

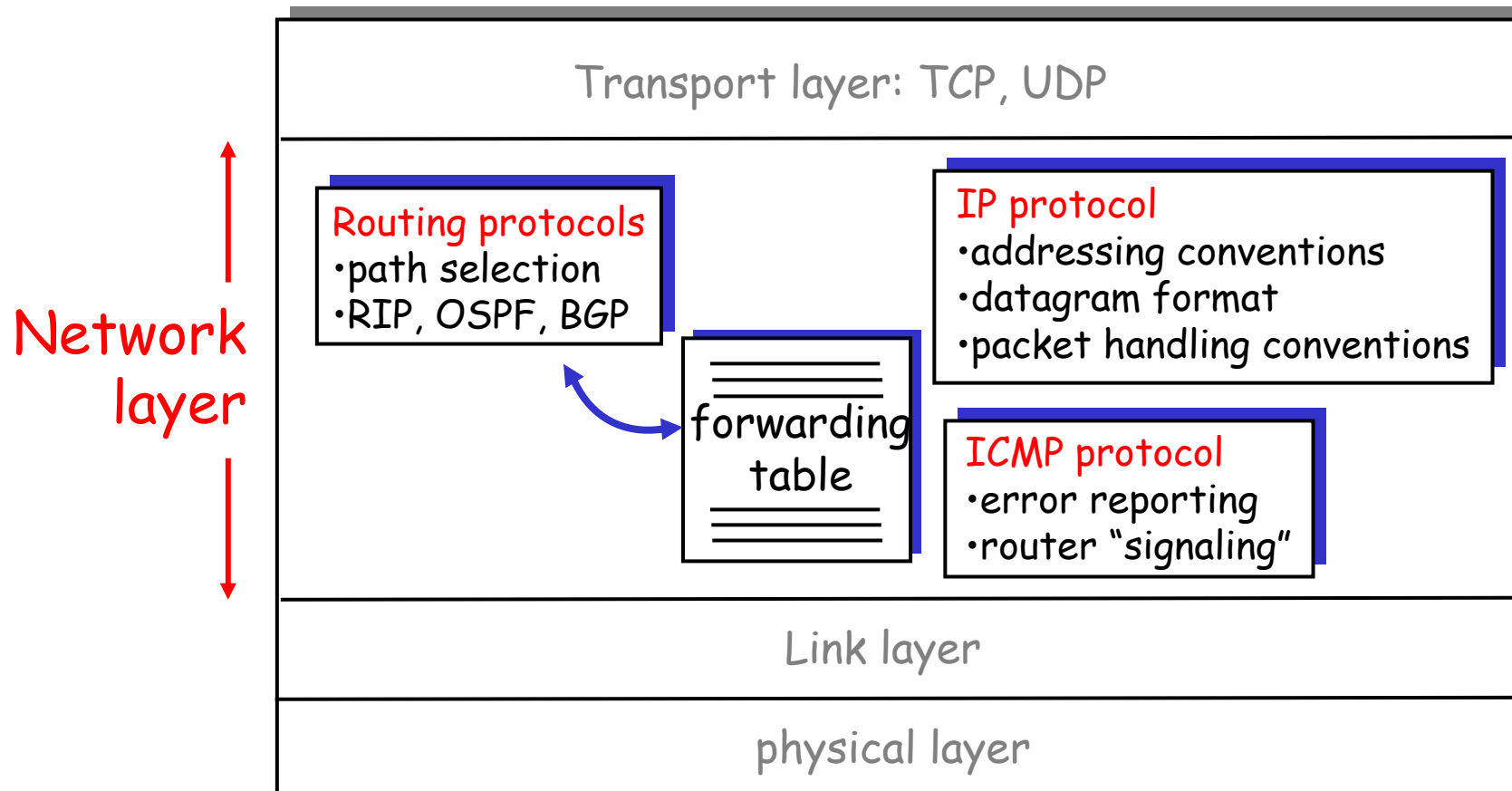
# Chapter 4: Network Layer

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

# The Internet Network layer

## 網際網路之網路層

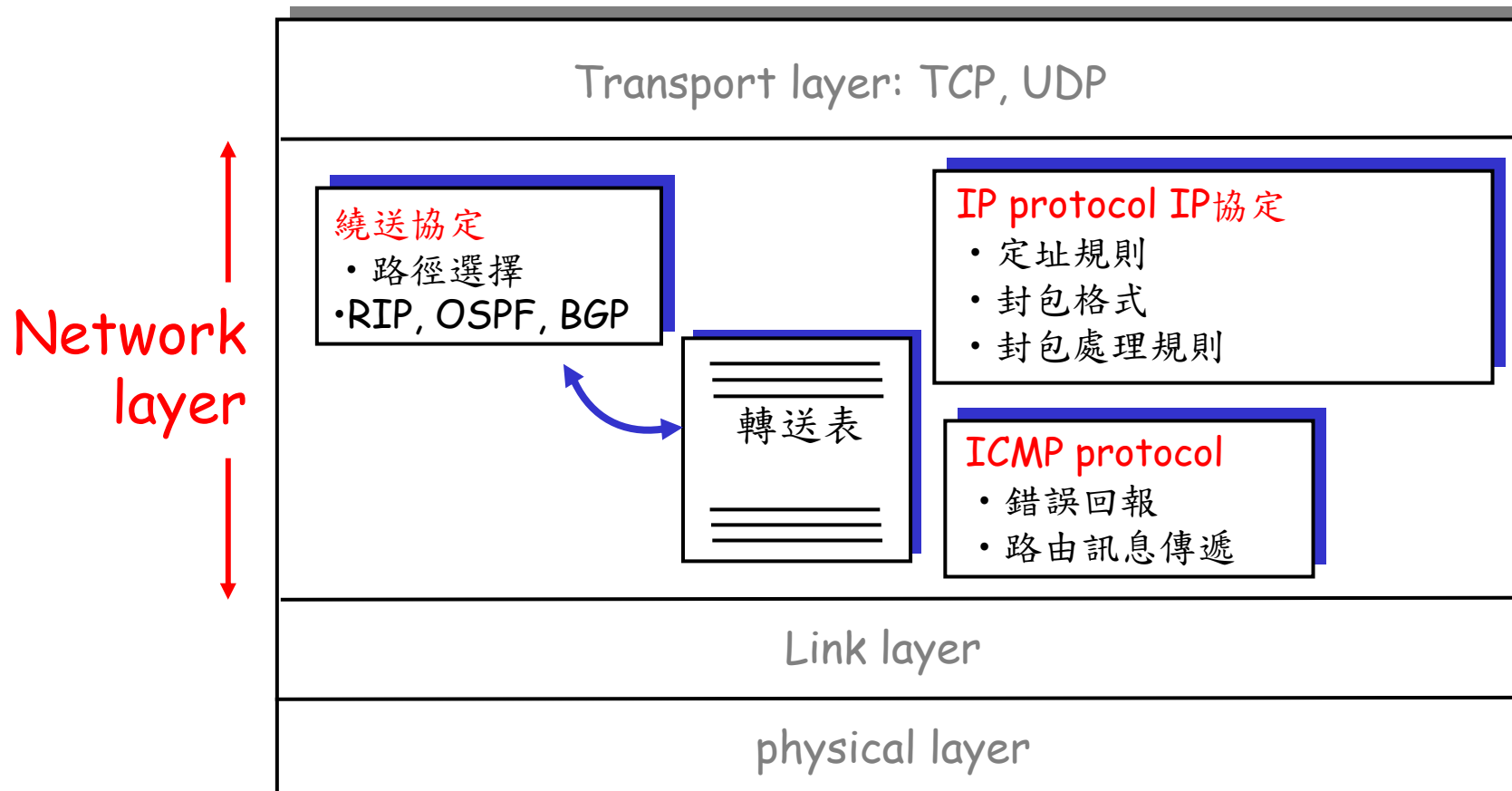
Host, router network layer functions:



# The Internet Network layer

## 網際網路之網路層

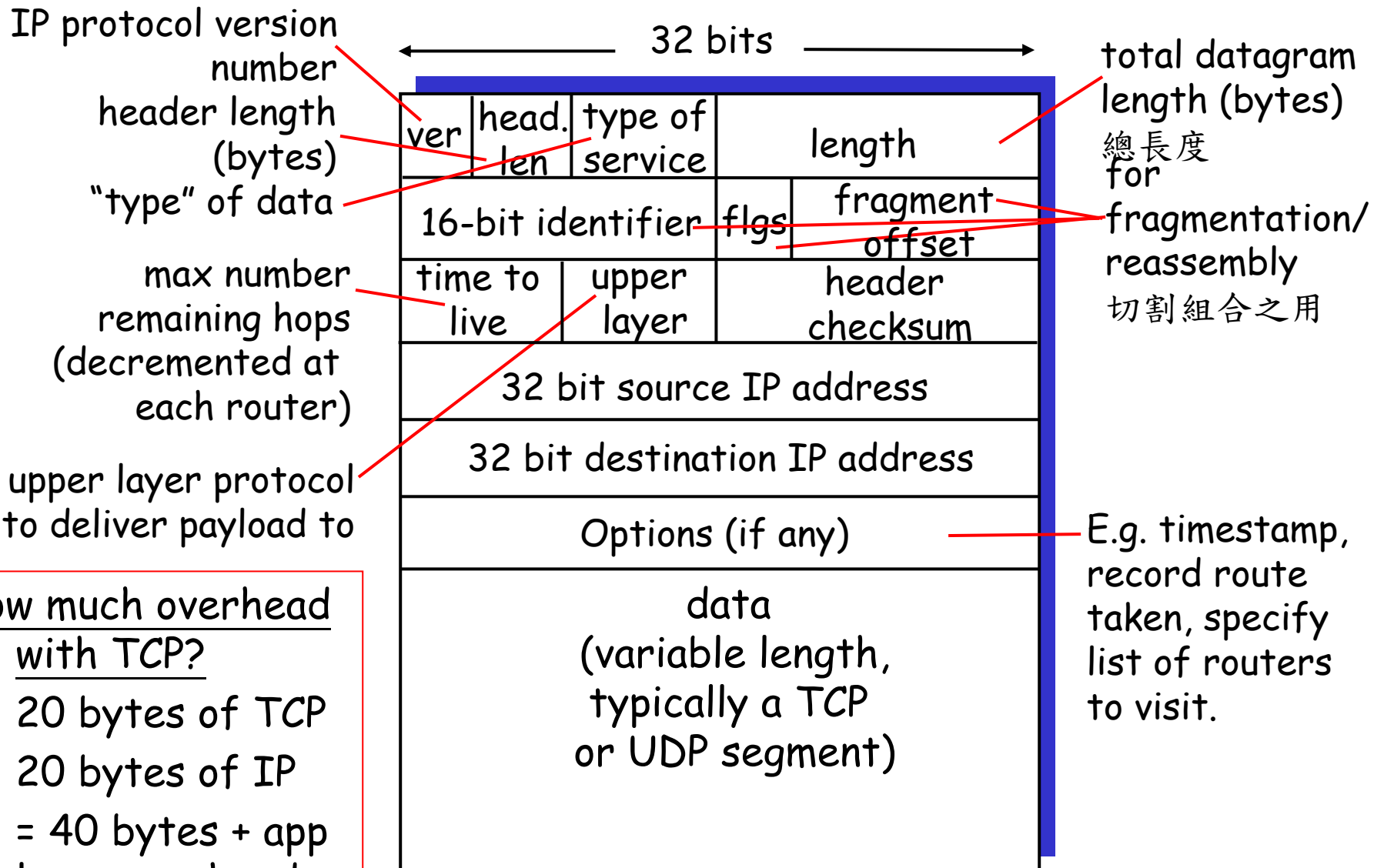
Host, router network layer functions:



# Chapter 4: Network Layer

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

# IP datagram format IP資料封包格式

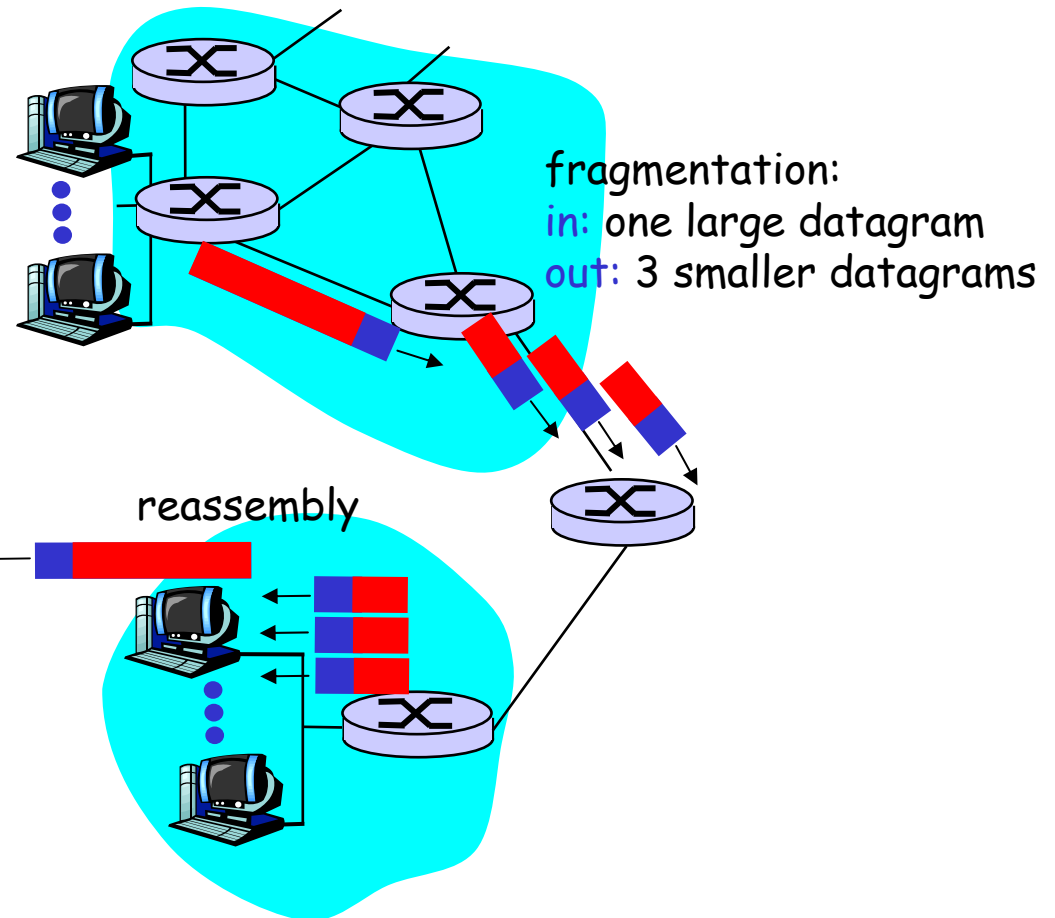


- how much overhead with TCP?
- ❑ 20 bytes of TCP
  - ❑ 20 bytes of IP
  - ❑ = 40 bytes + app layer overhead

# IP Fragmentation & Reassembly

## IP封包分割與組合

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



# IP Fragmentation and Reassembly

## Example

- ❑ 4000 byte datagram
- ❑ MTU = 1500 bytes
- ❑ offset以8 bytes為單位

1480 bytes in  
data field

offset =  
 $1480/8$

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

One large datagram becomes  
several smaller datagrams

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

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

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

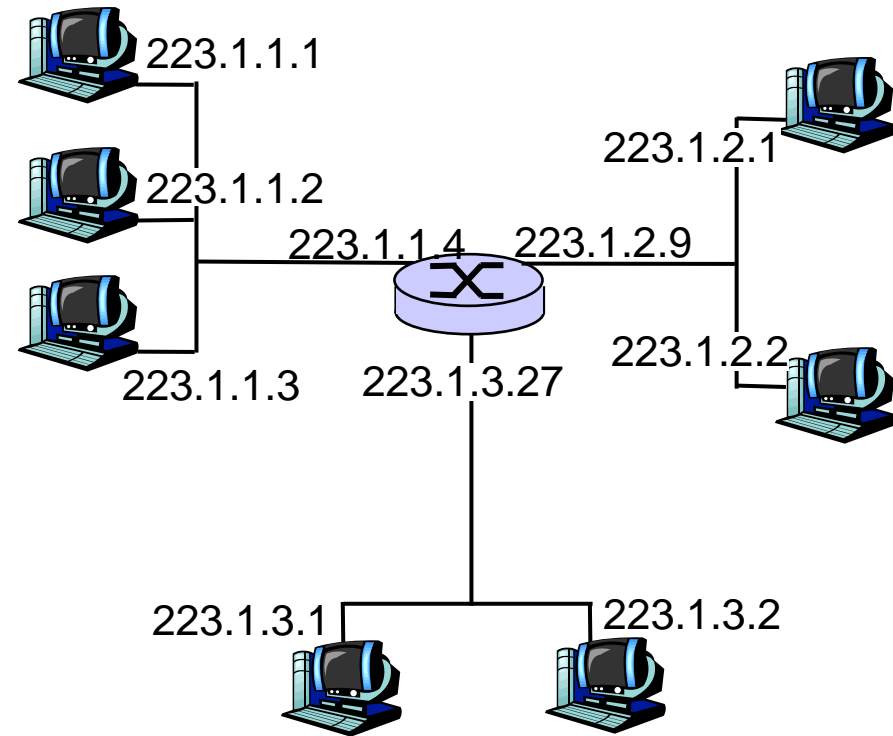
# Chapter 4: Network Layer

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



# IP Addressing: introduction 定址

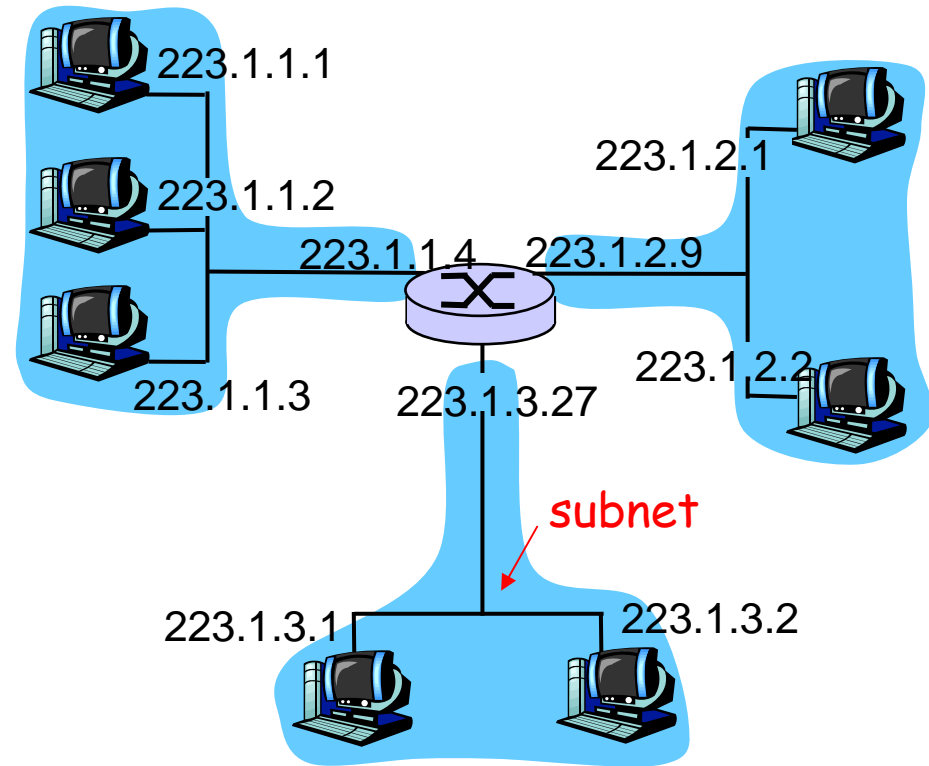
- IP address: 32-bit identifier for host, router *interface*
- *interface*: connection between host/router and physical link 介面
  - router's typically have multiple interfaces
  - host typically has one interface
  - IP addresses associated with each interface 每個介面有一個位址



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

# Subnets 子網路

- IP address: IP位址
  - subnet part (high order bits) 子網路部分
  - host part (low order bits) 主機部分
- *What's a subnet?*
  - device interfaces with same subnet part of IP address 相同IP位址中子網路部分的介面群
  - can physically reach each other without intervening router 可以不透過router即進行相互溝通

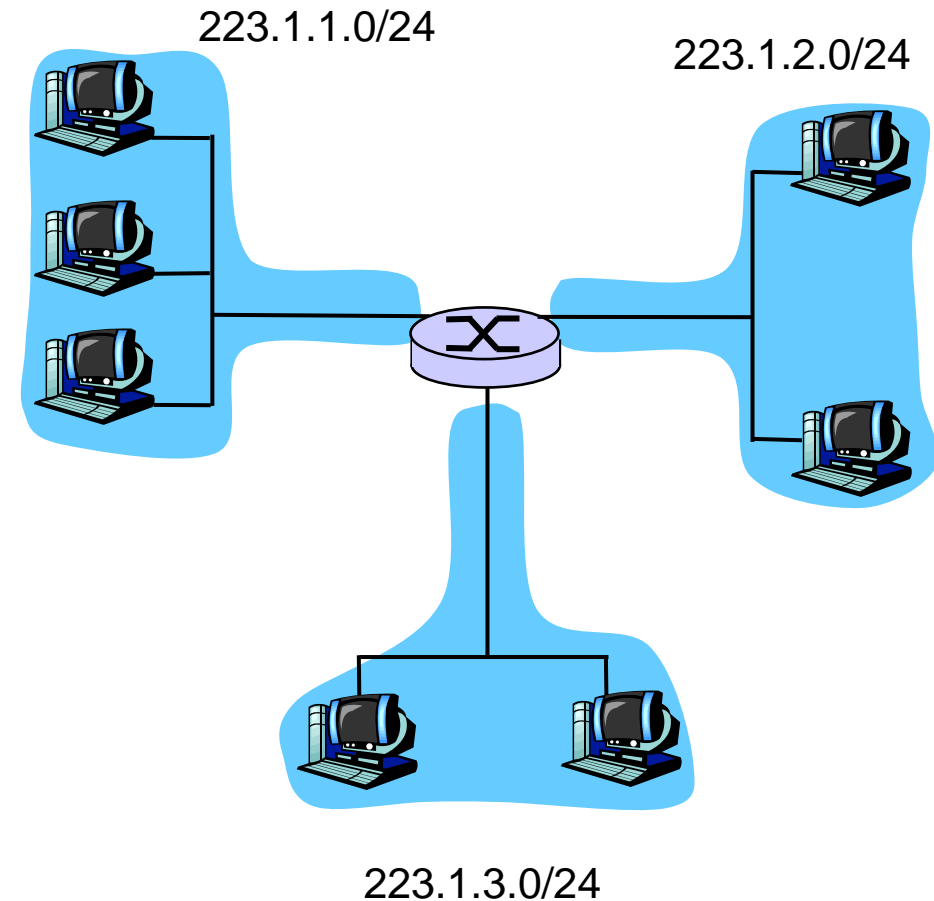


network consisting of 3 subnets

# Subnets 子網路

## 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**.

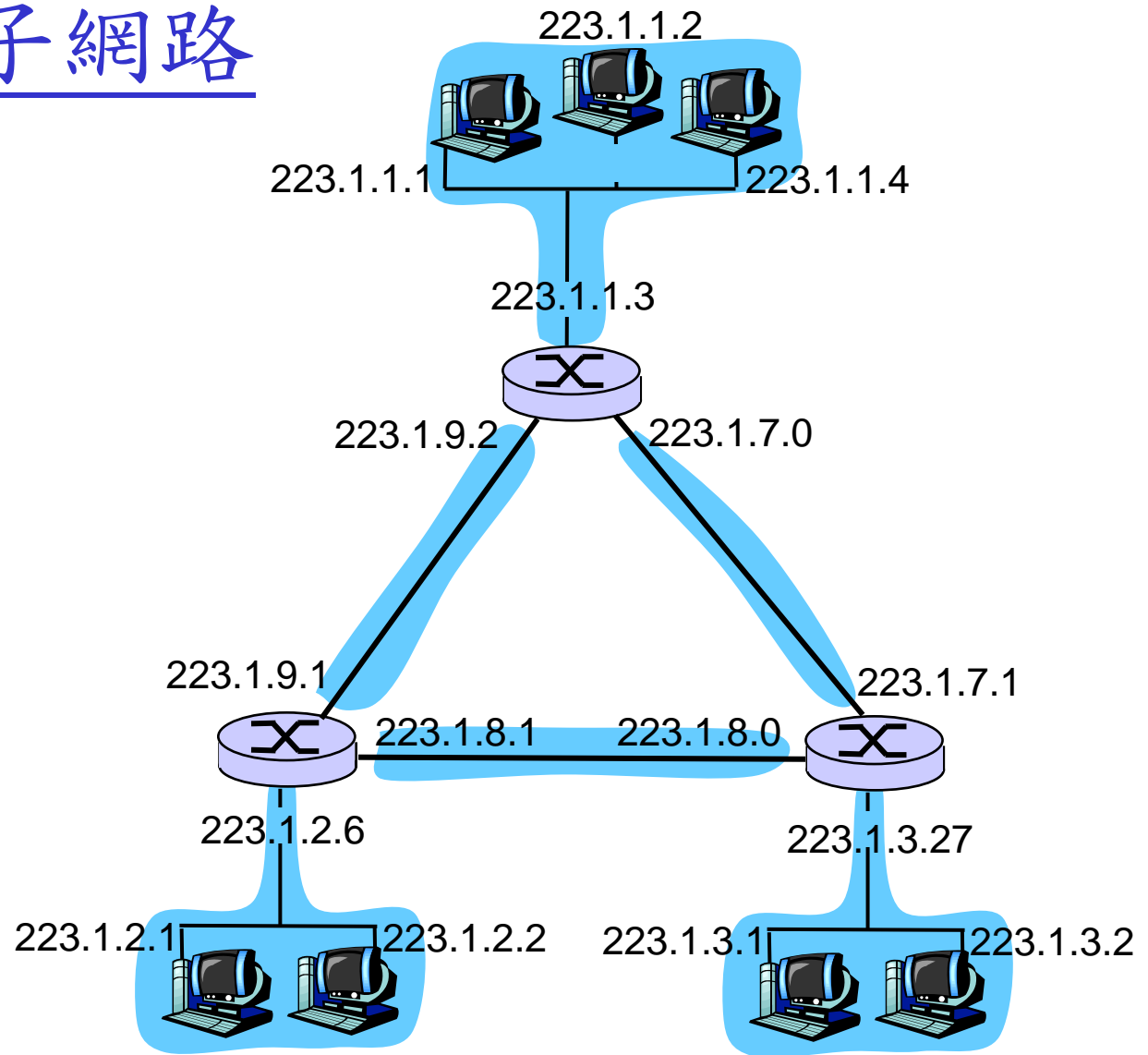


Subnet mask: /24

# Subnets 子網路

How many?

ANS :



# IP addressing: CIDR

## 無分級跨網路繞送

### **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

**x**代表子網路的bit數



# IP addresses: how to get one?

## 如何取得IP位址？

Q: How does *host* get IP address?

- hard-coded by system admin in a file
  - Wintel: 控制台 -> 網際網路 -> 區域網路 -> TCP/IP -> 內容
  - UNIX: /etc/rc.config
- **DHCP**: Dynamic Host Configuration Protocol:  
dynamically get address from a server
  - "plug-and-play" PNP

# DHCP: Dynamic Host Configuration Protocol

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

使電腦可以在連上網路時，動態取得IP位址

Can renew its lease on address in use

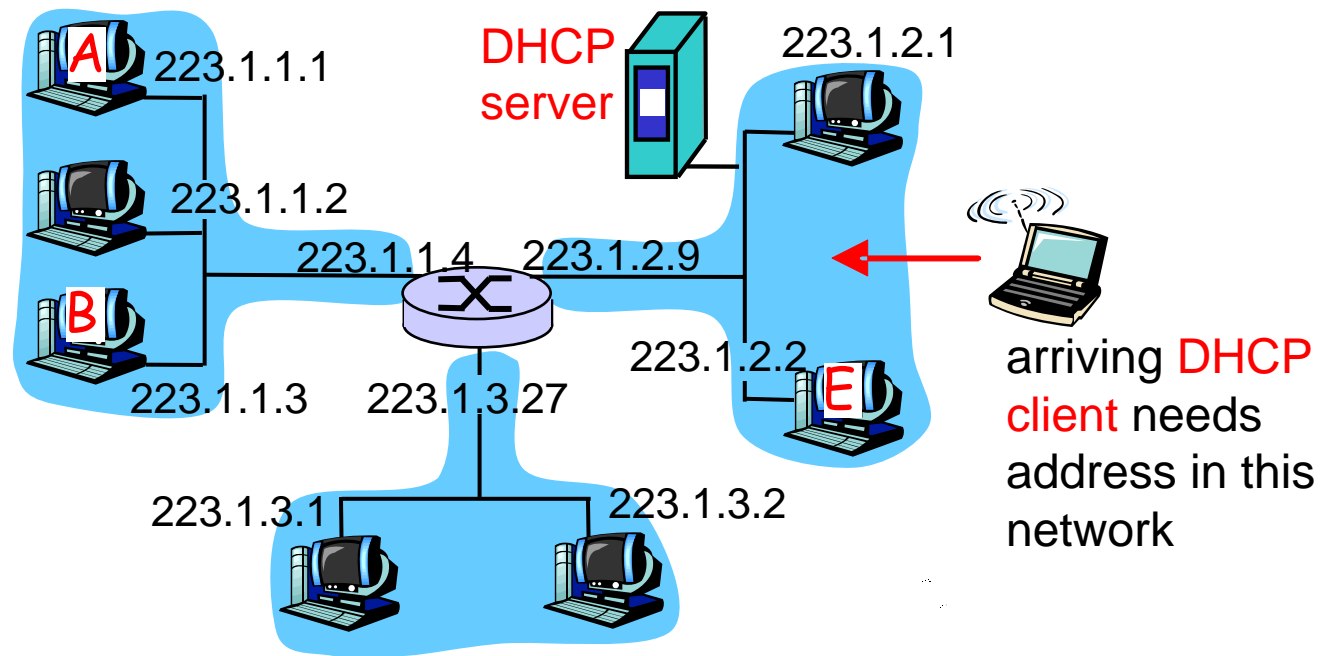
Allows reuse of addresses (only hold address while connected an "on")

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

DHCP overview:

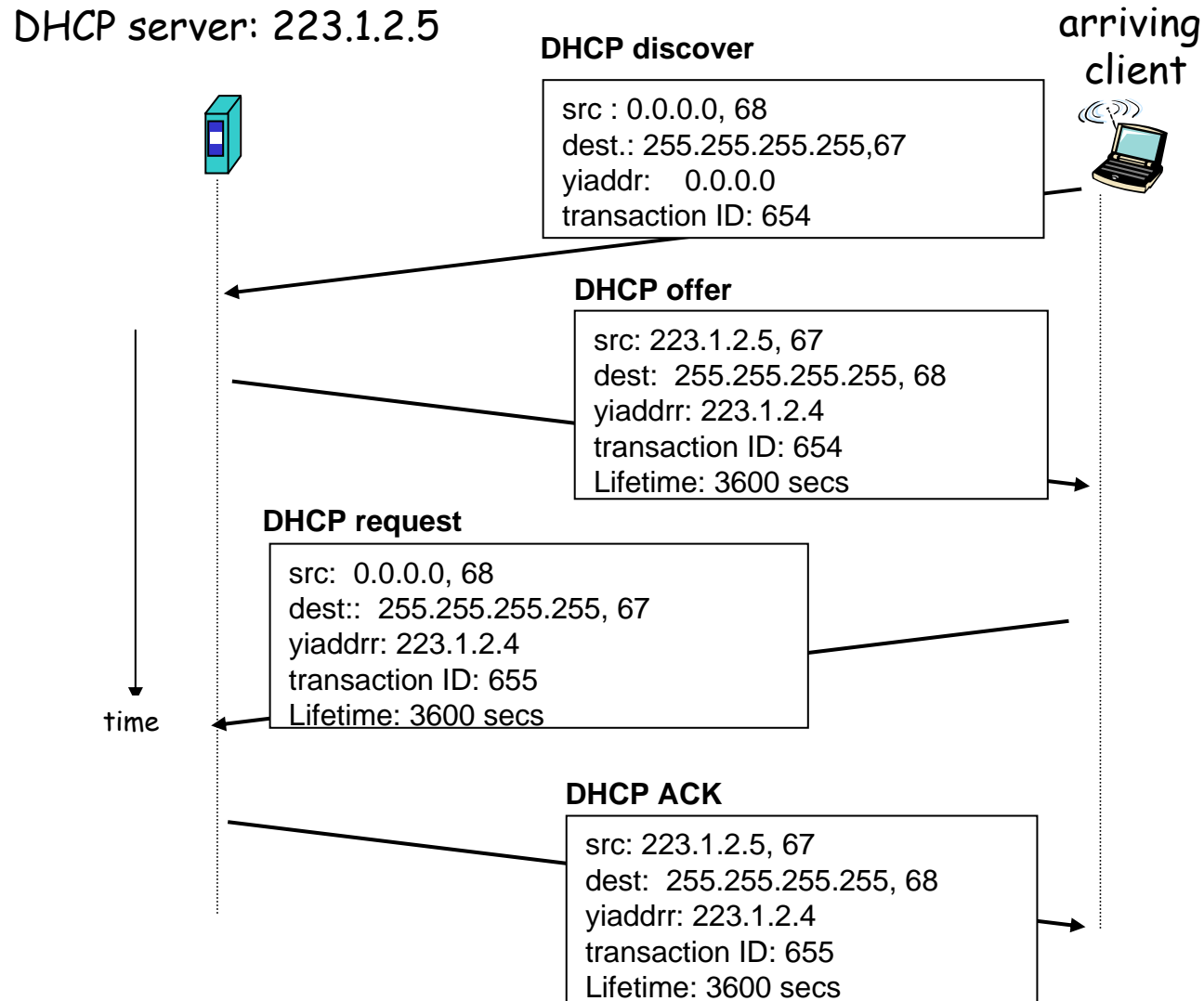
- host broadcasts "DHCP discover" msg
- DHCP server responds with "DHCP offer" msg
- host requests IP address: "DHCP request" msg
- DHCP server sends address: "DHCP ack" msg

# DHCP client-server scenario





# DHCP client-server scenario



# IP addresses: how to get one?

## 取得位址區塊

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

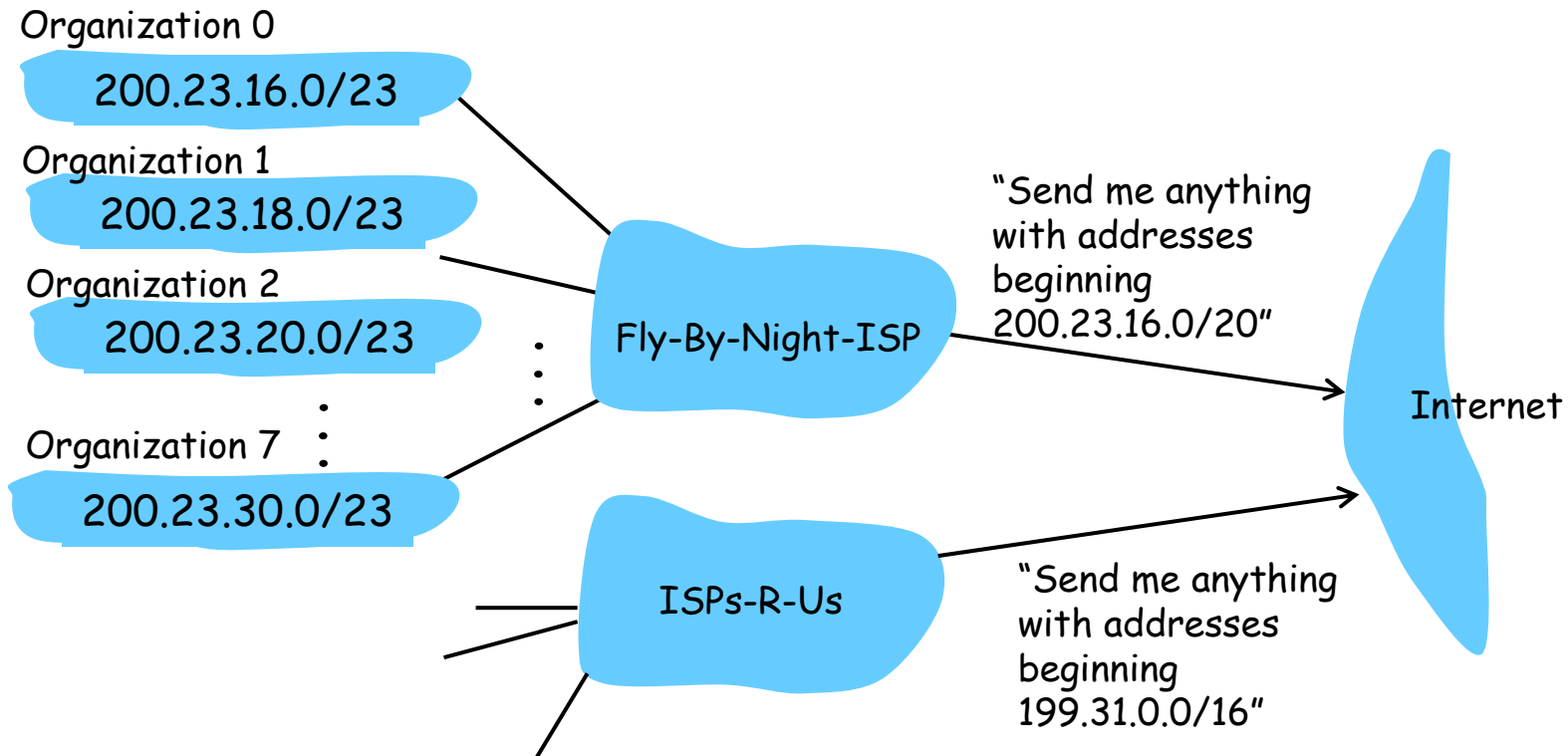
A: gets allocated portion of its provider ISP's address space

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

# Hierarchical addressing: route aggregation

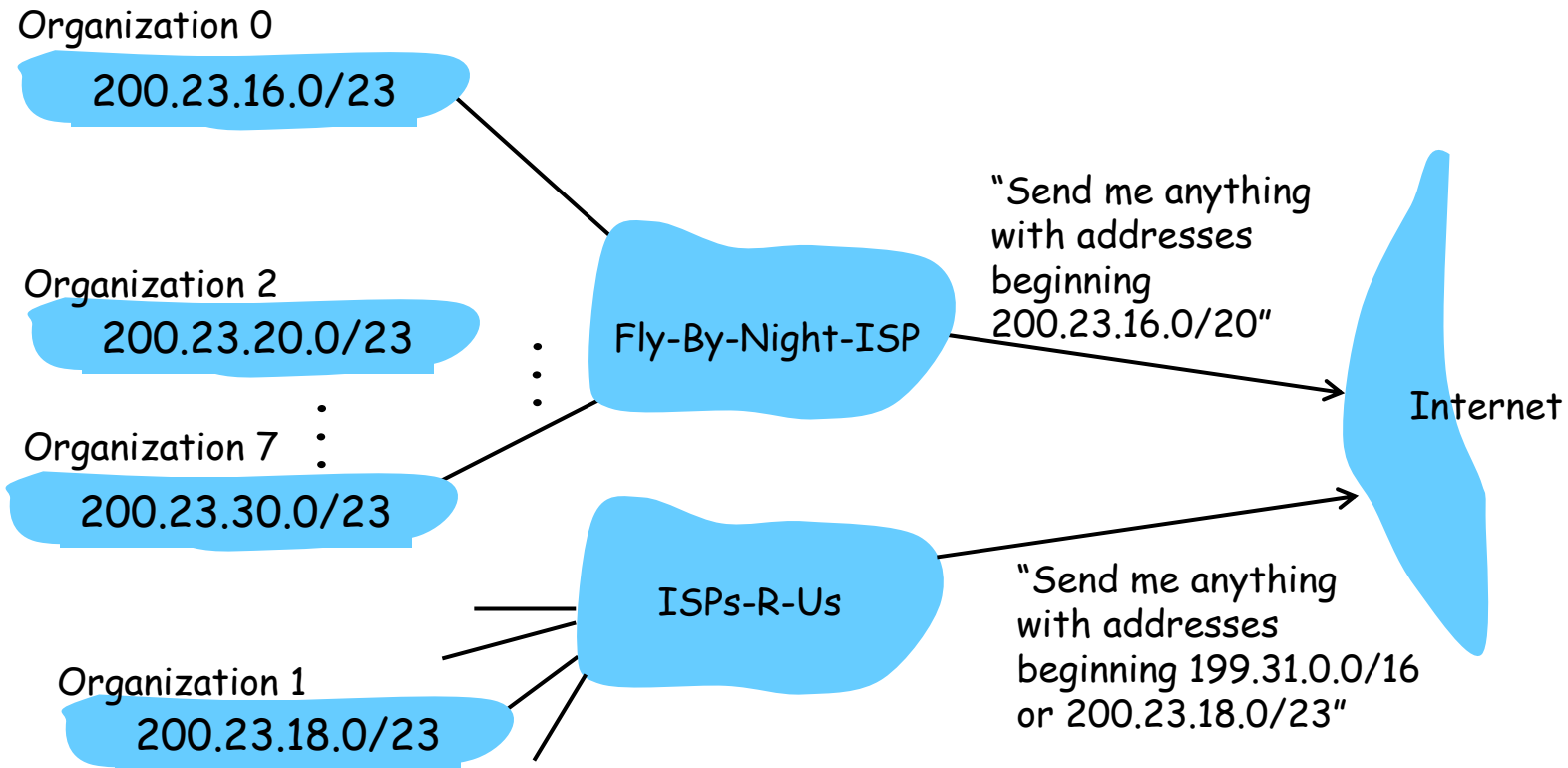
## 階層式定址：路由聚集

Hierarchical addressing allows efficient advertisement of routing information:



# Hierarchical addressing: more specific routes 特別的路由

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



## IP addressing: the last word... 最高組織

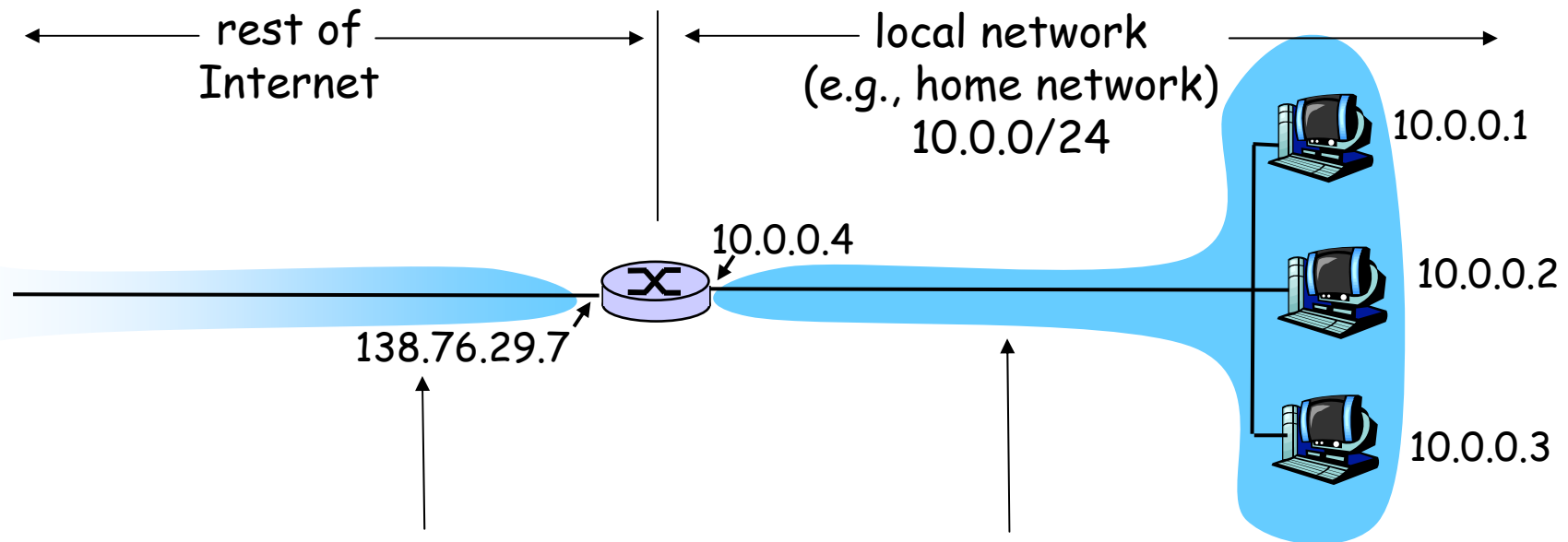
Q: How does an ISP get block of addresses?

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

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

# NAT: Network Address Translation

## 網路位址轉譯



*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).

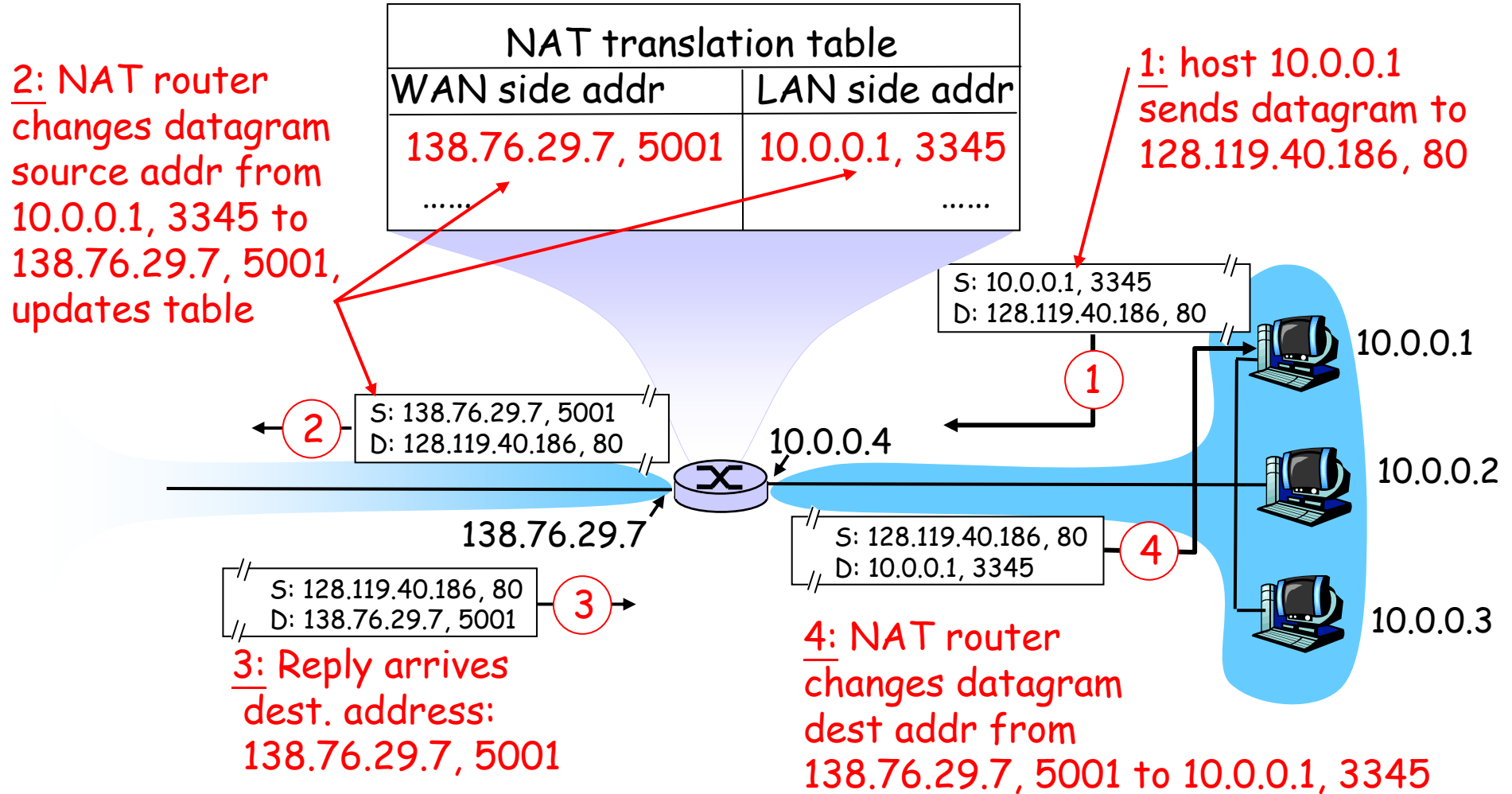
# NAT: Network Address Translation

**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

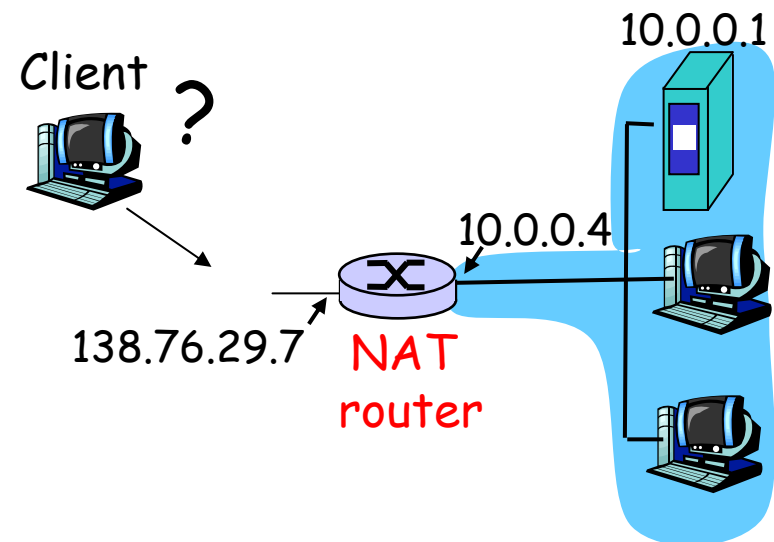


# NAT: Network Address Translation

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

# NAT traversal problem

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

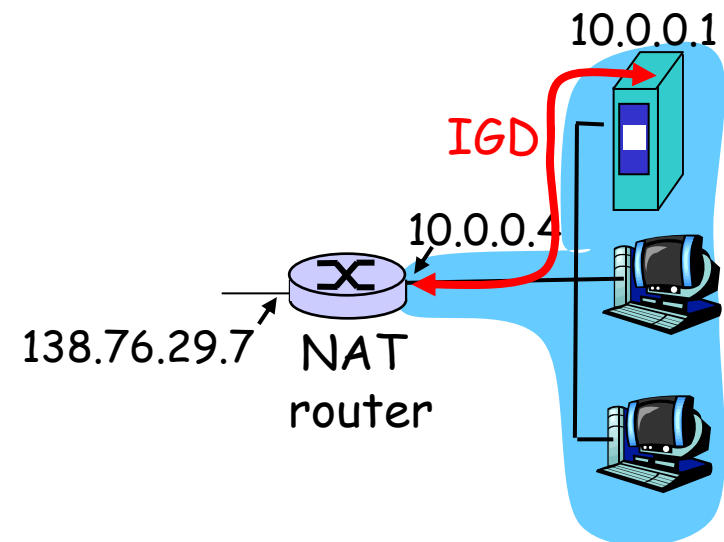


# NAT traversal problem

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

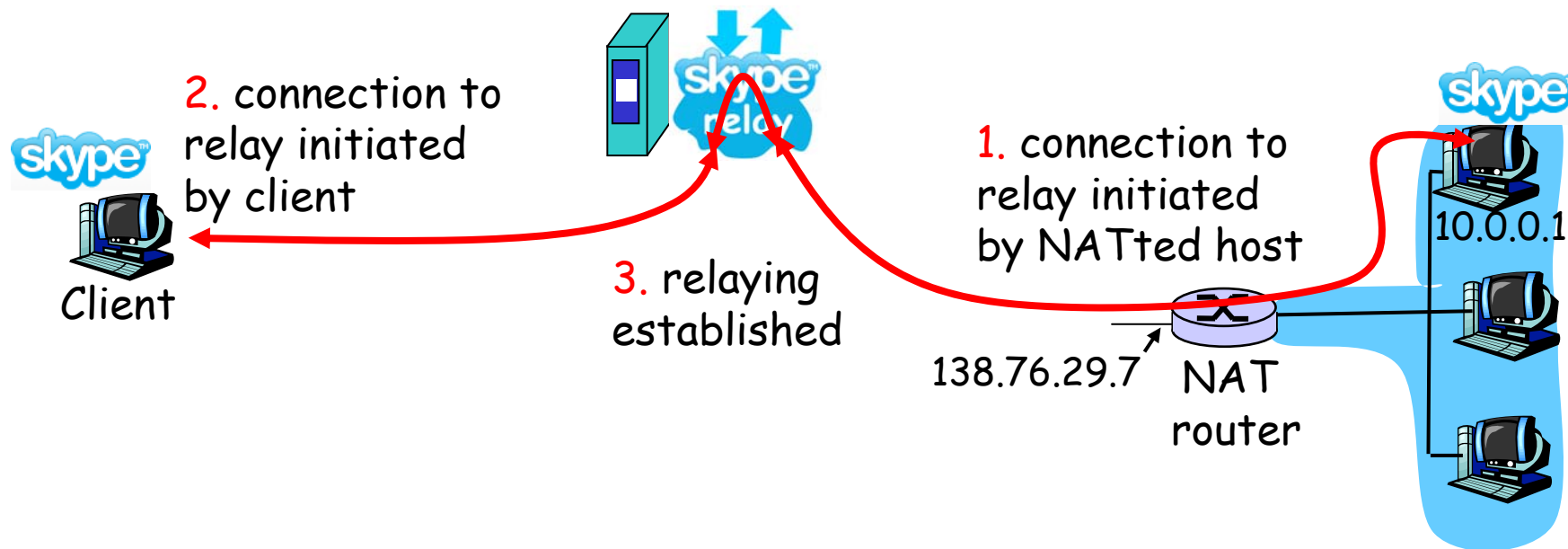
- ❖ learn public IP address (138.76.29.7)
- ❖ enumerate existing port mappings
- ❖ add/remove port mappings (with lease times)

i.e., automate static NAT port map configuration



# NAT traversal problem

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



# Chapter 4: Network Layer

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

# ICMP: Internet Control Message Protocol

## 網際網路控制訊息協定

- used by hosts & routers to communicate network-level information
  - error reporting: unreachable host, network, port, protocol
  - echo request/reply (used by ping)
- network-layer "above" IP:
  - ICMP msgs carried in IP datagrams
- **ICMP message:** type, code plus first 8 bytes of IP datagram causing error

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

# Traceroute and ICMP

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



# Chapter 4: Network Layer

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

# IPv6

- **Initial motivation:** 32-bit address space soon to be completely allocated. IP位址不足
- **Additional motivation:**
  - header format helps speed processing/forwarding
  - header changes to facilitate QoS

## **IPv6 datagram format:**

- fixed-length 40 byte header **40 bytes**的標頭區
- no fragmentation allowed 不切割封包

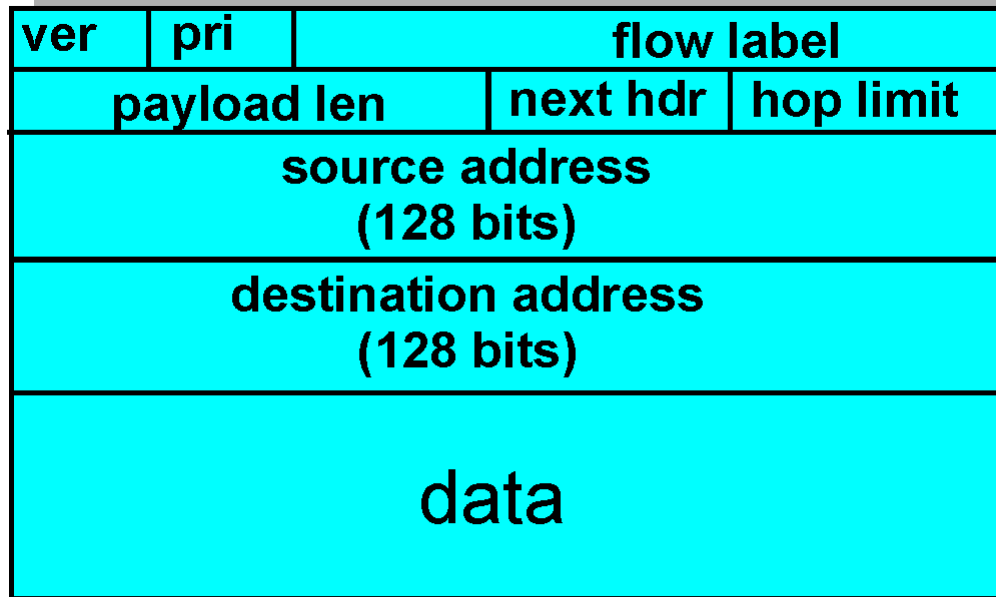
# IPv6 Header (Cont)

*Priority:* identify priority among datagrams in flow

*Flow Label:* identify datagrams in same "flow."

(concept of "flow" not well defined).

*Next header:* identify upper layer protocol for data



← 32 bits →

# Other Changes from IPv4

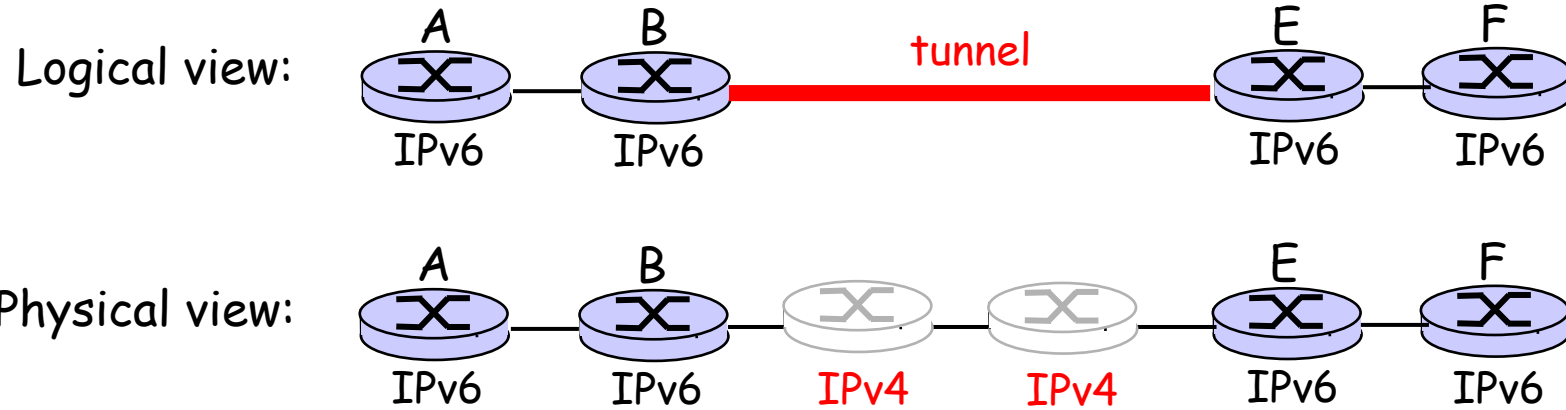
- ❑ *Checksum*: removed entirely to reduce processing time at each hop
- ❑ *Options*: allowed, but outside of header, indicated by "Next Header" field
- ❑ *ICMPv6*: new version of ICMP
  - additional message types, e.g. "Packet Too Big"
  - multicast group management functions

# Transition From IPv4 To IPv6

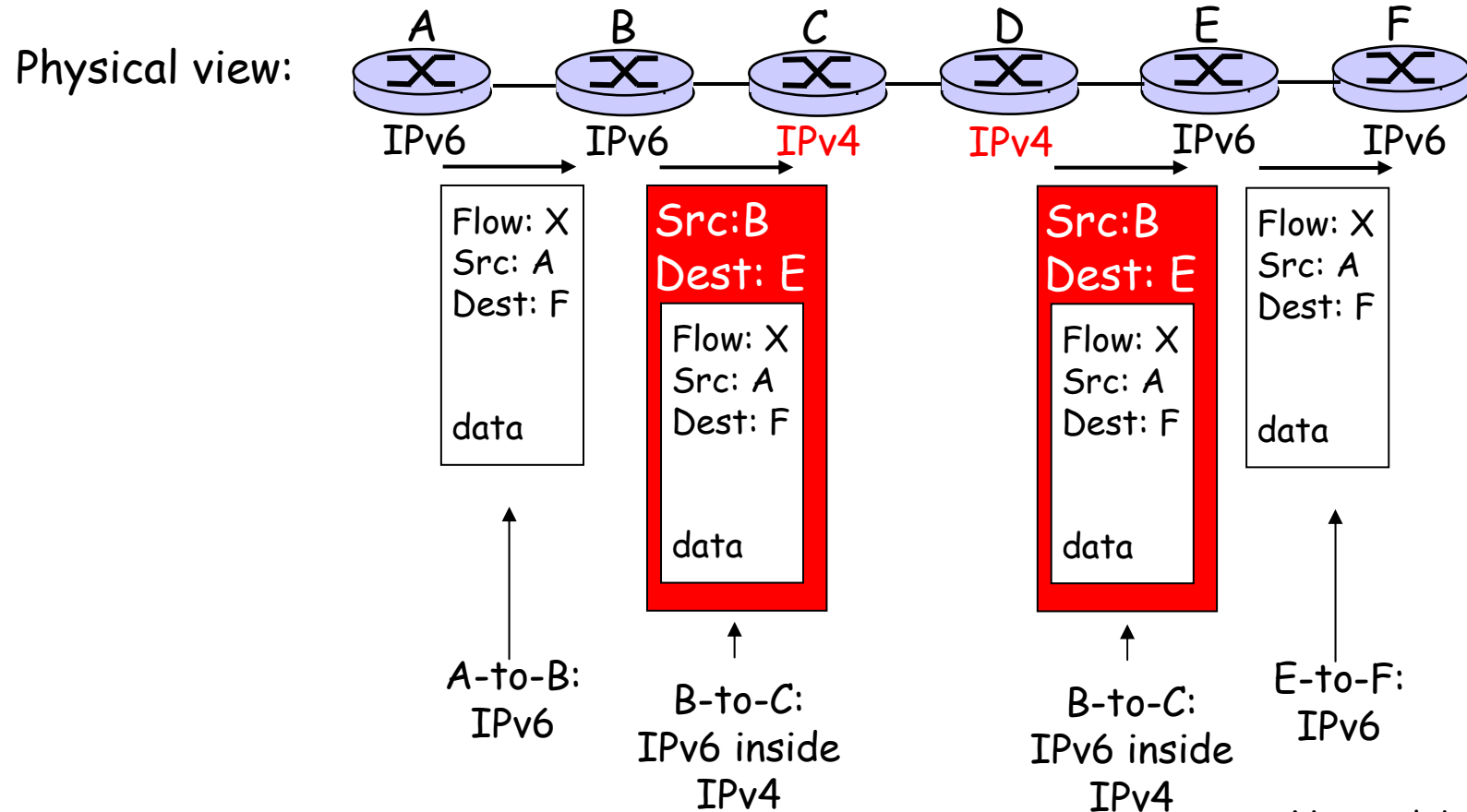
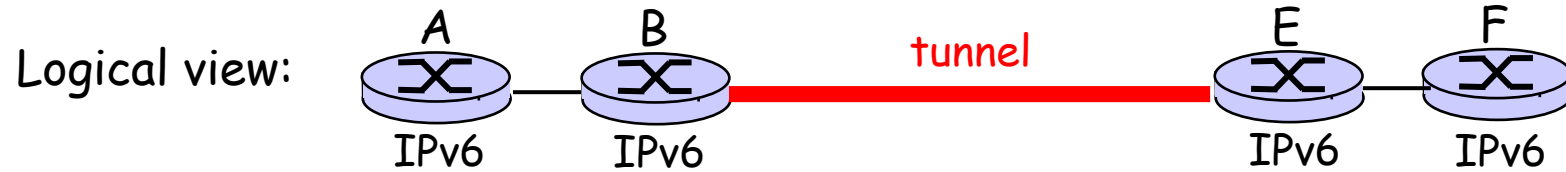
## 由IPv4到IPv6

- Not all routers can be upgraded simultaneous
  - no “flag days” 沒有固定的時程表
  - How will the network operate with mixed IPv4 and IPv6 routers? 如何整合
- *Tunneling*: IPv6 carried as payload in IPv4 datagram among IPv4 routers  
建隧道

# Tunneling 隧道化



# Tunneling 隧道化

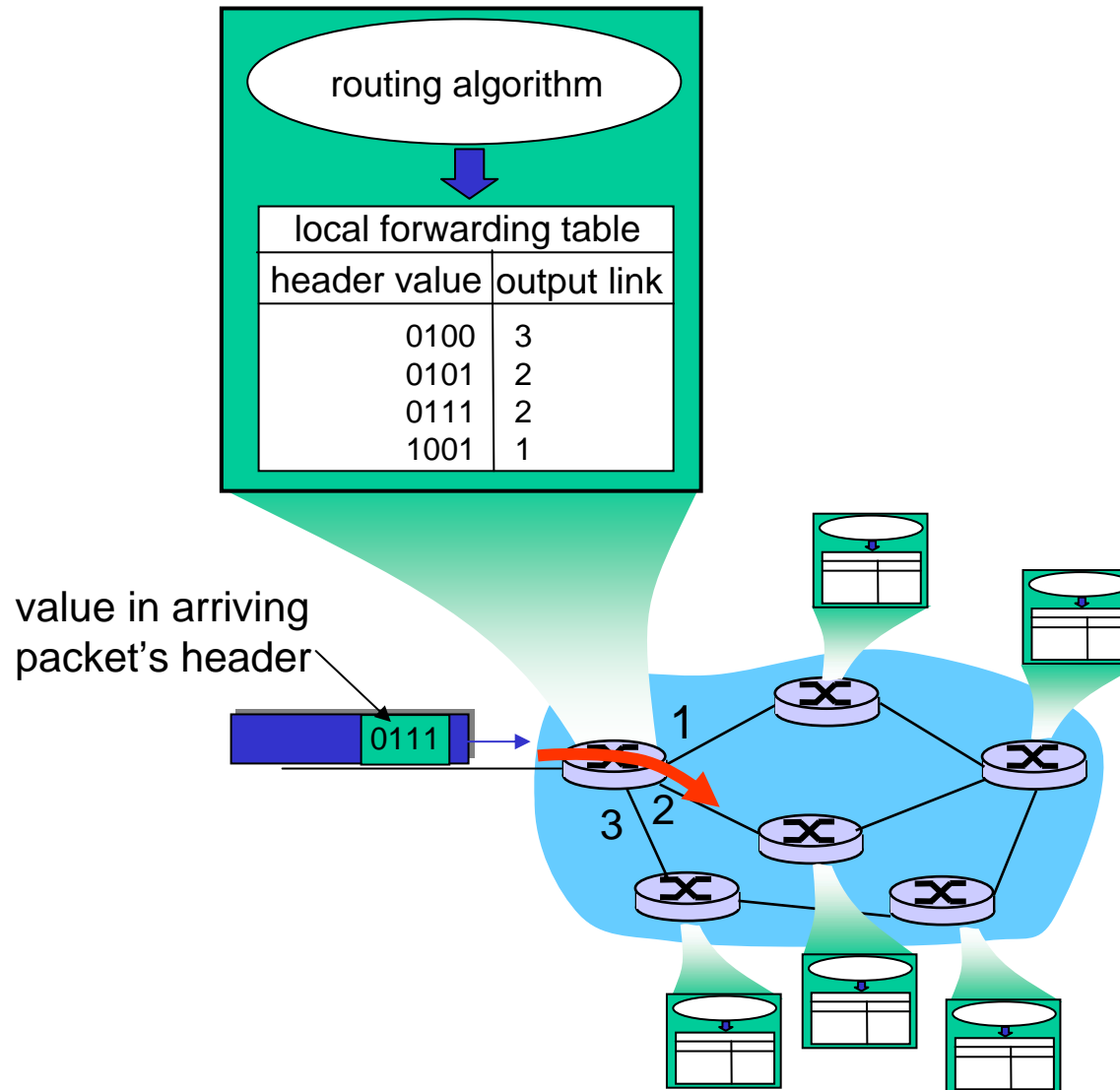


# Chapter 4: Network Layer

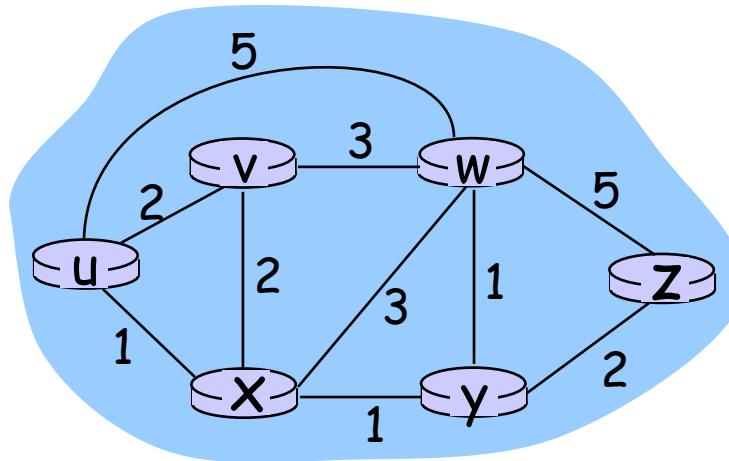
- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 **Routing algorithms**
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet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



# Interplay between routing, forwarding



# Graph abstraction 圖形化



Graph:  $G = (N, E)$

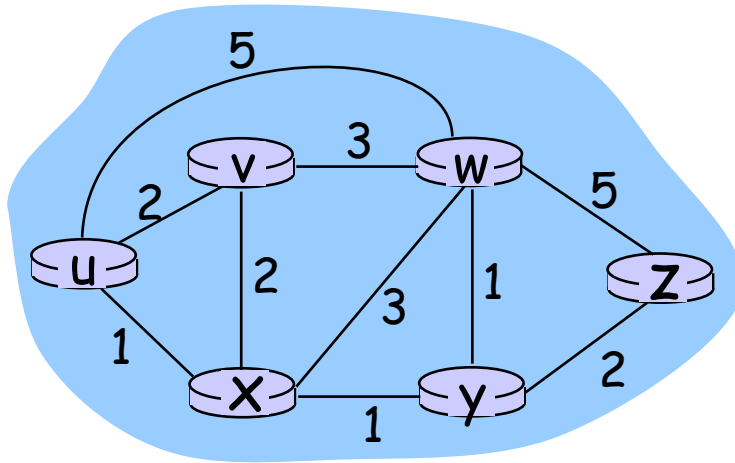
$N = \text{set of routers} = \{ u, v, w, x, y, z \}$

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

Remark: Graph abstraction is useful in other network contexts

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

# Graph abstraction: costs 成本



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

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

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

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

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

Routing algorithm: algorithm that finds least-cost path  
成本最低的路徑

# Routing Algorithm classification

## 繞徑演算法的分類

### Global or decentralized information?

分散式或集中式資訊

Global: 集中式

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

Decentralized: 分散式

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

### Static or dynamic?

Static:

- ❑ routes change slowly over time

Dynamic:

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

# Chapter 4: Network Layer

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

# A Link-State Routing Algorithm


## Dijkstra's algorithm

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

## Notation:

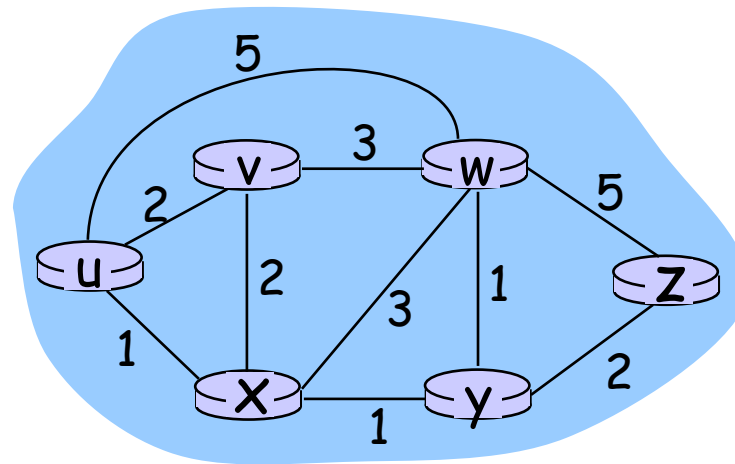
- ❑  $c(x,y)$ : link cost from node  $x$  to  $y$ ;  $= \infty$  if not direct neighbors
- ❑  $D(v)$ : current value of cost of path from source to dest.  $v$
- ❑  $p(v)$ : predecessor node along path from source to  $v$
- ❑  $N'$ : set of nodes whose least cost path definitively known

# Dijsktra's Algorithm 最短路徑演算法

- 1 **Initialization:**
  - 2  $N' = \{u\}$
  - 3 for all nodes  $v$
  - 4   if  $v$  adjacent to  $u$
  - 5     then  $D(v) = c(u,v)$
  - 6   else  $D(v) = \infty$
  - 7
  - 8 **Loop**
  - 9   find  $w$  not in  $N'$  such that  $D(w)$  is a minimum
  - 10   add  $w$  to  $N'$
  - 11   update  $D(v)$  for all  $v$  adjacent to  $w$  and not in  $N'$  :
  - 12      $D(v) = \min( D(v), D(w) + c(w,v) )$
  - 13   /\* new cost to  $v$  is either old cost to  $v$  or known
  - 14   shortest path cost to  $w$  plus cost from  $w$  to  $v$  \*/
  - 15 **until all nodes in  $N'$**
- 

# Dijkstra's algorithm: example

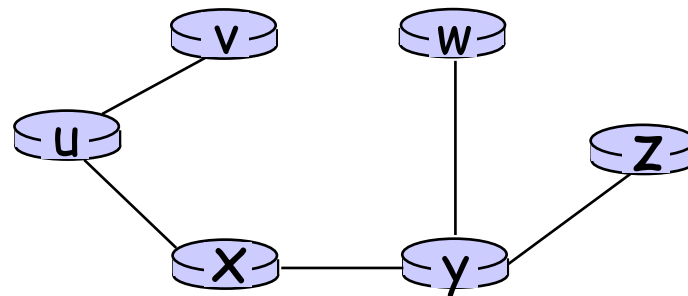
Step	N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0	u	2,u	5,u	1,u	$\infty$	$\infty$
1	ux	2,u	4,x		2,x	$\infty$
2	uxy	2,u	3,y			4,y
3	uxyv		3,y			4,y
4	uxyvw					4,y
5	uxyvwz					





# Dijkstra's algorithm: example (2)

Resulting shortest-path tree from u:



Resulting forwarding table in u:

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

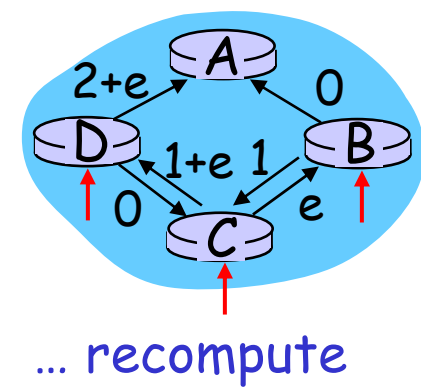
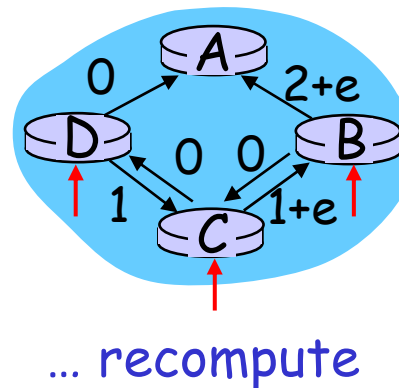
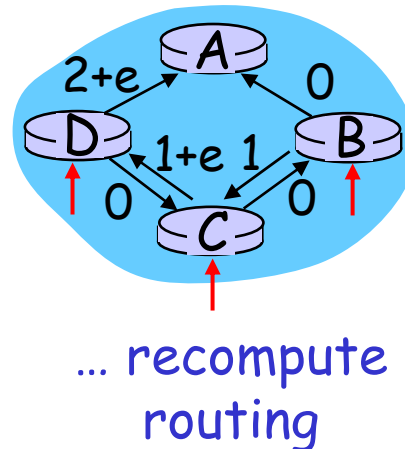
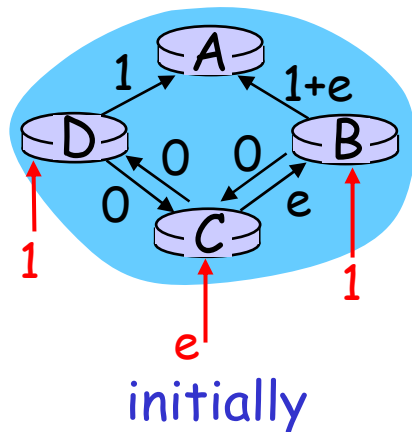
# Dijkstra's algorithm, discussion

**Algorithm complexity:**  $n$  nodes

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

**Oscillations possible:**

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



# Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 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

## 距離向量演算法

### Bellman-Ford Equation (dynamic programming)

Define

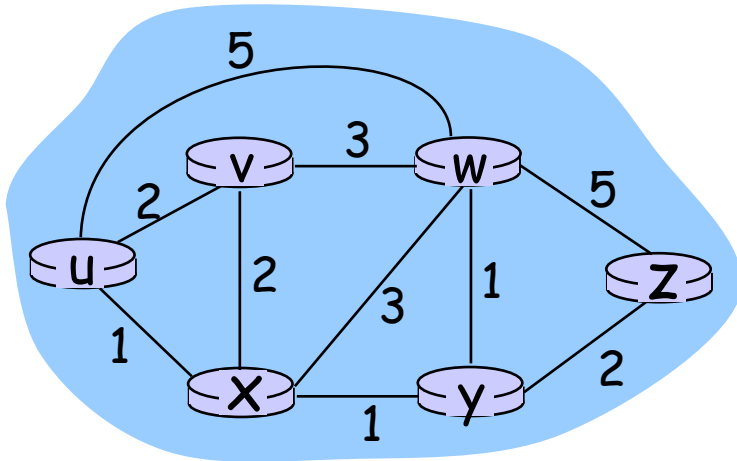
$d_x(y) :=$  cost of least-cost path from  $x$  to  $y$

Then

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

where min is taken over all neighbors  $v$  of  $x$

# Bellman-Ford example



Clearly,  $d_v(z) = 5$ ,  $d_x(z) = 3$ ,  $d_w(z) = 3$

B-F equation says:

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

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

# Distance Vector Algorithm

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

# Distance vector algorithm (4)

## Basic idea:

- Each node periodically sends its own distance vector estimate to neighbors
- When a node  $x$  receives new DV estimate from neighbor, it updates its own DV using B-F equation:

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

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

# Distance Vector Algorithm (5)

## Iterative, asynchronous:

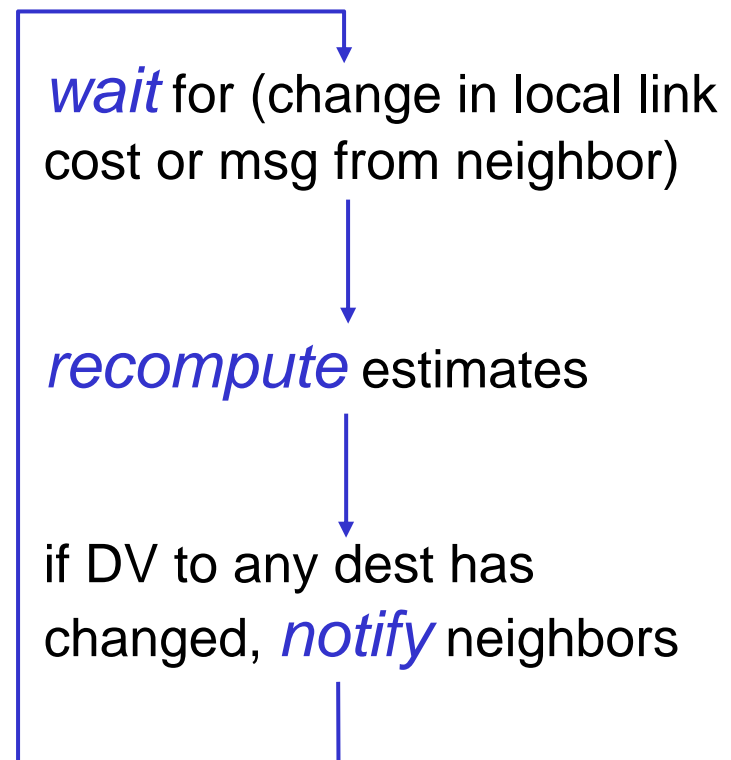
each local iteration caused by:

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

## Distributed:

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

## Each node:





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

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

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

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

**node x table**

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

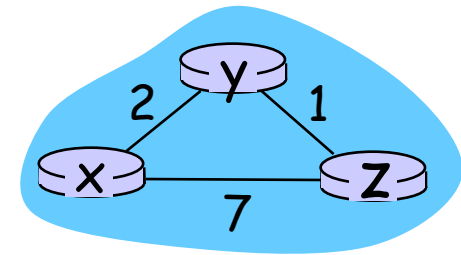
		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	7	1	0

**node y table**

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

**node z table**

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



time

$$D_x(y) = \min\{c(x,y) + D_y(y), c(x,z) + D_z(y)\} \\ = \min\{2+0, 7+1\} = 2$$

$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\} \\ = \min\{2+1, 7+0\} = 3$$

**node x table**

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

**node y table**

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

**node z table**

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

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

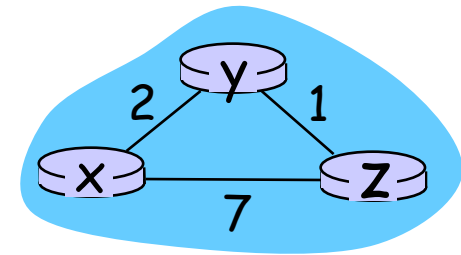
		cost to		
		x	y	z
from	x	0	2	7
	y	2	0	1
	z	7	1	0

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

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

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

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



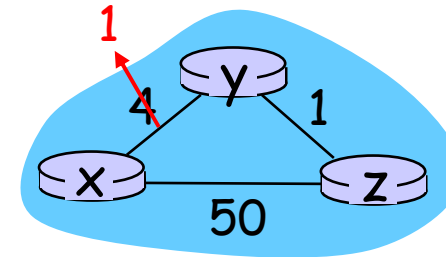
time →

# Distance Vector: link cost changes

## 成本改變的狀況

### Link cost changes:

- ❑ node detects local link cost change
- ❑ updates routing info, recalculates distance vector
- ❑ if DV changes, notify neighbors



“good  
news  
travels  
fast”

At time  $t_0$ ,  $y$  detects the link-cost change, updates its DV, and informs its neighbors.

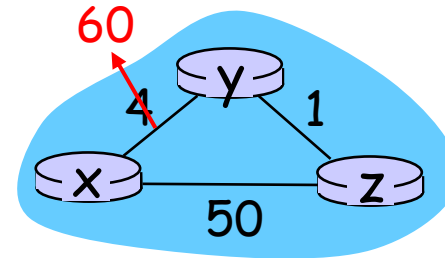
At time  $t_1$ ,  $z$  receives the update from  $y$  and updates its table. It computes a new least cost to  $x$  and sends its neighbors its DV.

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

# Distance Vector: link cost changes

## Link cost changes:

- ❑ good news travels fast  
好事傳得快
- ❑ 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

## 兩演算法的比較

### Message complexity

複雜度

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

### Speed of Convergence

收斂速度

- 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 強韌度**: what happens if router malfunctions?

### LS:

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

### DV:

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

# Chapter 4: Network Layer

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

# Hierarchical Routing 階層式繞徑

Our routing study thus far - idealization

- all routers identical
- network "flat"

... *not* true in practice

**scale:** with 200 million destinations:

- can't store all dest's in routing tables!
- routing table exchange would swamp links!

**administrative autonomy**

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

# Hierarchical Routing 階層式繞徑

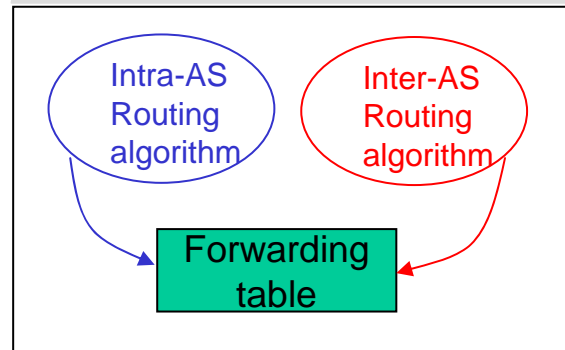
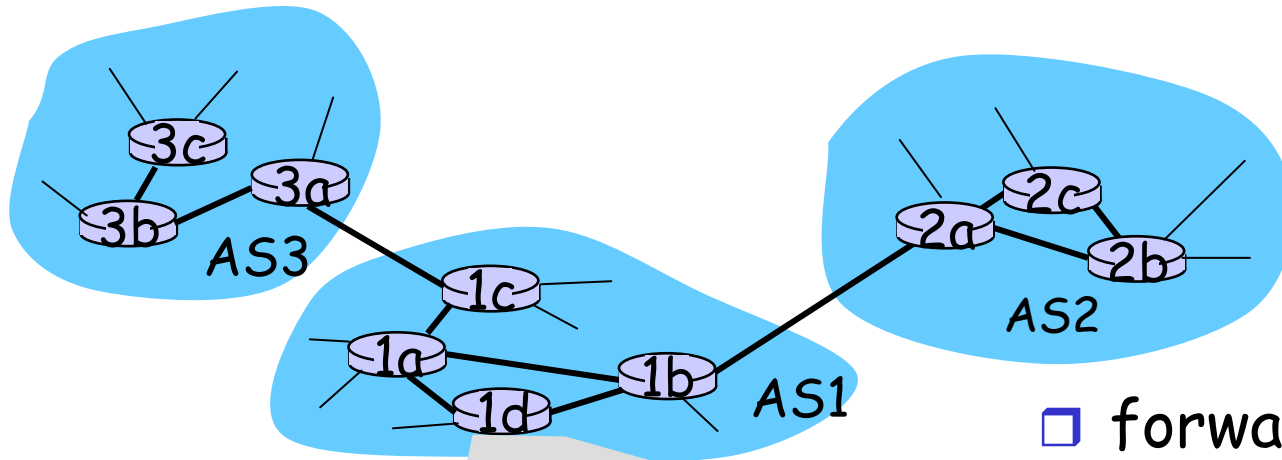
- ❑ aggregate routers into regions, "autonomous systems" (AS)
- ❑ routers in same AS run same routing protocol
  - "intra-AS" routing protocol
  - routers in different AS can run different intra-AS routing protocol

## Gateway router

- ❑ Direct link to router in another AS



# Interconnected ASes



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

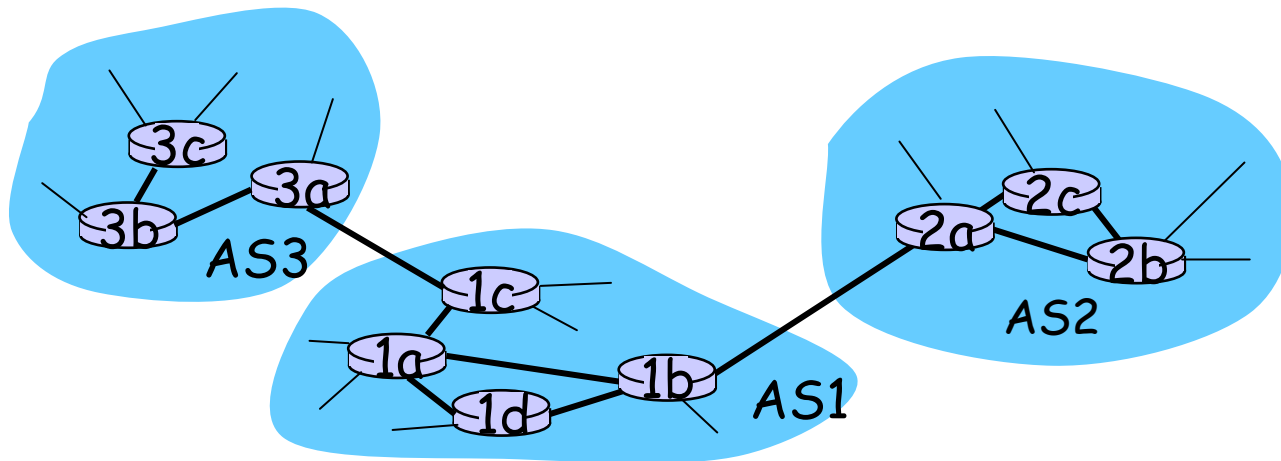
# Inter-AS tasks

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

## AS1 must:

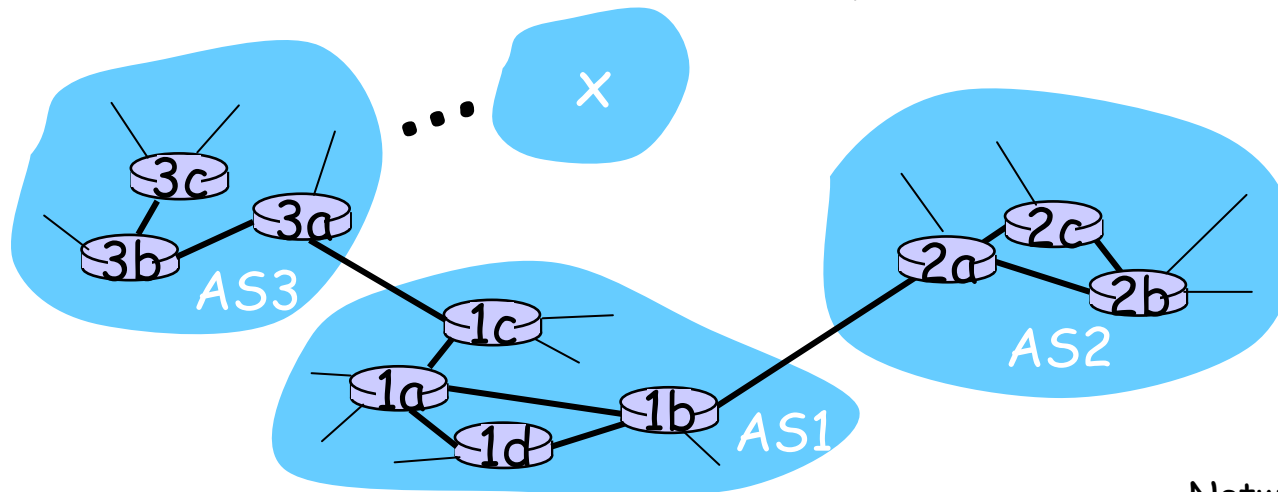
1. learn which dests reachable through AS2, which through AS3
2. propagate this reachability info to all routers in AS1

Job of inter-AS routing!



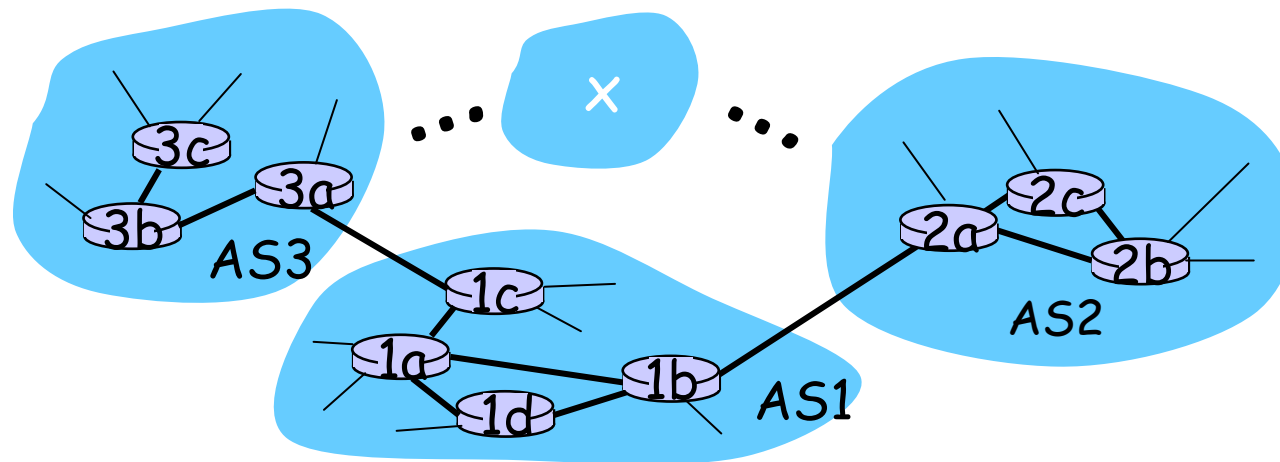
## Example: Setting forwarding table in router 1d

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



## Example: Choosing among multiple ASes

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



## Example: Choosing among multiple ASes

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

