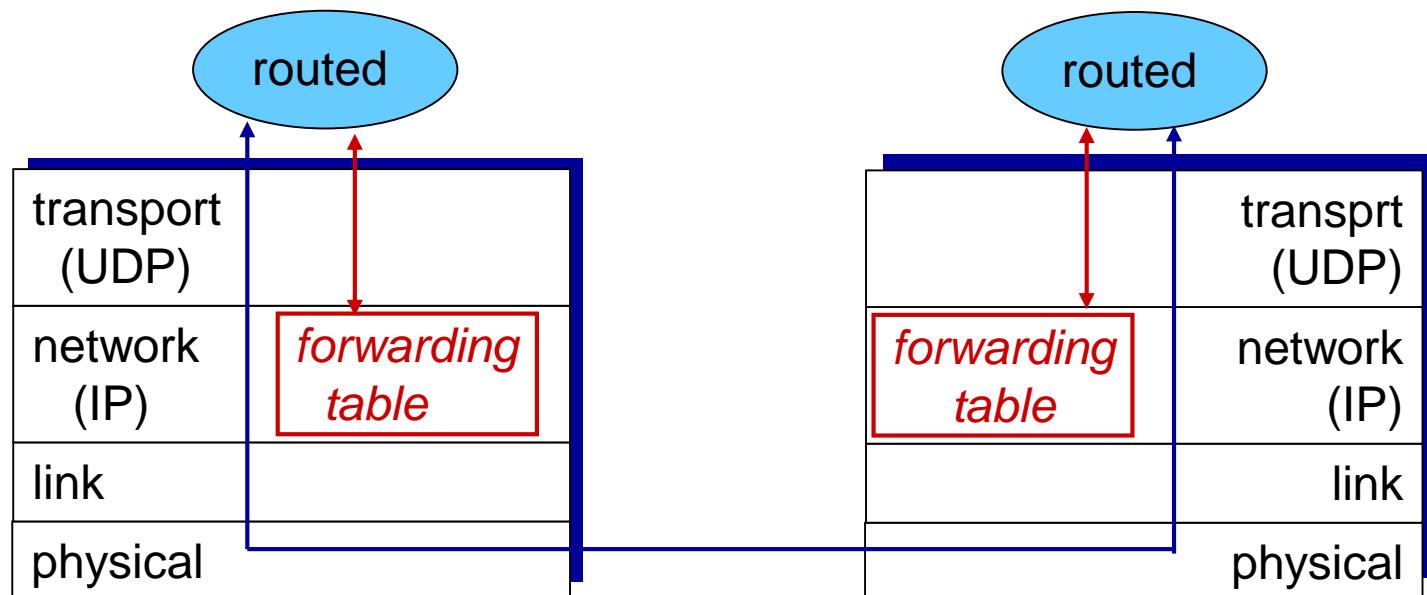
RIP-Request, Response msg.

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

# RIP table processing

- ❖ RIP routing tables managed by *application-level* process called route-d (daemon)
- ❖ advertisements sent in UDP packets port no 520, periodically repeated



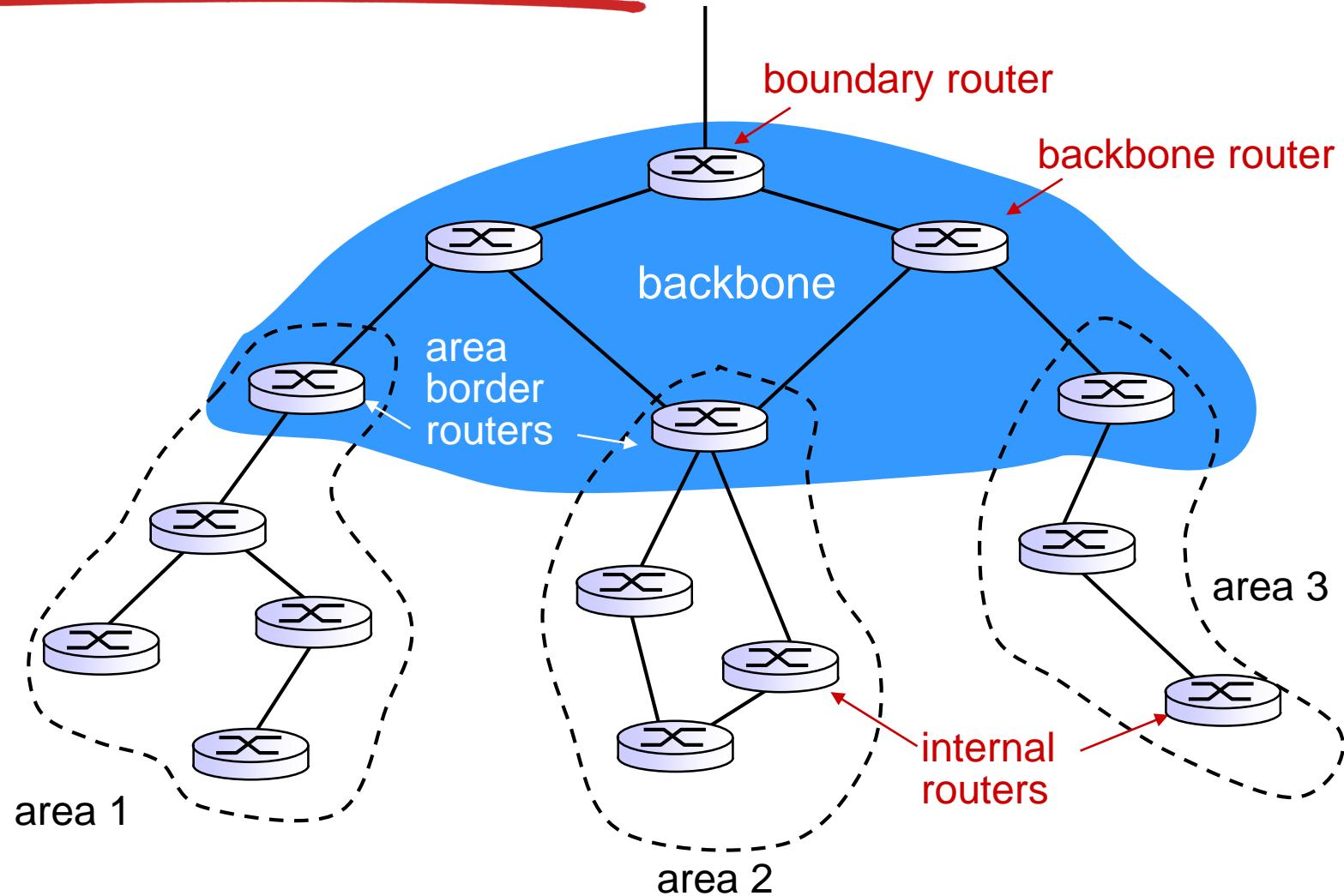
# OSPF (Open Shortest Path First)

- ❖ “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
- ❖ **RIP**- lower tier ISPs; OSPF-upper tier ISP

## 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

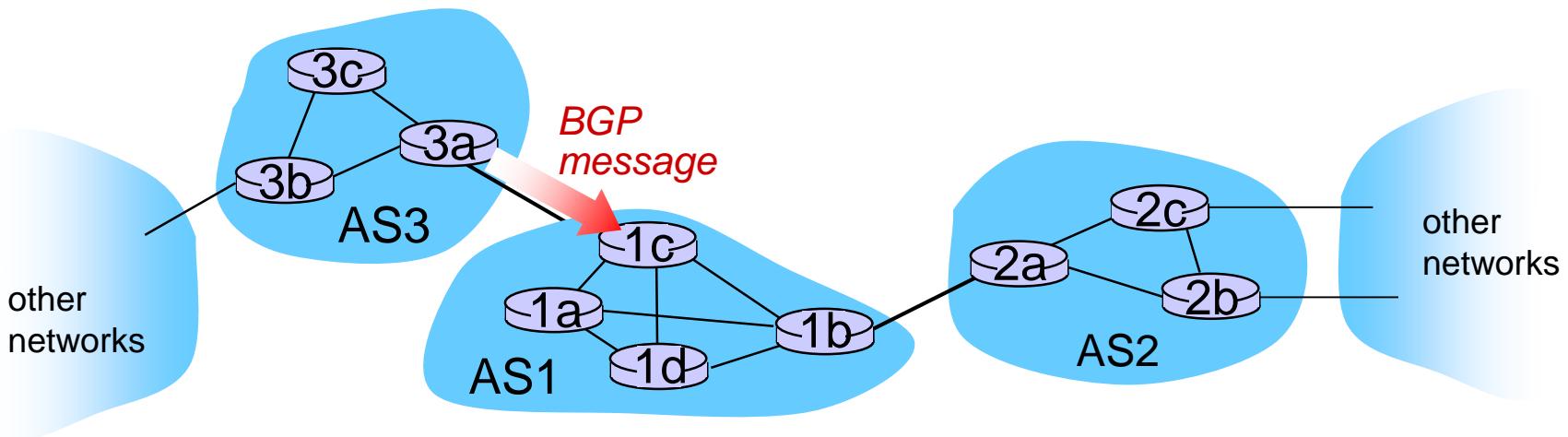
- ❖ **two-level hierarchy:** local area, backbone.
  - link-state advertisements only in area
  - each node has detailed area topology; only know direction (shortest path) to nets in other areas.
- ❖ **area border routers:** “summarize” distances to nets in own area, advertise to other Area Border routers.
- ❖ **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 standard 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**”

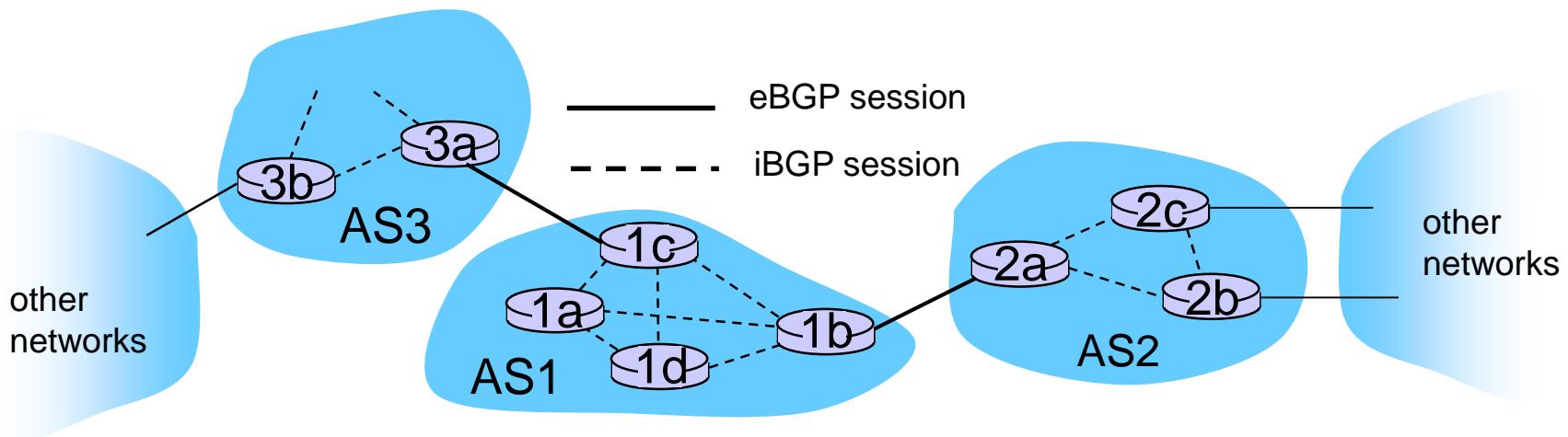
# BGP basics

- ❖ **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 do 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

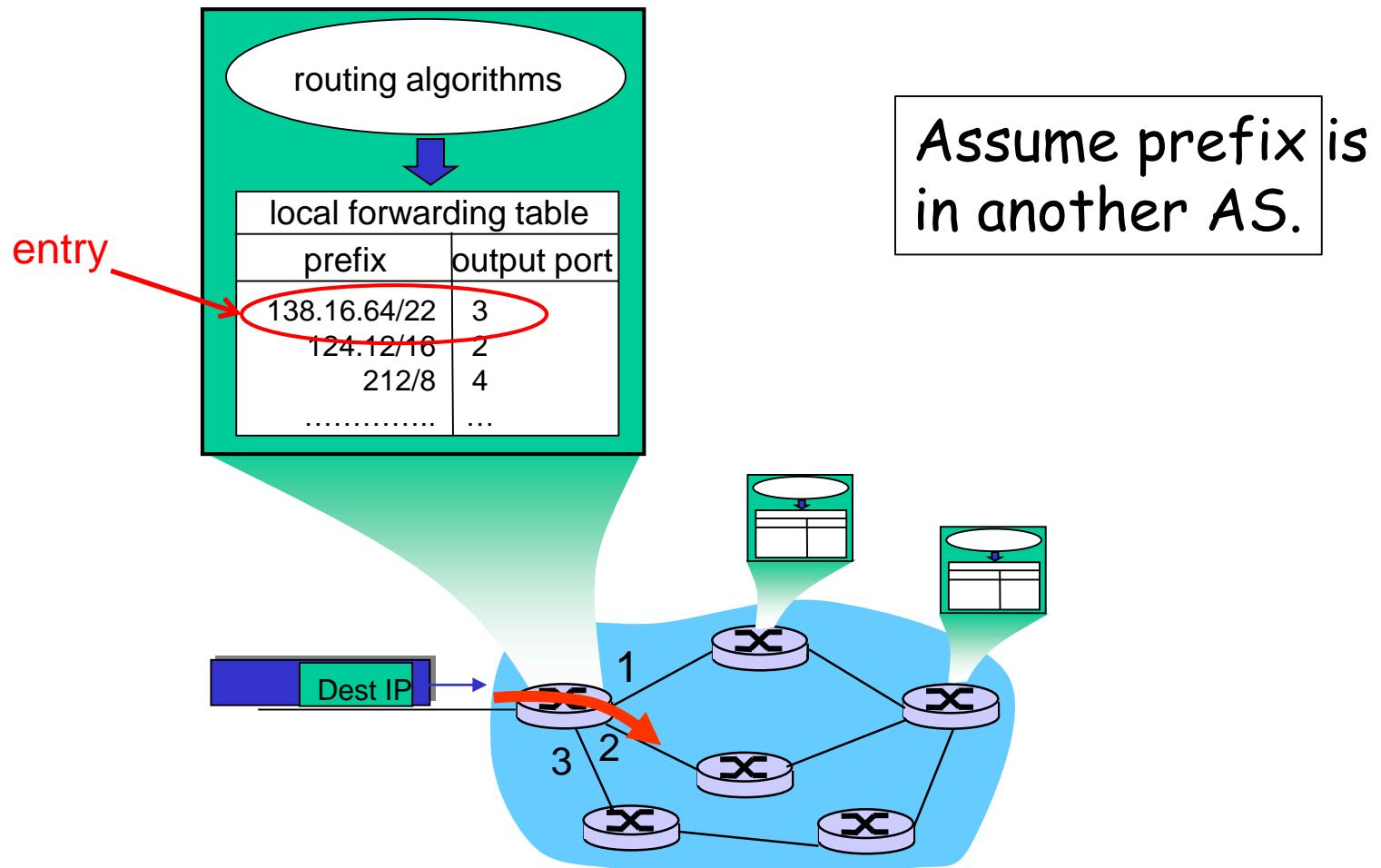
# BGP messages

- ❖ 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

# *Putting it Altogether: How Does an Entry Get Into a Router's Forwarding Table?*

- ❖ Answer is complicated!
- ❖ Ties together hierarchical routing (Section 4.5.3) with BGP (4.6.3) and OSPF (4.6.2).
- ❖ Provides nice overview of BGP!

# How does entry get in forwarding table?

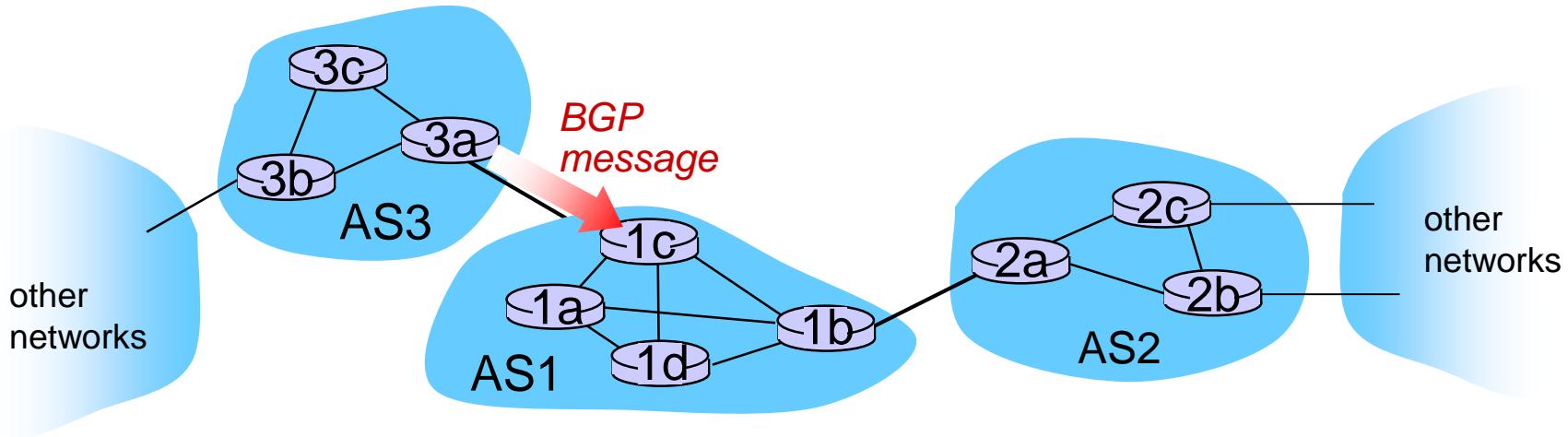


# How does entry get in forwarding table?

## High-level overview

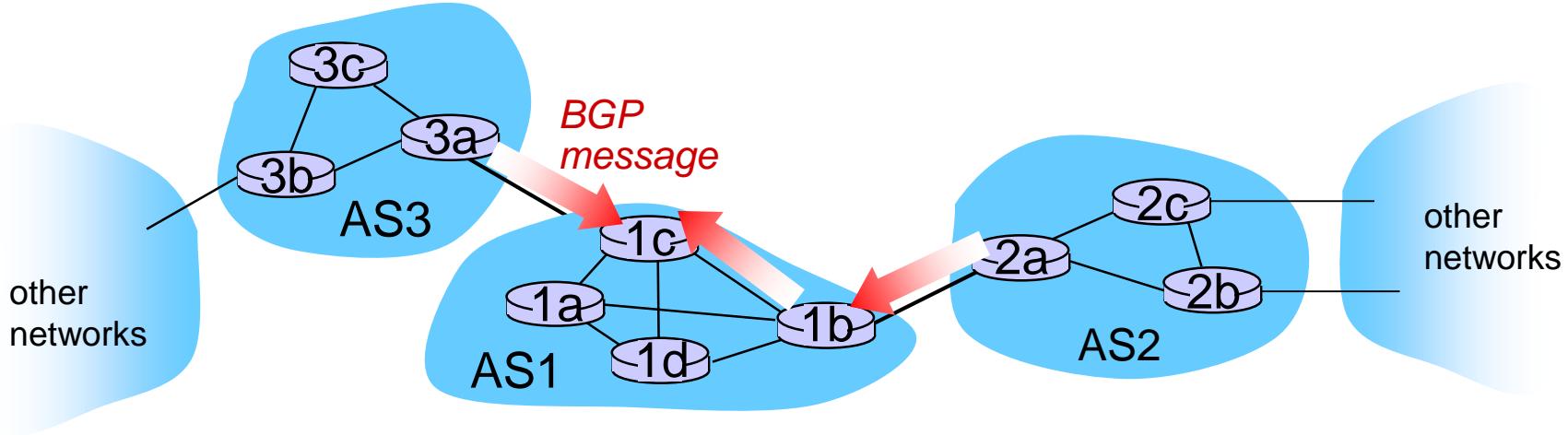
1. Router becomes aware of prefix
2. Router determines output port for prefix
3. Router enters prefix-port in forwarding table

# Router becomes aware of prefix



- ❖ BGP message contains “routes”
- ❖ “route” is a prefix and attributes: AS-PATH, NEXT-HOP,...
- ❖ Example: route:
  - ❖ Prefix:138.16.64/22 ; AS-PATH: AS3 AS131 ; NEXT-HOP: 201.44.13.125

# Router may receive multiple routes

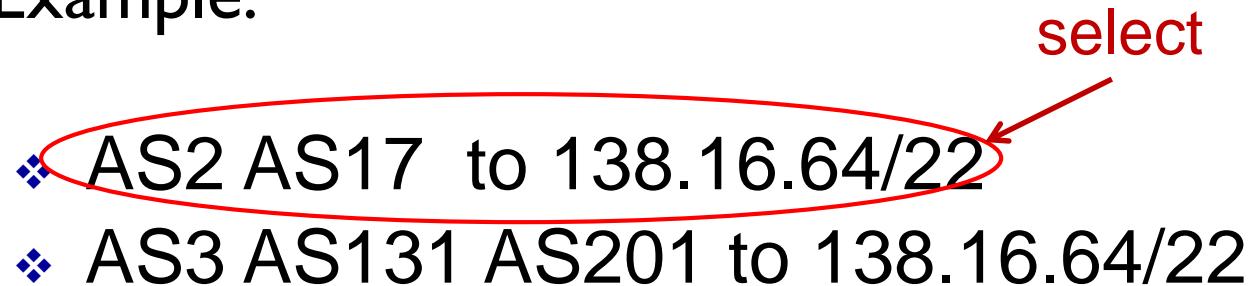


- ❖ Router may receive multiple routes for same prefix
- ❖ Has to select one route

# Select best BGP route to prefix

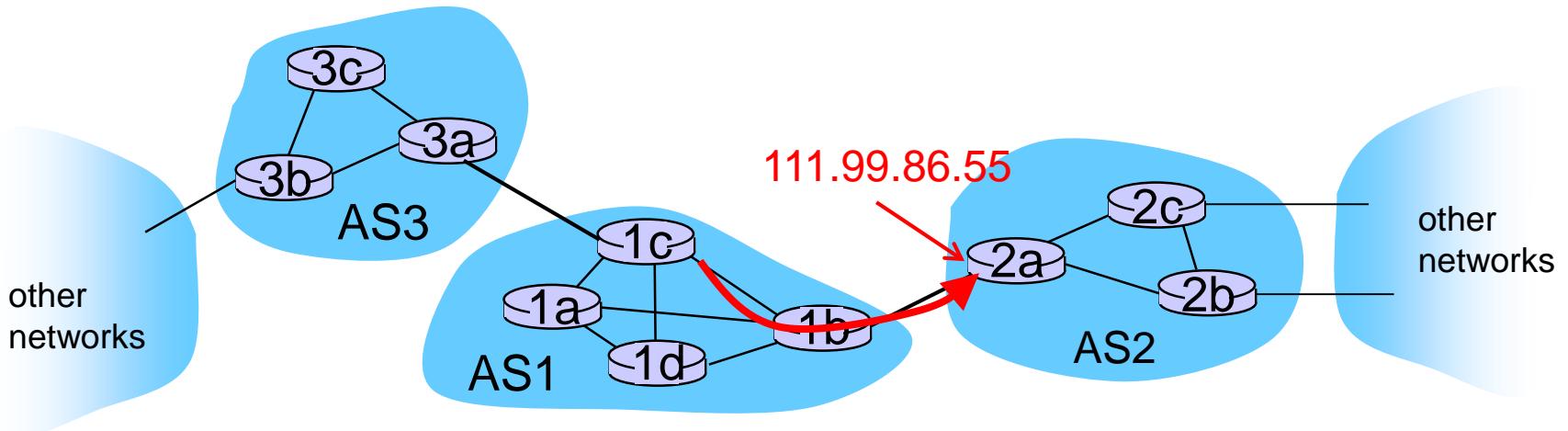
- ❖ Router selects route based on shortest AS-PATH

- ❖ Example:

- ❖ AS2 AS17 to 138.16.64/22
    - ❖ AS3 AS131 AS201 to 138.16.64/22
- 
- select

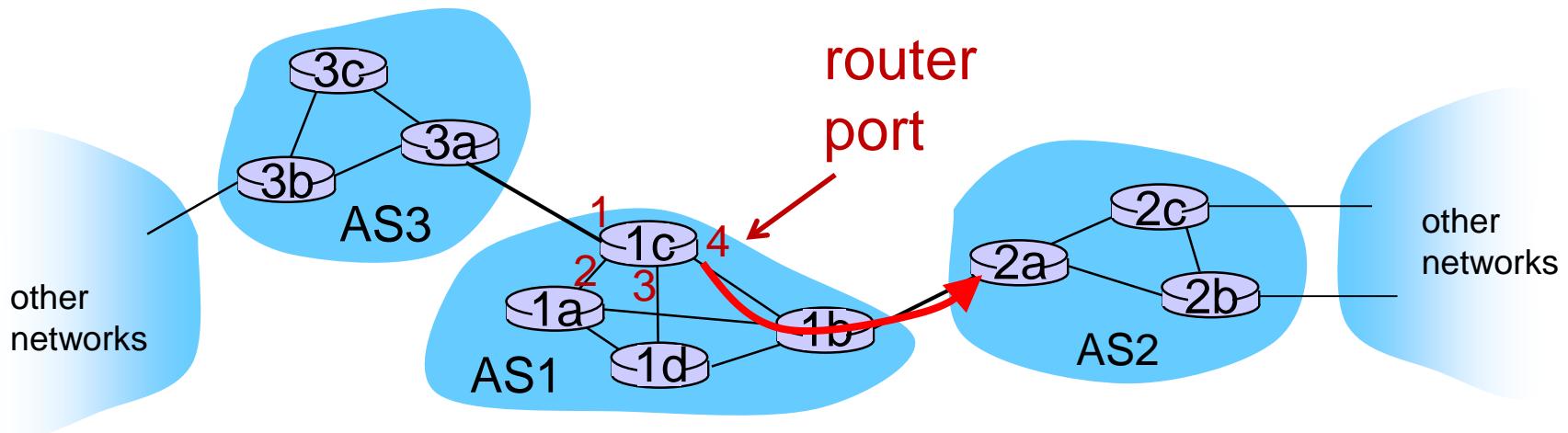
# Find best intra-route to BGP route

- ❖ Use selected route's NEXT-HOP attribute
  - Route's NEXT-HOP attribute is the IP address of the router interface that begins the AS PATH.
- ❖ Example:
  - ❖ AS-PATH: AS2 AS17 ; NEXT-HOP: 111.99.86.55
- ❖ Router uses OSPF to find shortest path from 1c to 111.99.86.55



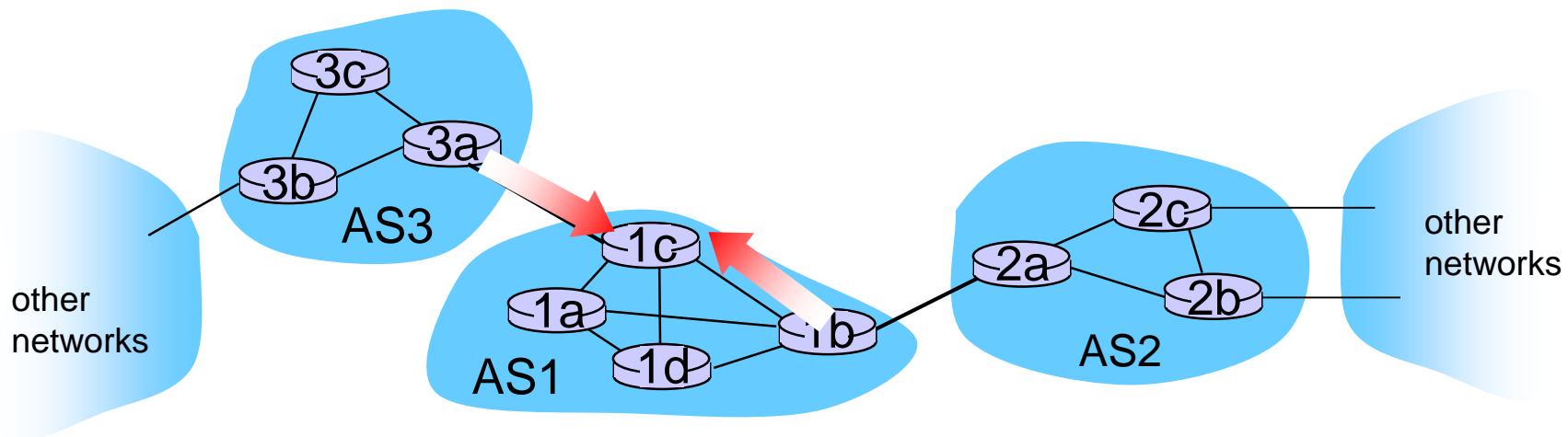
# Router identifies port for route

- ❖ Identifies port along the OSPF shortest path
  - ❖ Adds prefix-port entry to its forwarding table:
    - (138.16.64/22 , port 4)



# Hot Potato Routing

- ❖ Suppose there are two or more best inter-routes.
- ❖ Then choose route with closest NEXT-HOP
  - Use OSPF to determine which gateway is closest
  - Q: From 1c, chose AS3 AS131 or AS2 AS17?
  - A: route AS3 AS201 since it is closer



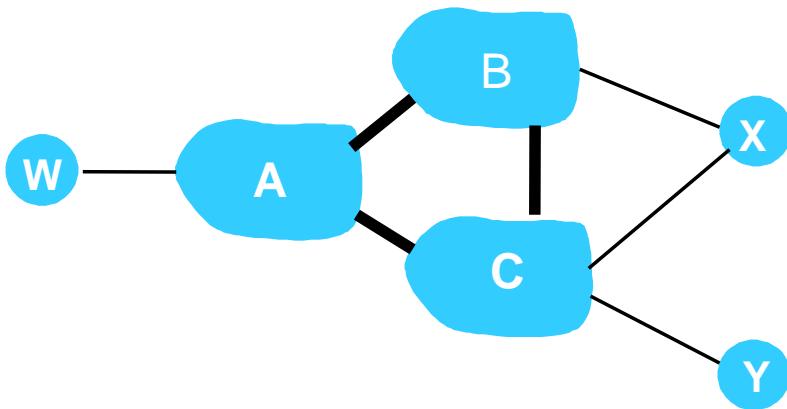
# How does entry get in forwarding ~~table~~?

---

## Summary

1. Router becomes aware of prefix
  - via BGP route advertisements from other routers
2. Determine router output port for prefix
  - Use BGP route selection to find best inter-AS route
  - Use OSPF to find best intra-AS route leading to best inter-AS route
  - Router identifies router port for that best route
3. Enter prefix-port entry in forwarding table

# BGP routing policy

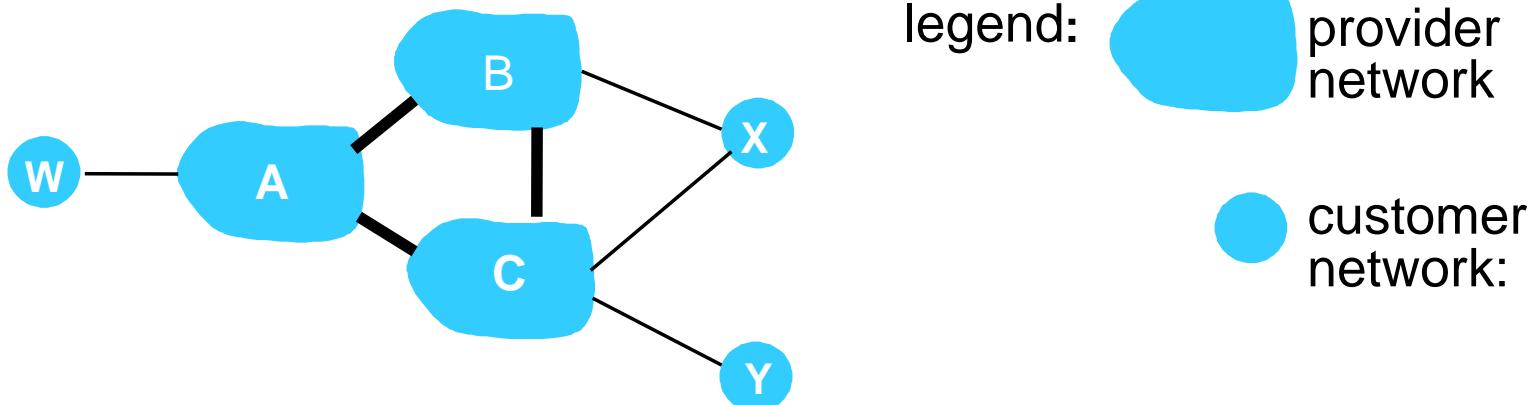


legend:

- provider network
- customer network:

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

## BGP routing policy (2)



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

# Why different Intra-, 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: 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

❖ 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

4.5 routing algorithms

- link state, distance  
vector, hierarchical  
routing

4.6 routing in the  
Internet

- RIP, OSPF, BGP