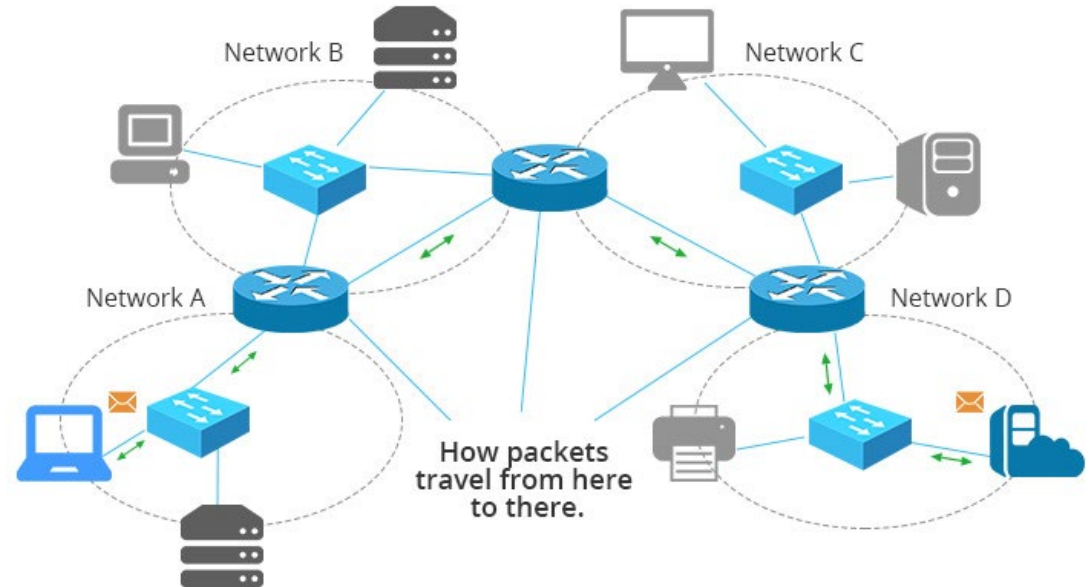


# Previous lectures:

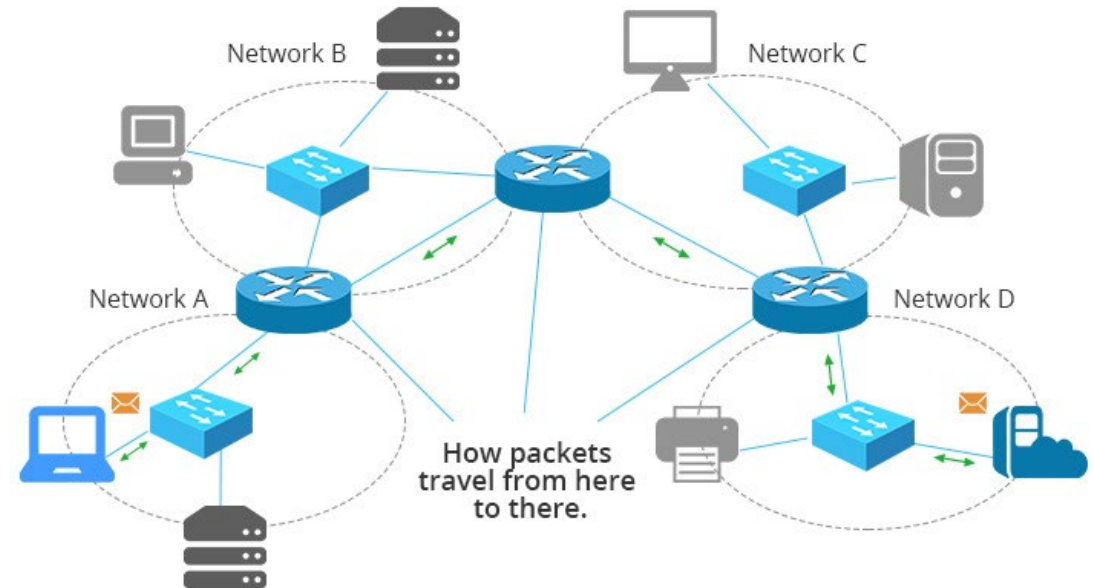
## Network Layer: “data plane”

- Forwarding in a router
- IP Protocol
- IP Addressing
- DHCP
- Network Address Translation (NAT)



# Network layer: “control plane” roadmap

- introduction
- routing protocols
  - link state
  - distance vector
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- Internet Control Message Protocol



# Network-layer functions

- **forwarding:** move packets from router's input to appropriate router output

*data plane*

- **routing:** determine route taken by packets from source to destination

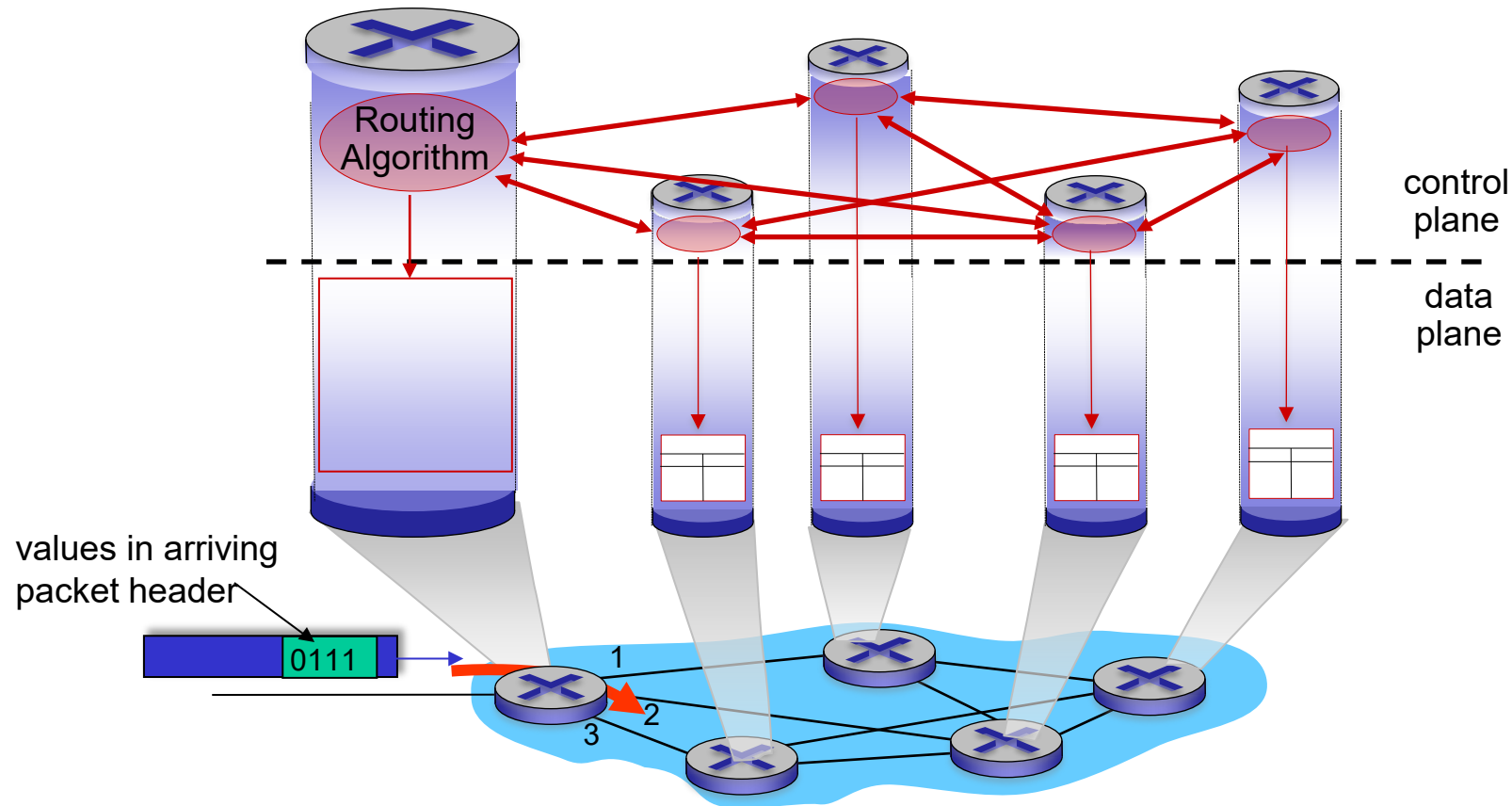
*control plane*

## Two approaches to structuring network control plane:

- per-router control (traditional)
- logically centralized control (software defined networking)

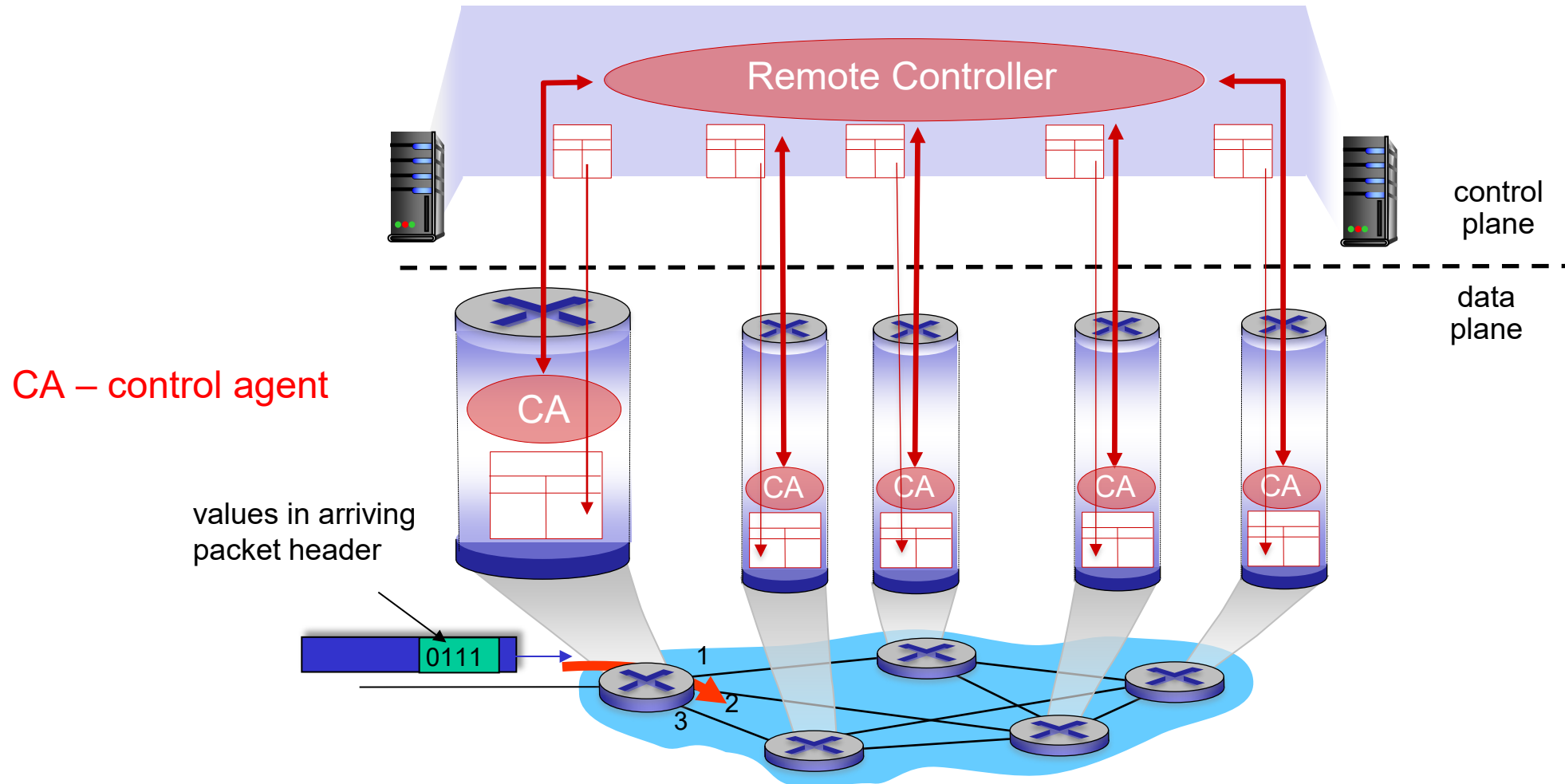
# Per-router control plane

Individual routing algorithm components *in each and every router* interact in the control plane

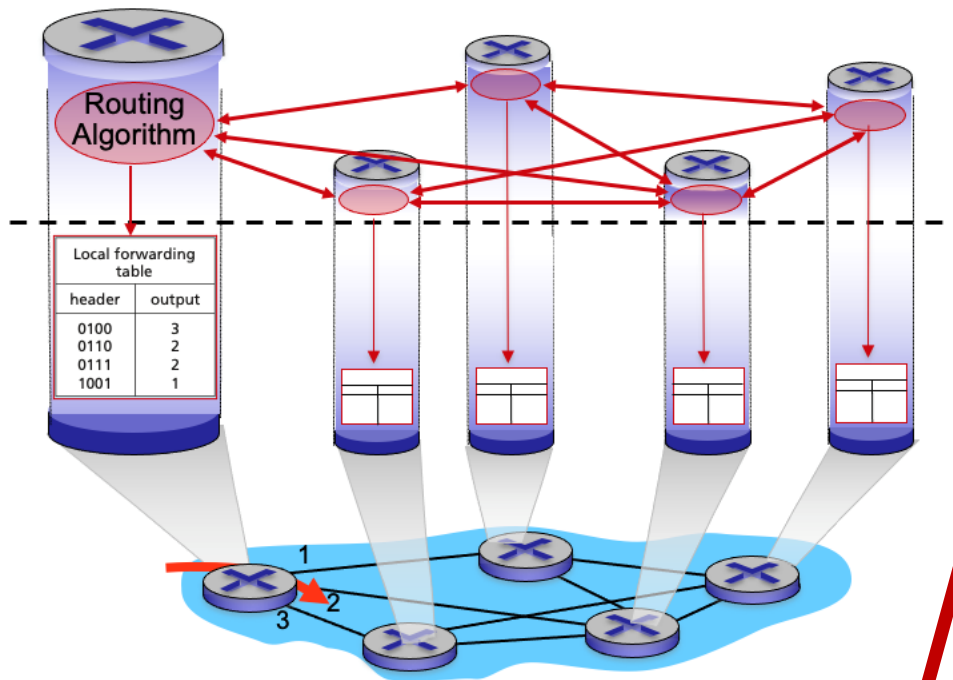


# Software-Defined Networking (SDN) control plane

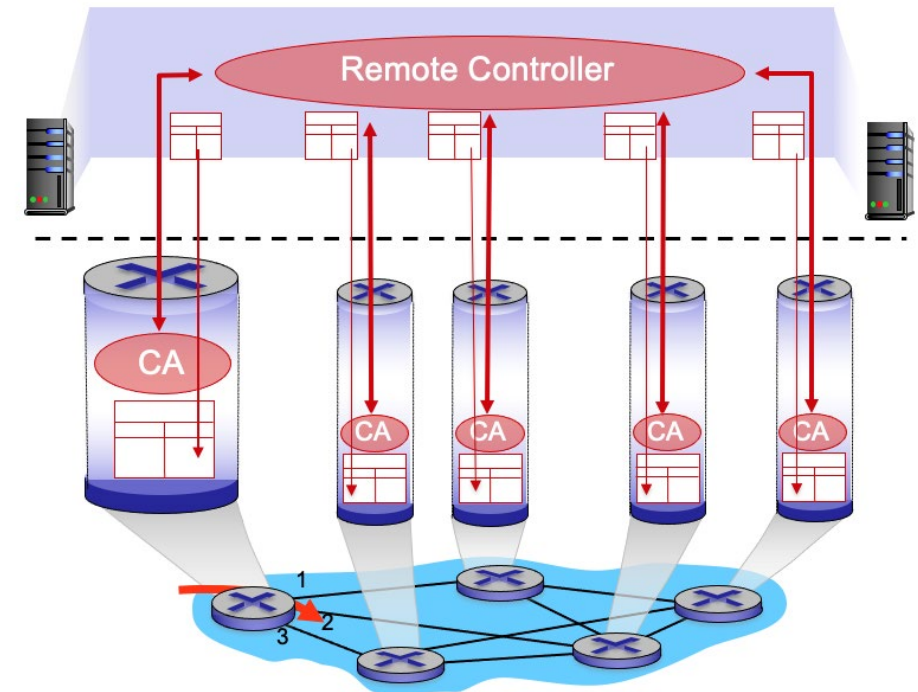
Remote controller computes, installs forwarding tables in routers



# Per-router control plane

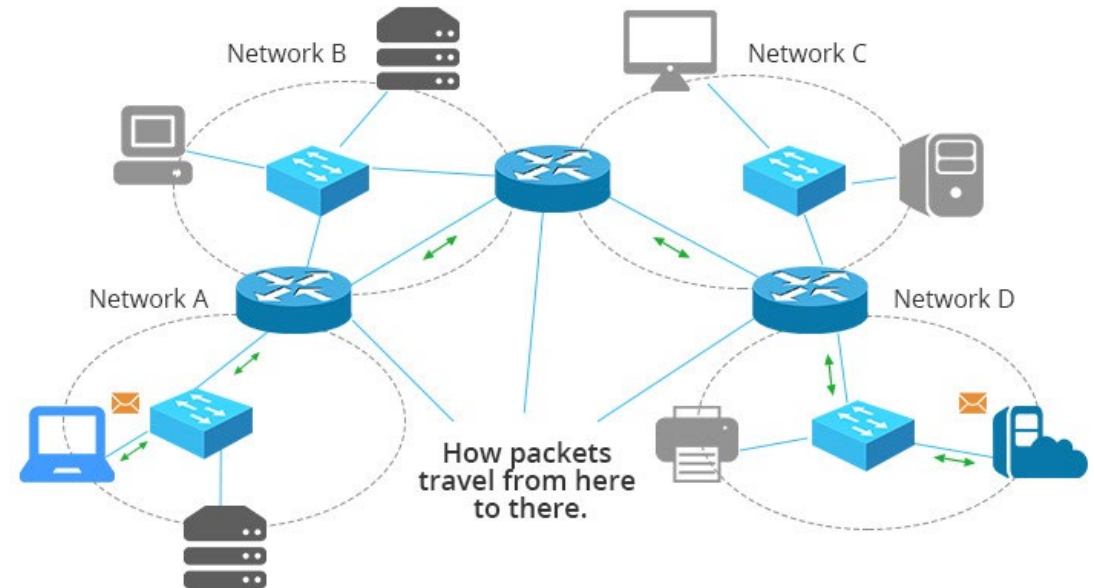


# SDN control plane



# Network layer: “control plane” roadmap

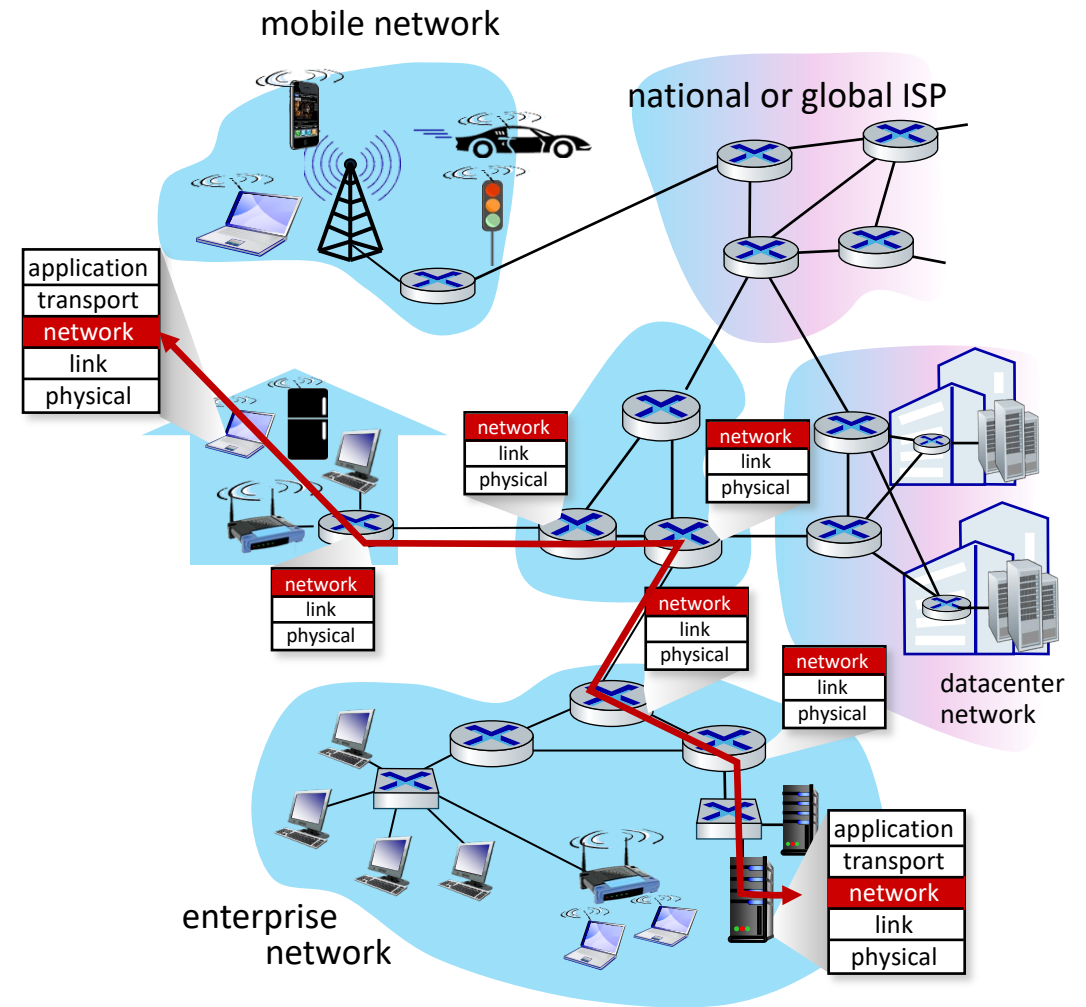
- introduction
- routing protocols
  - link state
  - distance vector
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- Internet Control Message Protocol



# Routing protocols

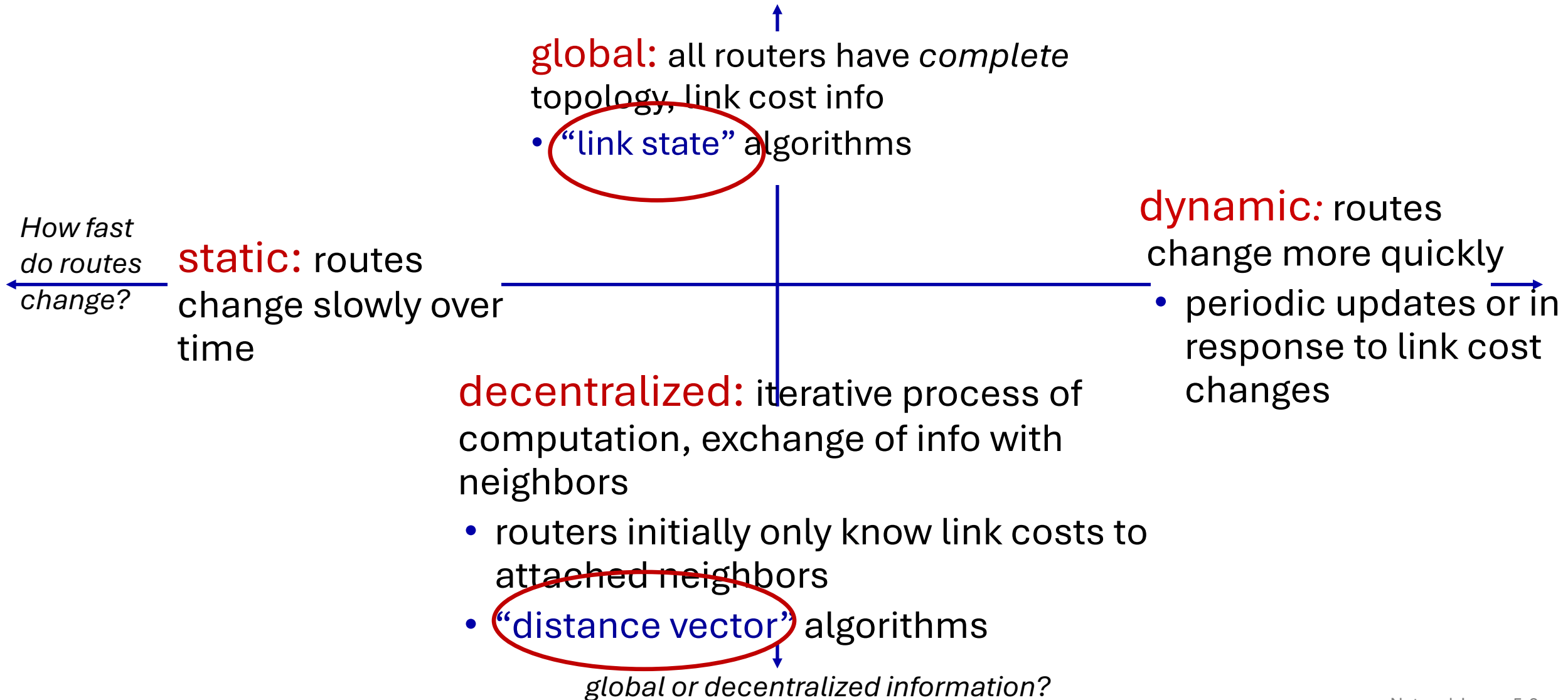
**Routing protocol goal:** determine “good” paths (equivalently, routes), from sending hosts to receiving host, through network of routers

- **path:** sequence of routers packets traverse from given initial source host to final destination host
- **“good”:** least “cost”, “fastest”, “least congested”
- routing: a “top-10” networking challenge!



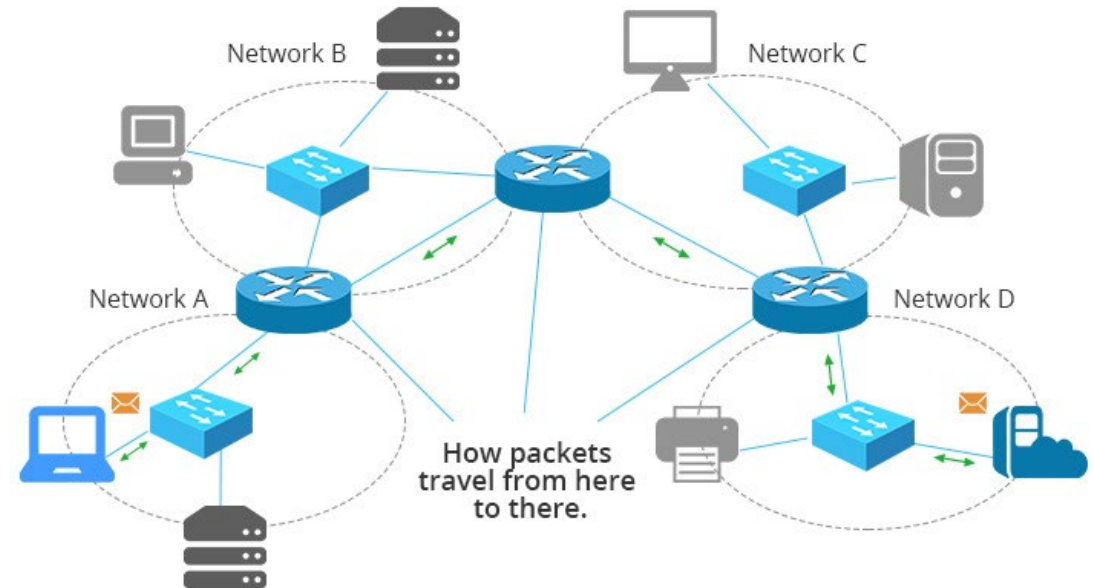


# Routing algorithm classification

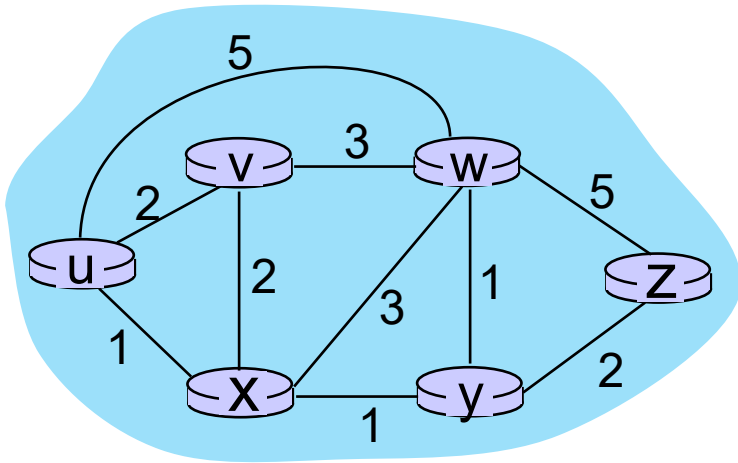


# Network layer: “control plane” roadmap

- introduction
- routing protocols
  - link state (Dijkstra's Algo)
  - distance vector
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- Internet Control Message Protocol



# Graph abstraction: link costs



$c_{a,b}$ : cost of *direct* link connecting  $a$  and  $b$   
e.g.,  $c_{w,z} = 5$ ,  $c_{u,z} = \infty$

cost defined by network operator:  
could always be 1, or inversely  
related to bandwidth, or proportional  
to congestion

graph:  $G = (N, E)$

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

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

# Dijkstra's link-state routing algorithm

- **centralized**: network 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  $k$  destinations

## notation

- $C_{x,y}$ : direct link cost from node  $x$  to  $y$ ;  $= \infty$  if not direct neighbors
- $D(v)$ : *current* estimate of cost of least-cost-path from source to destination  $v$
- $p(v)$ : predecessor node along path from source to  $v$
- $N'$ : set of nodes whose least-cost-path *definitively* known

# Dijkstra's link-state routing algorithm

1 *Initialization:*

2  $N' = \{u\}$  /\* compute least cost path from **u** to all other nodes \*/

3 for all nodes  $v$

4 if  $v$  adjacent to  $u$  /\*  $u$  initially knows direct-path-cost only to direct neighbors \*/

5 then  $D(v) = c_{u,v}$  /\* but may not be *minimum* cost! \*/

6 else  $D(v) = \infty$

7



8 *Loop*

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

11 update  $D(b)$  for all  $b$  adjacent to  $a$  and **not in  $N'$** :

12  **$D(b) = \min ( D(b), D(a) + c_{a,b} )$**

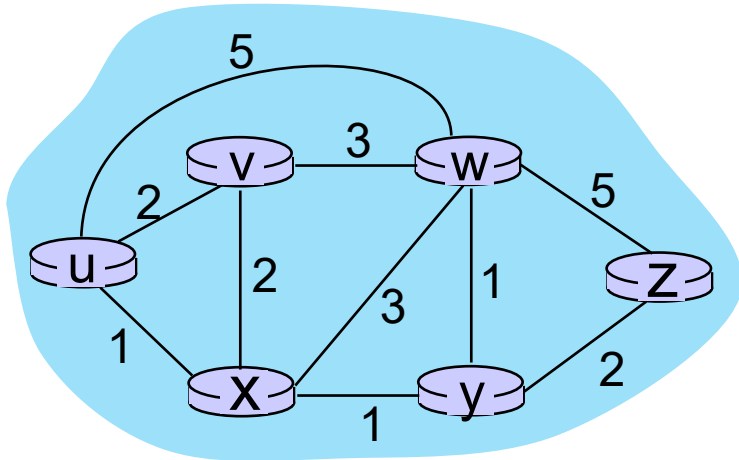
13 /\* new least-path-cost to  $v$  is either old least-cost-path to  $v$  or known

14 least-cost-path to  $w$  plus direct-cost from  $a$  to  $b$  \*/

15 *until all nodes in  $N'$*

# Dijkstra's algorithm: an example

		<b>v</b>	<b>w</b>	<b>x</b>	<b>y</b>	<b>z</b>
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						
2						
3						
4						
5						

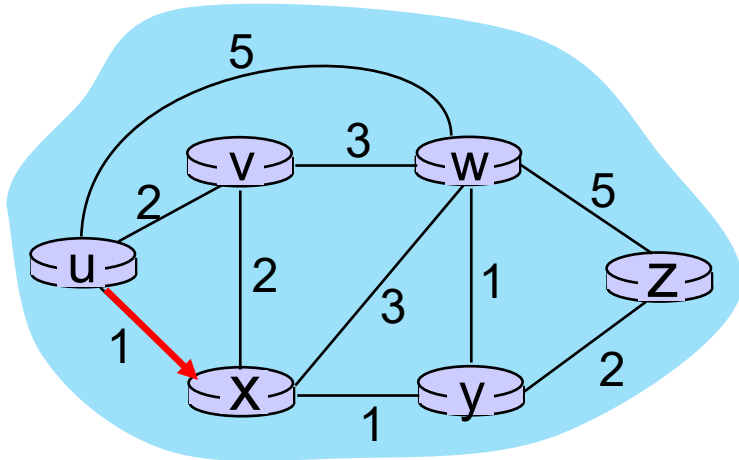


Initialization (step 0):

For all  $a$ : if  $a$  adjacent to  $u$  then  $D(a) = c_{u,a}$

# Dijkstra's algorithm: an 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						
3						
4						
5						



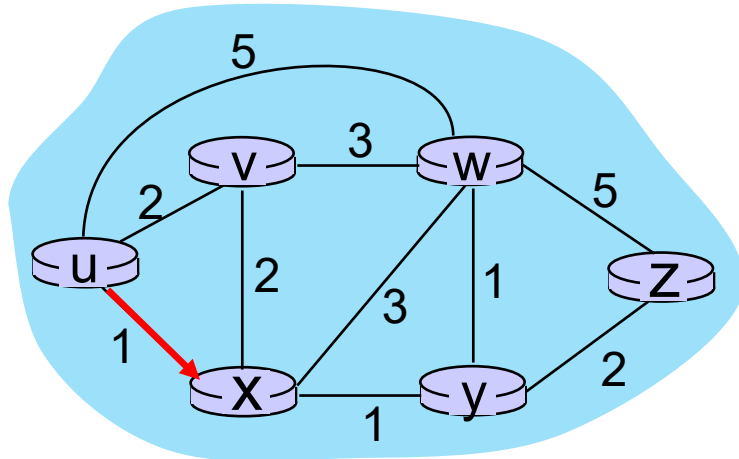
8 Loop

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

# Dijkstra's algorithm: an example

		<b>v</b>	<b>w</b>	<b>x</b>	<b>y</b>	<b>z</b>
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						
3						
4						
5						



8 Loop

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

11 update  $D(b)$  for all  $b$  adjacent to  $a$  and not in  $N'$ :

$$D(b) = \min ( D(b), D(a) + c_{a,b} )$$

$$D(v) = \min ( D(v), D(x) + c_{x,v} ) = \min(2, 1+2) = 2$$

$$D(w) = \min ( D(w), D(x) + c_{x,w} ) = \min(5, 1+3) = 4$$

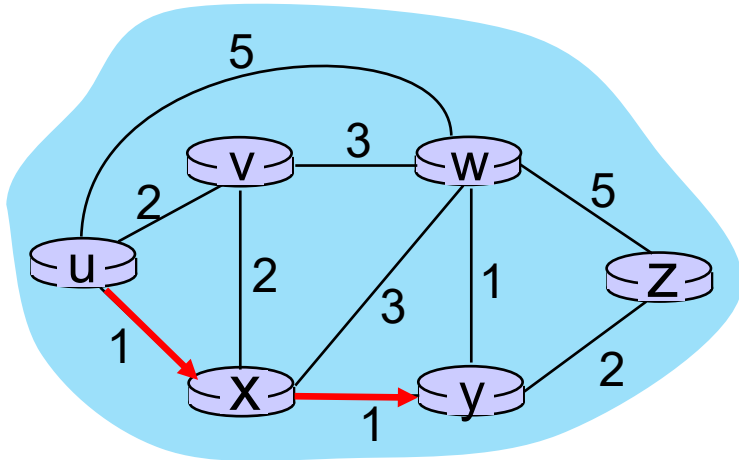
$$D(y) = \min ( D(y), D(x) + c_{x,y} ) = \min(\infty, 1+1) = 2$$





# Dijkstra's algorithm: an example

		<b>v</b>	<b>w</b>	<b>x</b>	<b>y</b>	<b>z</b>
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					
3						
4						
5						



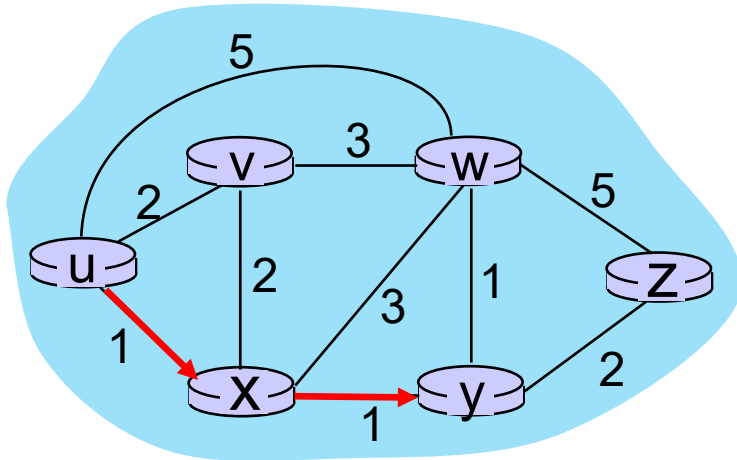
## 8 Loop

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

# Dijkstra's algorithm: an example

		<b>v</b>	<b>w</b>	<b>x</b>	<b>y</b>	<b>z</b>
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						
4						
5						



8 Loop

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

11 update  $D(b)$  for all  $b$  adjacent to  $a$  and not in  $N'$ :

$$D(b) = \min ( D(b), D(a) + c_{a,b} )$$

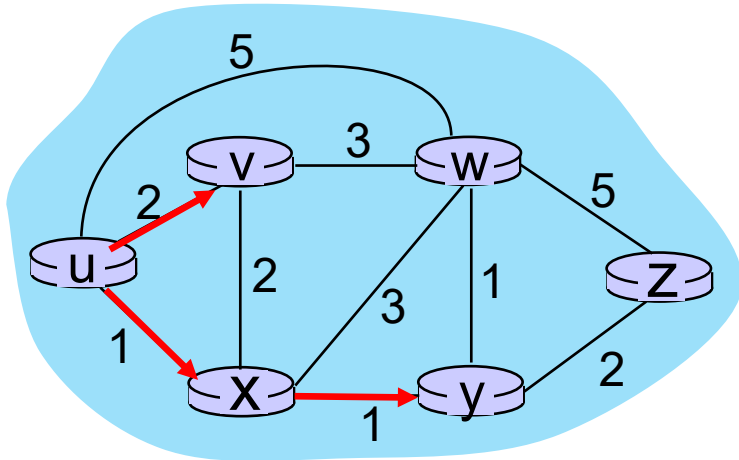
$$D(w) = \min ( D(w), D(y) + c_{y,w} ) = \min ( 4, 2+1 ) = 3$$

$$D(z) = \min ( D(z), D(y) + c_{y,z} ) = \min ( \infty, 2+2 ) = 4$$

NEW!  
NEW!

# Dijkstra's algorithm: an 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	$\infty$	$\infty$
2	uxy	2, u	3, y		4, y	
3	uxyv					
4						
5						



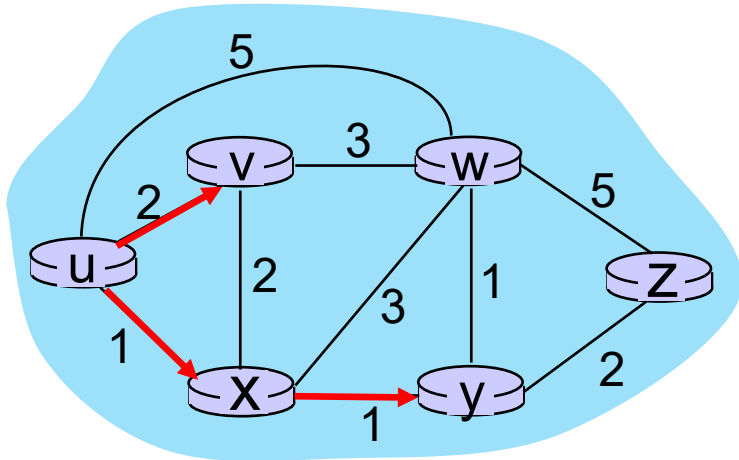
## 8 Loop

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

# Dijkstra's algorithm: an example

		<b>v</b>	<b>w</b>	<b>x</b>	<b>y</b>	<b>z</b>
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						
5						



8 *Loop*

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

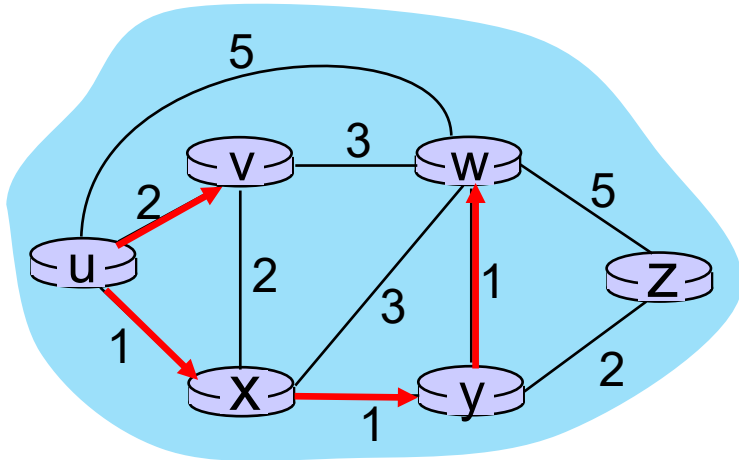
11 update  $D(b)$  for all  $b$  adjacent to  $a$  and not in  $N'$ :

$$D(b) = \min ( D(b), D(a) + c_{a,b} )$$

$$D(w) = \min ( D(w), D(v) + c_{v,w} ) = \min ( 3, 2+3 ) = 3$$

# Dijkstra's algorithm: an 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	$\infty$	$\infty$
2	uxy	2, u	3, y	4, y	$\infty$	$\infty$
3	uxyv	3, y	4, y	5, y	$\infty$	$\infty$
4	uxyvw	4, v	5, v	6, v	3, y	5, y
5	uxyvwz	5, z	6, z	7, z	4, y	5, y



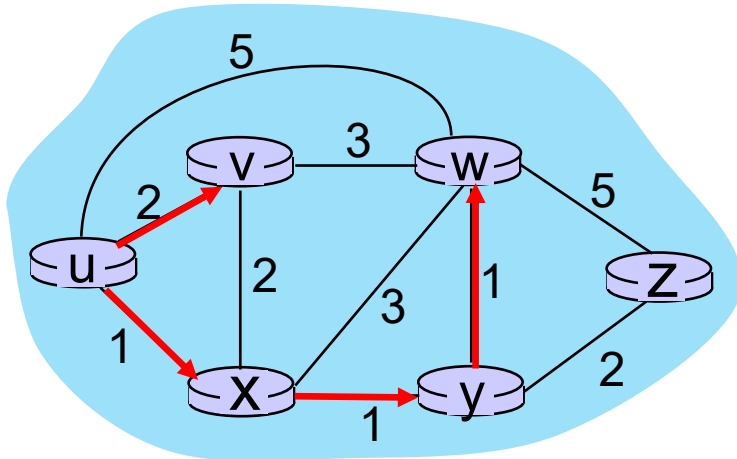
## 8 Loop

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

# Dijkstra's algorithm: an example

		<b>v</b>	<b>w</b>	<b>x</b>	<b>y</b>	<b>z</b>
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	$\infty$	$\infty$
2	uxy	2,u	3,y	2,x	4,y	$\infty$
3	uxyv	2,u	3,y	2,x	4,y	$\infty$
4	uxyvw	2,u	3,y	2,x	4,y	4,y
5						



8 Loop

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

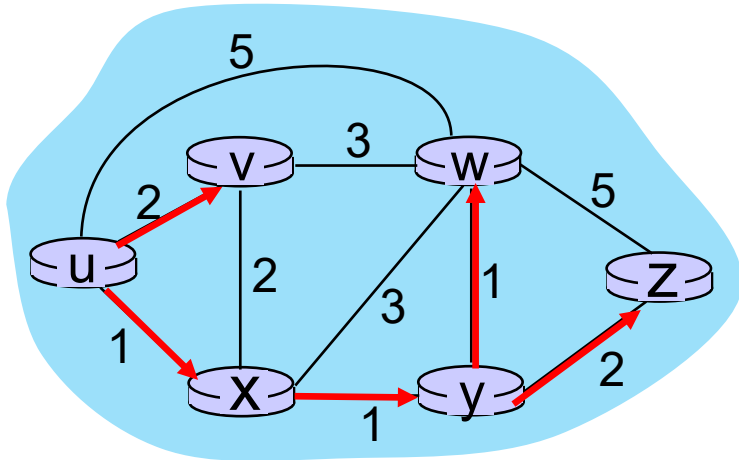
11 update  $D(b)$  for all  $b$  adjacent to  $a$  and not in  $N'$ :

$$D(b) = \min ( D(b), D(a) + c_{a,b} )$$

$$D(z) = \min ( D(z), D(w) + c_{w,z} ) = \min ( 4, 3+5 ) = 4$$

# Dijkstra's algorithm: an 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					



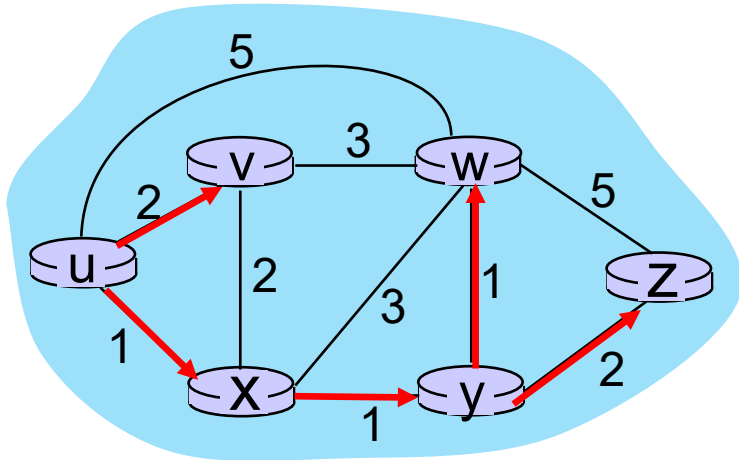
## 8 Loop

9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum

10 add  $a$  to  $N'$

# Dijkstra's algorithm: an example

		<b>v</b>	<b>w</b>	<b>x</b>	<b>y</b>	<b>z</b>
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					



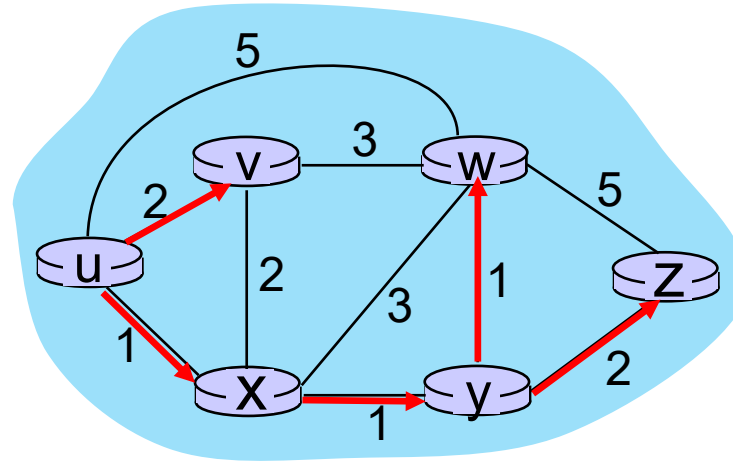
## 8 Loop

- 9 find  $a$  not in  $N'$  such that  $D(a)$  is a minimum
- 10 add  $a$  to  $N'$
- 11 update  $D(b)$  for all  $b$  adjacent to  $a$  and not in  $N'$  :  

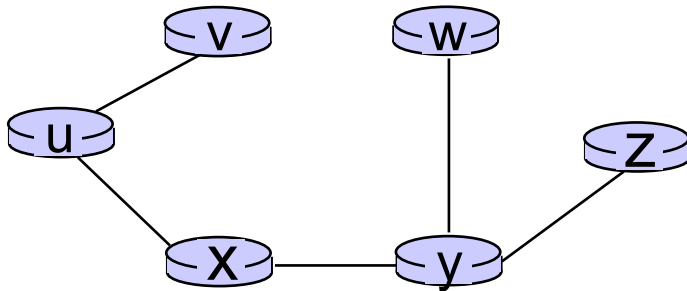
$$D(b) = \min ( D(b), D(a) + c_{a,b} )$$



# Dijkstra's algorithm: an example



resulting least-cost-path tree from u:



resulting forwarding table in u:

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

route from u to v directly

route from u to all other destinations via x

# Dijkstra's algorithm: discussion

algorithm complexity:  $n$  nodes

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

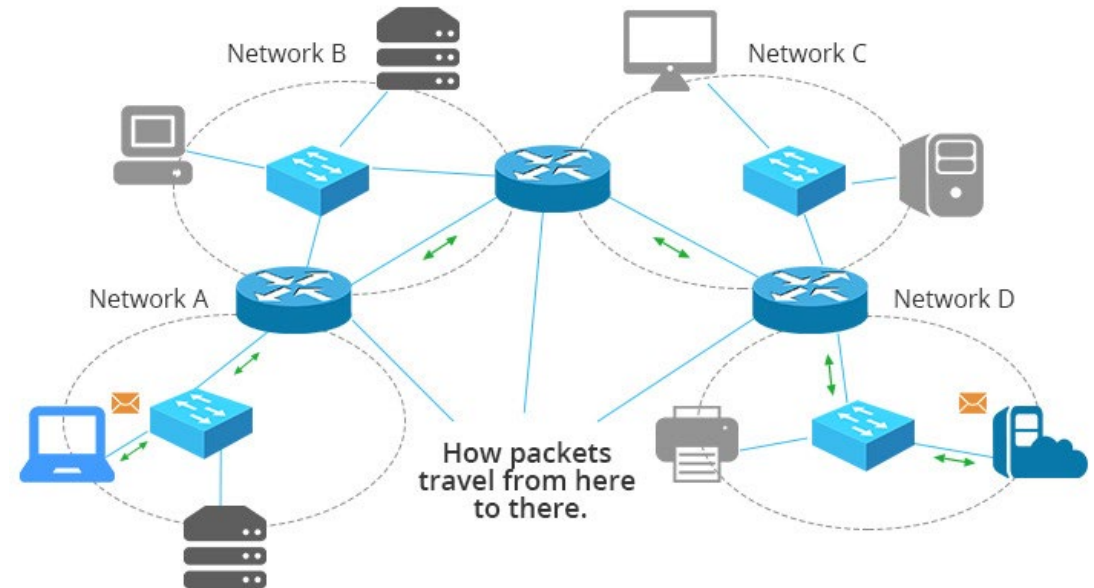
message complexity:

- each router must *broadcast* its link state information to other  $n$  routers
- efficient (and interesting!) broadcast algorithms:  $O(n)$  link crossings to disseminate a broadcast message from one source
- each router's message crosses  $O(n)$  links: overall message complexity:  $O(n^2)$

<https://www.cs.usfca.edu/~galles/visualization/Dijkstra.html>

# Network layer: “control plane” roadmap

- introduction
- routing protocols
  - link state
  - **distance vector**
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- Internet Control Message Protocol



# Distance vector algorithm

$D_x(y)$ : estimated distance from node  $x$  to  $y$

Distance vector of node  $x$ :  $\{D_x(y): \text{for all nodes } y \text{ in the network}\}$

Based on *Bellman-Ford* (BF) equation (dynamic programming):

## Bellman-Ford equation

Let  $D_x(y)$ : cost of least-cost path from  $x$  to  $y$ .

Then:

$$D_x(y) = \min_v \{ c_{x,v} + D_v(y) \}$$

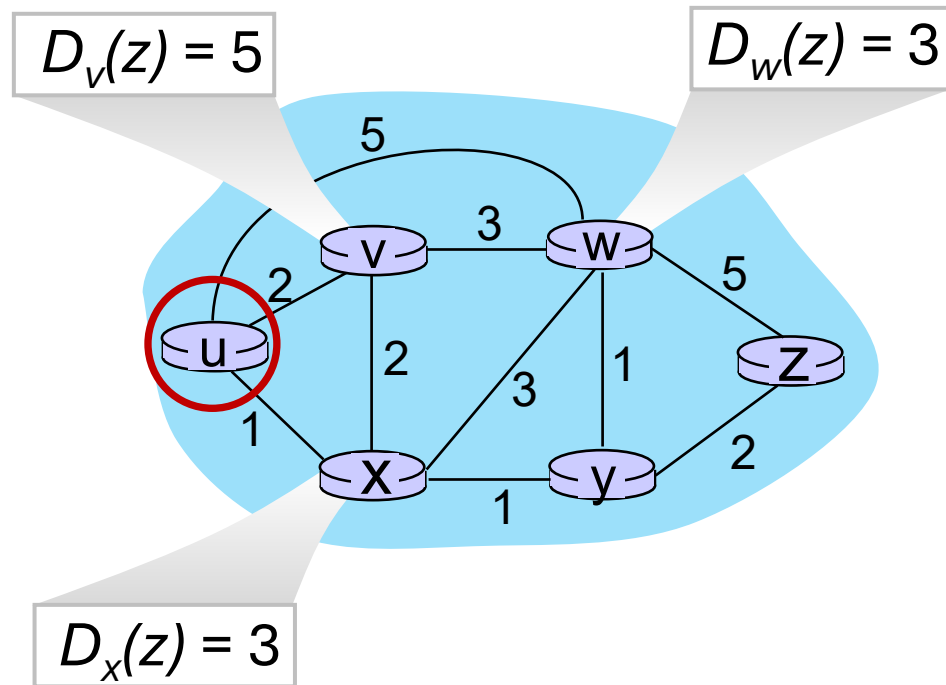
*min* taken over all neighbors  $v$  of  $x$

direct cost of link from  $x$  to  $v$

$v$ 's estimated least-cost-path cost to  $y$

# Bellman-Ford Example

Suppose that  $u$ 's neighboring nodes,  $x, v, w$ , know that for destination  $z$ :



Bellman-Ford 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 ( $x$ )  
is next hop on estimated  
least-cost path to destination  
( $z$ )*

# 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 any 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 examp



t=1

- b receives DVs from a, c, e

## DV in a:

$D_a(a)=0$   
 $D_a(b)=8$   
 $D_a(c)=\infty$   
 $D_a(d)=1$   
 $D_a(e)=\infty$   
 $D_a(f)=\infty$   
 $D_a(g)=\infty$   
 $D_a(h)=\infty$   
 $D_a(i)=\infty$

## DV in b:

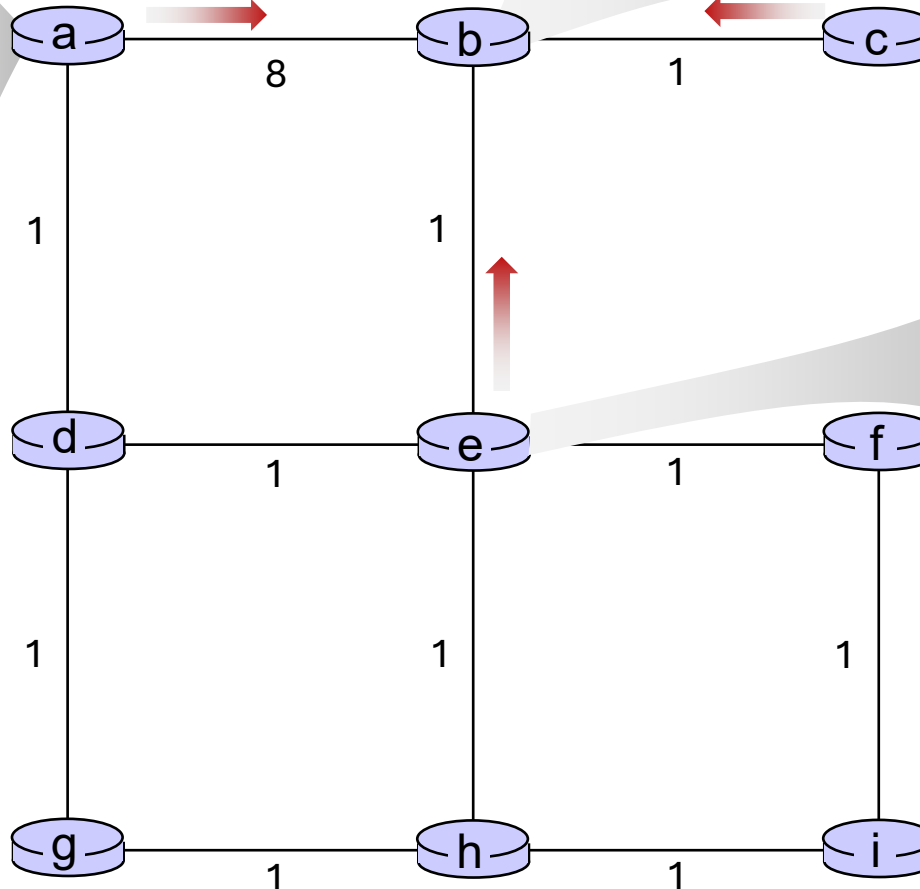
$D_b(a)=8$     $D_b(f)=\infty$   
 $D_b(c)=1$     $D_b(g)=\infty$   
 $D_b(d)=\infty$     $D_b(h)=\infty$   
 $D_b(e)=1$     $D_b(i)=\infty$

## DV in c:

$D_c(a)=\infty$   
 $D_c(b)=1$   
 $D_c(c)=0$   
 $D_c(d)=\infty$   
 $D_c(e)=\infty$   
 $D_c(f)=\infty$   
 $D_c(g)=\infty$   
 $D_c(h)=\infty$   
 $D_c(i)=\infty$

## DV in e:

$D_e(a)=\infty$   
 $D_e(b)=1$   
 $D_e(c)=\infty$   
 $D_e(d)=1$   
 $D_e(e)=0$   
 $D_e(f)=1$   
 $D_e(g)=\infty$   
 $D_e(h)=1$   
 $D_e(i)=\infty$



# DV Example1:

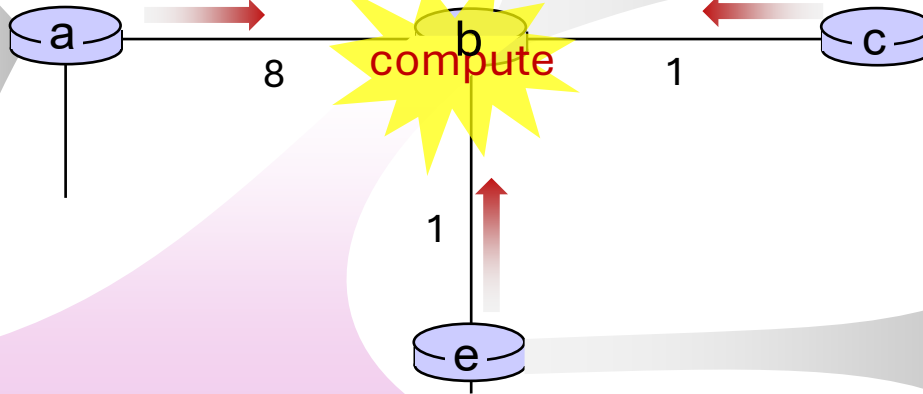


t=1

- b receives DVs from a, c, e, computes:

$$\begin{aligned}
 D_b(a) &= \min\{c_{b,a} + D_a(a), c_{b,c} + D_c(a), c_{b,e} + D_e(a)\} = \min\{8, \infty, \infty\} = 8 \\
 D_b(c) &= \min\{c_{b,a} + D_a(c), c_{b,c} + D_c(c), c_{b,e} + D_e(c)\} = \min\{\infty, 1, \infty\} = 1 \\
 D_b(d) &= \min\{c_{b,a} + D_a(d), c_{b,c} + D_c(d), c_{b,e} + D_e(d)\} = \min\{9, 2, \infty\} = 2 \\
 D_b(e) &= \min\{c_{b,a} + D_a(e), c_{b,c} + D_c(e), c_{b,e} + D_e(e)\} = \min\{\infty, \infty, 1\} = 1 \\
 D_b(f) &= \min\{c_{b,a} + D_a(f), c_{b,c} + D_c(f), c_{b,e} + D_e(f)\} = \min\{\infty, \infty, 2\} = 2 \\
 D_b(g) &= \min\{c_{b,a} + D_a(g), c_{b,c} + D_c(g), c_{b,e} + D_e(g)\} = \min\{\infty, \infty, \infty\} = \infty \\
 D_b(h) &= \min\{c_{b,a} + D_a(h), c_{b,c} + D_c(h), c_{b,e} + D_e(h)\} = \min\{\infty, \infty, 2\} = 2 \\
 D_b(i) &= \min\{c_{b,a} + D_a(i), c_{b,c} + D_c(i), c_{b,e} + D_e(i)\} = \min\{\infty, \infty, \infty\} = \infty
 \end{aligned}$$

DV in a:
$D_a(a) = 0$
$D_a(b) = 8$
$D_a(c) = \infty$
$D_a(d) = 1$
$D_a(e) = \infty$
$D_a(f) = \infty$
$D_a(g) = \infty$
$D_a(h) = \infty$
$D_a(i) = \infty$



DV in b:	
$D_b(a) = 8$	$D_b(f) = \infty$
$D_b(c) = 1$	$D_b(g) = \infty$
$D_b(d) = \infty$	$D_b(h) = \infty$
$D_b(e) = 1$	$D_b(i) = \infty$

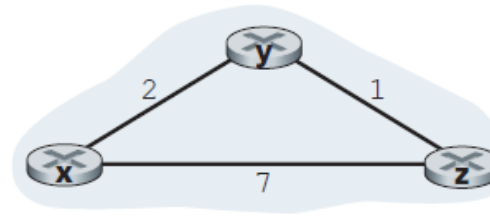
DV in c:
$D_c(a) = \infty$
$D_c(b) = 1$
$D_c(c) = 0$
$D_c(d) = \infty$
$D_c(e) = \infty$
$D_c(f) = \infty$
$D_c(g) = \infty$
$D_c(h) = \infty$
$D_c(i) = \infty$

DV in e:
$D_e(a) = \infty$
$D_e(b) = 1$
$D_e(c) = \infty$
$D_e(d) = 1$
$D_e(e) = 0$
$D_e(f) = 1$
$D_e(g) = \infty$
$D_e(h) = 1$
$D_e(i) = \infty$

DV in b:	
$D_b(a) = 8$	$D_b(f) = 2$
$D_b(c) = 1$	$D_b(g) = \infty$
$D_b(d) = 2$	$D_b(h) = 2$
$D_b(e) = 1$	$D_b(i) = \infty$



# DV Example 2:



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

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

Node y table

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

		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	3
	y	2	0	1
	z	3	1	0

Node z table

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

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

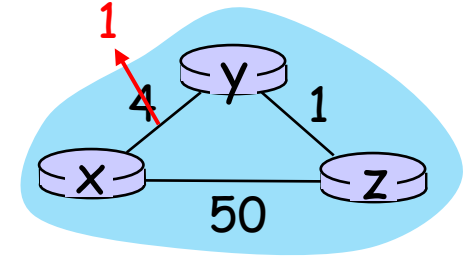
		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	3	1	0

.....> Time

# Distance vector: link cost changes

## link cost changes:

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



“good news  
travels fast”

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

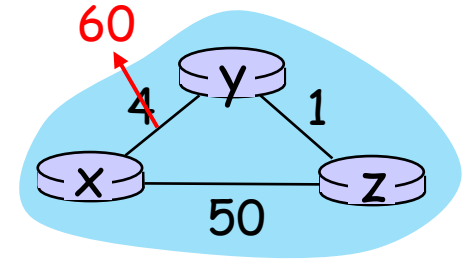
$t_1$ : z receives update from y, updates its DV, computes new least cost to x, sends its neighbors its DV.

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

# Distance vector: count-to-infinity problem

## link cost changes:

- node detects local link cost change
- “**bad news travels slow**” – count-to-infinity problem:
  - y sees direct link to x has new cost 60, but z has said it has a path at cost of 5. So y computes “my new cost to x will be 6, via z); notifies z of new cost of 6 to x.
  - z learns that path to x via y has new cost 6, so z computes “my new cost to x will be 7 via y), notifies y of new cost of 7 to x.
  - y learns that path to x via z has new cost 7, so y computes “my new cost to x will be 8 via y), notifies z of new cost of 8 to x.
  - z learns that path to x via y has new cost 8, so z computes “my new cost to x will be 9 via y), notifies y of new cost of 9 to x.
  - ...
- **Poisoned reverse:** Z tells Y, its distance to X is infinity



# Making routing scalable

our routing study thus far - idealized

- all routers identical
- network “flat”

... not true in practice

**scale:** billions of destinations:

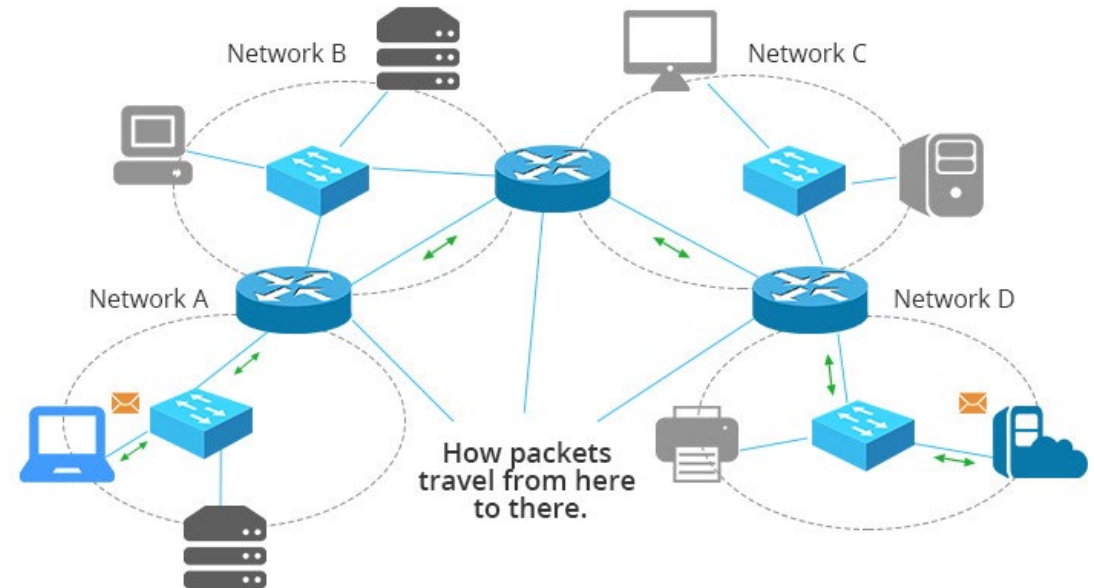
- can't store all destinations in routing tables!
- routing table exchange would swamp links!

**administrative autonomy:**

- Internet: a network of networks
- each network admin may want to control routing in its own network

# Network layer: “control plane” roadmap

- introduction
- routing protocols
  - link state
  - distance vector
- **intra-ISP routing: OSPF**
- routing among ISPs: BGP
- Internet Control Message Protocol



# Internet approach to scalable routing

aggregate routers into regions known as “autonomous systems” (AS) (a.k.a. “domains”)

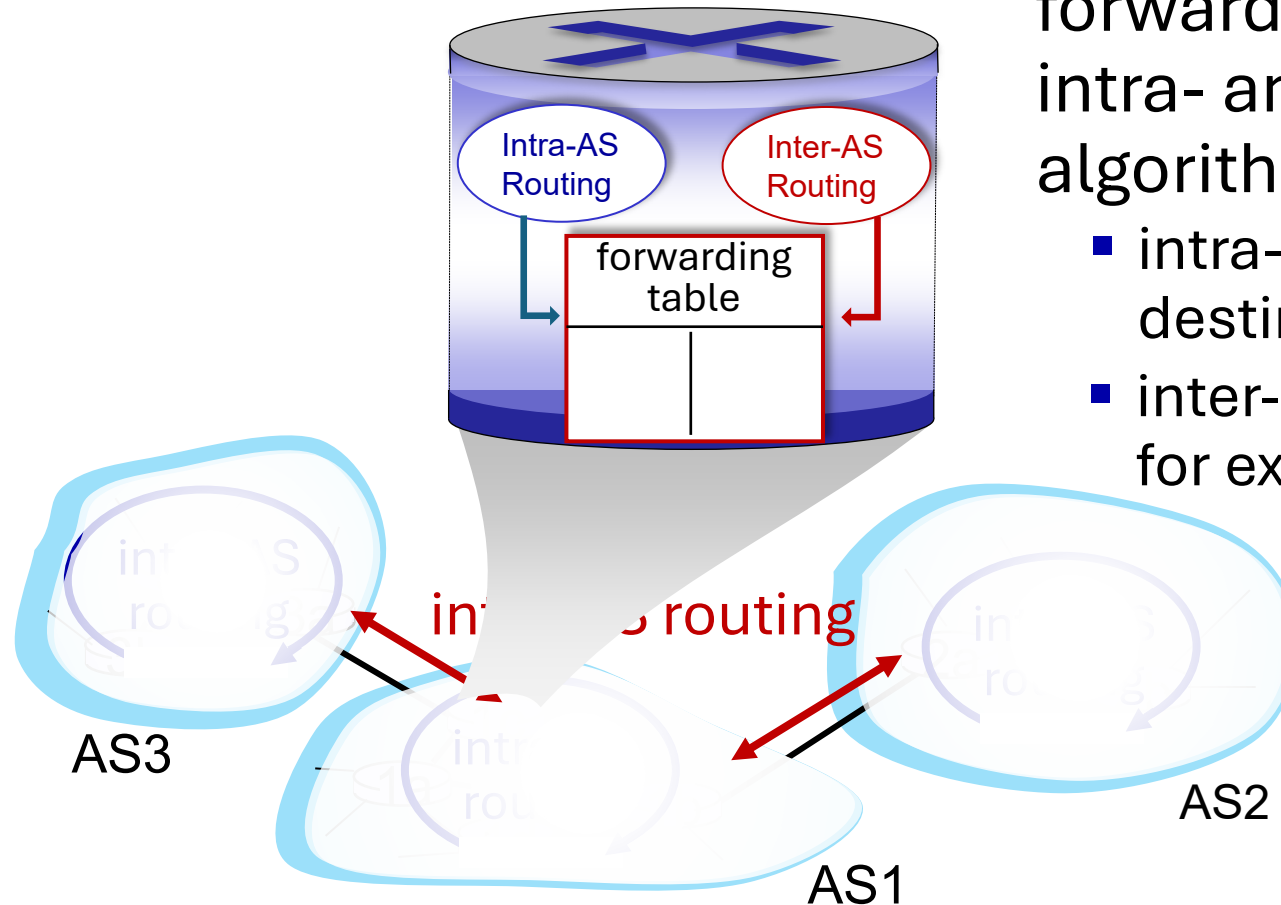
**intra-AS (aka “intra-domain”):**  
routing among routers *within same AS (“network”)*

- all routers in AS must run same intra-domain protocol
- routers in different AS can run different intra-domain routing protocols
- **gateway router:** at “edge” of its own AS, has link(s) to router(s) in other AS'es

**inter-AS (aka “inter-domain”):** routing *among AS'es*

- gateways perform inter-domain routing (as well as intra-domain routing)

# Interconnected ASes



forwarding table configured by intra- and inter-AS routing algorithms

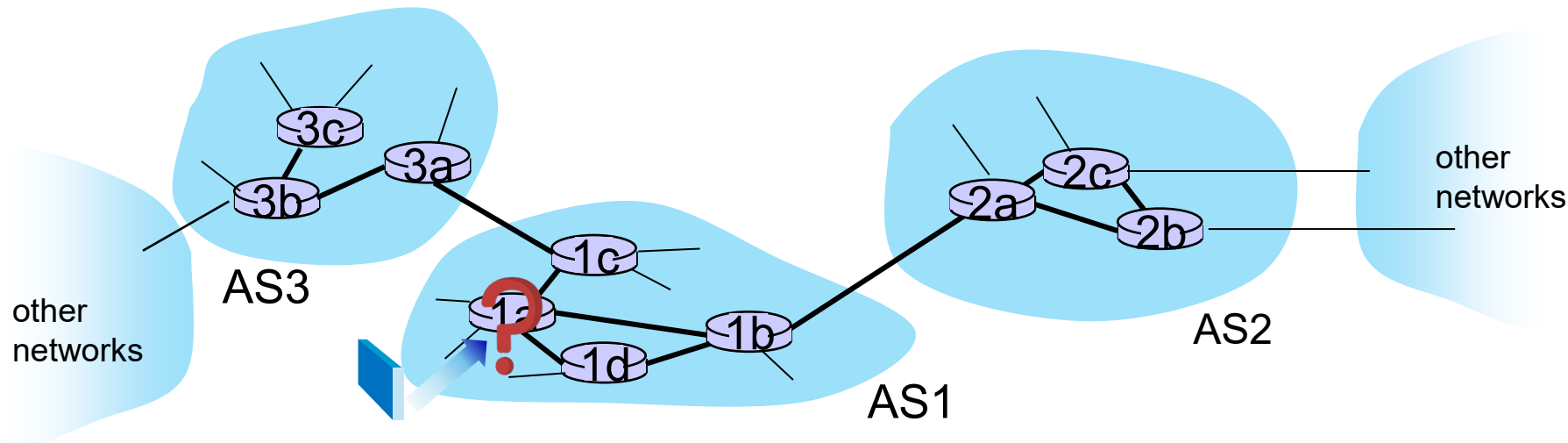
- intra-AS routing determine entries for destinations within AS
- inter-AS & intra-AS determine entries for external destinations

# Inter-AS routing: a role in intradomain forwarding

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

## AS1 inter-domain routing must:

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





# Intra-AS routing: routing within an AS

most common intra-AS routing protocols:

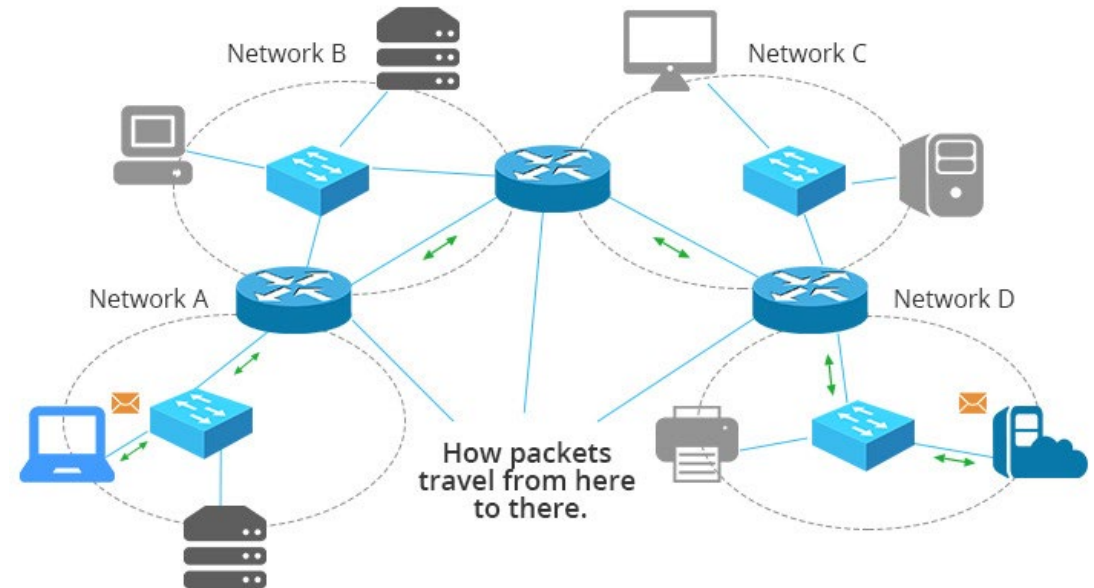
- **RIP: Routing Information Protocol** [RFC 1723]
  - classic DV: DVs exchanged every 30 secs
  - no longer widely used
- **EIGRP: Enhanced Interior Gateway Routing Protocol**
  - DV based
  - formerly Cisco-proprietary for decades (became open in 2013 [RFC 7868])
- **OSPF: Open Shortest Path First** [RFC 2328]
  - link-state routing
  - IS-IS protocol (ISO standard, not RFC standard) essentially same as OSPF

# OSPF (Open Shortest Path First) routing

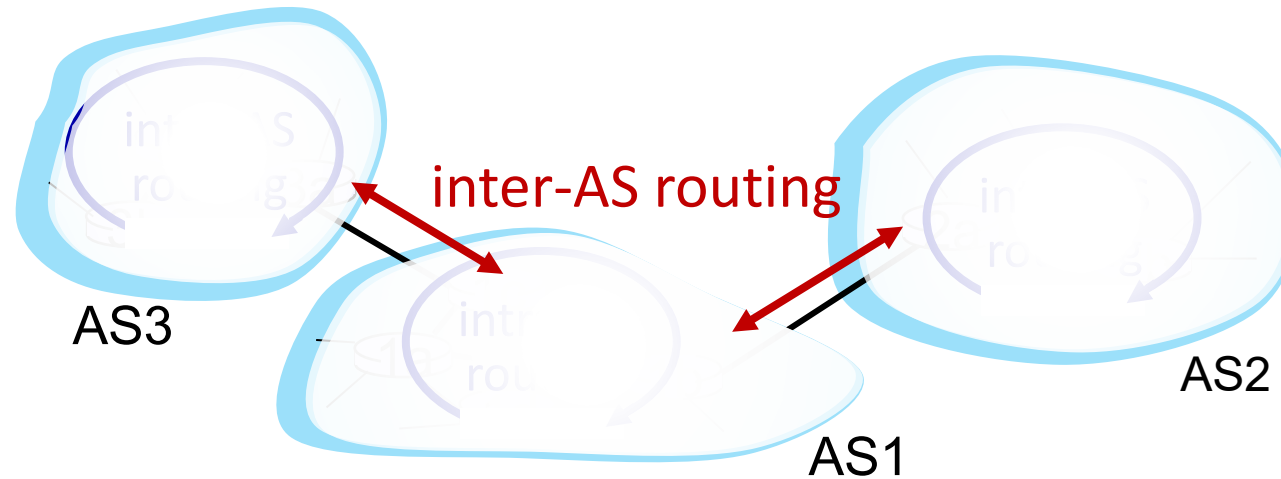
- “open”: publicly available
- classic link-state
  - each router floods OSPF link-state advertisements (directly over IP rather than using TCP/UDP) to all other routers in entire AS
  - multiple link costs metrics possible: bandwidth, delay
  - each router has full topology, uses Dijkstra’s algorithm to compute forwarding table
- *security*: all OSPF messages authenticated (to prevent malicious intrusion)

# Network layer: “control plane” roadmap

- introduction
- routing protocols
  - link state
  - distance vector
- intra-ISP routing: OSPF
- **routing among ISPs: BGP**
- Internet Control Message Protocol



# Interconnected ASes



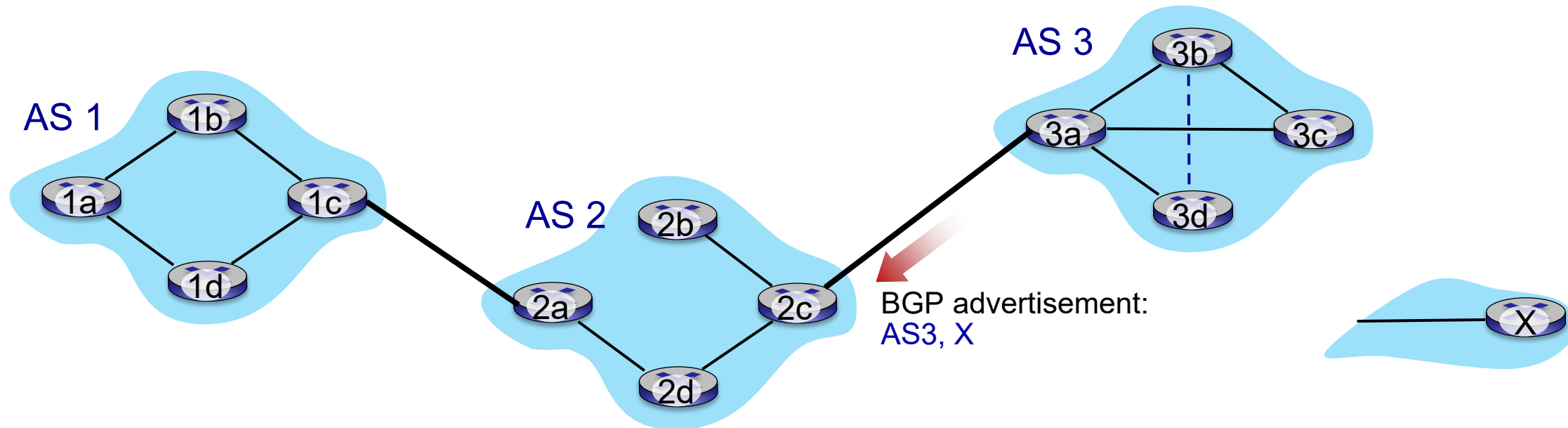
- ✓ **intra-AS (aka “intra-domain”)**: routing among routers *within same* AS (“*network*”)
- ➡ **inter-AS (aka “inter-domain”)**: routing *among* AS'es

# Internet inter-AS routing: BGP

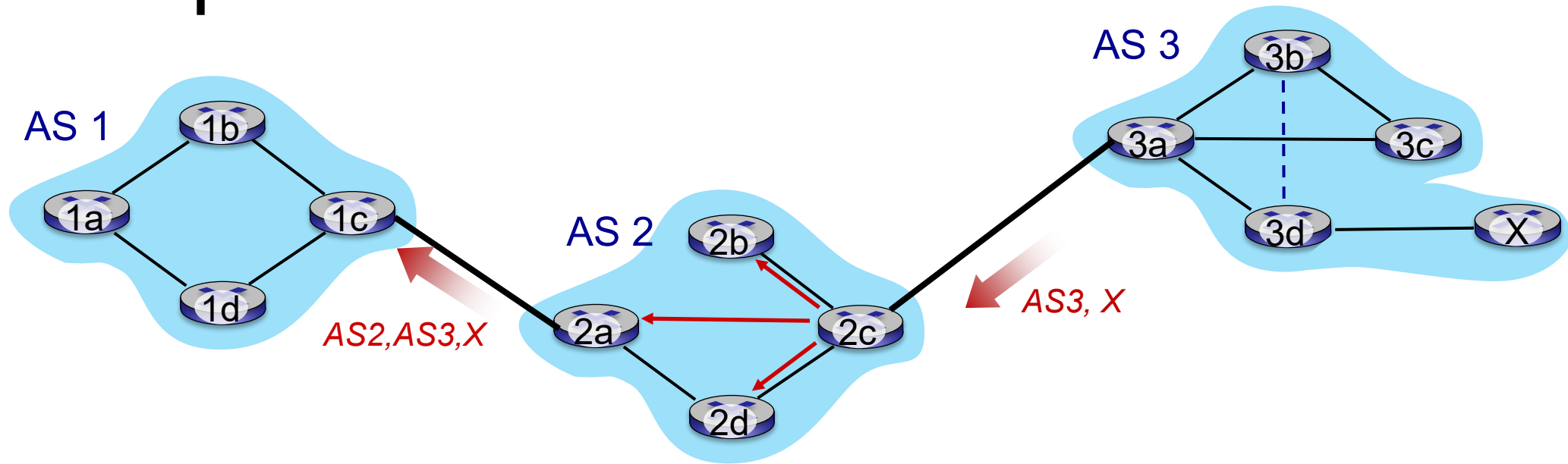
- **BGP (Border Gateway Protocol):** *the de facto inter-domain routing protocol*
  - “glue that holds the Internet together”
- allows subnet to advertise its existence, and the destinations it can reach, to rest of Internet: *“I am here, here is who I can reach, and how”*
- BGP provides each AS a means to:
  - obtain destination network reachability info from neighboring ASes (**eBGP**)
  - determine routes to other networks based on reachability information and *policy*
  - propagate reachability information to all AS-internal routers (**iBGP**)
  - **advertise** (to neighboring networks) destination reachability info

# BGP basics

- **BGP session:** two BGP routers (“peers”) exchange BGP messages over semi-permanent TCP connection:
  - advertising *paths* to different destination network prefixes (BGP is a “path vector” protocol)
- when AS3 gateway 3a advertises *path AS3,X* to AS2 gateway 2c:
  - AS3 *promises* to AS2 it will forward datagrams towards X

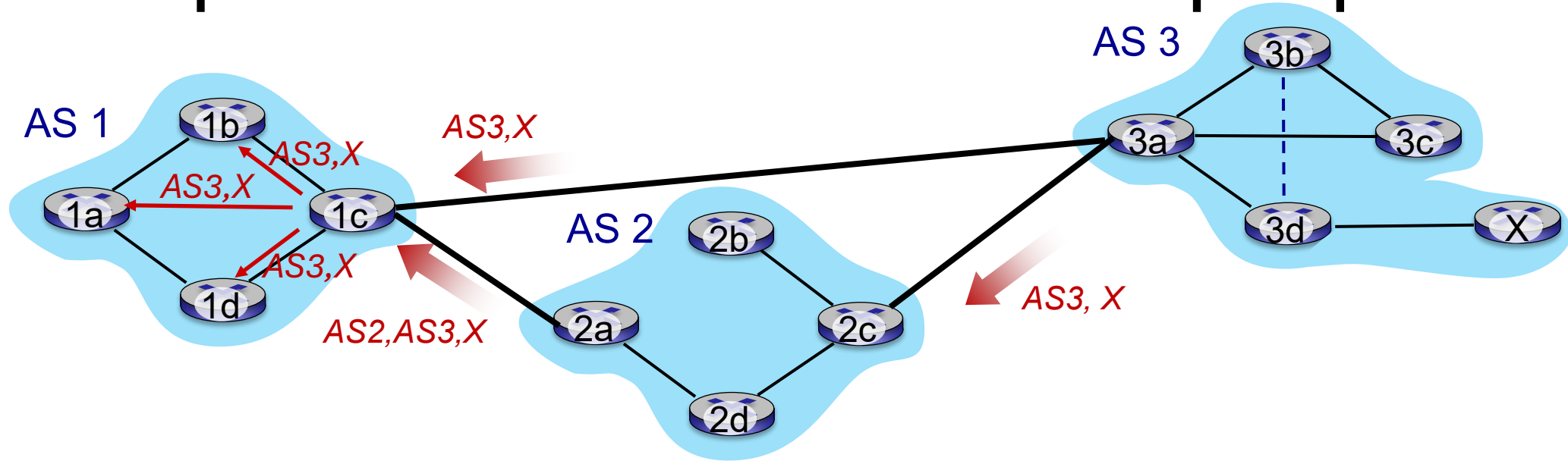


# BGP path advertisement



- AS2 router 2c receives path advertisement **AS3, X** (via eBGP) from AS3 router 3a
- based on AS2 policy, AS2 router 2c accepts path AS3, X, propagates (via iBGP) to all AS2 routers
- based on AS2 policy, AS2 router 2a advertises (via eBGP) path **AS2, AS3, X** to AS1 router 1c

# BGP path advertisement: multiple paths

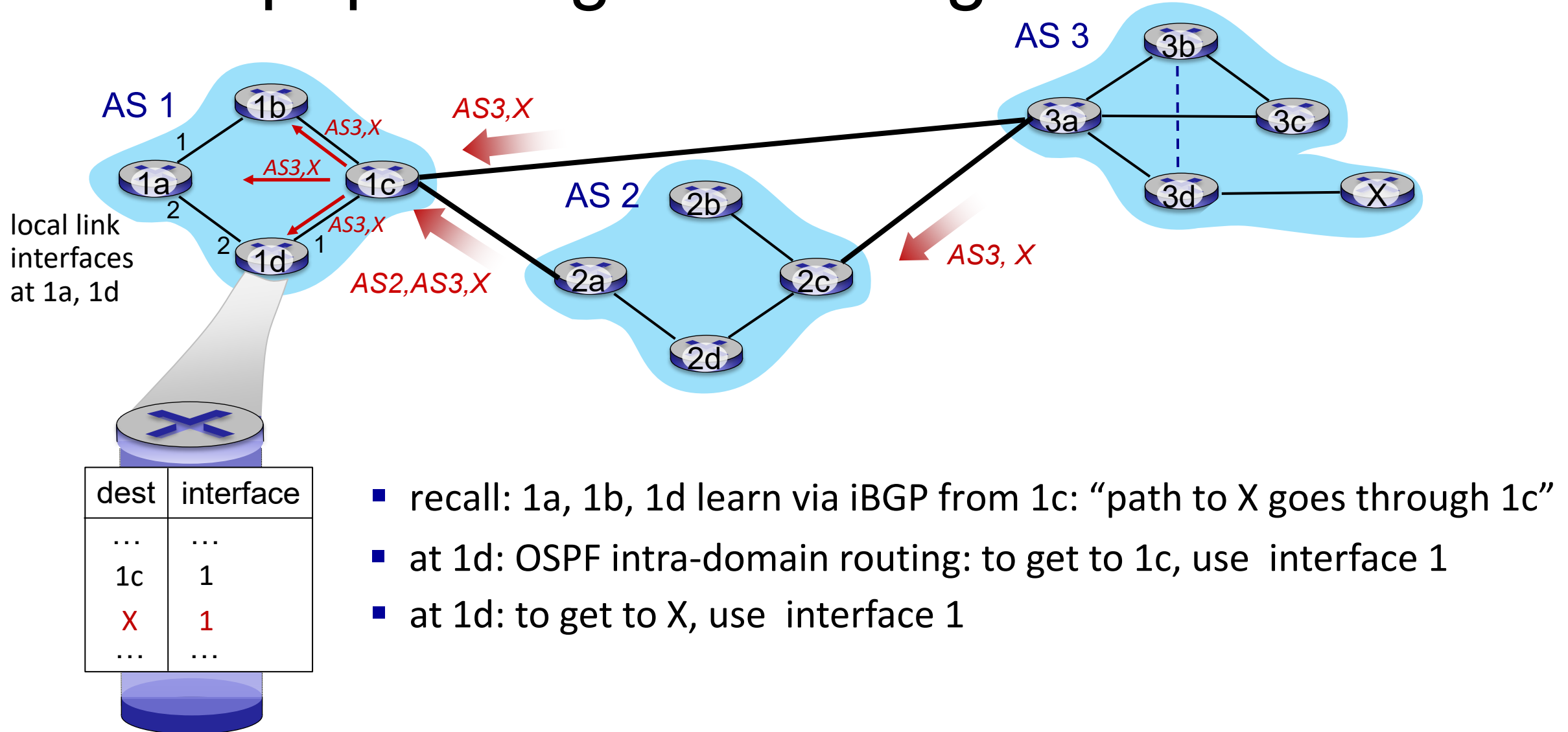


gateway router may learn about **multiple** paths to destination:

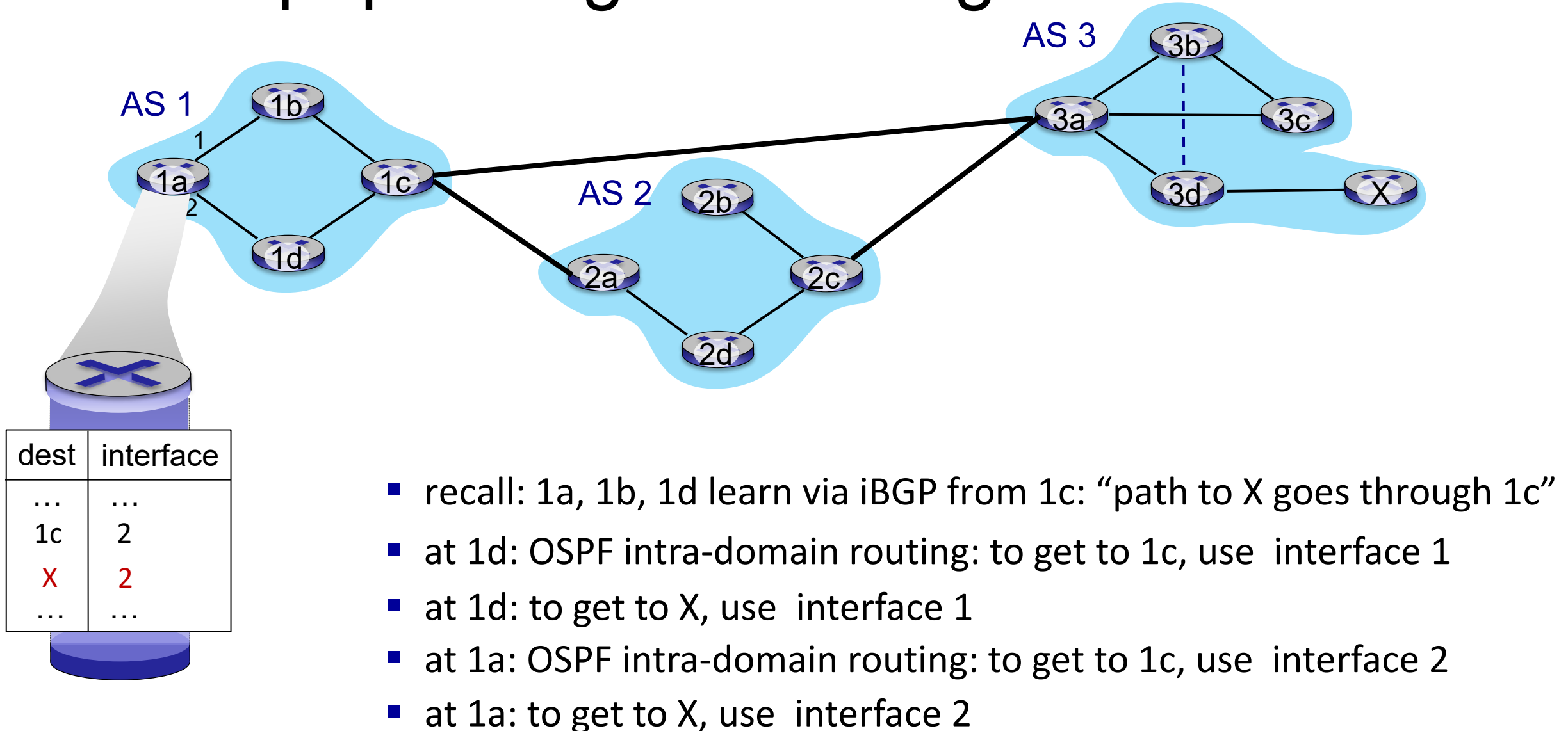
- AS1 gateway router 1c learns path **AS2,AS3,X** from 2a
- AS1 gateway router 1c learns path **AS3,X** from 3a
- based on **policy**, AS1 gateway router 1c chooses path **AS3,X** and advertises path within AS1 via iBGP



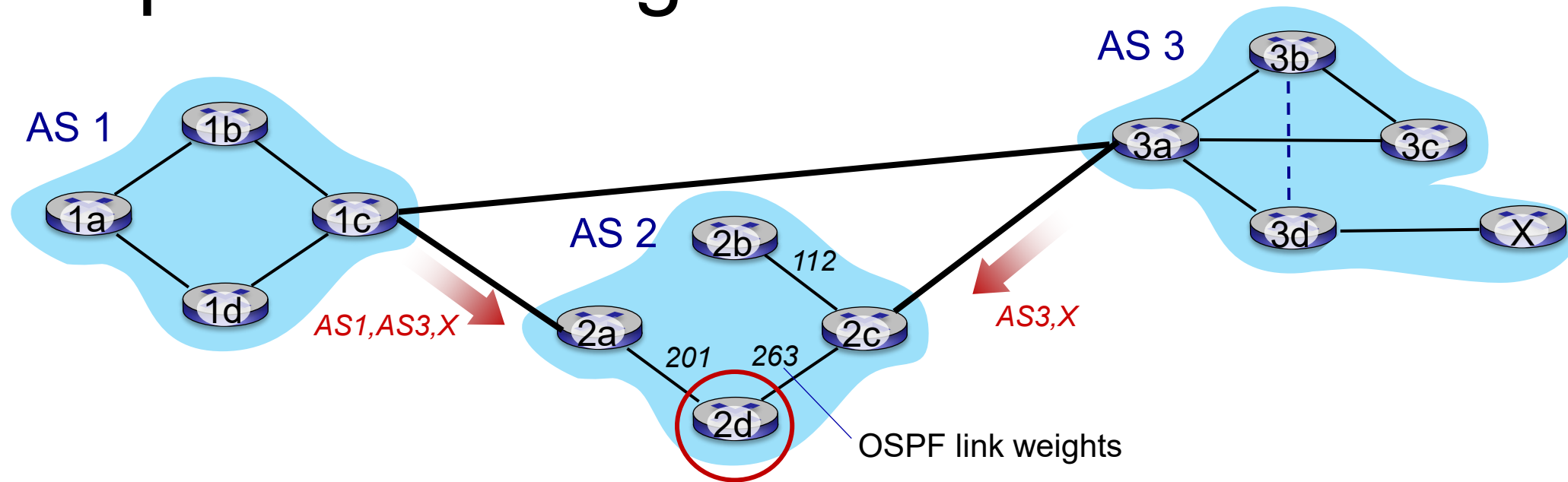
# BGP: populating forwarding tables



# BGP: populating forwarding tables



# Hot potato routing



- 2d learns (via iBGP) it can route to X via 2a or 2c
- **hot potato routing**: choose local gateway that has least *intra-domain* cost (e.g., 2d chooses 2a, even though more AS hops to X): don't worry about inter-domain cost!

# BGP route selection

- router may learn about more than one route to destination AS, selects route based on:
  1. closest NEXT-HOP router: hot potato routing
  2. local preference value attribute: policy decision
  3. shortest AS-PATH
  4. additional criteria

# Why different Intra-, Inter-AS routing ?

## policy:

- inter-AS: admin wants control over how its traffic routed, who routes through its network
- intra-AS: single admin, so policy less of an issue

## scale:

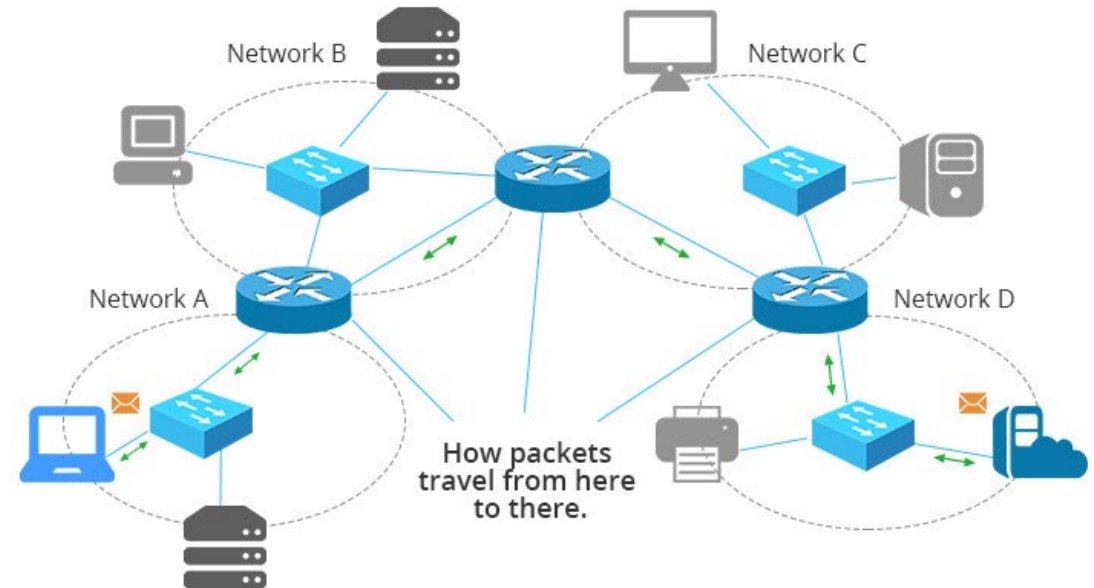
- hierarchical routing saves table size, reduced update traffic

## performance:

- intra-AS: can focus on performance
- inter-AS: policy dominates over performance

# Network layer: “control plane” roadmap

- introduction
- routing protocols
  - link state
  - distance vector
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- Internet Control Message Protocol

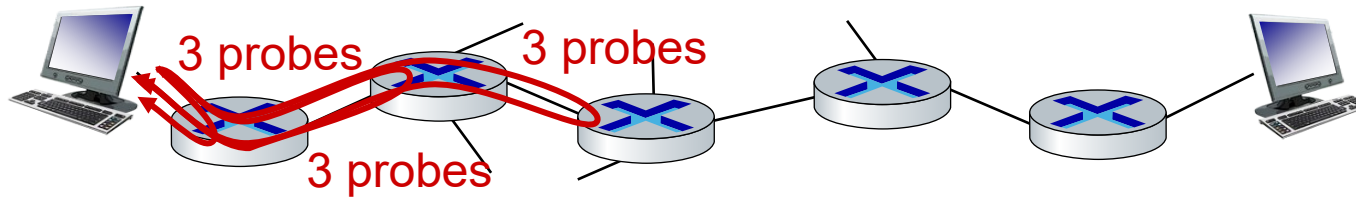


# ICMP: internet control message protocol

- used by hosts and routers to communicate network-level information
  - error reporting: unreachable host, network, port, protocol
  - echo request/reply (used by ping)
- network-layer “above” IP:
  - ICMP messages 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 sets of UDP segments to destination
    - 1<sup>st</sup> set has TTL =1, 2<sup>nd</sup> set has TTL=2, etc.
  - datagram in  $n$ th set arrives to  $n$ th router:
    - router discards datagram and sends source ICMP message (type 11, code 0)
    - ICMP message possibly includes name of router & IP address
  - when ICMP message arrives at source: record RTTs
- stopping criteria:
- UDP segment eventually arrives at destination host
  - destination returns ICMP “port unreachable” message (type 3, code 3)
  - source stops