

# Chapter 5 - Network layer: “control plane” roadmap

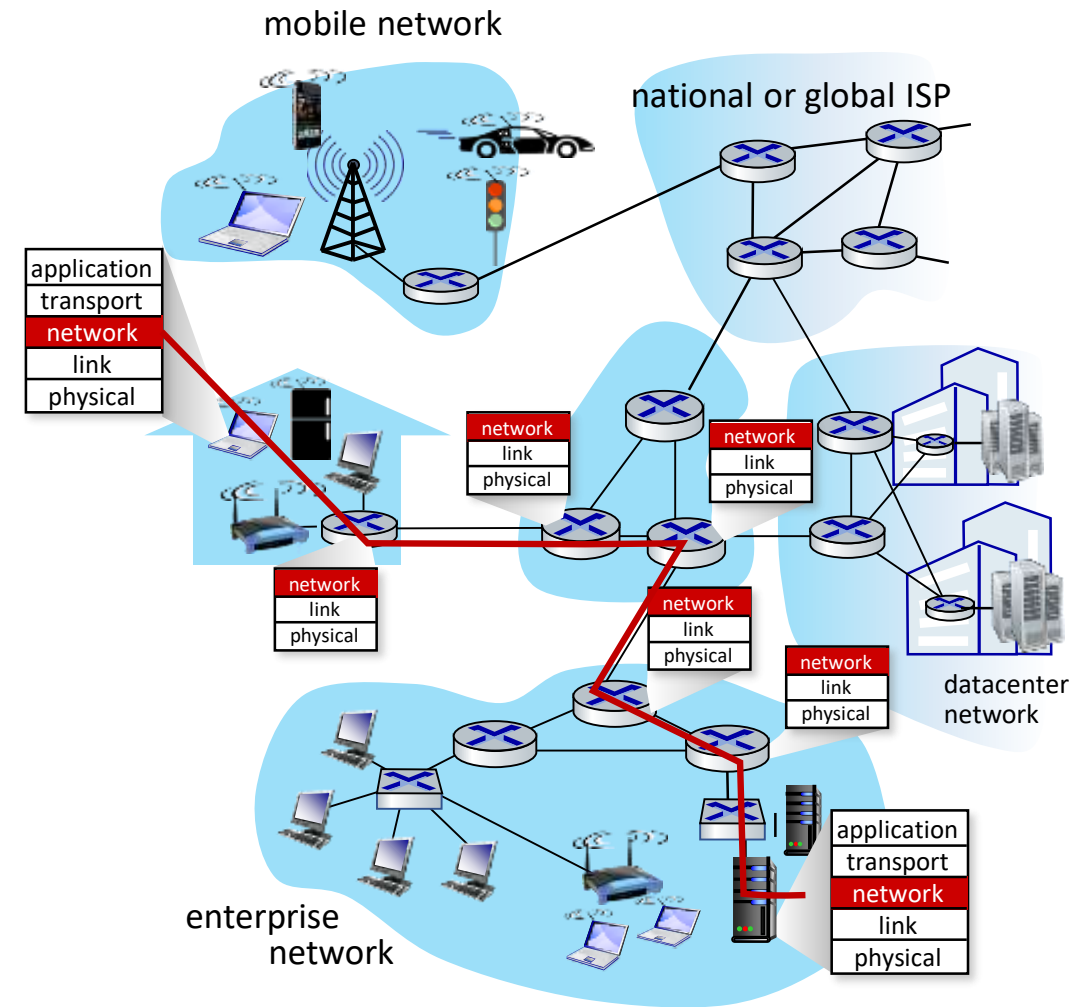
- Introduction
- routing protocols
  - link state
  - distance vector
- intra-ISP routing: OSPF
- routing among ISPs: BGP



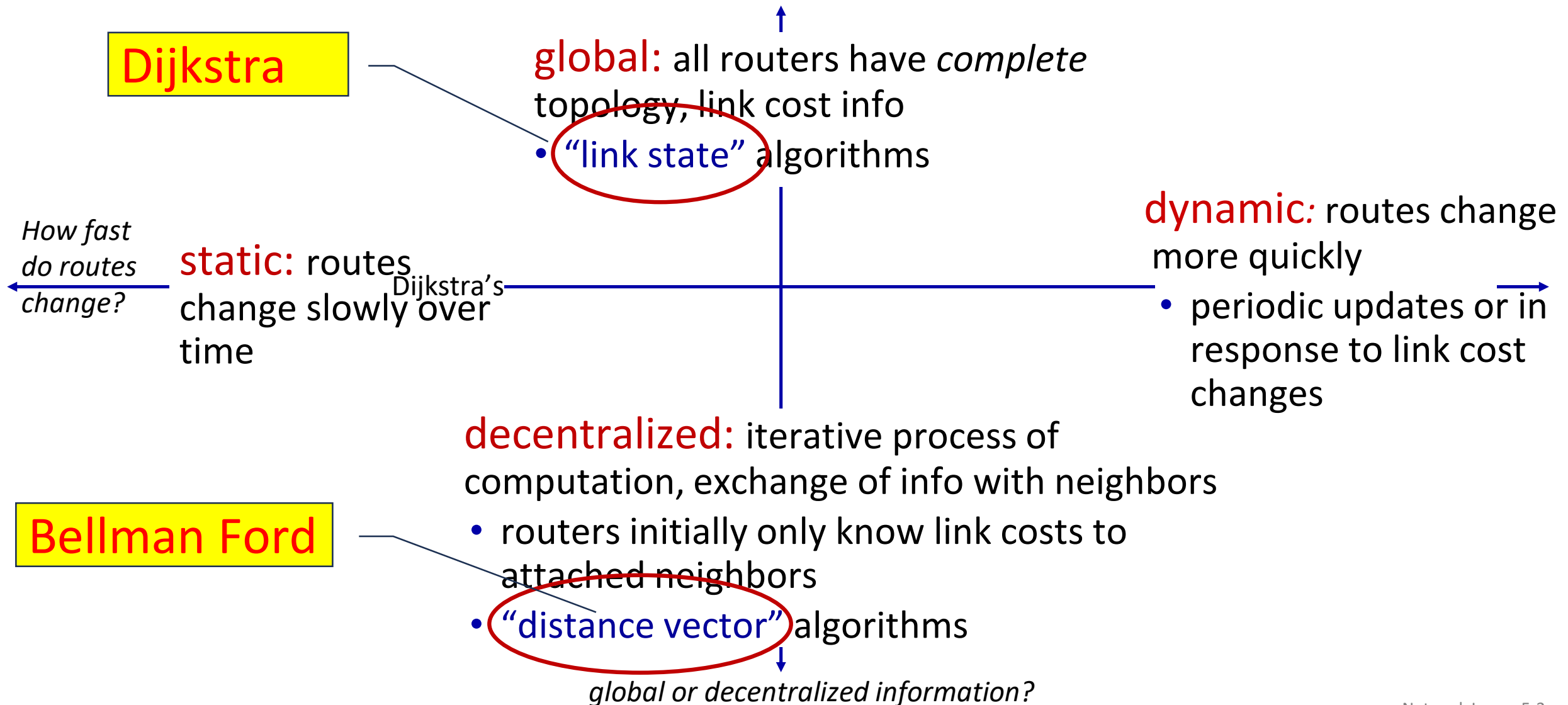
# Routing protocols

**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



# link-state routing - Dijkstra's algorithm

- Greedy Algorithm
- computes least cost paths from one node ("source") to all other nodes
  - gives *forwarding table* for that node

## centralized

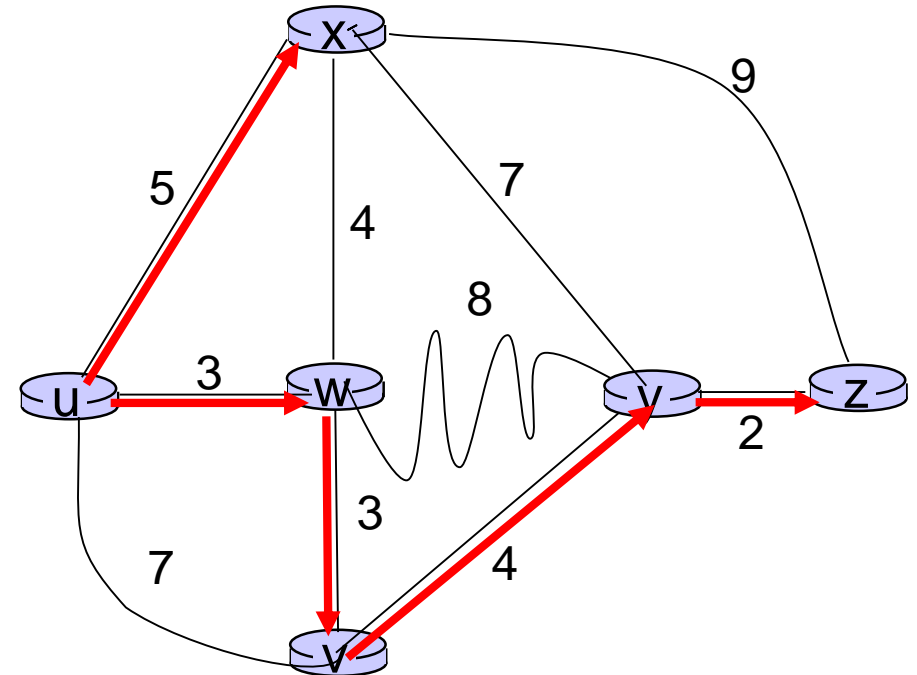
- network topology, link costs known to *all* nodes
  - accomplished via "link state broadcast"
  - all nodes have same info

## iterative

- after  $k$  iterations, know least cost path to  $k$  destinations

# Dijkstra's algorithm: example

Step	$N'$	$D(v), p(v)$	$D(w), p(w)$	$D(x), p(x)$	$D(y), p(y)$	$D(z), p(z)$
0	u	7, u	3, u	5, u	$\infty$	$\infty$
1	uw	6, w		5, u	11, w	$\infty$
2	uwx	6, w			11, w	14, x
3	uwxv				10, v	14, x
4	uwxvy					12, y
5	uwxvyz					



## notes:

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

# 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:  $O(n^2)$

# Distance vector routing – Bellman Ford algorithm

dynamic programming

Bellman-Ford equation

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

Then:

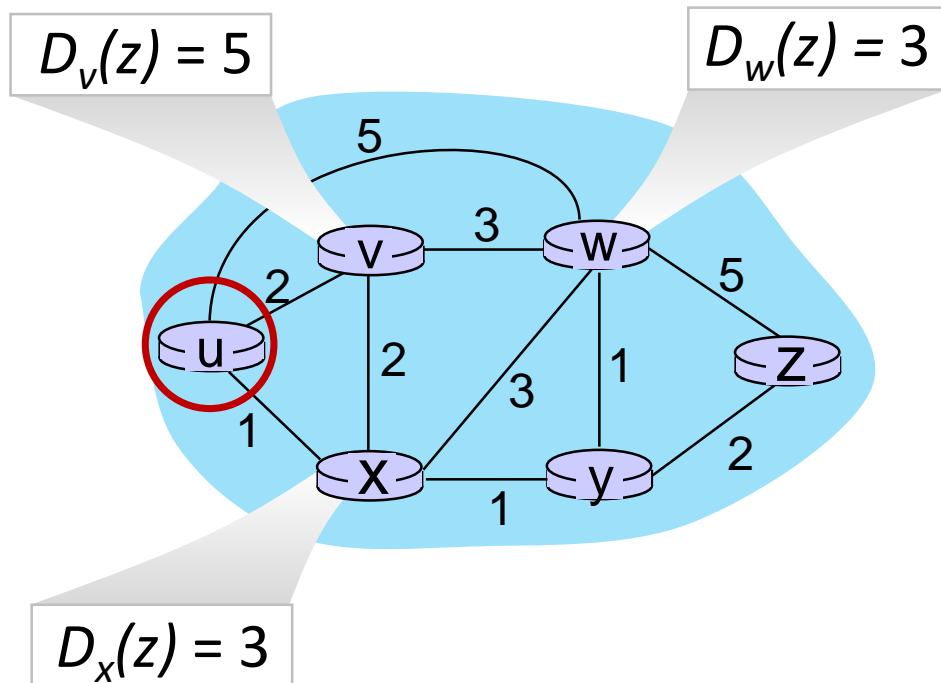
$$D_x(y) = \min_v \{ c_{x,v} + D_v(y) \}$$

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

$v$ 's estimated least-cost-path cost to  $y$   
direct cost of link from  $x$  to  $v$

# Bellman-Ford Example

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



Bellman-Ford equation says:

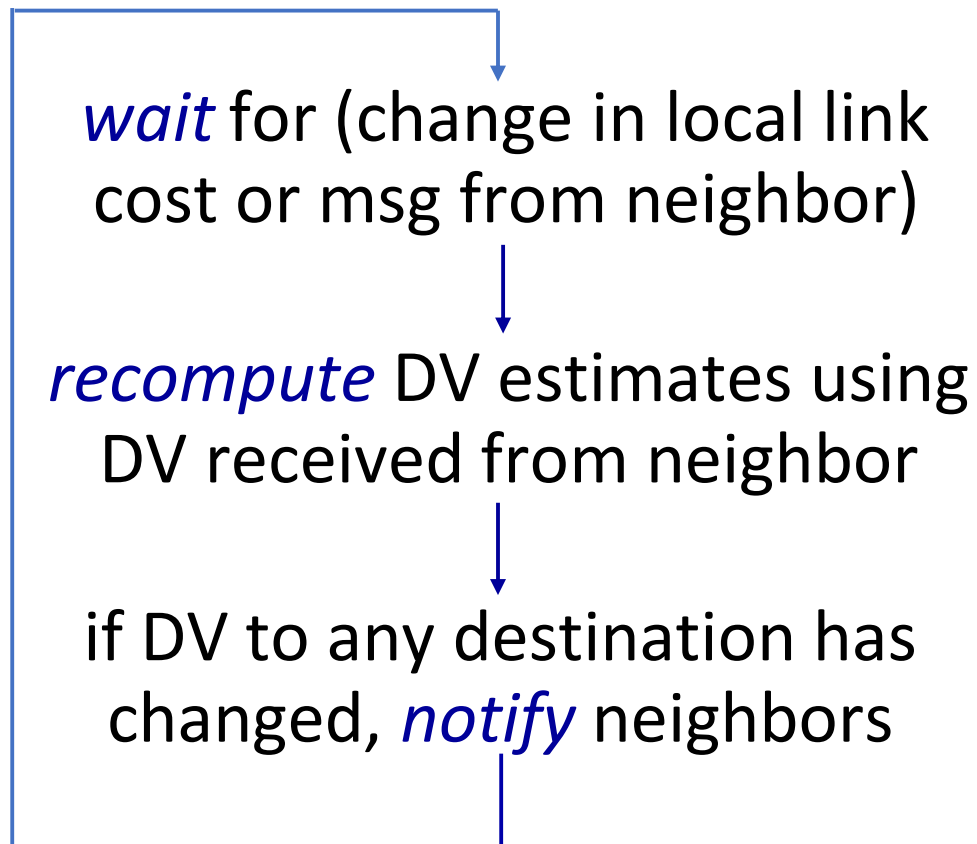
$$\begin{aligned} D_u(z) &= \min \{ c_{u,v} + D_v(z), \\ &\quad c_{u,x} + D_x(z), \\ &\quad c_{u,w} + D_w(z) \} \\ &= \min \{ 2 + 5, \\ &\quad 1 + 3, \\ &\quad 5 + 3 \} = 4 \end{aligned}$$

*node achieving minimum (x) is next hop on estimated least-cost path to destination (z)*



# Distance vector algorithm:

each node:



**iterative, asynchronous:** each local iteration caused by:

- local link cost change
- DV update message from neighbor

**distributed, self-stopping:** each node notifies neighbors *only* when its DV changes

- neighbors then notify their neighbors – *only if necessary*
- no notification received, no actions taken!

# Distance vector example



**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:

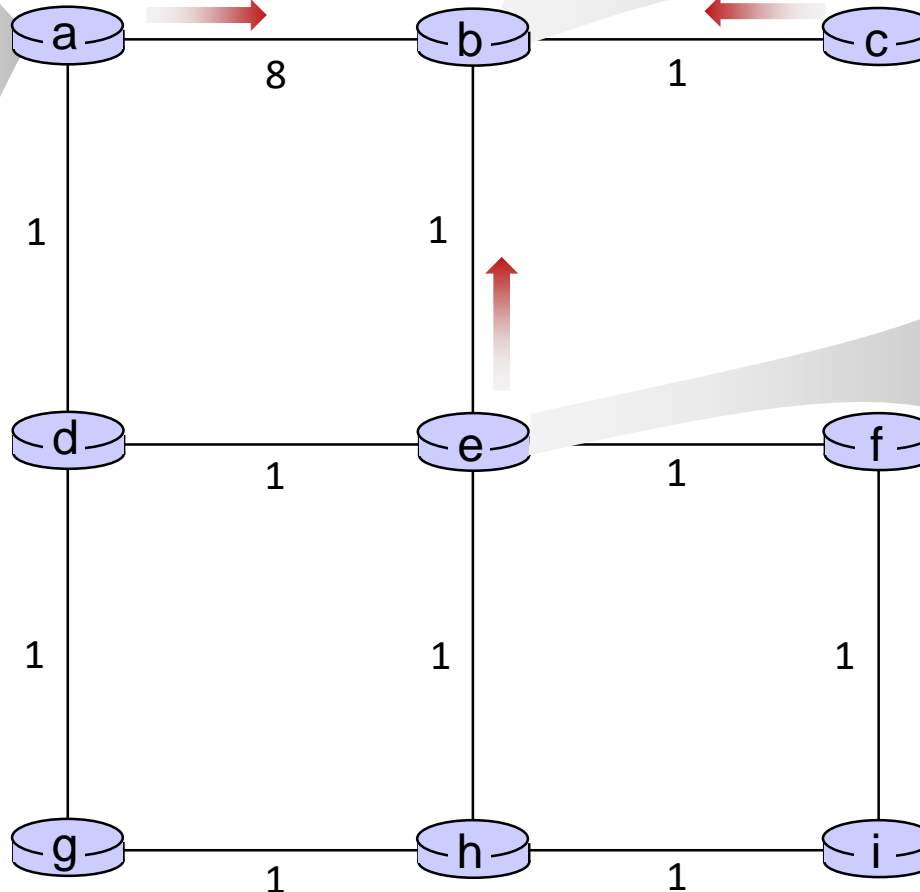
$D_b(a)=8$     $D_b(f)=\infty$   
 $D_b(c)=1$     $D_b(g)=\infty$   
 $D_b(d)=\infty$     $D_b(h)=\infty$   
 $D_b(e)=1$     $D_b(i)=\infty$

## 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:

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



# Distance vector example

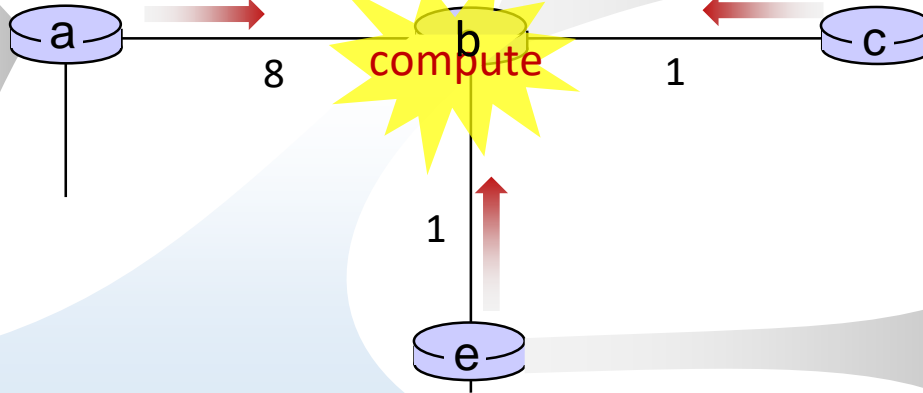


**t=1**

- 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:

$$D_b(a) = 8$$

$$D_b(c) = 1$$

$$D_b(d) = \infty$$

$$D_b(e) = 1$$

$$D_b(f) = \infty$$

$$D_b(g) = \infty$$

$$D_b(h) = \infty$$

$$D_b(i) = \infty$$

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:

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

# Distance vector example



**t=1**

- c receives DVs from b

## 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:

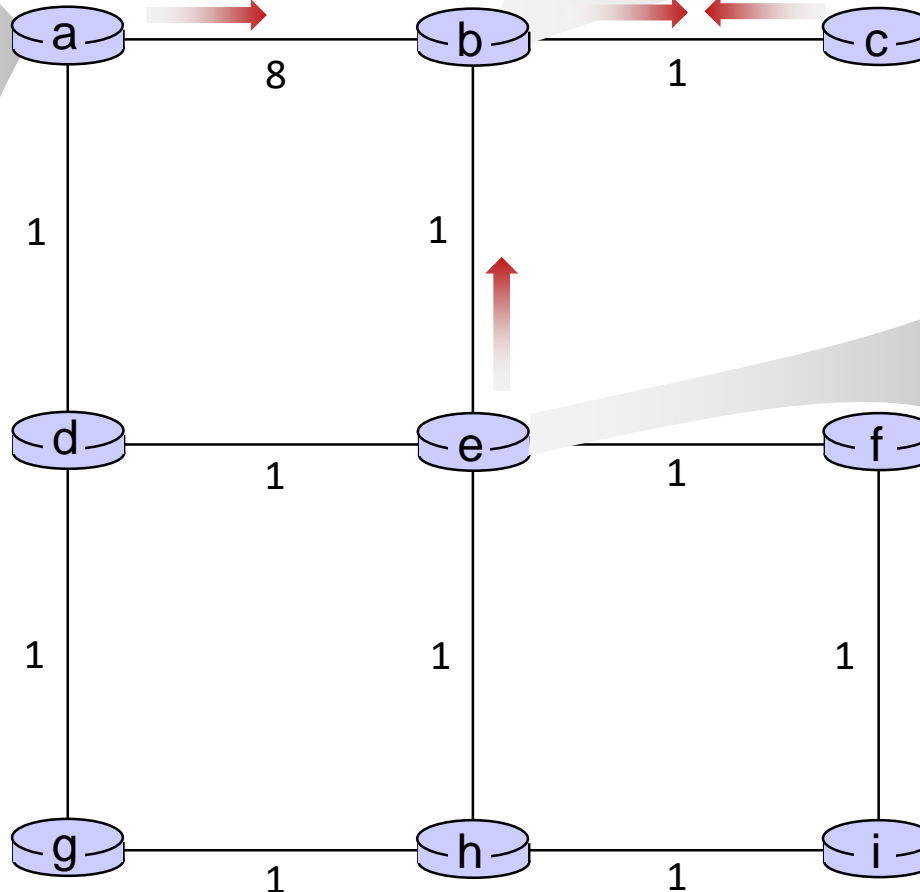
$D_b(a)=8$     $D_b(f)=\infty$   
 $D_b(c)=1$     $D_b(g)=\infty$   
 $D_b(d)=\infty$     $D_b(h)=\infty$   
 $D_b(e)=1$     $D_b(i)=\infty$

## 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:

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



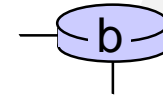
# Distance vector example



**t=1**

- c receives DVs from b computes:

$$\begin{aligned}
 D_c(a) &= \min\{c_{c,b} + D_b(a)\} = 1 + 8 = 9 \\
 D_c(b) &= \min\{c_{c,b} + D_b(b)\} = 1 + 0 = 1 \\
 D_c(d) &= \min\{c_{c,b} + D_b(d)\} = 1 + \infty = \infty \\
 D_c(e) &= \min\{c_{c,b} + D_b(e)\} = 1 + 1 = 2 \\
 D_c(f) &= \min\{c_{c,b} + D_b(f)\} = 1 + \infty = \infty \\
 D_c(g) &= \min\{c_{c,b} + D_b(g)\} = 1 + \infty = \infty \\
 D_c(h) &= \min\{c_{c,b} + D_b(h)\} = 1 + \infty = \infty \\
 D_c(i) &= \min\{c_{c,b} + D_b(i)\} = 1 + \infty = \infty
 \end{aligned}$$



1

compute

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 c:

$D_c(a) = 9$
$D_c(b) = 1$
$D_c(c) = 0$
$D_c(d) = 2$
$D_c(e) = \infty$
$D_c(f) = \infty$
$D_c(g) = \infty$
$D_c(h) = \infty$
$D_c(i) = \infty$

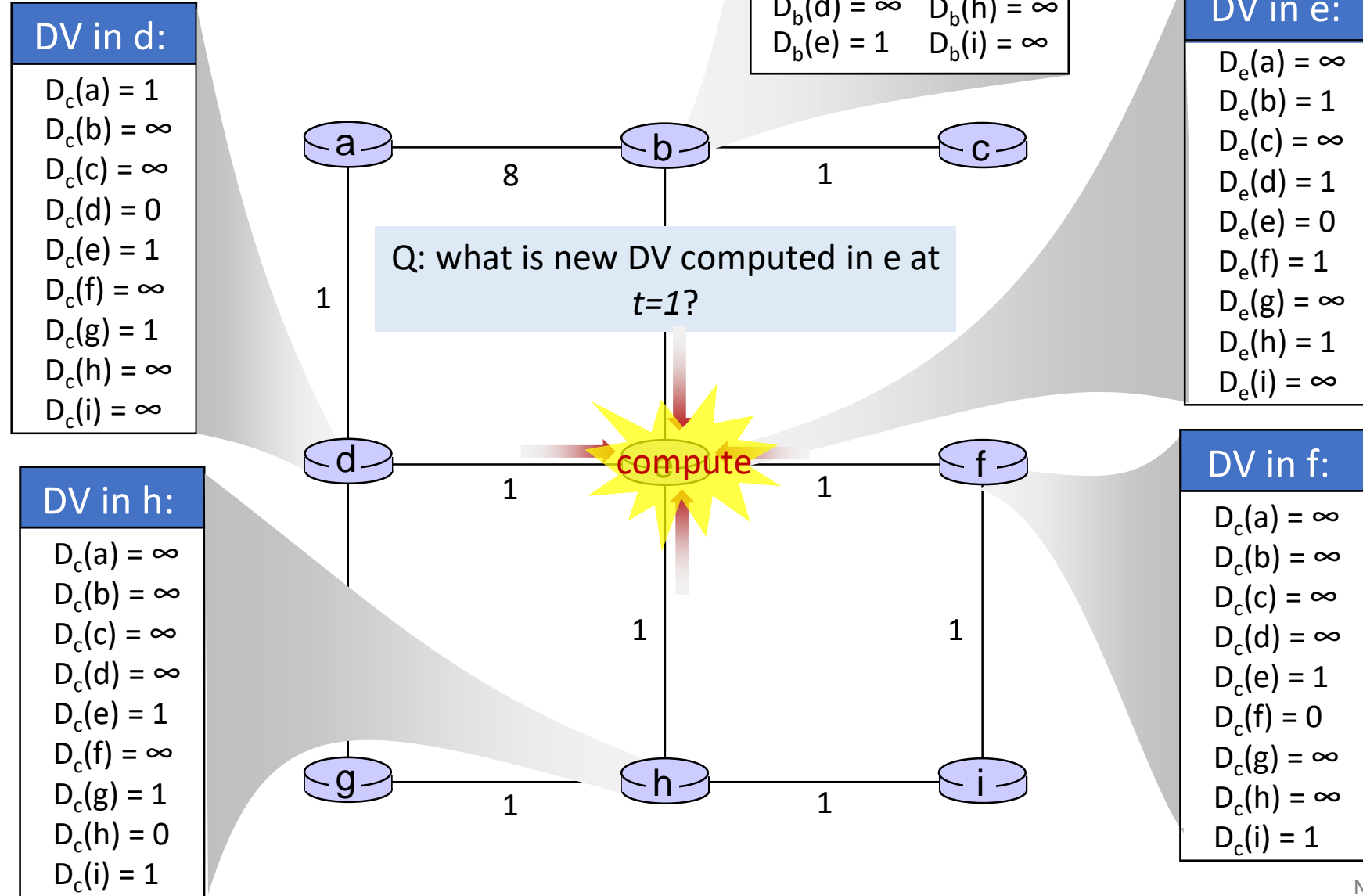
\* Check out the online interactive exercises for more examples:  
[http://gaia.cs.umass.edu/kurose\\_ross/interactive/](http://gaia.cs.umass.edu/kurose_ross/interactive/)

# Distance vector example








**t=1**

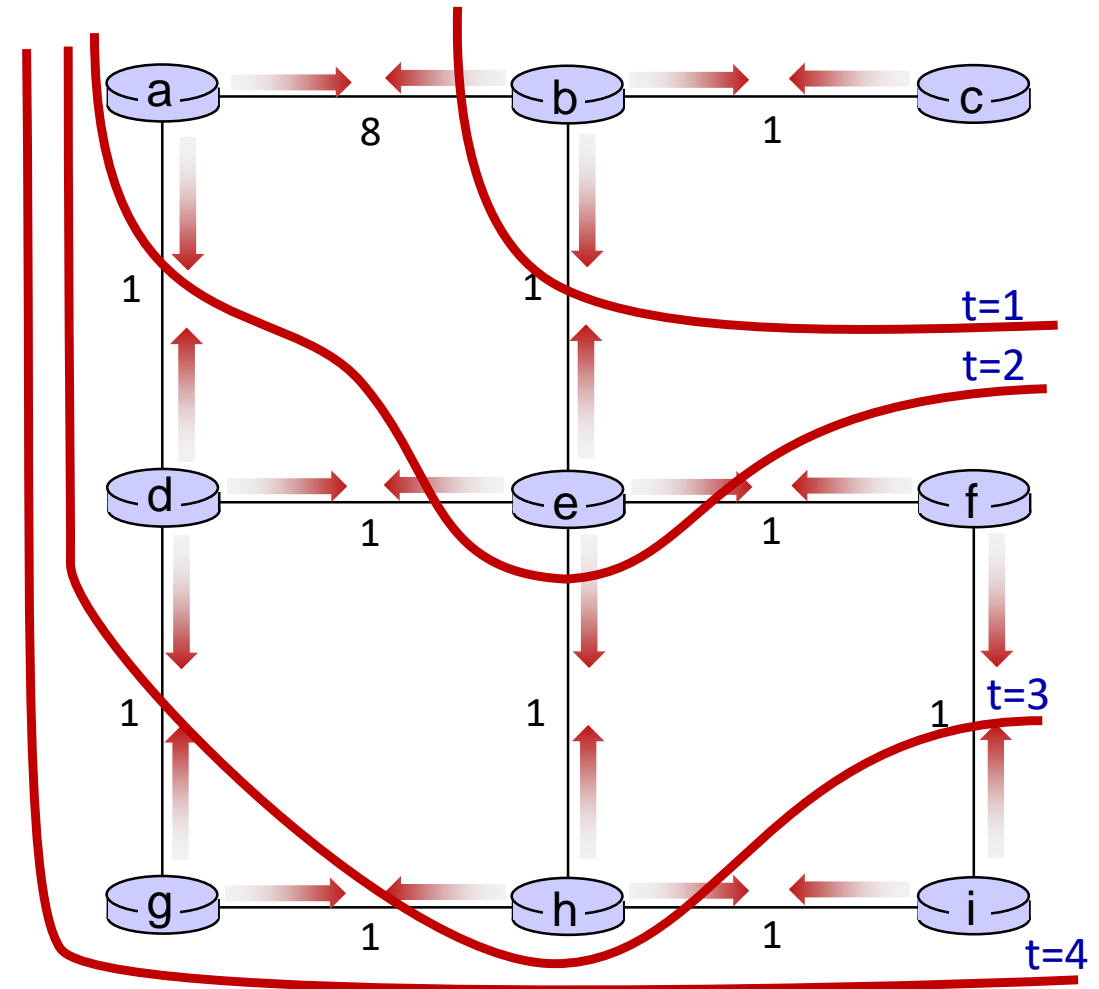
- e receives DVs from b, d, f, h



# Distance vector: state information diffusion

Iterative communication, computation steps diffuses information through network:

-   $t=0$   $c$ 's state at  $t=0$  is at  $c$  only
-   $t=1$   $c$ 's state at  $t=0$  has propagated to  $b$ , and may influence distance vector computations up to **1** hop away, i.e., at  $b$
-   $t=2$   $c$ 's state at  $t=0$  may now influence distance vector computations up to **2** hops away, i.e., at  $b$  and now at  $a$ ,  $e$  as well
-   $t=3$   $c$ 's state at  $t=0$  may influence distance vector computations up to **3** hops away, i.e., at  $d$ ,  $f$ ,  $h$
-   $t=4$   $c$ 's state at  $t=0$  may influence distance vector computations up to **4** hops away, i.e., at  $g$ ,  $i$



# Comparison of LS and DV algorithms

## message complexity

LS:  $n$  routers,  $O(n^2)$  messages sent

DV: exchange between neighbors;  
convergence time varies

## speed of convergence

LS:  $O(n^2)$  algorithm,  $O(n^2)$  messages

- may have oscillations

DV: convergence time varies

- may have routing loops
- count-to-infinity problem

robustness: what happens if router malfunctions, or is compromised?

LS:

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

DV:

- DV router can advertise incorrect *path* cost (“I have a *really* low-cost path to everywhere”): *black-holing*
- each router’s DV is used by others: error propagate thru network



# Internet approach to scalable routing

## intra-AS and Inter-AS routing

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

# 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)

# 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

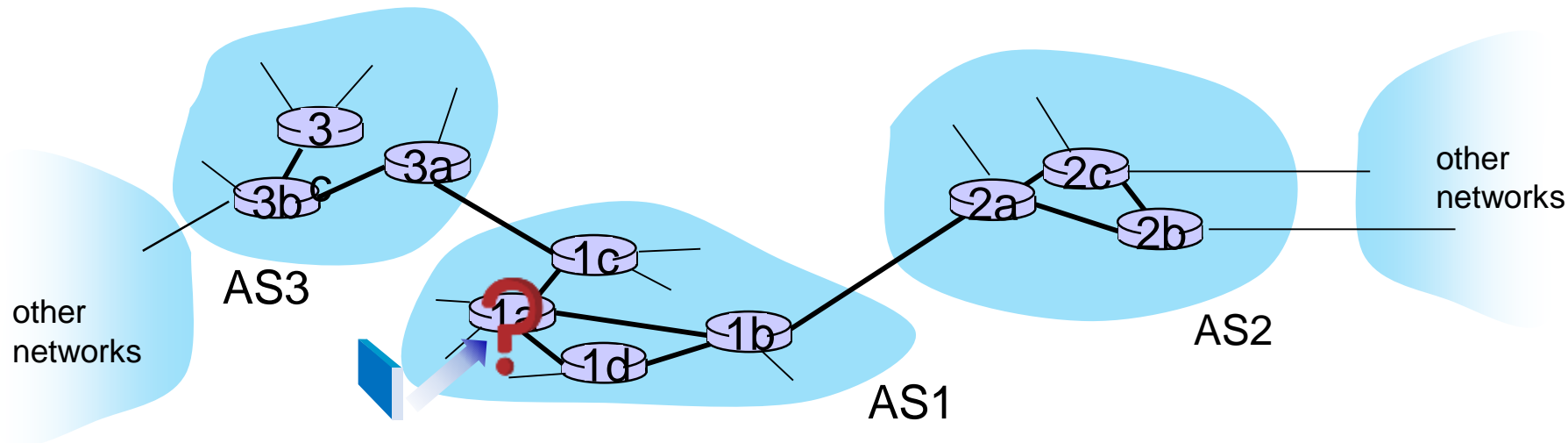
# Inter-AS routing: a role in intradomain forwarding

- 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:**

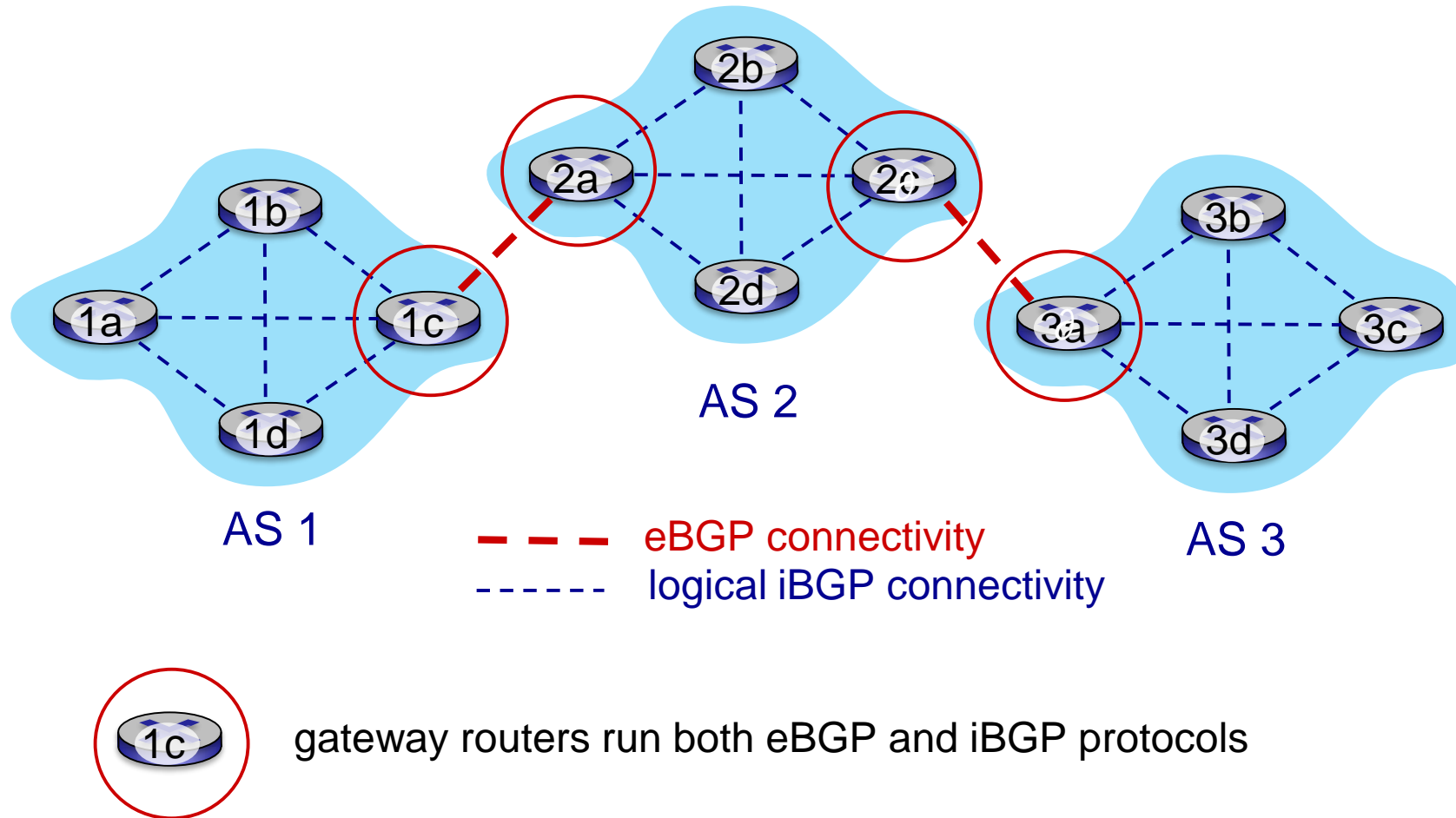
1. learn which destinations reachable through AS2, which through AS3
2. propagate this reachability info to all routers in AS1



# 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

# eBGP, iBGP connections



# BGP essentials

- **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

