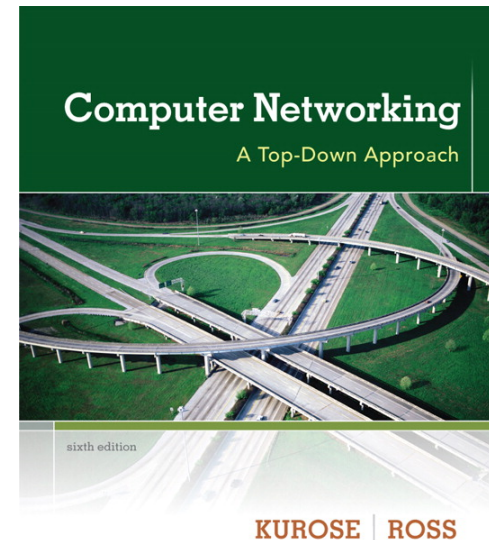


Chapter 4

Network Layer

Reti degli Elaboratori
Canale AL
Prof.ssa Chiara Petrioli
a.a. 2013/2014

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

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

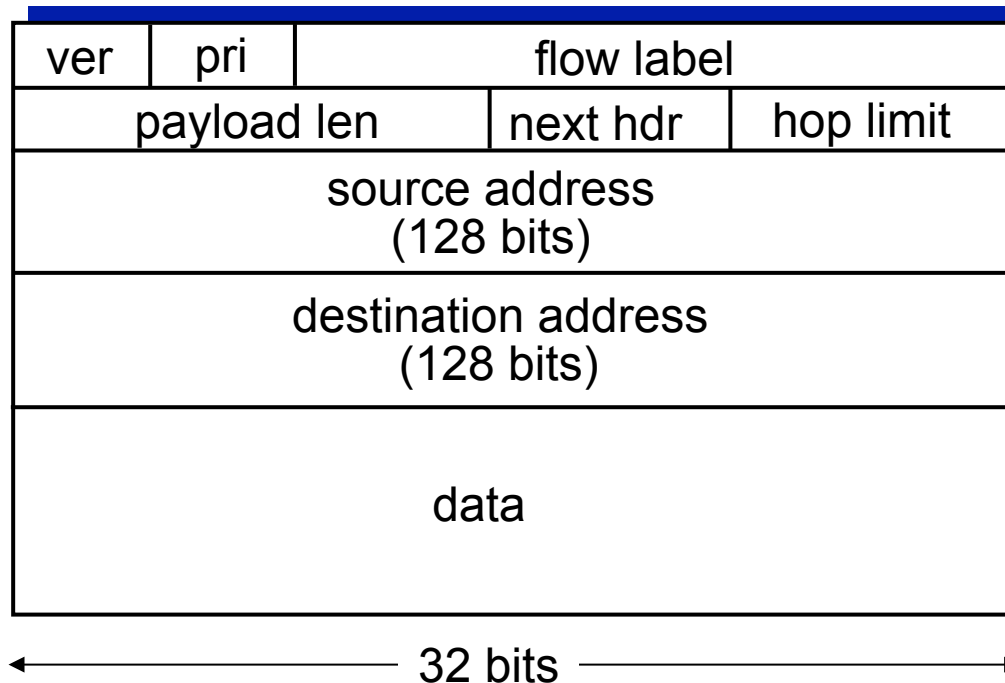
IPv6 datagram format

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

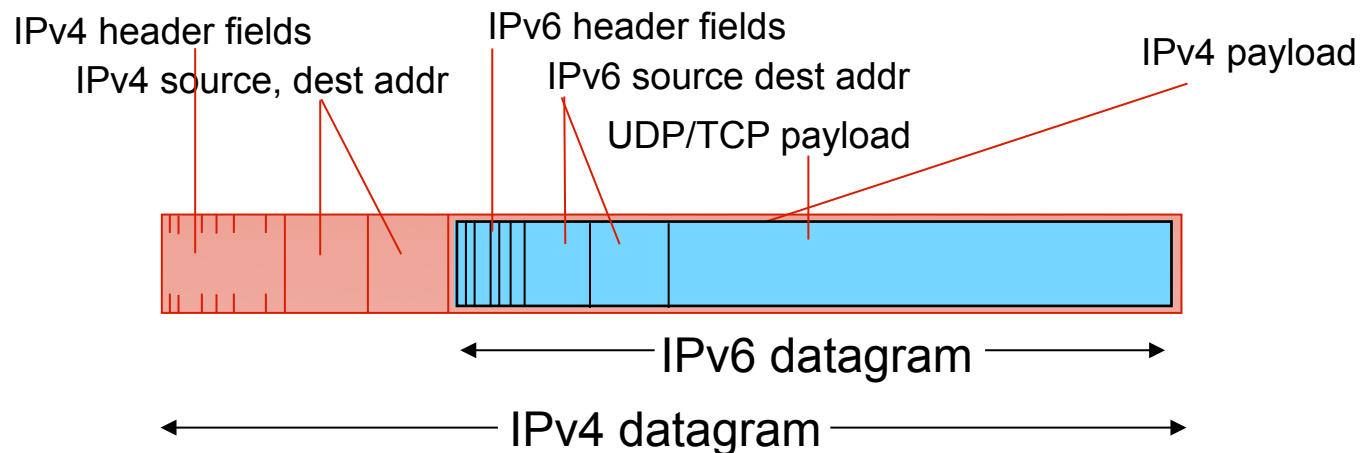


Other changes from IPv4

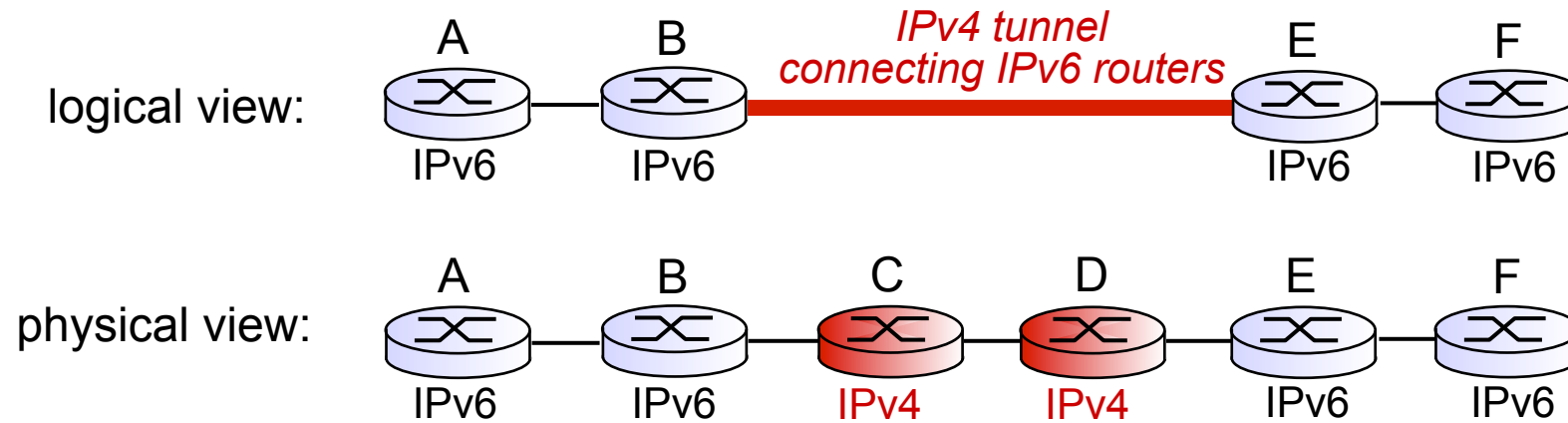
- ❑ *checksum*: removed entirely to reduce processing time at each hop
- ❑ *options*: allowed, but outside of header, indicated by “Next Header” field
- ❑ *ICMPv6*: new version of ICMP
 - additional message types, e.g. “Packet Too Big”
 - multicast group management functions

Transition from IPv4 to IPv6

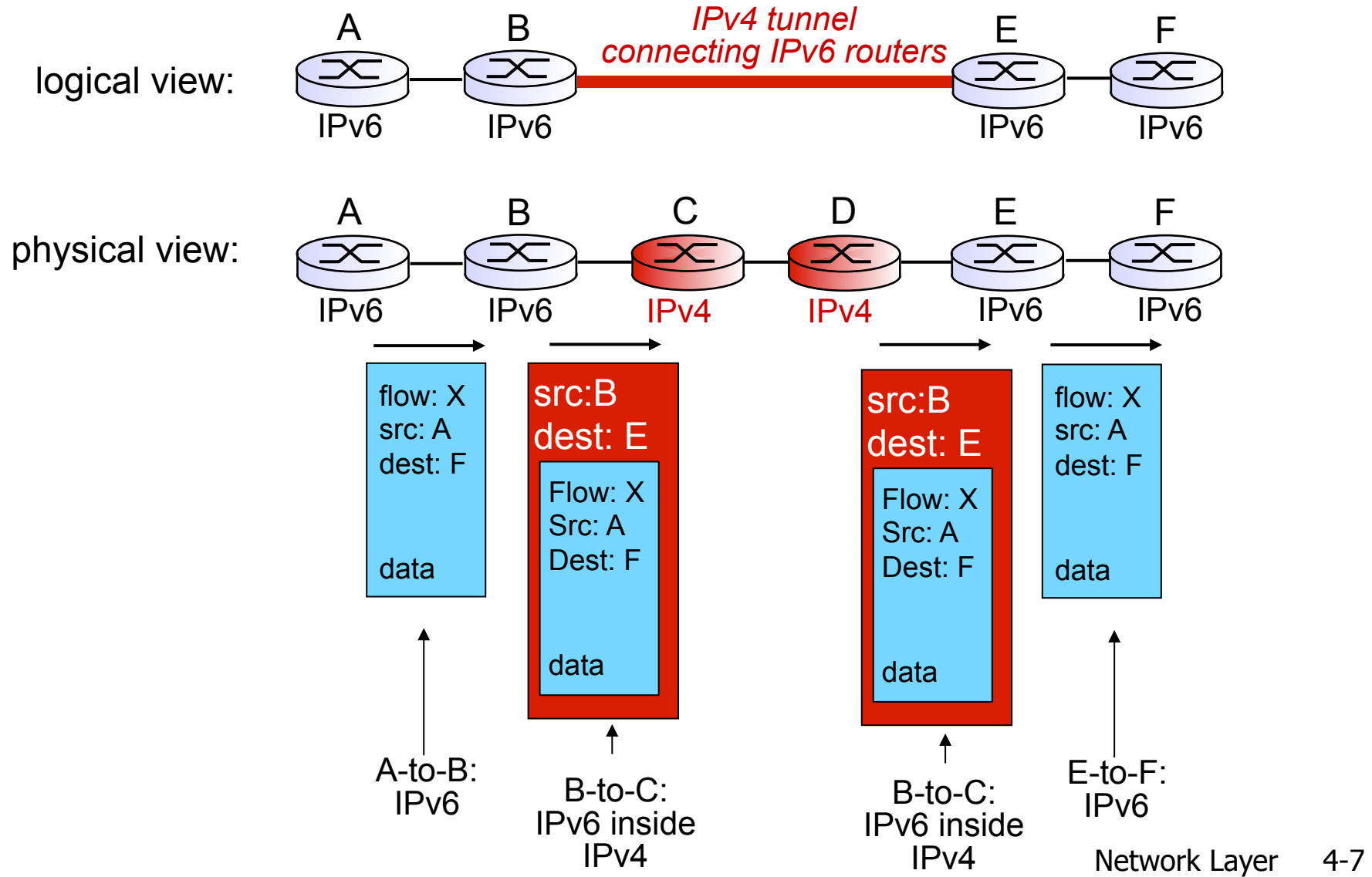
- not all routers can be upgraded simultaneously
 - no “flag days”
 - how will network operate with mixed IPv4 and IPv6 routers?
- **tunneling**: IPv6 datagram carried as *payload* in IPv4 datagram among IPv4 routers



Tunneling



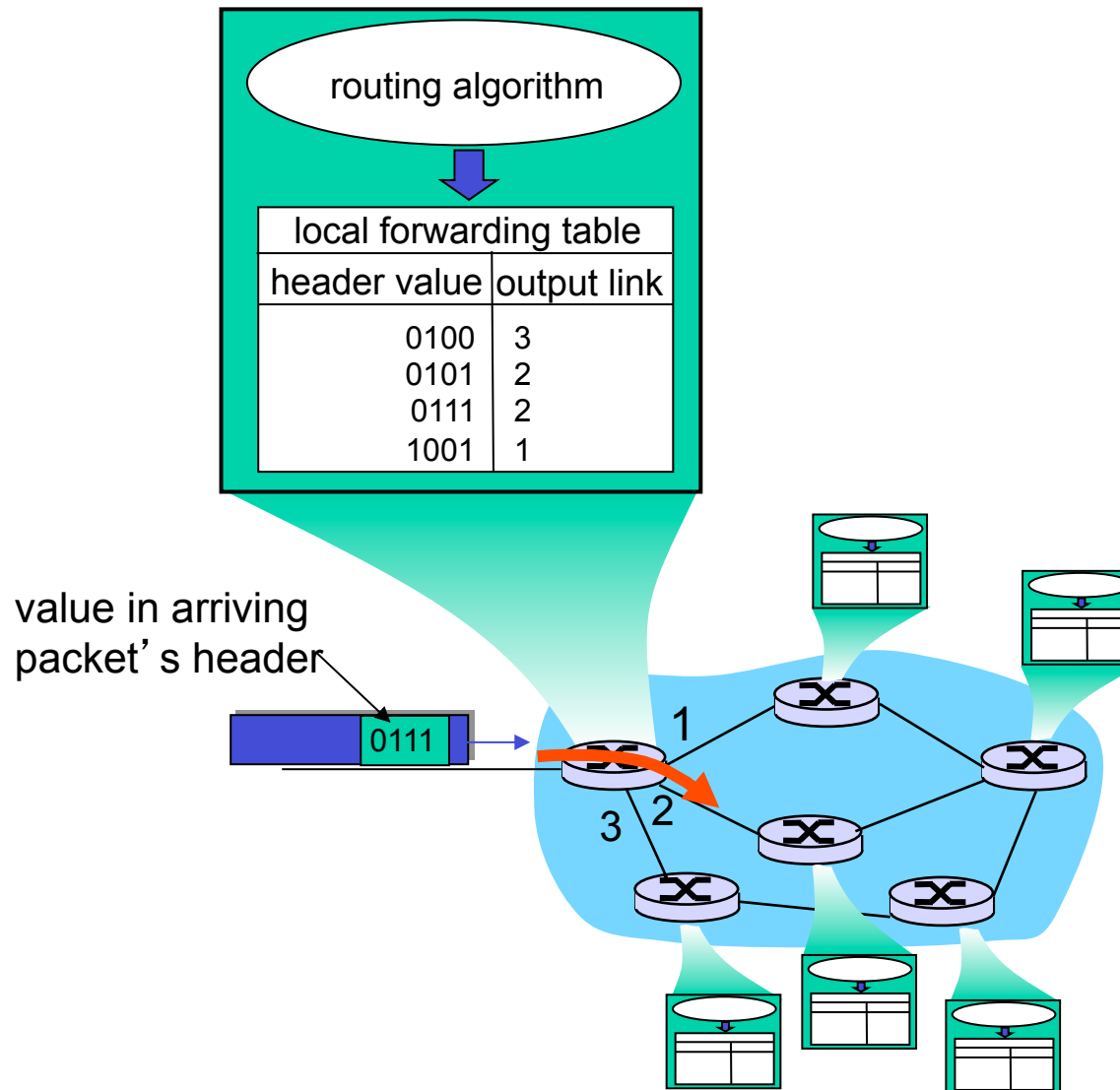
Tunneling



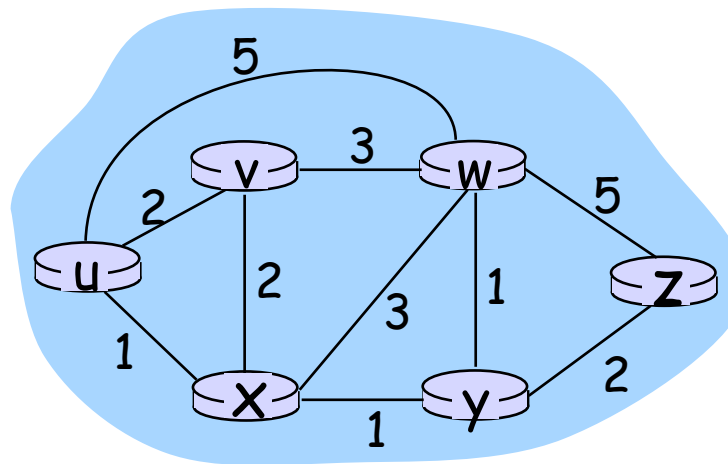
Chapter 4: Network Layer

- ❑ 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



Graph: $G = (N, E)$

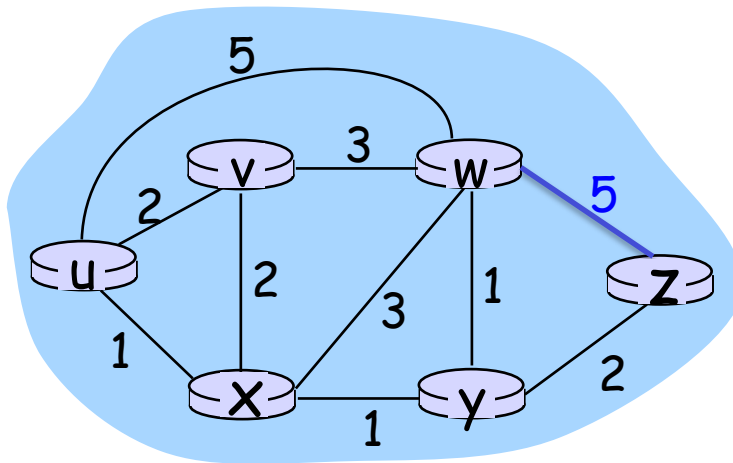
$N = \text{set of routers} = \{ u, v, w, x, y, z \}$

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

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



- $c(x,x')$ = 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

Routing Algorithm classification

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

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:

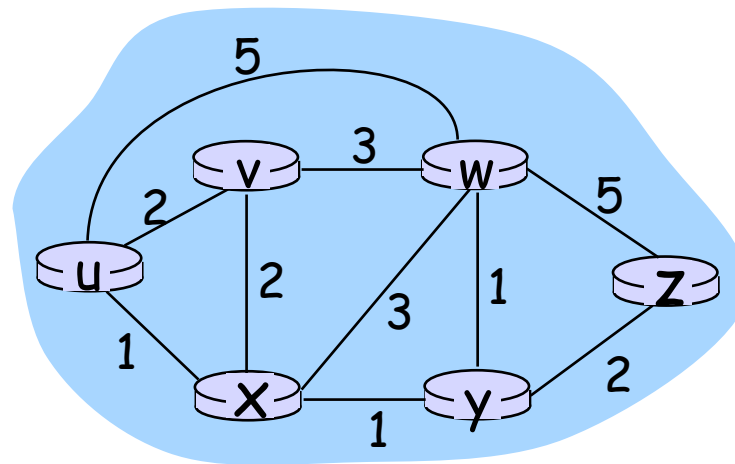
- ❑ $c(x,y)$: link cost from node x to y ; $= \infty$ if not direct neighbors
- ❑ $D(v)$: current value of cost of path from source to dest. v
- ❑ $p(v)$: predecessor node along path from source to v
- ❑ N' : set of nodes whose least cost path definitively known
- ❑ $c(x,x)=0$.

Dijkstra's Algorithm

```
1 Initialization:  
2  $N' = \{u\}$  u e' la sorgente  
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 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  $N'$ 
```

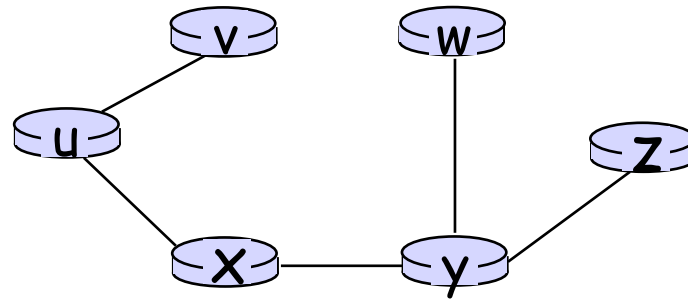
Dijkstra's algorithm: example

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	∞	∞
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)

Dijkstra Algorithm-correctness

Teorema

Se l'algoritmo di Dijkstra è eseguito su un grafo $G=(N,E)$ diretto e pesato i cui pesi sugli archi sono tutti non negativi, e con una sorgente s , allora Dijkstra termina con tutti i vertici w in N con valore $D(w)$ pari alla lunghezza del cammino minimo da s a w .

Approccio:

- Terminazione banale (perchè?)
- Mostriamo che ogni volta che un vertice w è inserito in N' allora $D(w)$ è pari alla lunghezza del cammino minimo da s a w

(dato che il valore di $D(w)$ non viene ad essere mai più modificato dopo che w è inserito in N' -v. linea 11 dell'algoritmo- questo consente di dimostrare l'assunto).

Dijkstra's Algorithm

```
1 Initialization:  
2  $N' = \{u\}$  u e' la sorgente  
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 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 Algorithm-correctness

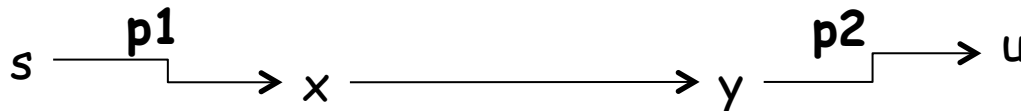
Mostriamo che ogni volta che un vertice w è inserito in N' allora $D(w)$ è pari alla lunghezza del cammino minimo da s a w

Si dimostra per assurdo.

Sia u il primo nodo inserito in N' che non rispetta la condizione, ovvero sia u il primo nodo che al momento del suo inserimento in N' abbia $D(u) \neq \delta(s,u)$, dove $\delta(s,u)$ denota la lunghezza del percorso minimo da s a u .

$u \neq s$ dato che s è il primo nodo inserito in N' e $D(s) = \delta(s,s) = 0$.

Ci deve essere un percorso da s a u in G (altrimenti $D(u) = \text{infinity}$, e $D(u) = \delta(s,u)$) e quindi anche un cammino minimo p da s a u . Prima di aggiungere u a N' p univa un vertice in N' (il vertice s) and un vertice in $N-N'$ (il vertice u). Sia y il primo vertice in p non in N' , e x il suo predecessore.



Dijkstra Algorithm-correctness

Mostriamo che ogni volta che un vertice w è inserito in N' allora $D(w)$ è pari alla lunghezza del cammino minimo da s a w

Si dimostra per assurdo.

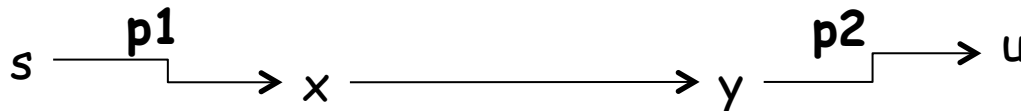
Sia y il primo vertice in p non in N' , e x il suo predecessore.

Vogliamo far vedere che $D(y) = \delta(s, y)$ quando u è stato aggiunto in N' .

Dato che il primo nodo che è aggiunto in N' senza rispettare l'assunto è u , e che x è aggiunto prima in N' , allora al momento del suo inserimento $D(x) = \delta(s, x)$

Dato che p è uno shortest path da s ad u allora anche il percorso p_1 è uno shortest path da s a y (perchè??)

Quindi quando è stato aggiunto x in N' e si sono ricalcolate le distanze dei percorsi per raggiungere s dai vicini di x passando tramite x , il valore $D(y)$ è stato aggiornato in modo che $D(y) = \delta(s, y)$



Dijkstra Algorithm-correctness

Mostriamo che ogni volta che un vertice w è inserito in N' allora $D(w)$ è pari alla lunghezza del cammino minimo da s a w

Si dimostra per assurdo.

Vogliamo ora dimostrare che $D(u) = \delta(s, u)$

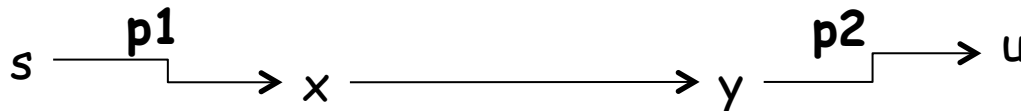
Dato che y viene prima di u in un percorso da s a u e tutti gli archi hanno pesi non negativi vale che $\delta(s, y) \leq \delta(s, u)$ e quindi

$$\begin{aligned} D(y) &= \delta(s, y) \\ &\leq \delta(s, u) \\ &\leq D(u) \end{aligned}$$

D'altra parte dato che sia u che y erano in $N - N'$ quando u è stato scelto per essere inserito in N' , al momento del suo inserimento

$D(u) \leq D(y)$ (per la regola di selezione del vertice da inserire in N')

Quindi: $D(y) = \delta(s, y) = \delta(s, u) = D(u)$ **ASSURDO C.V.D**



Dijkstra's Algorithm

```
1 Initialization:  
2  $N' = \{u\}$  u e' la sorgente  
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 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  $N'$ 
```

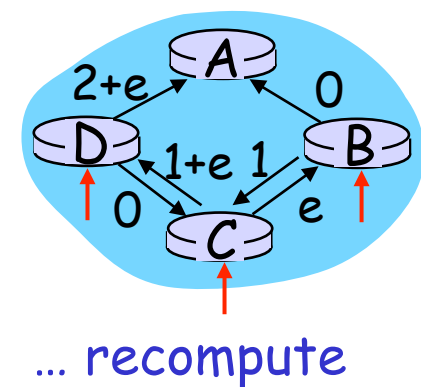
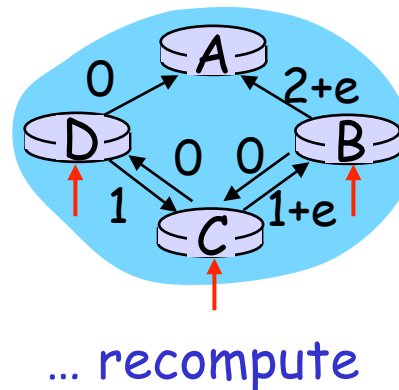
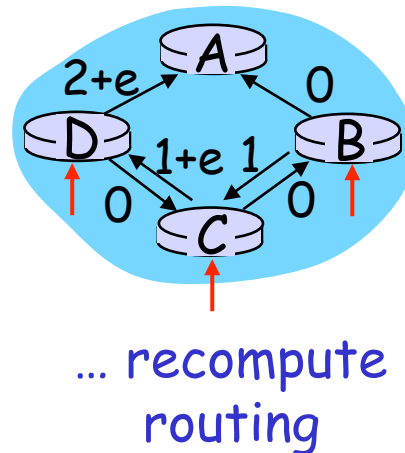
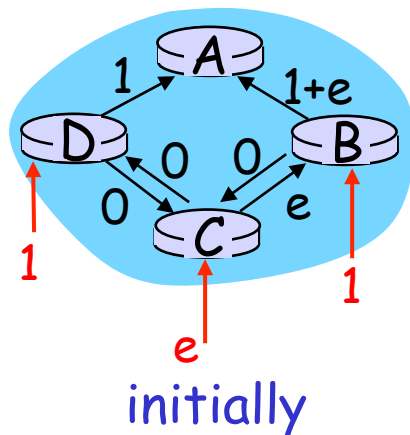
Dijkstra's algorithm, discussion

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: $O(n \log n + |E|)$

Oscillations possible:

- ❑ e.g., link cost = 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^h_i .

Bellman-Ford rule:

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

Iteration:

$$D^{h+1}_i = \min_k [c_{i,k} + D^h_k]$$

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

Bellman-Ford

$$D^{h+1}_i = \min_k [c_{i,k} + D^h_k]$$

Can be computed locally.

What do I need?

For each neighbor k , I need to know

-the cost of the link to it (known info)

-The cost of the best route from the neighbor k to the destination
(←this is an info that each of my neighbor has to send to me via messages)

In the real world: I need to know the best routes among each pair of nodes → we apply distributed Bellman Ford to get the best route for each of the possible destinations

Distance Vector Routing Algorithm -Distributed Bellman Ford

iterative:

- ❑ continues until no nodes exchange info.
- ❑ *self-terminating*: no “signal” to stop

asynchronous:

- ❑ nodes need *not* exchange info/iterate in lock step!

Distributed, based on local info:

- ❑ each node communicates *only* with directly-attached neighbors

Distance Table data structure

each node has its own

- ❑ row for each possible destination
- ❑ column for each directly-attached neighbor to node
- ❑ example: in node X, for dest. Y via neighbor V: what is the cost?

Distance vector algorithm

Bellman-Ford equation (dynamic programming)

let

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

then

Info maintained at v . Min must be communicated

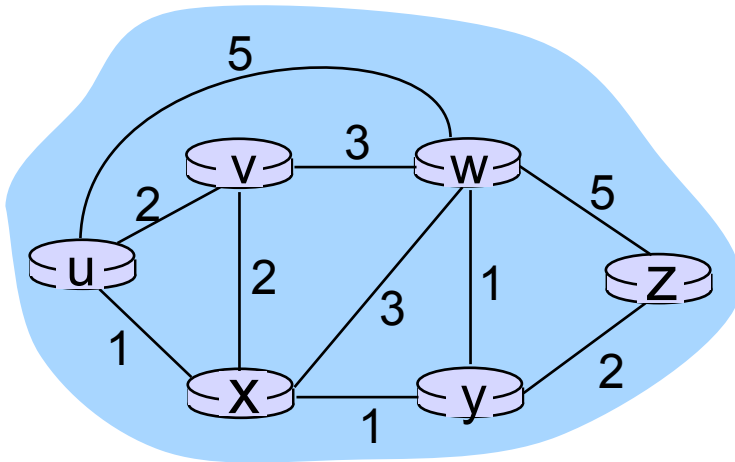
$$d_x(y) = \min_v \{ c(x,v) + d_v(y) \}$$

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 $\mathbf{D}_x = [D_x(y): y \in N]$
- node x:
 - knows cost to each neighbor v: $c(x,v)$
 - maintains its neighbors' distance vectors. For each neighbor v, x maintains $\mathbf{D}_v = [D_v(y): y \in N]$

Distance vector algorithm

key idea:

- from time-to-time, each node sends its own distance vector estimate to neighbors
- when x receives new DV estimate from neighbor, it updates its own DV using B-F equation:

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

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

Distance vector algorithm

iterative, asynchronous:

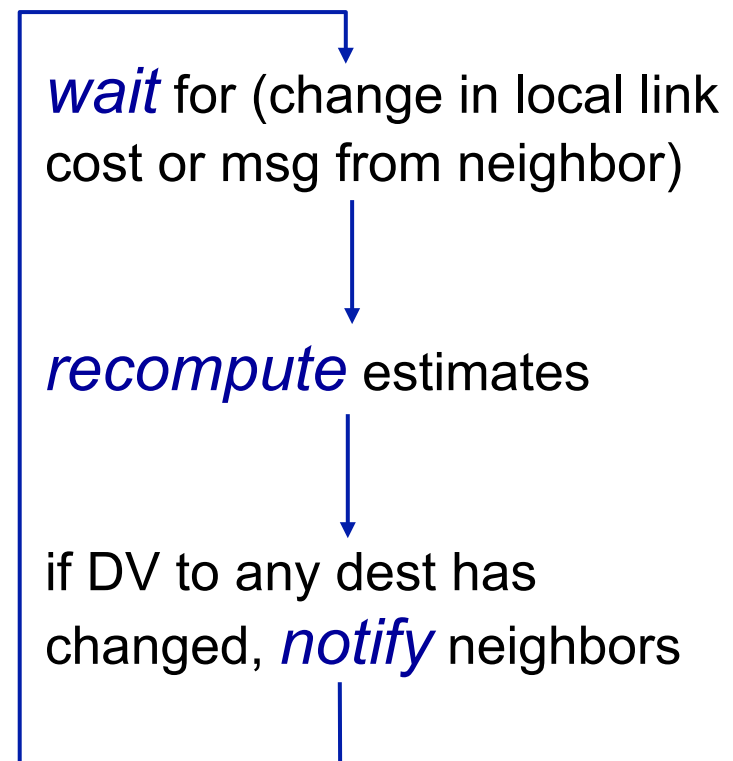
each local iteration
caused by:

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

distributed:

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

each node:



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

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

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

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

node x table

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

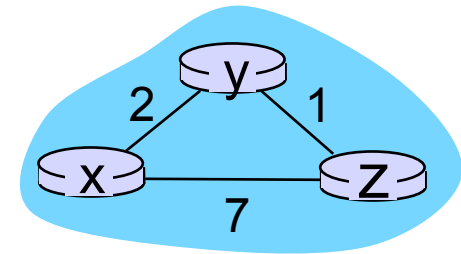
		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	7	1	0

node y table

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

node z table

		cost to		
		x	y	z
from	x	∞	∞	∞
	y	∞	∞	∞
	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	∞	∞	∞
	z	∞	∞	∞

node y table

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

node z table

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

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

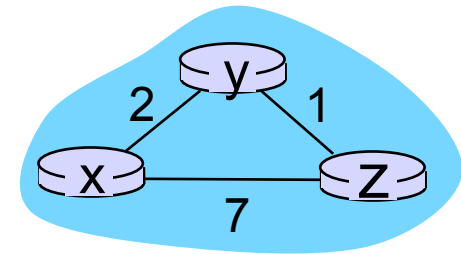
		cost to		
		x	y	z
from	x	0	2	7
	y	2	0	1
	z	7	1	0

		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

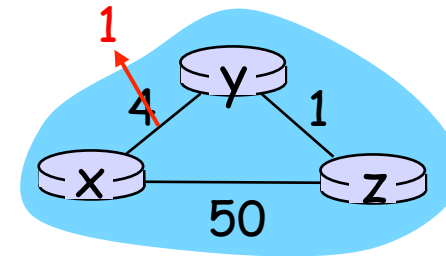


time →

Distance vector: link cost changes

link cost changes:

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



“good
news
travels
fast”

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

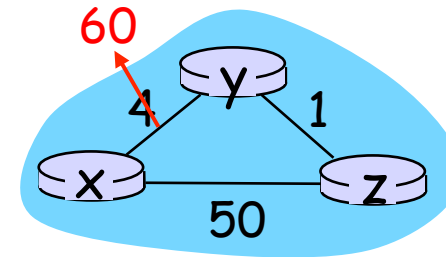
t_1 : z receives update from y , updates its table, computes new least cost to x , sends its neighbors its DV.

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

Distance vector: link cost changes

link cost changes:

- ❖ node detects local link cost change
- ❖ *bad news travels slow* - “count to infinity” problem!
- ❖ 44 iterations before algorithm stabilizes: 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?

Distributed Bellman Ford-Count to Infinity (we will now use a slightly different notation- lightweigh)

Distance Table data structure

each node has its own

- row for each possible destination
- column for each directly-attached neighbor to node
- example: in node X, for dest. Y via neighbor Z:

Cost associated to the (X,Z) link

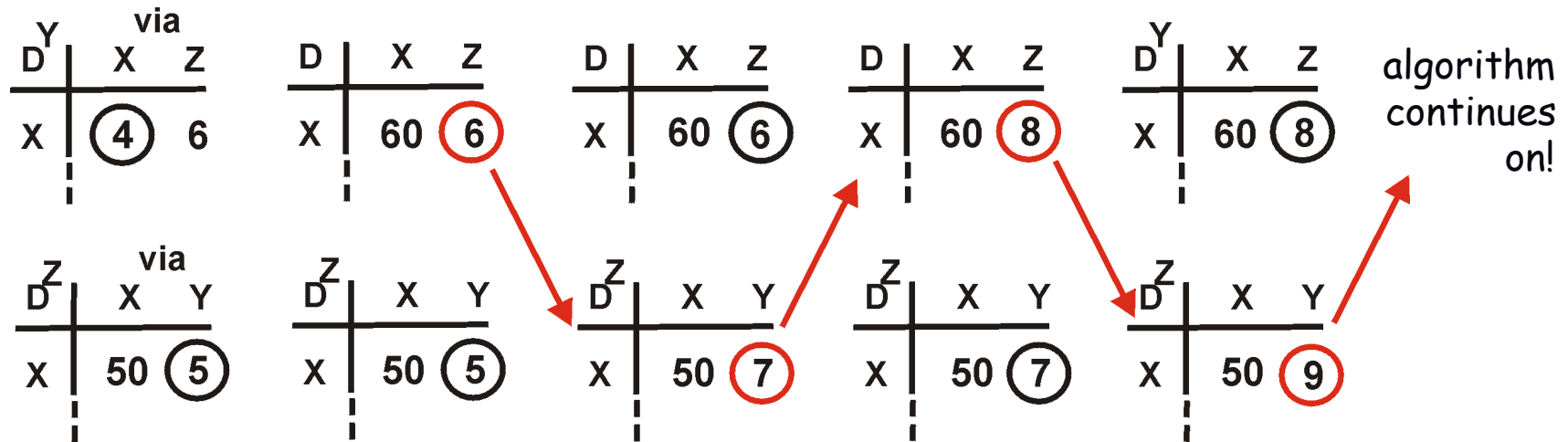
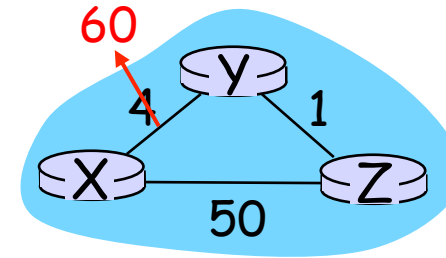
$$D_{(Y,Z)}^X = \text{distance from X to Y, via Z as next hop}$$
$$= c(X,Z) + \min_w \{D^Z(Y,w)\}$$

Info maintained at Z. Min must be communicated

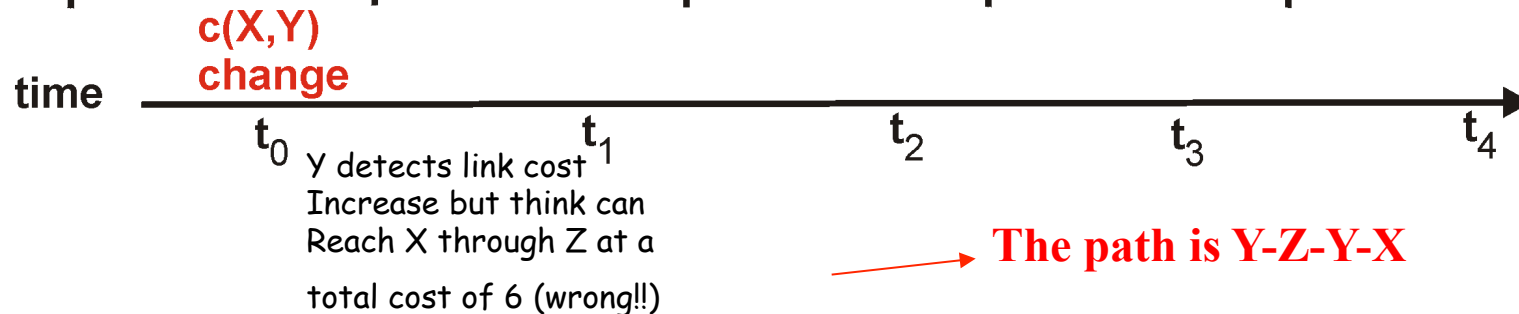
Distance Vector: link cost changes

Link cost changes:

- good news travels fast
- ***bad news travels slow*** - "count to infinity" problem!



algorithm continues on!



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

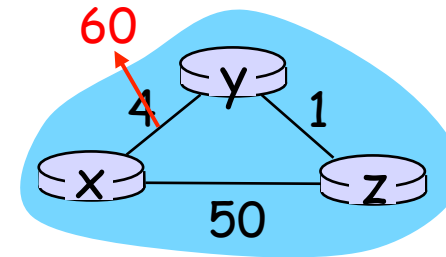
Colosseo that updates ‘OK I can go to P.Popolo via Al. Patria in 4Km’

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)!!

Distance vector: link cost changes

link cost changes:

- ❖ node detects local link cost change
- ❖ *bad news travels slow* - “count to infinity” problem!
- ❖ 44 iterations before algorithm stabilizes: 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

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