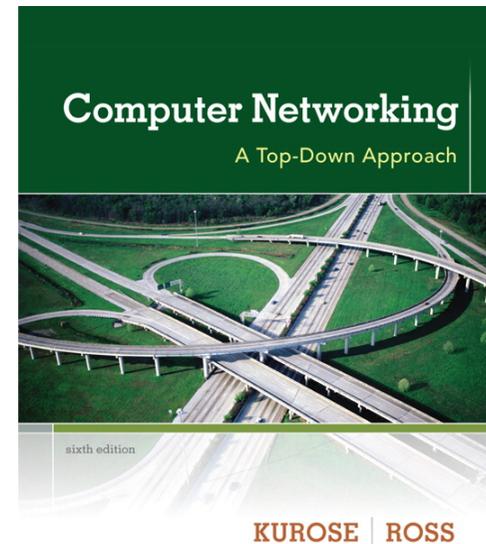


Chapter 4

Network Layer

Reti degli Elaboratori
Canale AL
Prof.ssa Chiara Petrioli
a.a. 2014/2015

We thank for the support material Prof. Kurose-Ross
All material copyright 1996-2012
© J.F Kurose and K.W. Ross, All Rights Reserved



*Computer
Networking: A Top
Down Approach*
6th edition
Jim Kurose, Keith Ross
Addison-Wesley
March 2012

Chapter 4: 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

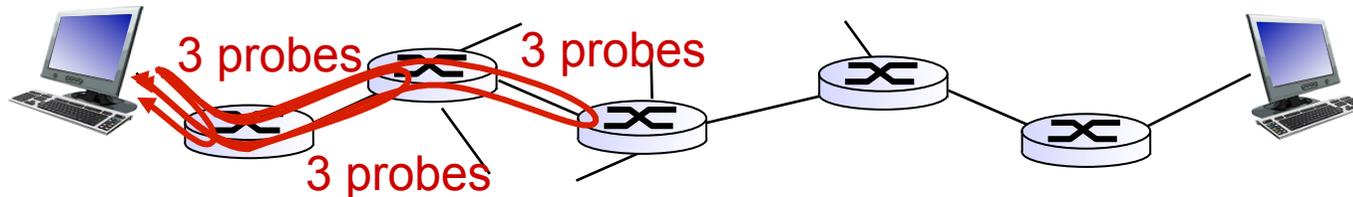
ICMP: internet control message protocol

- ❖ 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

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

Traceroute and ICMP

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



IPv6: motivation

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

IPv6 datagram format:

- fixed-length 40 byte header
- no fragmentation allowed

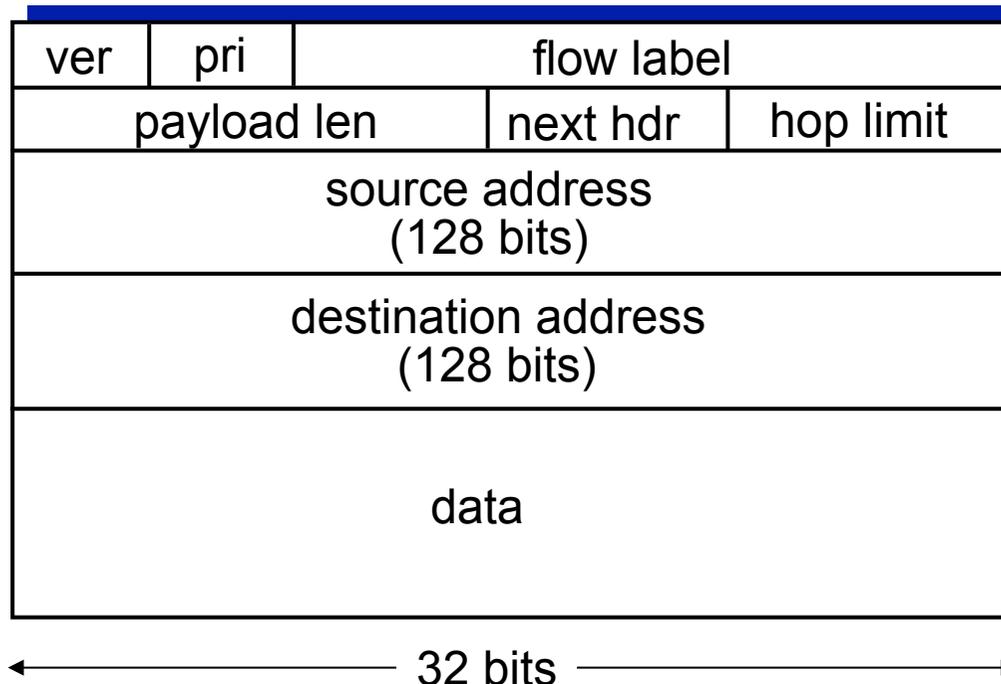
IPv6 datagram format

priority: identify priority among datagrams in flow

flow Label: identify datagrams in same “flow.”

(concept of “flow” not well defined).

next header: identify upper layer protocol for data

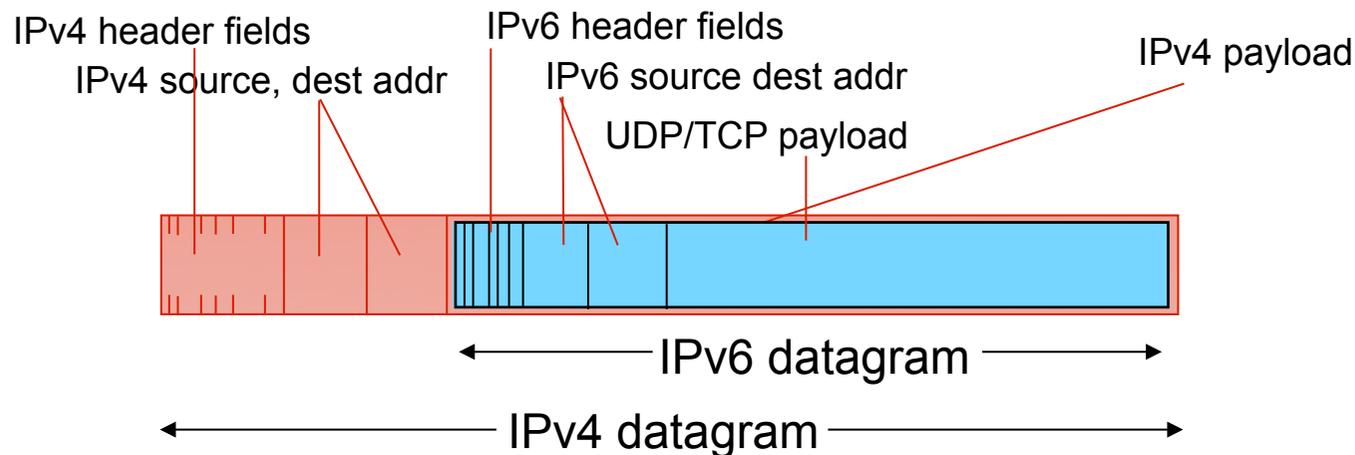


Other changes from IPv4

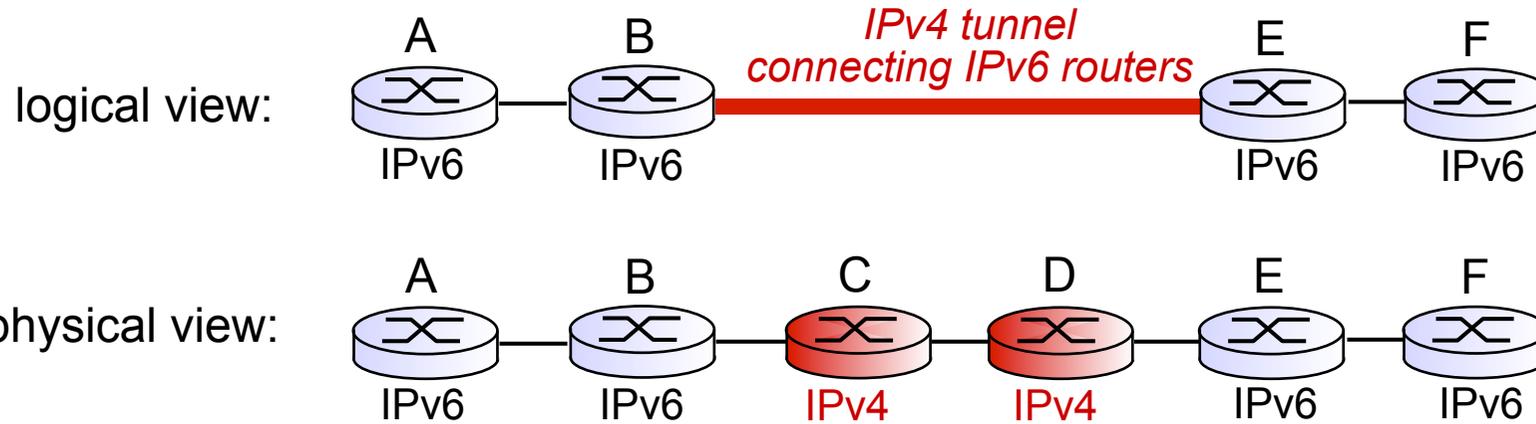
- ❖ *checksum*: removed entirely to reduce processing time at each hop
- ❖ *options*: allowed, but outside of header, indicated by “Next Header” field
- ❖ *ICMPv6*: new version of ICMP
 - additional message types, e.g. “Packet Too Big”
 - multicast group management functions

Transition from IPv4 to IPv6

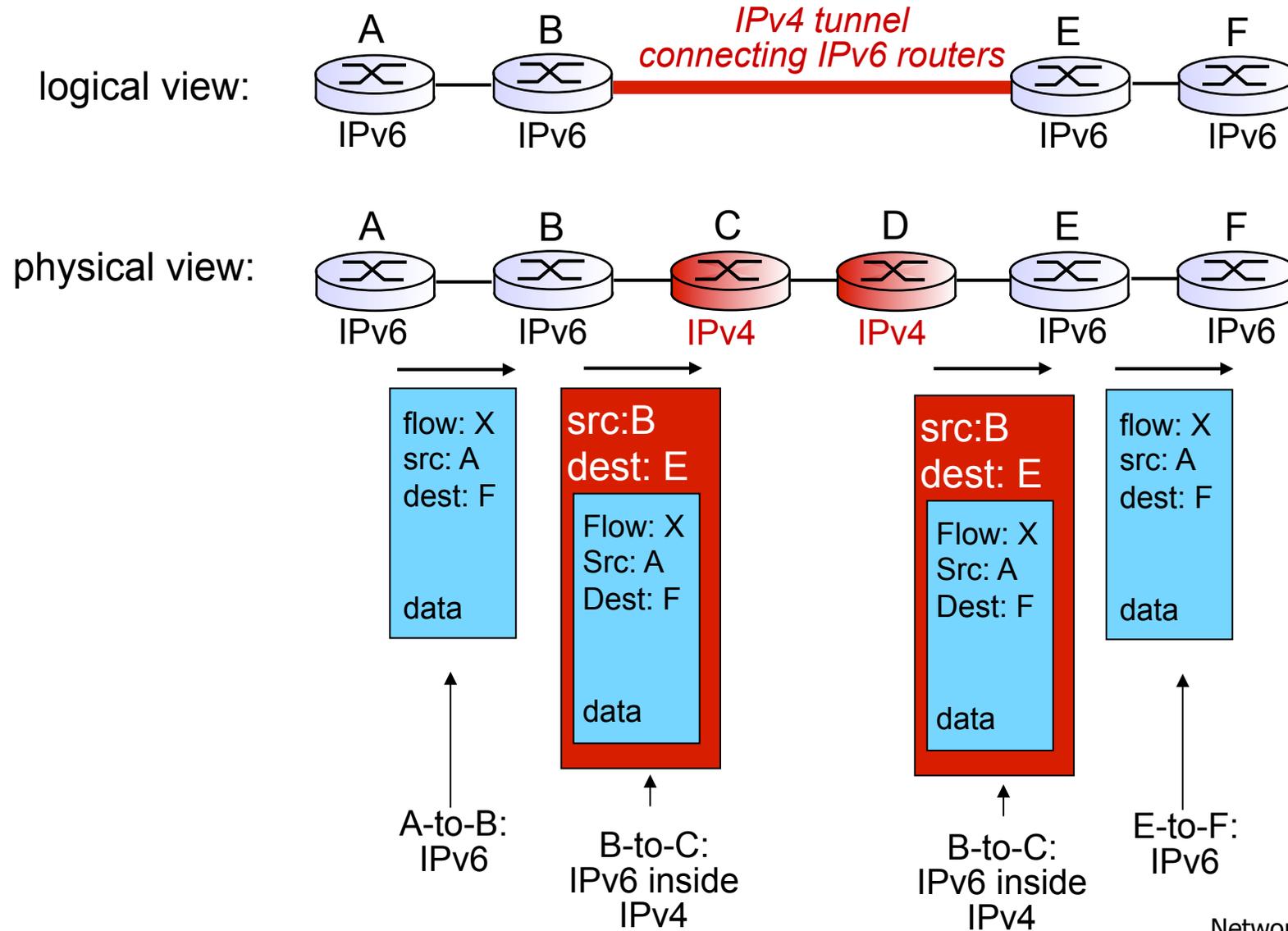
- ❖ not all routers can be upgraded simultaneously
 - no “flag days”
 - how will network operate with mixed IPv4 and IPv6 routers?
- ❖ **tunneling**: IPv6 datagram carried as *payload* in IPv4 datagram among IPv4 routers



Tunneling



Tunneling



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

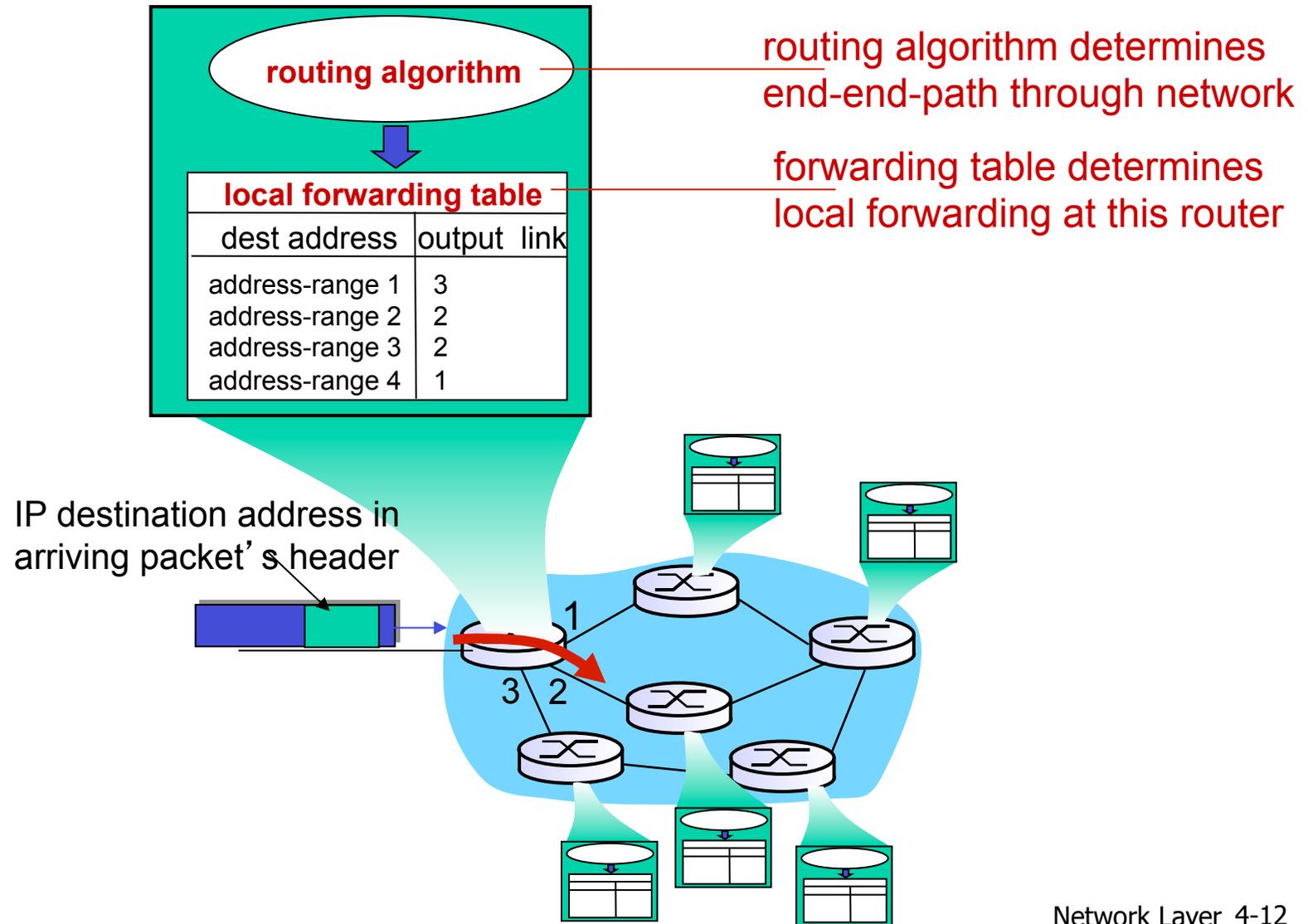
- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet

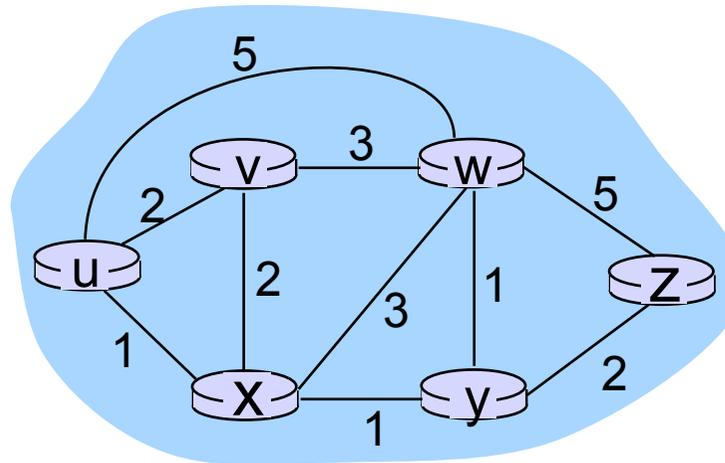
- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Interplay between routing, forwarding



Graph abstraction



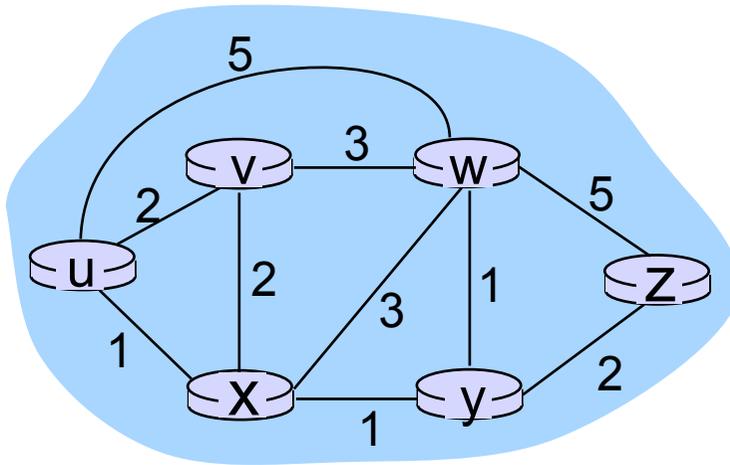
graph: $G = (N,E)$

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

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

aside: graph abstraction is useful in other network contexts, e.g., P2P, where N is set of peers and E is set of TCP connections

Graph abstraction: costs



$c(x, x')$ = cost of link (x, x')
e.g., $c(w, z) = 5$

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

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

key question: what is the least-cost path between u and z ?
routing algorithm: algorithm that finds that least cost path

Routing algorithm classification

Q: global or decentralized information?

global:

- ❖ all routers have complete topology, link cost info
- ❖ “link state” algorithms

decentralized:

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

Q: static or dynamic?

static:

- ❖ routes change slowly over time

dynamic:

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

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 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

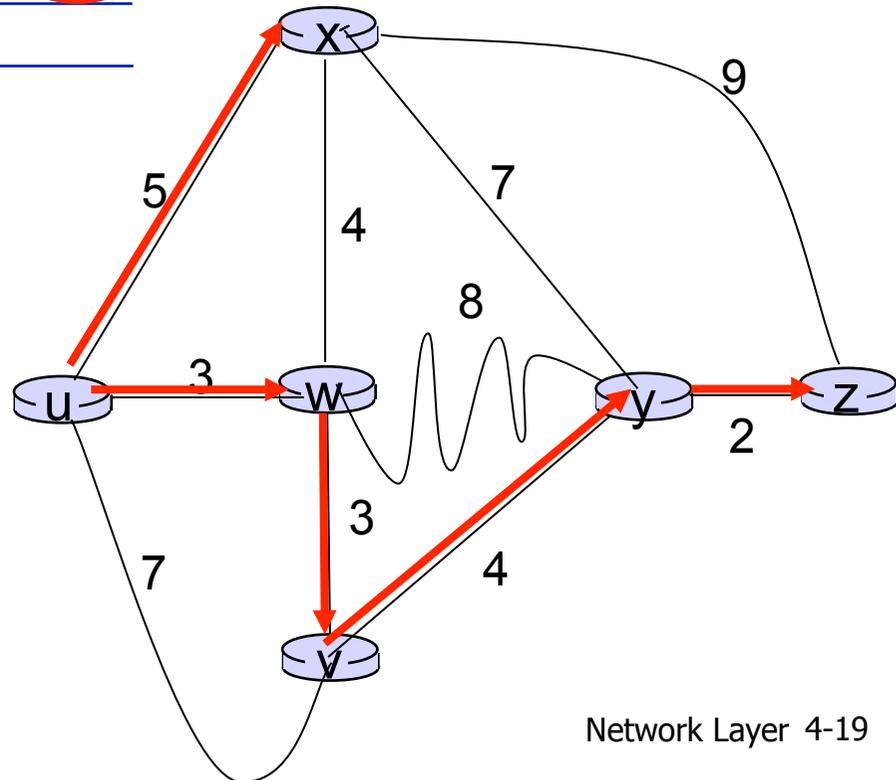
- 1 **Initialization:**
- 2 $N' = \{u\}$
- 3 for all nodes v
- 4 if v adjacent to u
- 5 then $D(v) = c(u,v)$
- 6 else $D(v) = \infty$
- 7
- 8 **Loop**
- 9 find w not in N' such that $D(w)$ is a minimum
- 10 add w to N'
- 11 update $D(v)$ for all v adjacent to w and not in N' :
- 12 **$D(v) = \min(D(v), D(w) + c(w,v))$**
- 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'**

Dijkstra's algorithm: example

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

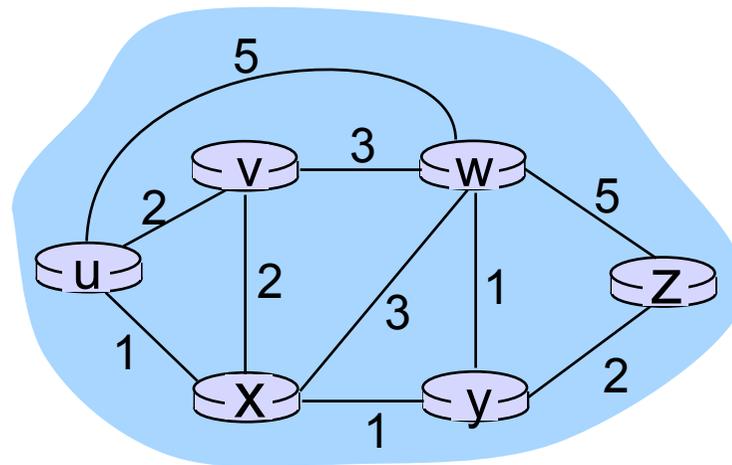
notes:

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



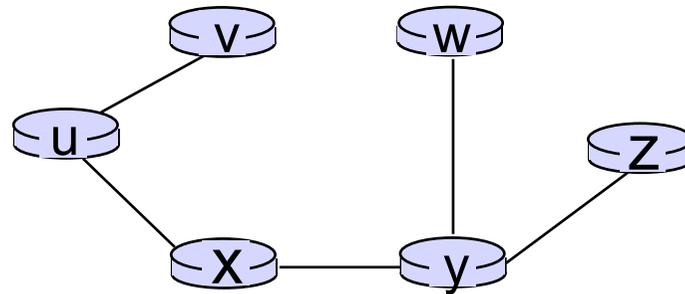
Dijkstra's algorithm: another example

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



Dijkstra's algorithm: example (2)

resulting shortest-path tree from u:



resulting forwarding table in u:

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

Correttezza

Se eseguiamo l'algoritmo di Dijkstra su un grafo pesato diretto $G=(N,E)$ con pesi non negativi, sorgente u , e funzione peso c allora alla terminazione $D(v)=\delta(u,v)$, per ogni nodo v in N . (dove $\delta(u,v)$ indica la lunghezza del cammino di peso minimo tra u e v).

Dim.

$D(v)$ non è più aggiornato nel momento in cui v è inserito in N' . Dovremmo quindi mostrare che $D(v)=\delta(u,v)$ nel momento in cui v è inserito in N' , per ogni v .

Ragioniamo per assurdo. Sia x il primo nodo (nell'ordine di inserimento in N') per cui vale $D(x) \neq \delta(u,x)$ al momento in cui x è inserito nell'insieme N' (linea 10 dell'algoritmo). $x \neq u$ dato che u , nodo sorgente, è inserito nella fase di inizializzazione e per lui vale $D(u)=\delta(u,u)=0$. Inoltre deve esistere un percorso di costo non infinito da u a x dato che altrimenti varrebbe che il valore a cui $D(x)$ è inizializzato (infinito) sarebbe uguale a $\delta(u,x)$. Quindi esiste un percorso di costo minimo $p=u \dots v \rightarrow y \dots x$ dove y è il primo nodo sul percorso di costo minimo NON in N' (quindi $u \dots v$ sono TUTTI in N'). Il percorso p può quindi essere diviso in due percorsi: p_1 che va da u a y e p_2 che va da y a x .

Da notare che il percorso p_1 è anch'esso il percorso di costo minimo che unisce u a y (se non lo fosse e ci fosse un percorso p_3 che unisce u a y di costo $<$ del costo di p_1 , allora la concatenazione di p_2 e p_3 sarebbe un percorso p' da u a x di costo $<$ di p , contro l'assunto che p sia un percorso di costo minimo).

...Correttezza

Se eseguiamo l'algoritmo di Dijkstra su un grafo pesato diretto $G=(N,E)$ con pesi non negativi, sorgente u , e funzione peso c allora alla terminazione $D(v)=\delta(u,v)$, per ogni nodo v in N . (dove $\delta(u,v)$ indica la lunghezza del cammino di peso minimo tra u e v).

Dim (...continua).

Quando x è inserito in N' $D(y)=\delta(u,y)$. Infatti in quel momento v è stato già inserito in N' e dopo il suo inserimento y ha ricalcolato $D(y)=D(v)+c(v,y)=\delta(u,v)+c(v,y)$ (dato che per ipotesi x è il primo nodo per cui all'inserimento in N' la stima dei costi non corrisponde al percorso di costo minimo) $=\delta(u,y)$.

Dato che y precede x sul percorso minimo ed i pesi sugli archi sono non negativi vale che:

$$\delta(u,x) \geq \delta(u,y) = D(y)$$

e quindi anche che

$$D(x) \geq \delta(u,x) \geq \delta(u,y) = D(y)$$

D'altro canto dato che x viene inserito in N' prima di y vale che

$$\delta(u,x) \leq D(x) \leq D(y) = \delta(u,y)$$

Quindi

$$\delta(u,x) = D(x) = D(y) = \delta(u,y)$$

Cosa che porta alla contraddizione.

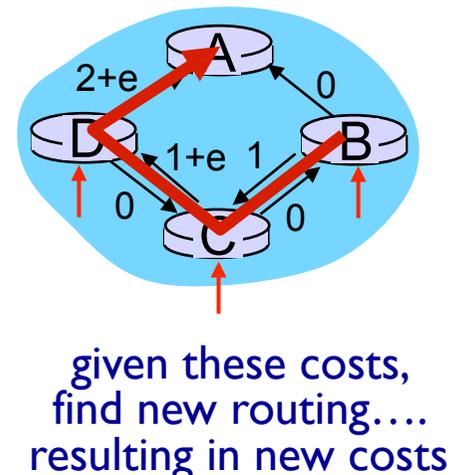
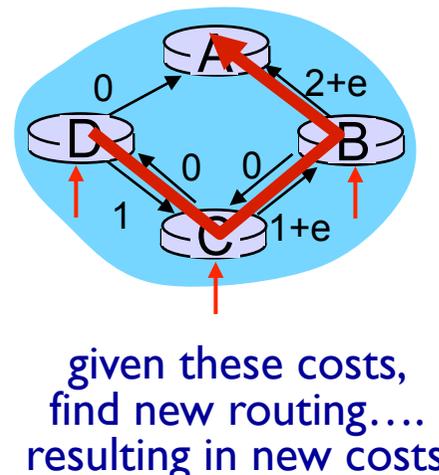
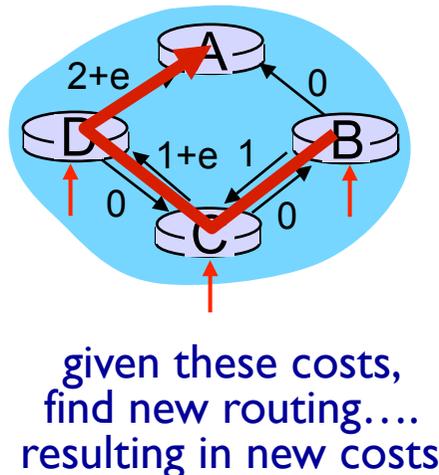
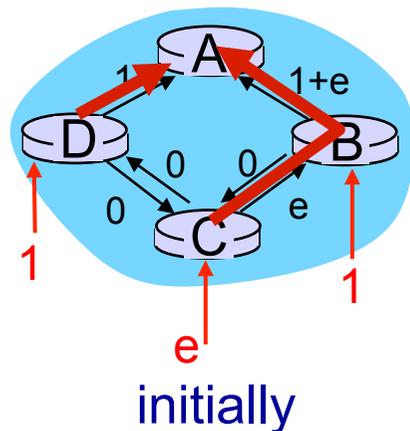
Dijkstra's algorithm, discussion

algorithm complexity: n nodes

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

oscillations possible:

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



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

Bellman-Ford

Given a graph $G=(N,E)$ and a node s finds the shortest path from s to every node in N .

A shortest walk from s to i subject to the constraint that the walk contains at most h arcs and goes through node s only once, is denoted shortest($\leq h$) walk and its length is D_i^h .

Bellman-Ford rule:

Initialization $D_s^h=0$, for all h ; $c_{i,k} = \text{infinity}$ if (i,k) NOT in E ; $c_{k,k} = 0$;
 $D_i^0 = \text{infinity}$ for all $i \neq s$.

Iteration:

$$D_i^{h+1} = \min_k [c_{i,k} + D_k^h]$$

Assumption: non negative cycles (**this is the case in a network!!**)

The Bellman-Ford algorithm first finds the one-arc shortest walk lengths, then the two-arc shortest walk length, then the three-arc...etc. \rightarrow distributed version used for routing

Distance vector algorithm

Bellman-Ford equation (dynamic programming)

let

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

then

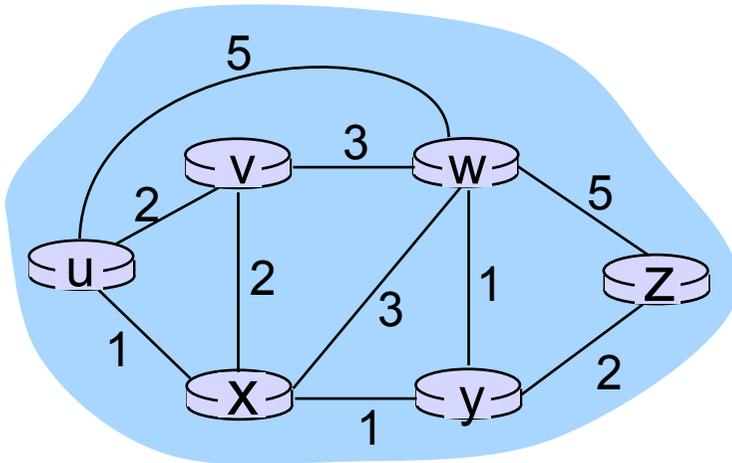
$$d_x(y) = \min_v \{ c(x,v) + d_v(y) \}$$

cost from neighbor v to destination y

cost to neighbor v

\min taken over all neighbors v of x

Bellman-Ford example



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

B-F equation says:

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

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

Distance vector algorithm

- ❖ $D_x(y)$ = estimate of least cost from x to y
 - x maintains distance vector $D_x = [D_x(y): y \in N]$
- ❖ 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:

At all nodes, X:

- 1 Initialization:
- 2 for all adjacent nodes v:
- 3 $D^X(*,v) = \text{infinity}$ /* the * operator means "for all rows" */
- 4 $D^X(v,v) = c(X,v)$
- 5 for all destinations, y
- 6 send $\min_w D^X(y,w)$ to each neighbor /* w over all X's neighbors */

From the node to whatever destination going through v

Distance Vector Algorithm (cont.):

```
8 loop
9  wait (until I see a link cost change to neighbor V
10      or until I receive update from neighbor V)
11
12  if (c(X,V) changes by d)
13      /* change cost to all dest's via neighbor v by d */
14      /* note: d could be positive or negative */
15      for all destinations y:  $D^X(y,V) = D^X(y,V) + d$ 
16
17  else if (update received from V wrt destination Y)
18      /* shortest path from V to some Y has changed */
19      /* V has sent a new value for its  $\min_w DV(Y,w)$  */
20      /* call this received new value is "newval" */
21      for the single destination y:  $D^X(Y,V) = c(X,V) + \text{newval}$ 
22
23  if we have a new  $\min_w D^X(Y,w)$  for any destination Y
24      send new value of  $\min_w D^X(Y,w) = \text{of } D_x(Y)$  to all neighbors
25
26 forever
```

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:

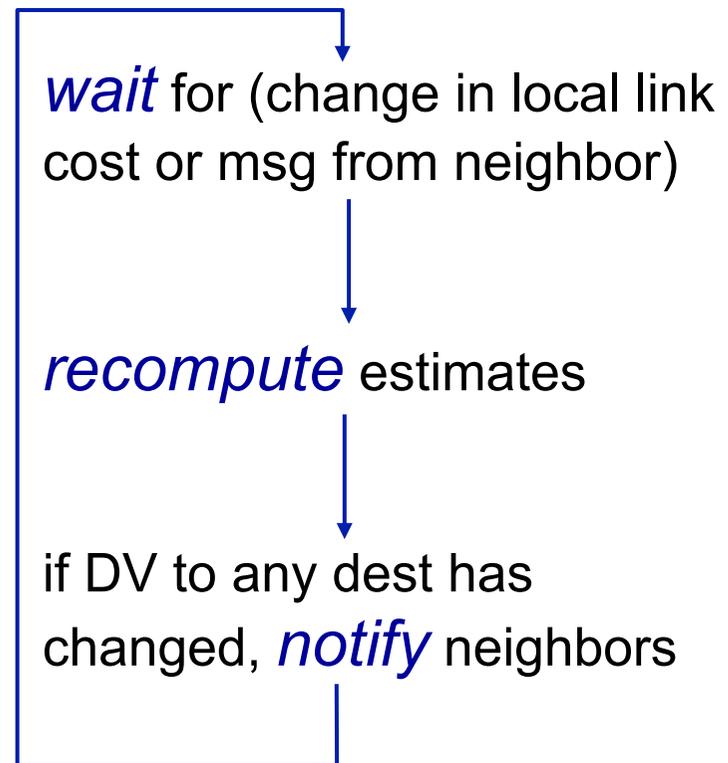
each local iteration
caused by:

- ❖ 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:



Previous lecture. Summary:

Distributed Belman Ford

- Based on Distributed Bellman Ford Equation

Cost associated to the (X,Z) link

$$D^X(Y,Z) = \text{distance from } X \text{ to } Y, \text{ via } Z \text{ as next hop}$$
$$= c(X,Z) + \min_w \{D^Z(Y,w)\}$$

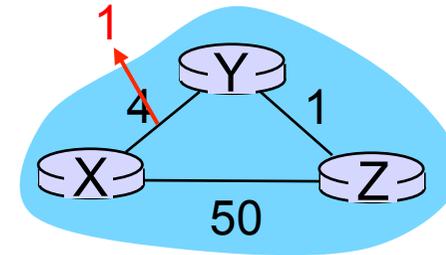
- $D^X(Y,Z)$ re-computed:
 - Upon reception of updates from the neighbors
 - Upon link cost change
- $\min_z D^X(Y,Z)$ communicated to the neighbors whenever its value changes, or periodically
- How long does it take for the algorithm to converge? ‘good news travel fast, bad news may not → count to infinity’

Distance Vector: link cost changes

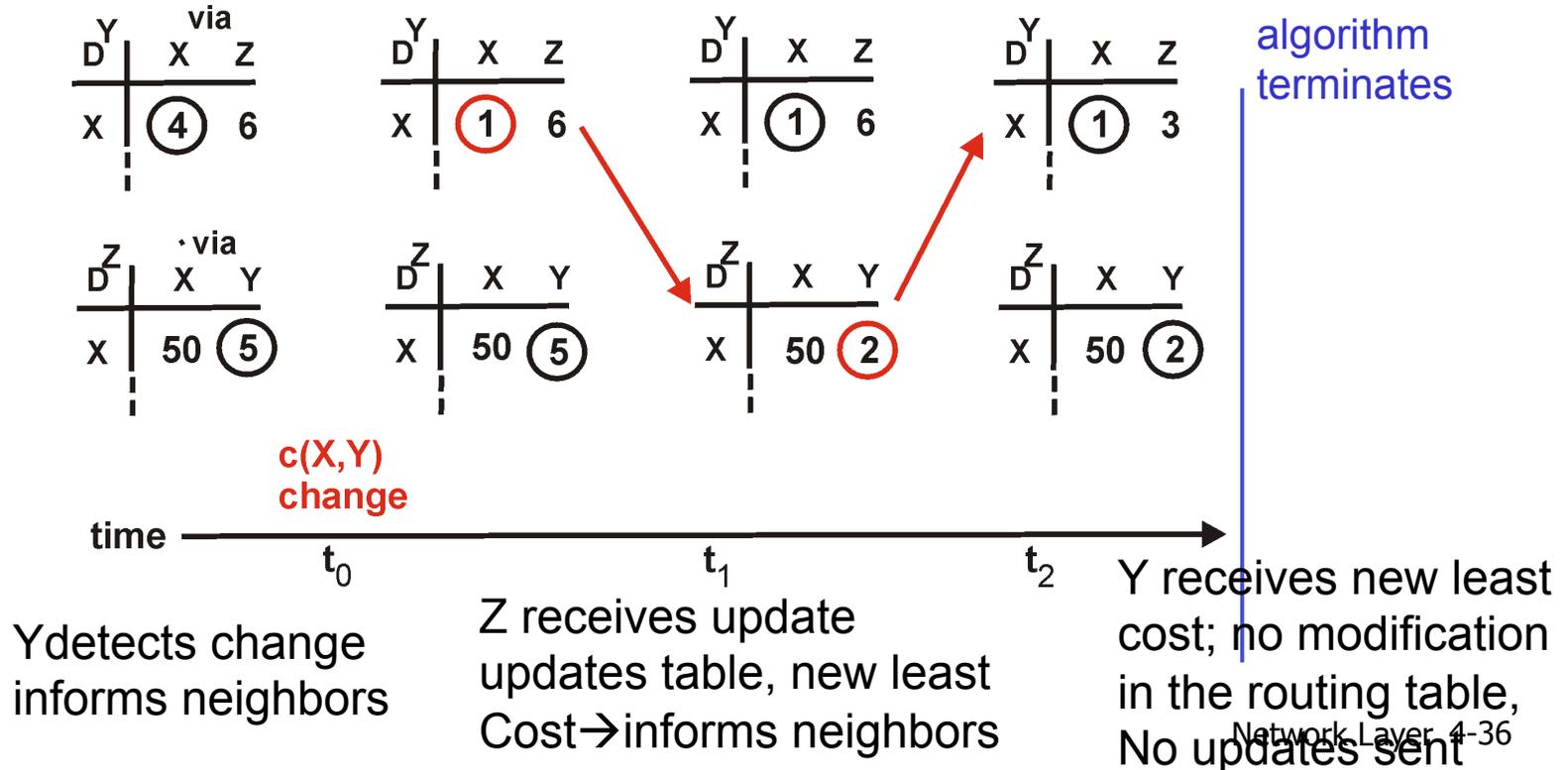
Subtitle: Distributed Bellman Ford converges but how fast?

Link cost changes:

- ❑ node detects local link cost change
- ❑ updates distance table (line 15)
- ❑ if cost change in least cost path, notify neighbors (lines 23,24)



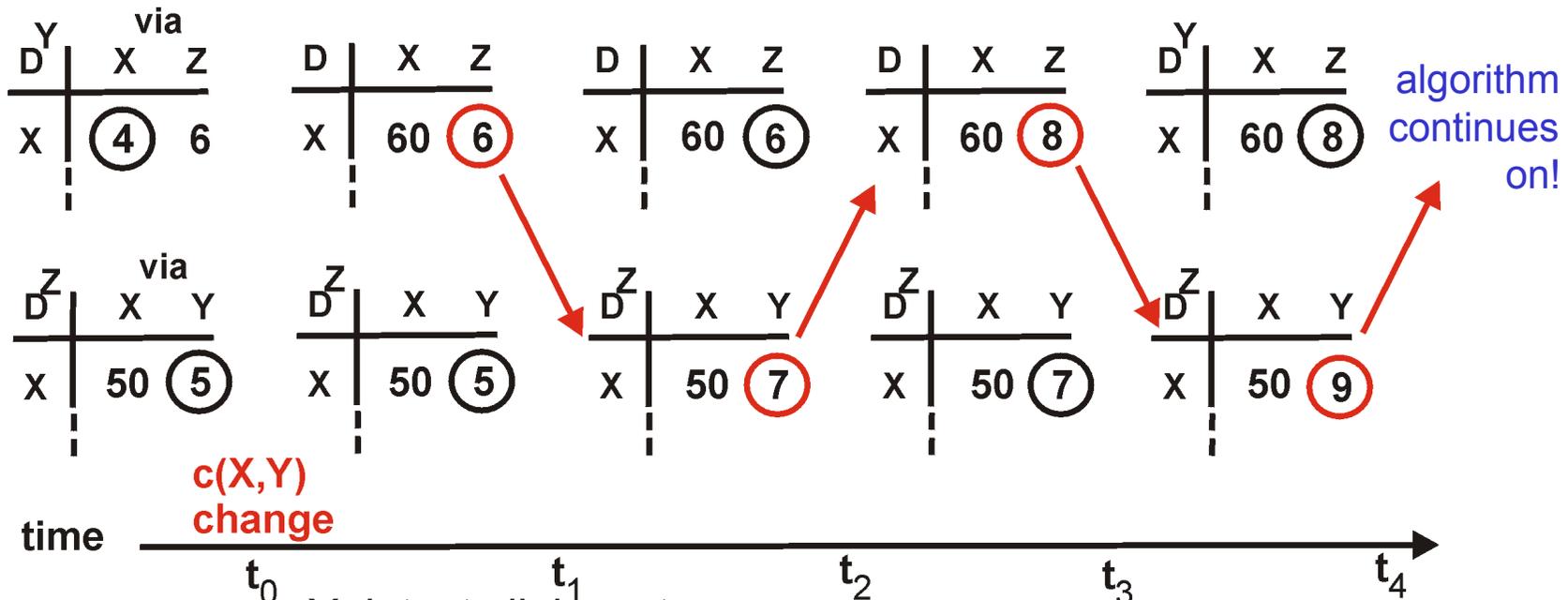
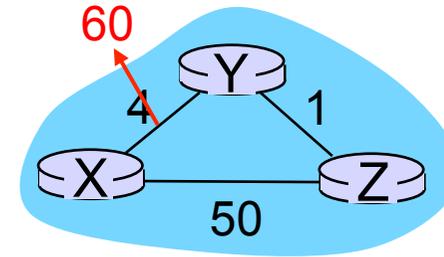
“good news travels fast”



Distance Vector: link cost changes

Link cost changes:

- good news travels fast
- bad news travels slow - “count to infinity” problem!



Y detects link cost Increase but think can Reach X through Z at a total cost of 6 (wrong!!)

The path is Y-Z-Y-X

Count-to-infinity –an everyday life example

Which is the problem here?

the info exchanged by the protocol!! ‘the best route to X I have has the following cost...’ (no additional info on the route)

A Roman example...

-assumption: there is only one route going from Colosseo to Altare della Patria: Via dei Fori Imperiali. Let us now consider a network, whose nodes are Colosseo., Altare della Patria, Piazza del Popolo



Count-to-infinity –everyday life example (2/2)



The Colosseo. and Alt. Patria nodes exchange the following info

- Colosseo says ‘the shortest route from me to P. Popolo is 2 Km’
- Alt. Patria says ‘the shortest path from me to P. Popolo is 1Km’

Based on this exchange from Colosseo you go to Al. Patria, and from there to Piazza del Popolo OK Now due to the big dig they close Via del Corso (Al. Patria—P.Popolo)

- Al. Patria thinks ‘I have to find another route from me to P.Popolo.

Look there is a route from Colosseo to P.Popolo that

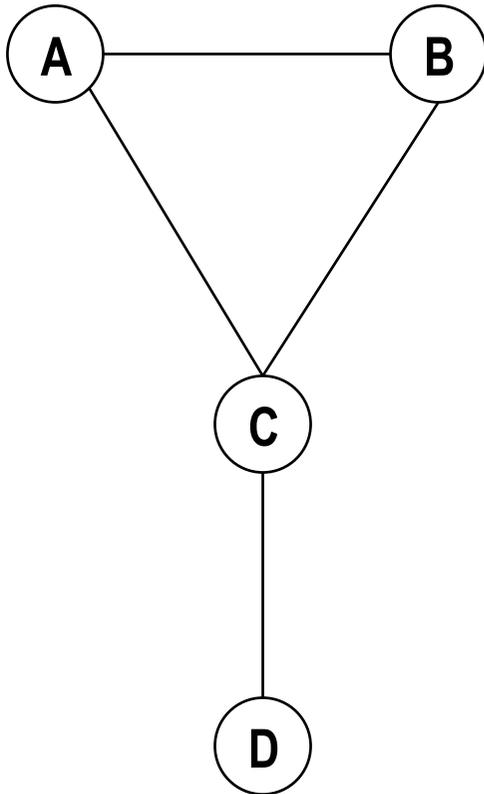
takes 2Km, I can be at Colosseo in 1Km → I have found

a 3Km route from me to P.Popolo!!’ Communicates the new cost to

Colosseo that updates ‘OK I can go to P.Popolo via Al. Patria in 4Km’

VERY WRONG!! Why is it so? I didn’t know that the route from Colosseo to P.Popolo was going through Via del Corso from Al.Patria to P.Popolo (which is closed)!!

Split horizon poison reverse failure



Line CD goes down...

- 1) because of split horizon rule, A and B tell C that $\text{dist}(D)=\text{inf}$
- 2) C concludes that D is unreachable and reports this to A and B
- 3) but A knows from B that $\text{dist}(D)=2$, and sets its $\text{dist}=3$
- 4) similarly, B knows from A distance from D... C estimates new value 4; A and B again through C estimate a value of 5....then again 1) ... etc until distance = infinite

Regardless the hack used, there is always a network topology that makes the trick fail!

Comparison of LS and DV algorithms

message complexity

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

speed of convergence

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

robustness: what happens if router malfunctions?

LS:

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

DV:

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

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 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

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

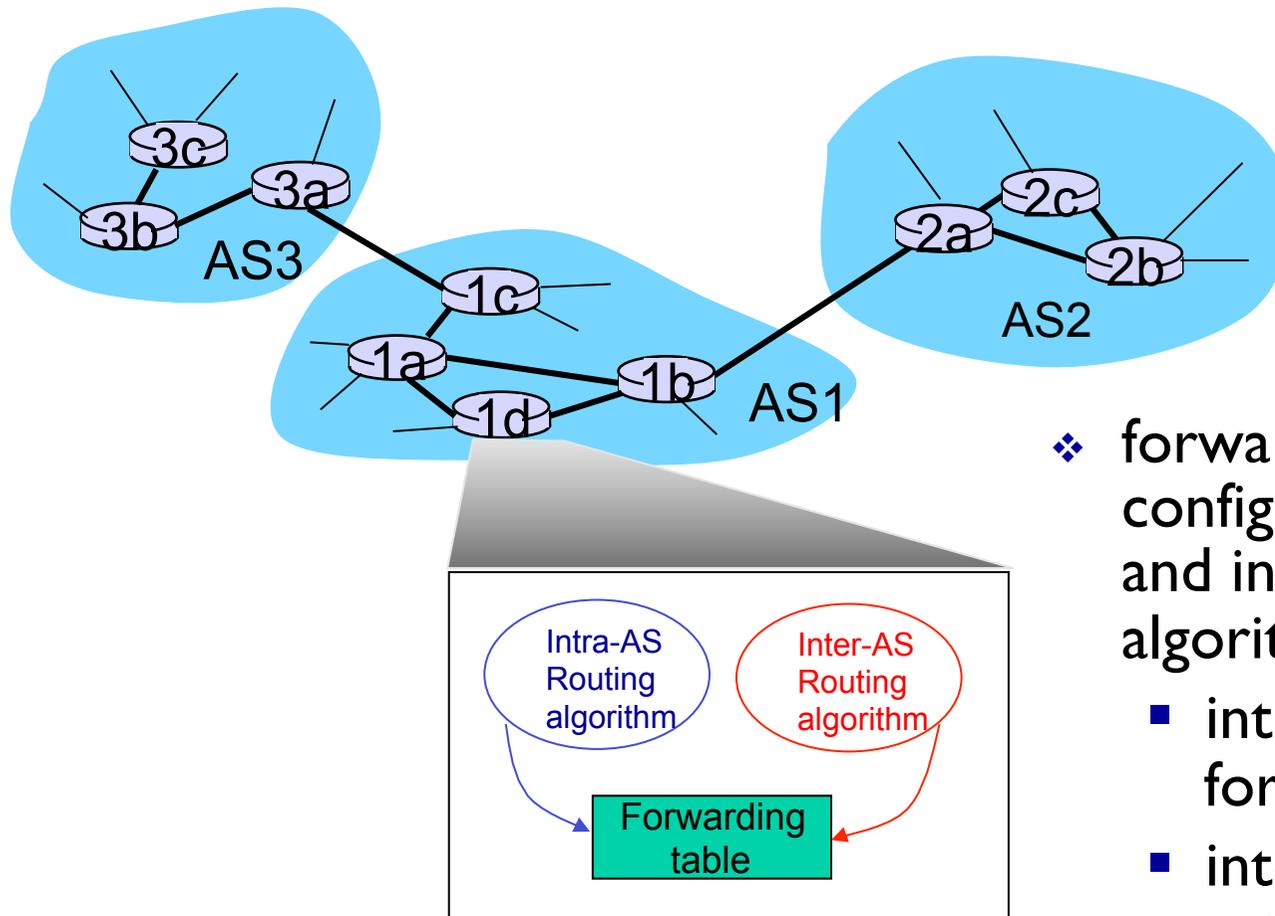
Hierarchical routing

- ❖ 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

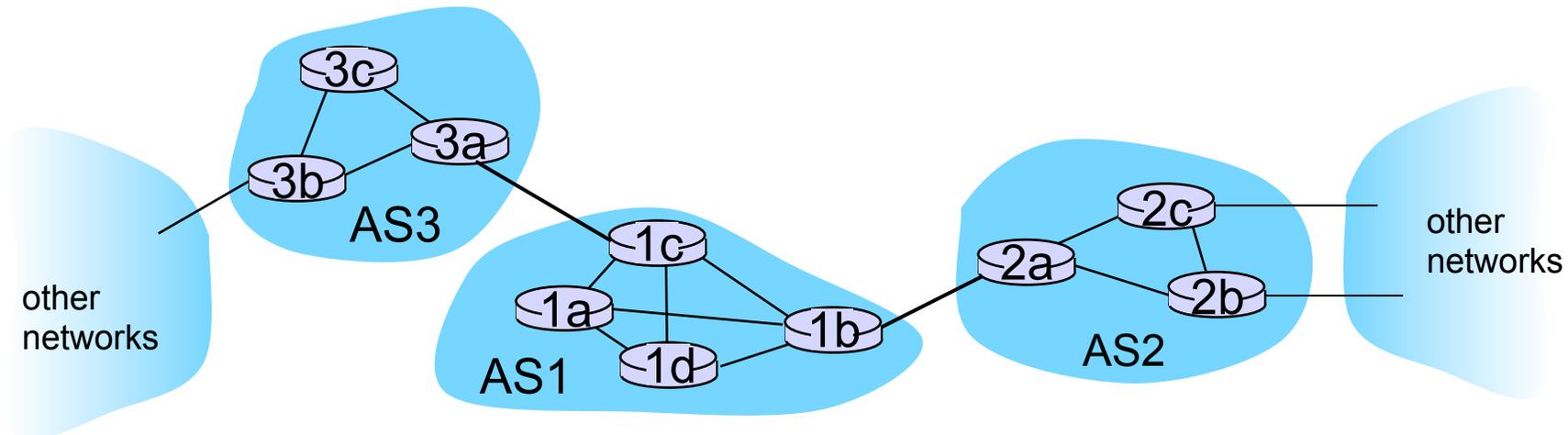
Inter-AS tasks

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

AS1 must:

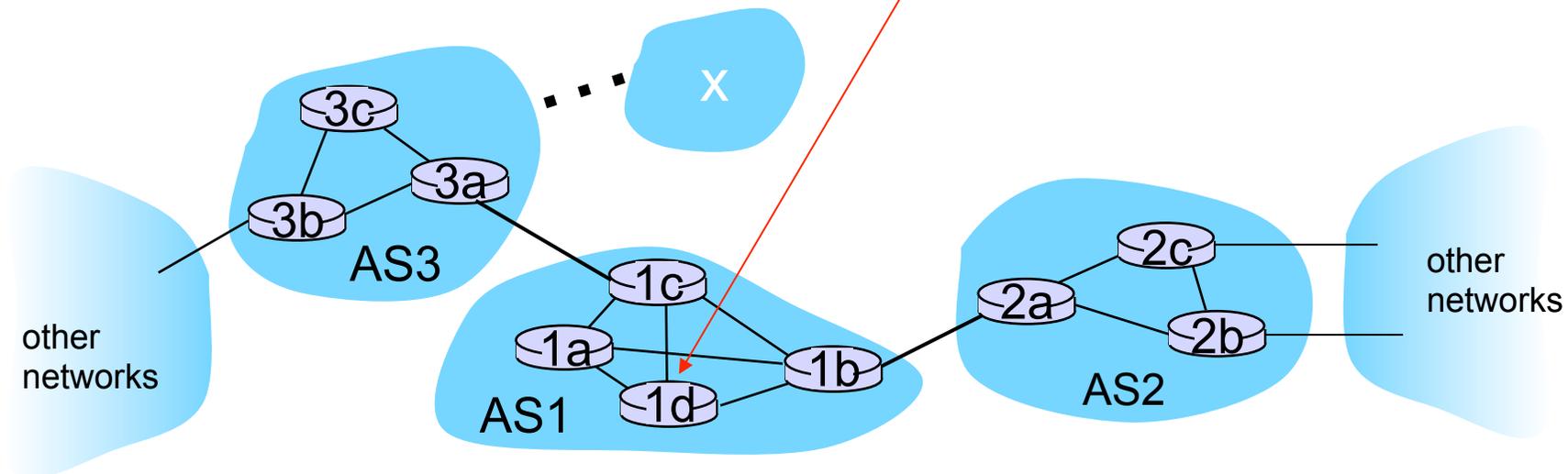
1. learn which destds are reachable through AS2, which through AS3
2. propagate this reachability info to all routers in AS1

job of inter-AS routing!



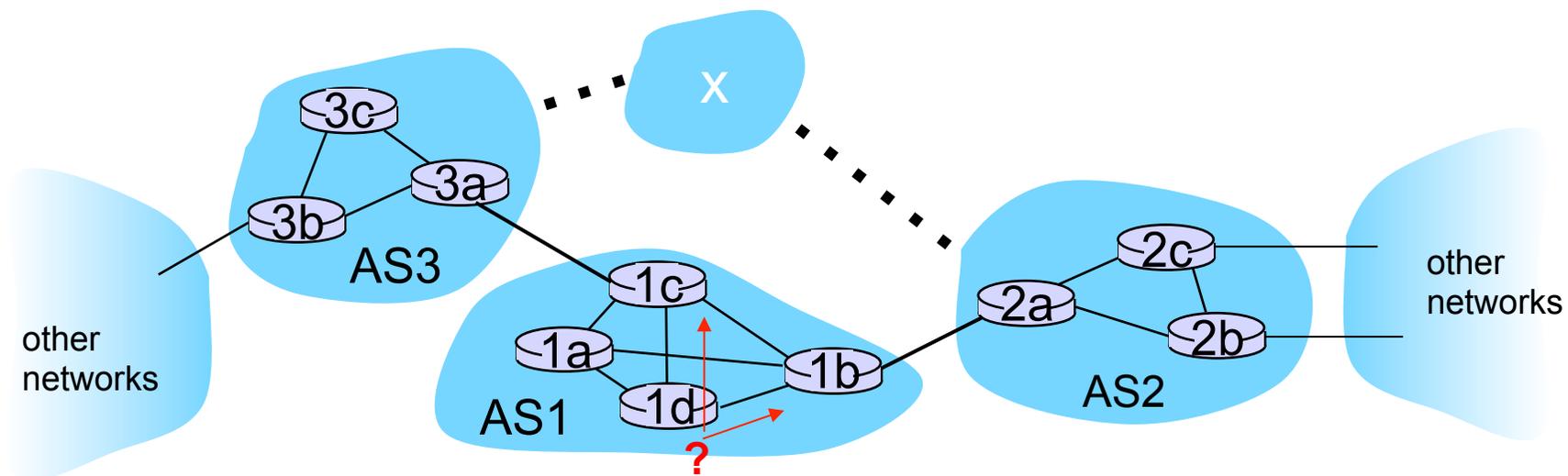
Example: setting forwarding table in router 1d

- ❖ 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 l is on the least cost path to 1c
 - installs forwarding table entry (x, l)



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 which gateway it should forward packets towards 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.

