# Routing Algorithms

- IP v6
- Routing algorithms
  - Link state
  - Distance Vector

# IPv6

❖ **Initial motivation:** 32-bit address space soon to be completely allocated.

❖ Additional motivation:

  ▪ header format helps speed processing/forwarding

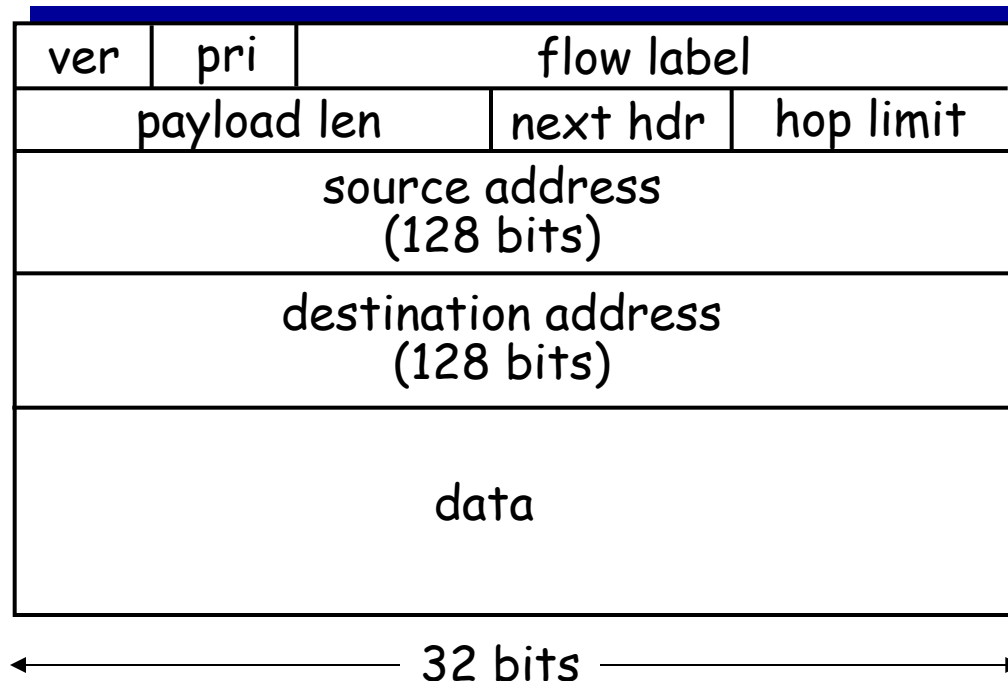  ▪ header changes to facilitate QoS

  IPv6 datagram format:

  ▪ fixed-length 40 byte header

  ▪ no fragmentation allowed

# IPv6 Header (Cont)

*Priority:* identify priority among datagrams in flow
*Flow Label:* identify datagrams in same "flow."
        (concept of "flow" not well defined).
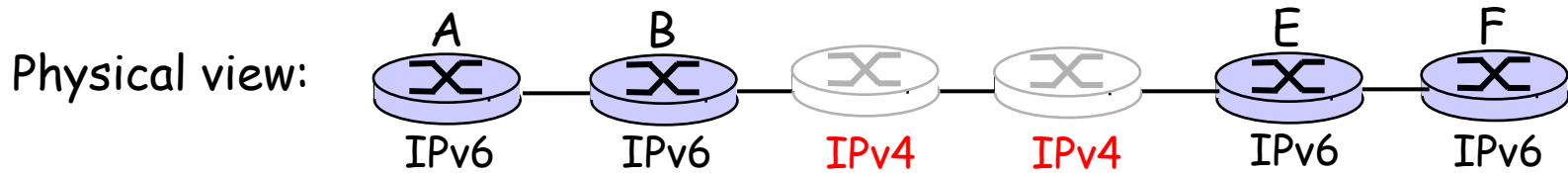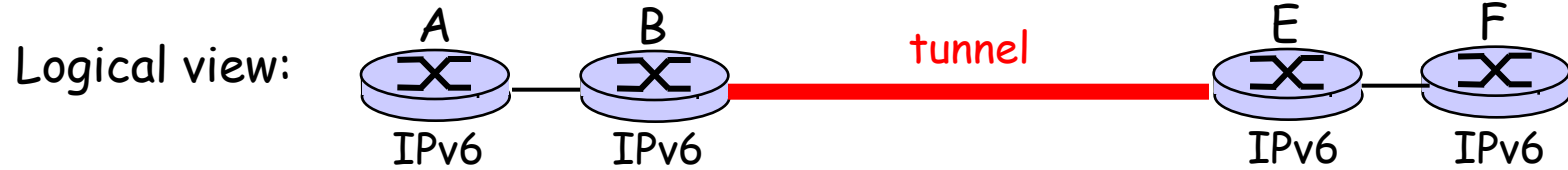*Next header:* identify upper layer protocol for data

| ver | pri | flow label | | |
|---|---|---|---|---|
| payload len | | | next hdr | hop limit |
| source address<br>(128 bits) | | | | |
| destination address<br>(128 bits) | | | | |
| 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"
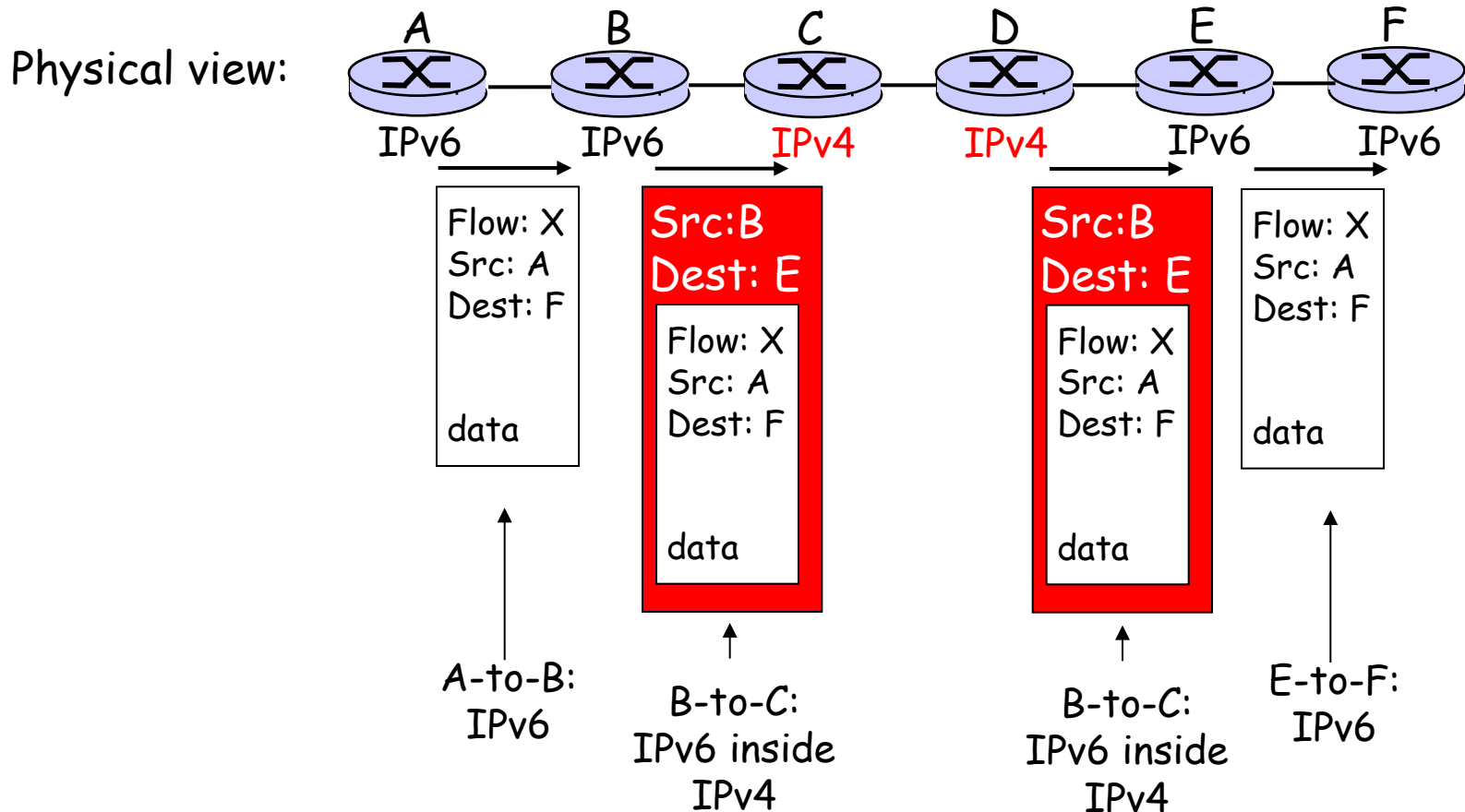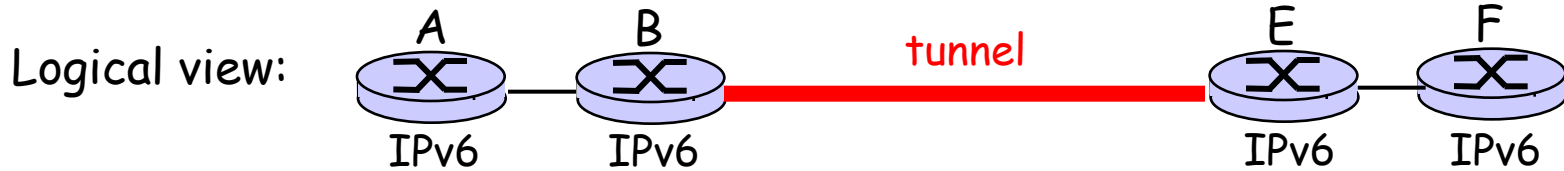  - ▪ multicast group management functions

# Transition From IPv4 To IPv6

❖ Not all routers can be upgraded simultaneous
  ▪ no "flag days"
  ▪ How will the network operate with mixed IPv4 and IPv6 routers?

❖ *Tunneling:* IPv6 carried as payload in IPv4 datagram among IPv4 routers

# Tunneling

Logical view:
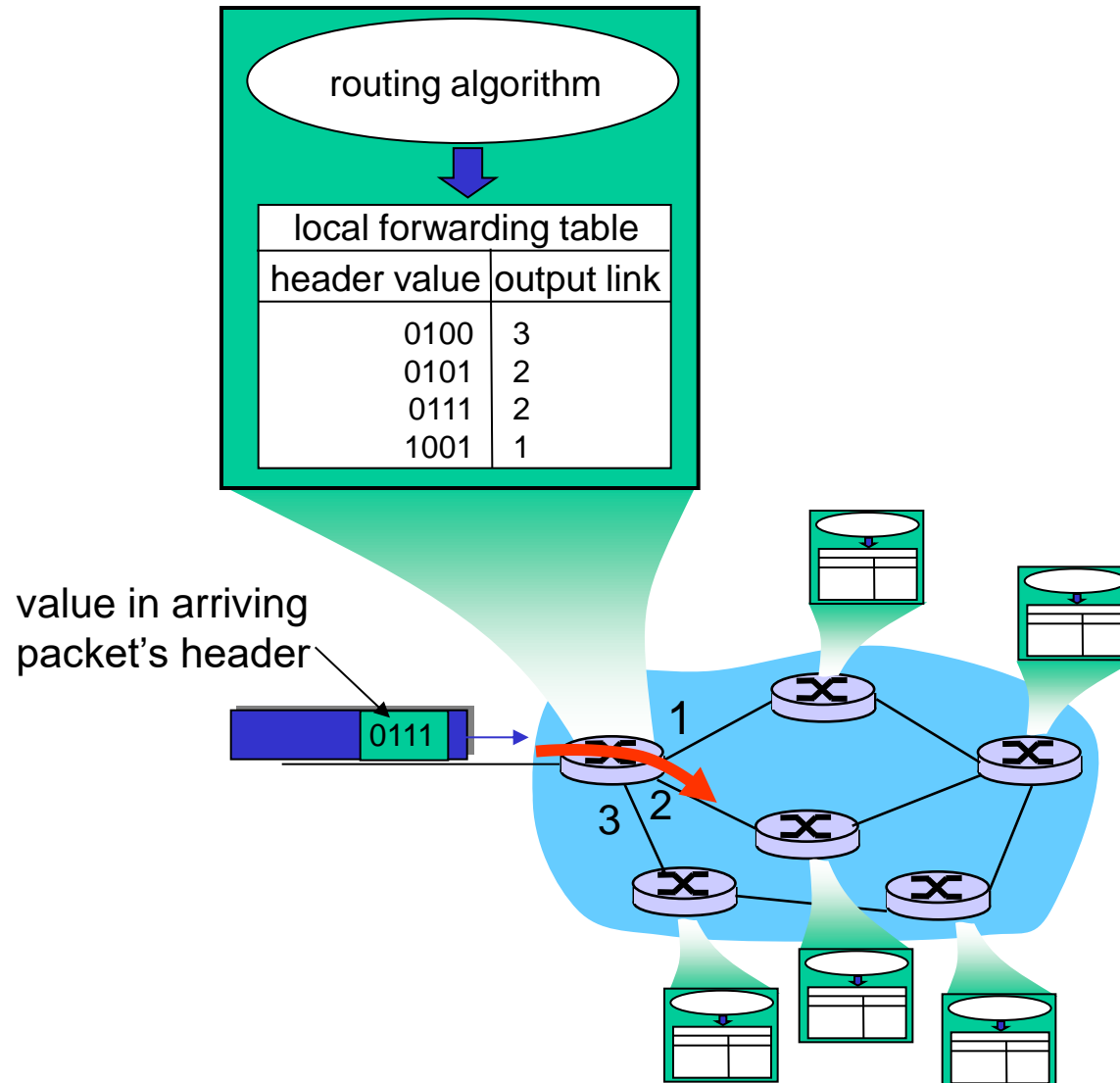
A — B ════ tunnel ════ E — F
IPv6  IPv6              IPv6  IPv6

Physical view:

A — B — (IPv4) — (IPv4) — E — F
IPv6  IPv6   IPv4    IPv4    IPv6  IPv6

# Tunneling

# Network Layer

Routing algorithms
- Link state
- Distance Vector

# Interplay between routing, forwarding



routing algorithm

| local forwarding table | |
|---|---|
| header value | output link |
| 0100 | 3 |
| 0101 | 2 |
| 0111 | 2 |
| 1001 | 1 |

value in arriving
packet's header

0111

1

3   2

# Graph abstraction



Graph: G = (N,E)
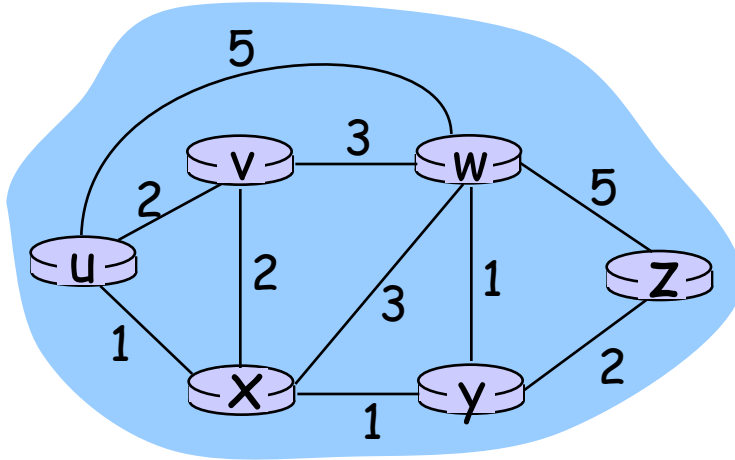
N = set of routers = { u, v, w, x, y, z }

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

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 $(x_1, x_2, x_3,..., x_p) = c(x_1,x_2) + c(x_2,x_3) + ... + c(x_{p-1},x_p)$

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

# Network Layer

Routing algorithms
- Link state
- Distance Vector

# 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

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) = ∞
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
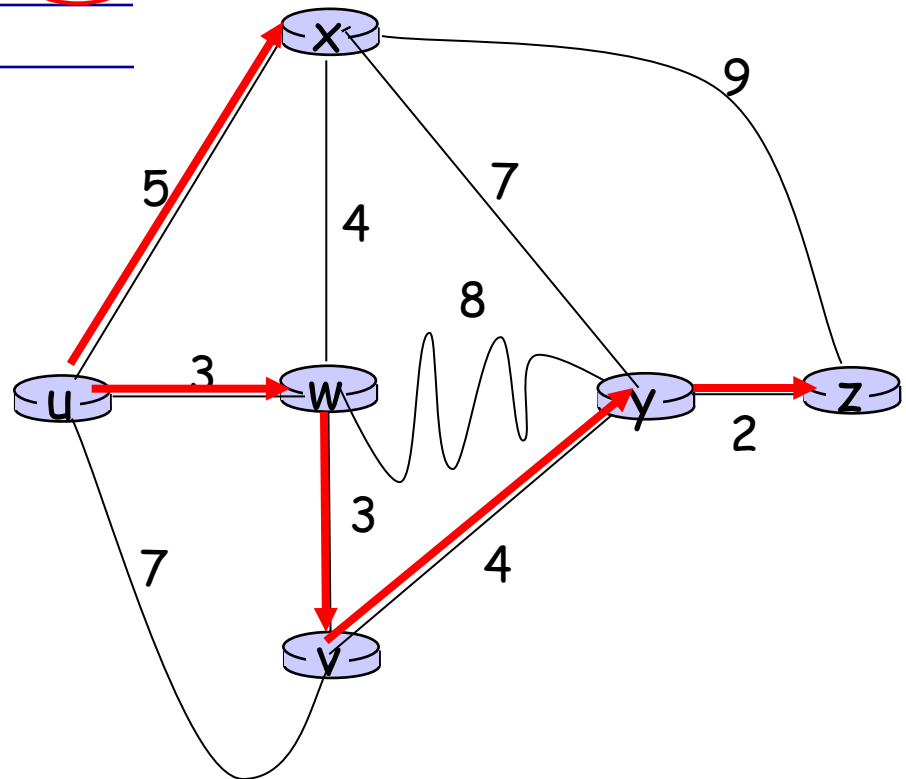14      shortest path cost to w plus cost from w to v */
15   *until all nodes in N'*

# Dijkstra's algorithm: example

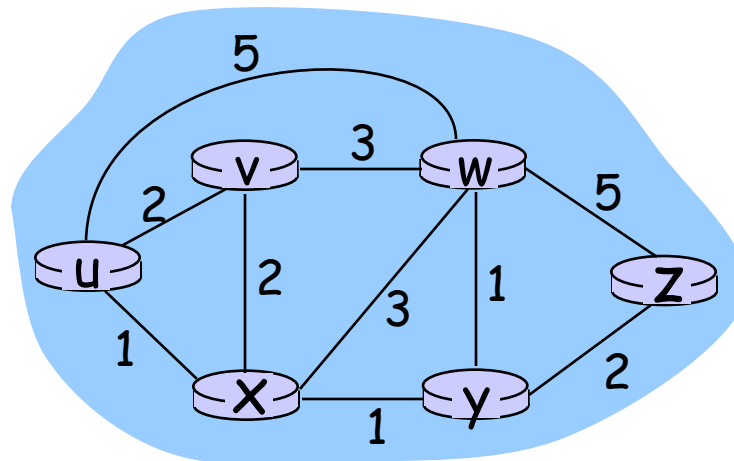| Step | N' | D(v) p(v) | D(w) p(w) | D(x) p(x) | D(y) p(y) | D(z) p(z) |
|------|------|-----------|-----------|-----------|-----------|-----------|
| 0 | u | 7,u | ③,u | 5,u | ∞ | ∞ |
| 1 | uw | 6,w | | ⑤,u | 11,w | ∞ |
| 2 | uwx | ⑥,w | | | 11,w | 14,x |
| 3 | uwxv | | | | ⑩,v | 14,x |
| 4 | uwxvy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

## Notes:

❖ construct shortest path tree by tracing predecessor nodes

❖ ties can exist (can be broken arbitrarily)

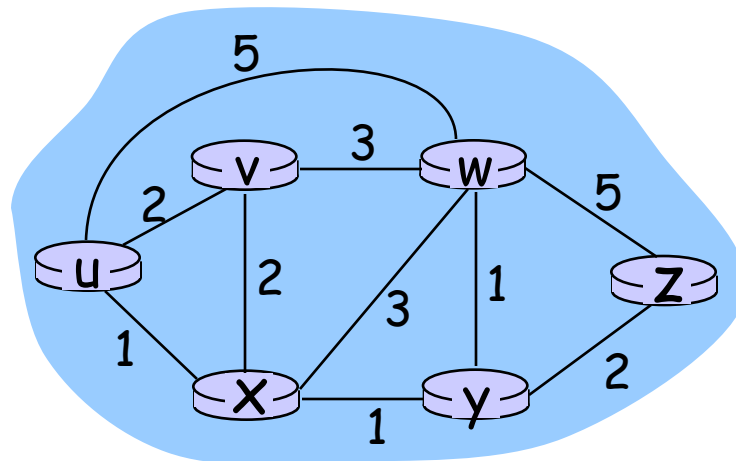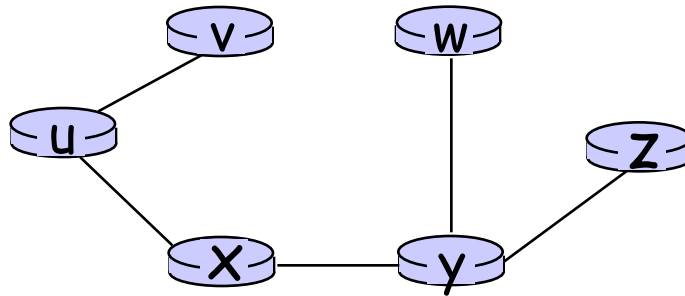# Dijkstra's algorithm: another example

| 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 | | | 3,y | | 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) |

# Network Layer

## Routing algorithms

- Link state
- Distance Vector

# Distance Vector Algorithm

Bellman-Ford Equation (dynamic programming)
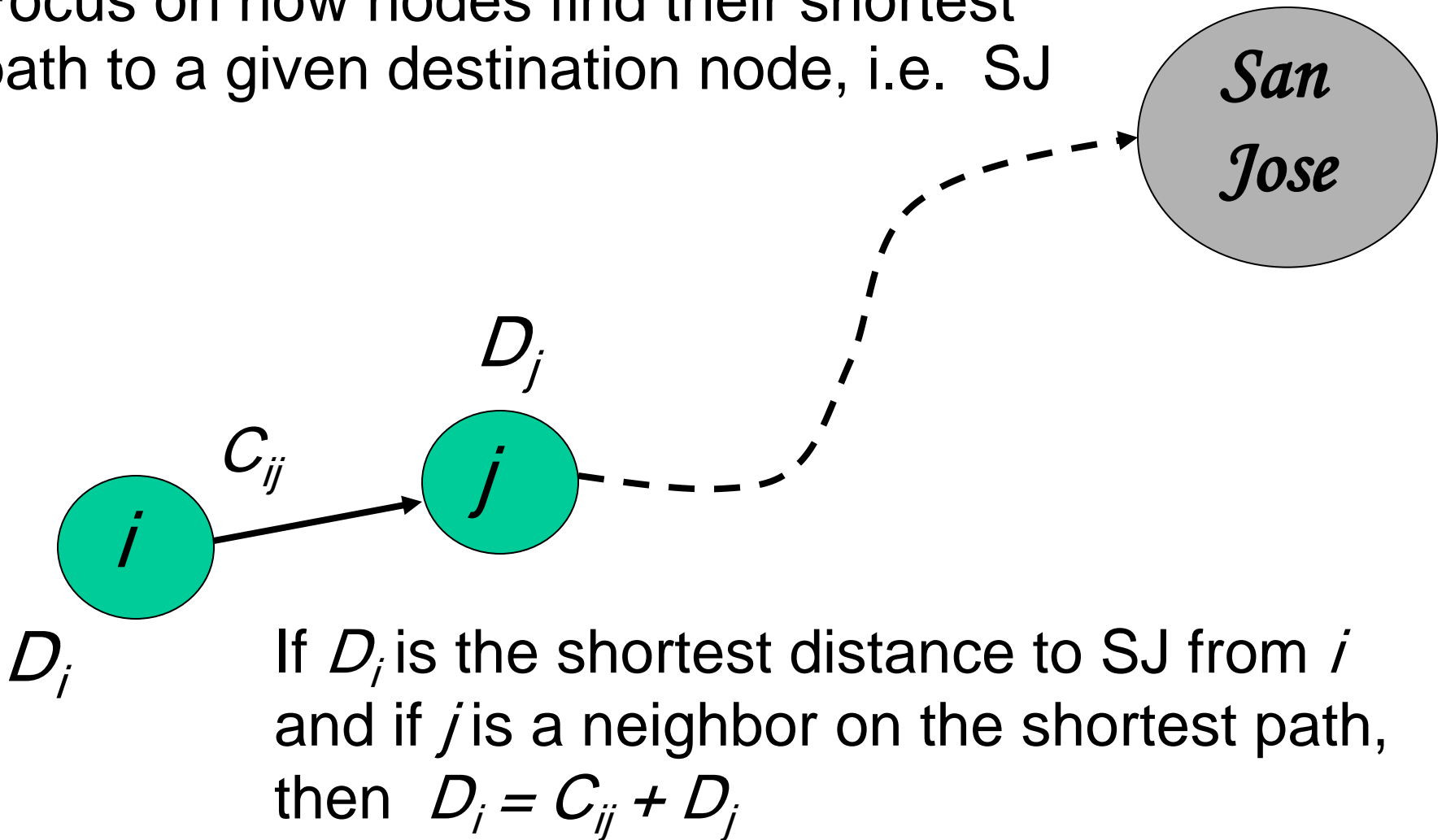Define
$d_x(y)$ := cost of least-cost path from x to y

Then

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

where min is taken over all neighbors v of x

# Shortest Path to SJ

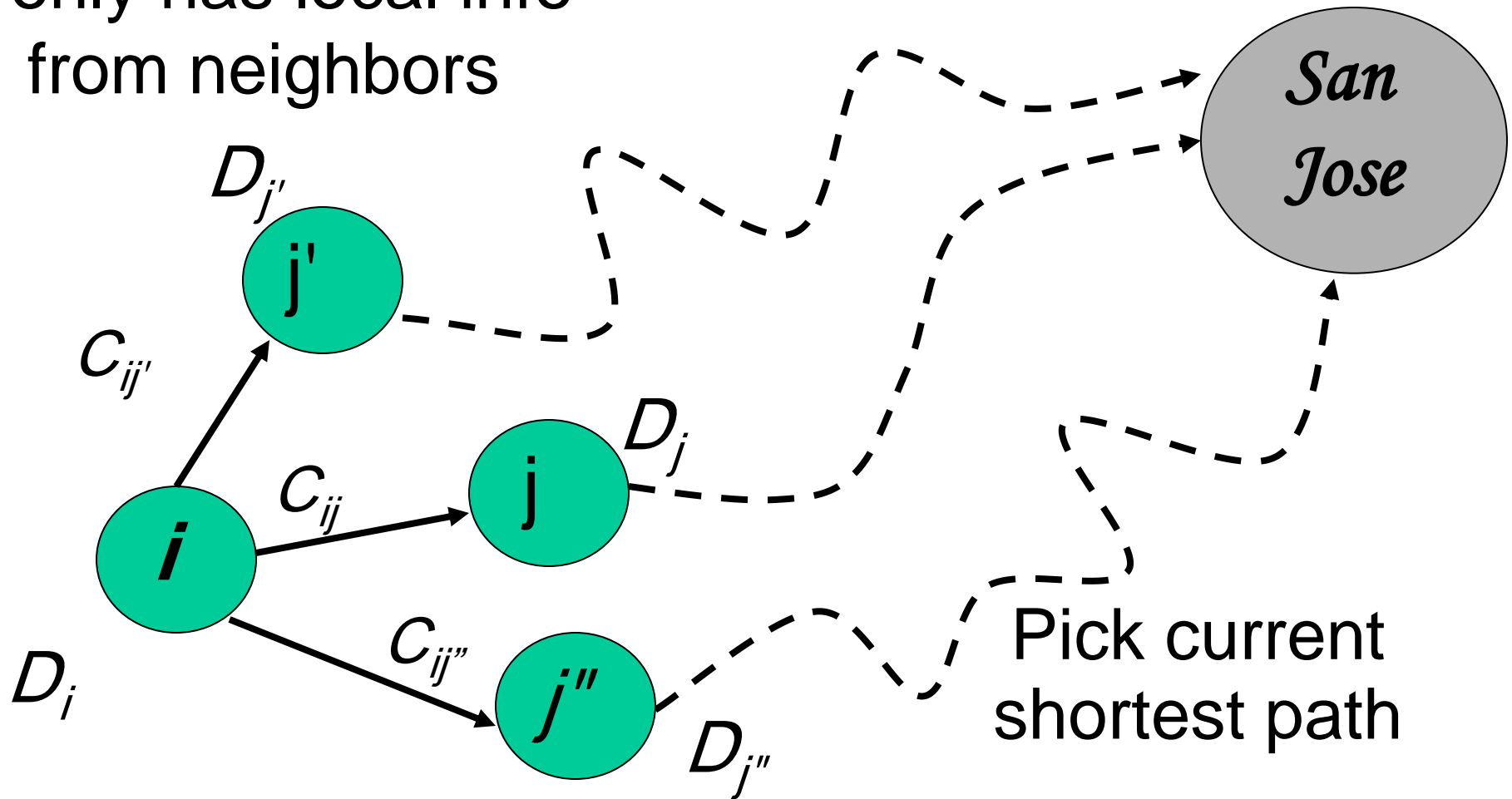Focus on how nodes find their shortest path to a given destination node, i.e. SJ

*San Jose*

$D_j$

$C_{ij}$

$i$ → $j$

$D_i$

If $D_i$ is the shortest distance to SJ from $i$ and if $j$ is a neighbor on the shortest path, then $D_i = C_{ij} + D_j$

# But we don't know the shortest paths

$i$ only has local info from neighbors



$D_{j'}$

$C_{ij'}$

$C_{ij}$

$C_{ij''}$

$D_i$

$D_j$

$D_{j''}$

*San Jose*

Pick current shortest path

# Why Distance Vector Works

SJ sends accurate info

San Jose

1 Hop From SJ

2 Hops From SJ

3 Hops From SJ

Hop-1 nodes calculate current (next hop, dist), & send to neighbors

Accurate info about SJ ripples across network, Shortest Path Converges

# Bellman-Ford example



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

B-F equation says:

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

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

   ▪ 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 (4)

**Basic 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:

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

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

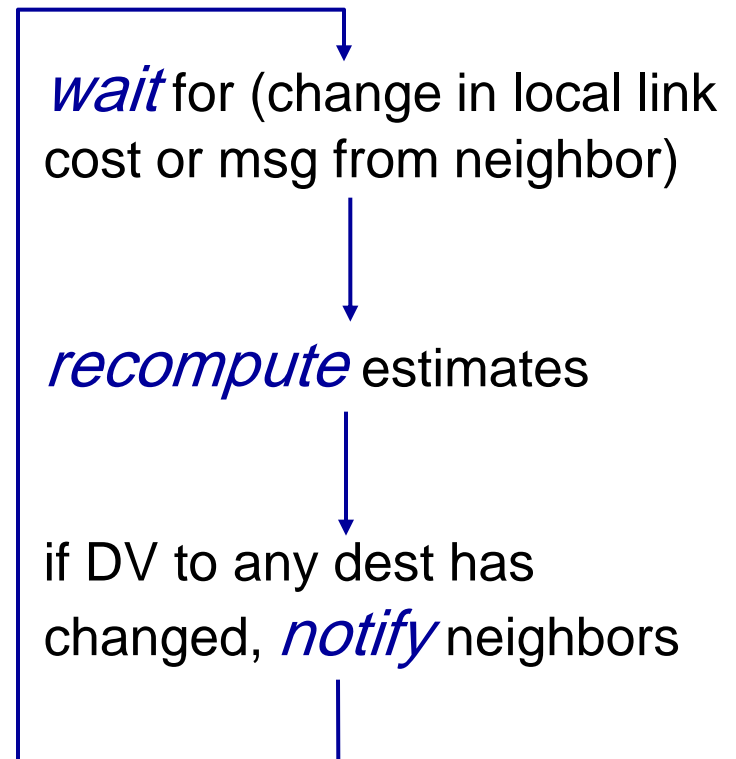# Distance Vector Algorithm (5)

**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:**

*wait* for (change in local link cost or msg from neighbor)

↓

*recompute* estimates

↓

if DV to any dest has changed, *notify* neighbors

# Distance Vector Algorithm

- $c(x,v)$ = cost for direct link from x to v
  - Node x maintains costs of direct links $c(x,v)$
- $D_x(y)$ = estimate of least cost from x to y
  - Node x maintains distance vector $\mathbf{D_x} = [D_x(y): y \in N]$
- Node x maintains its neighbors' distance vectors
  - For each neighbor v, x maintains $\mathbf{D_v} = [D_v(y): y \in N]$
- Each node v periodically sends $D_v$ to its neighbors
  - And neighbors update their own distance vectors
  - $D_x(y) \leftarrow \min_v\{c(x,v) + D_v(y)\}$    for each node $y \in N$
- Over time, the distance vector $D_x$ converges

```
1  Initialization:
2     for all destinations y in N:
3         Dx(y) = c(x,y)     /* if y is not a neighbor then c(x,y) = ∞ */
4     for each neighbor w
5         Dw(y) = ? for all destinations y in N
6     for each neighbor w
7         send distance vector Dx = [Dx(y): y in N] to w
8
9  loop
10    wait (until I see a link cost change to some neighbor w or
11          until I receive a distance vector from some neighbor w)
12
13    for each y in N:
14        Dx(y) = minv{c(x,v) + Dv(y)}
15
16    if Dx(y) changed for any destination y
17        send distance vector Dx = [Dx(y): y in N] to all neighbors
18
19 forever
```

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

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

**node x table**

cost to

|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 7 |
| y | ∞ | ∞ | ∞ |
| z | ∞ | ∞ | ∞ |

from

cost to

|  | x | y | z |
|---|---|---|---|
| x | 0 | 2 | 3 |
| y | 2 | 0 | 1 |
| z | 7 | 1 | 0 |

from

**node y table**

cost to

|  | x | y | z |
|---|---|---|---|
| x | ∞ | ∞ | ∞ |
| y | 2 | 0 | 1 |
| z | ∞ | ∞ | ∞ |

from

**node z table**

cost to

|  | x | y | z |
|---|---|---|---|
| x | ∞ | ∞ | ∞ |
| y | ∞ | ∞ | ∞ |
| z | 7 | 1 | 0 |

from



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 | ∞ | ∞ | ∞ |

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 | ∞ | ∞ | ∞ |
| y | 2 | 0 | 1 |
| z | ∞ | ∞ | ∞ |

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 | ∞ | ∞ | ∞ |
| y | ∞ | ∞ | ∞ |
| z | 7 | 1 | 0 |

cost to

|      | x | y | z |
|------|---|---|---|
| from x | 0 | 2 | 7 |
| y | 2 | 0 | 1 |
| z | 3 | 1 | 0 |

cost to

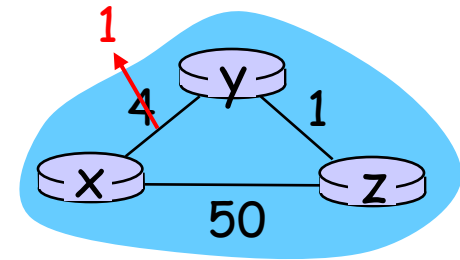|      | x | y | z |
|------|---|---|---|
| from x | 0 | 2 | 3 |
| y | 2 | 0 | 1 |
| z | 3 | 1 | 0 |

time

# Distance Vector: link cost changes

Link cost changes:

❖ node detects local link cost change

❖ updates routing info, recalculates distance vector

❖ if DV changes, notify neighbors

"good news travels fast"

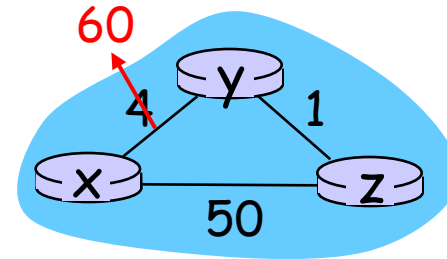$t_0$ : y detects link-cost change, updates its DV, informs its neighbors.

$t_1$ : 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:

❖ good news travels fast

❖ bad news travels slow - "count to infinity" problem!

❖ 44 iterations before algorithm stabilizes.

## 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?

# Routing Protocol

# 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
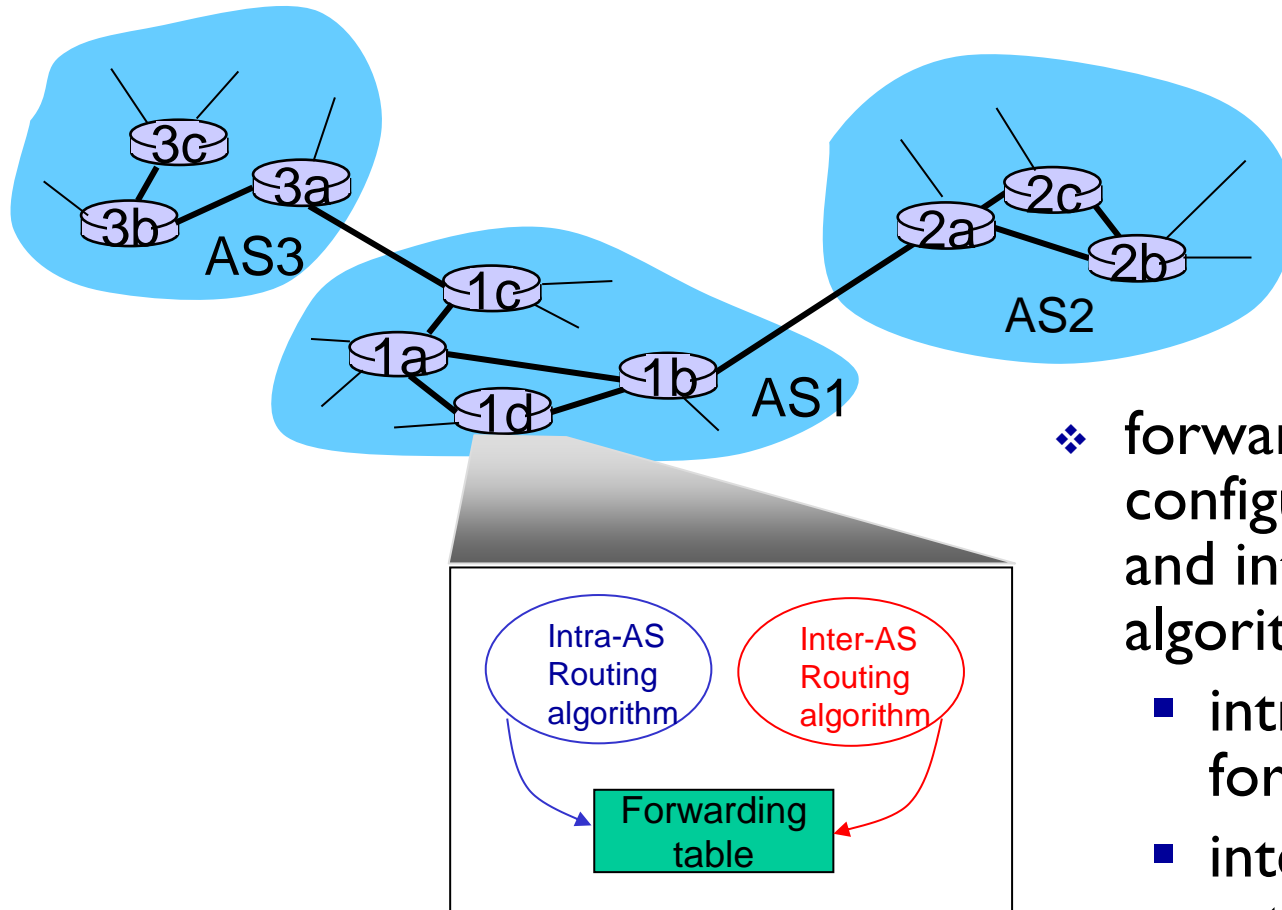- ❖ each network admin may want to control routing in its own network

# 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
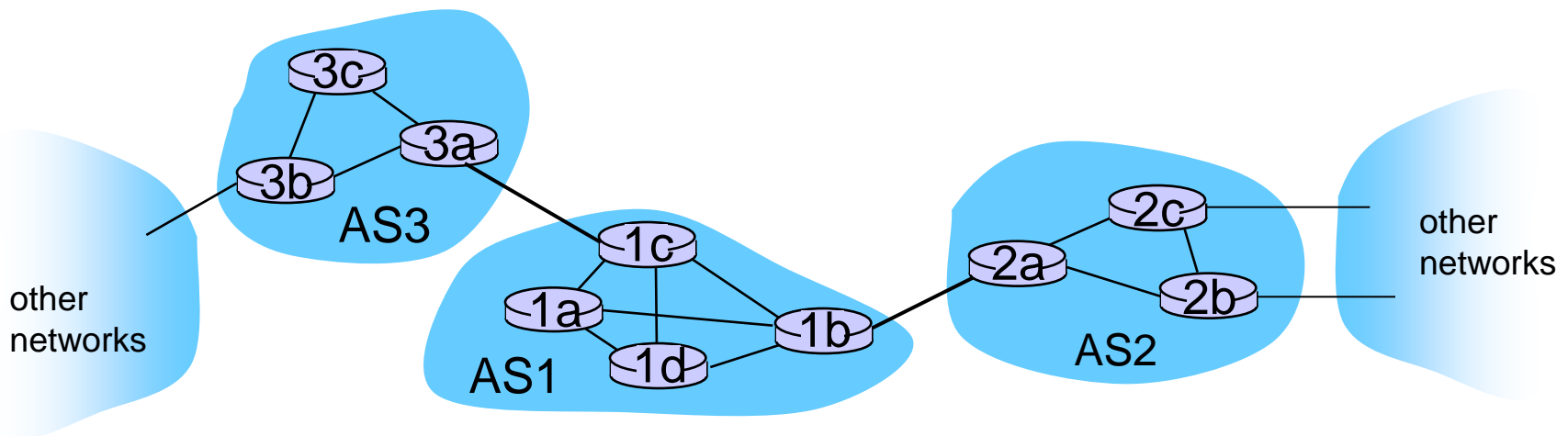    * inter-AS & intra-AS sets entries for external dests

# Inter-AS tasks

❖ suppose router in AS1 receives datagram destined outside of AS1:

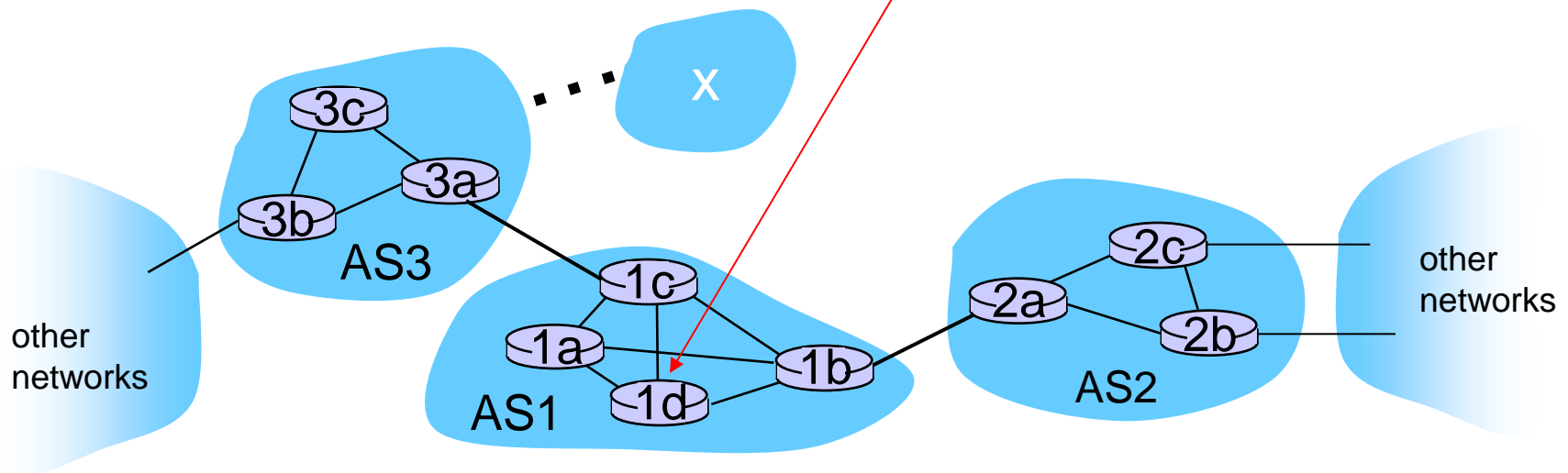  ▪ router should forward packet to gateway router, but which one?

*AS1 must:*

1. learn which dests are reachable through AS2, which through AS3

2. propagate this reachability info to all routers in AS1

*job of inter-AS routing!*

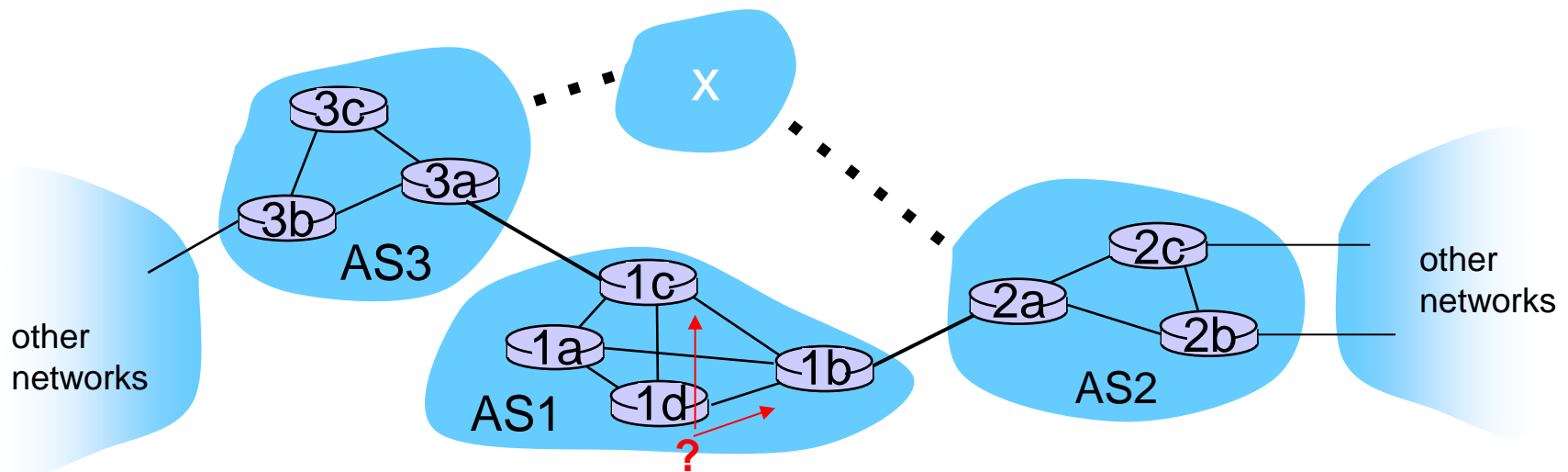# Example: setting forwarding table in router 1d

❖ 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 *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 which gateway it should forward packets towards for dest *x*
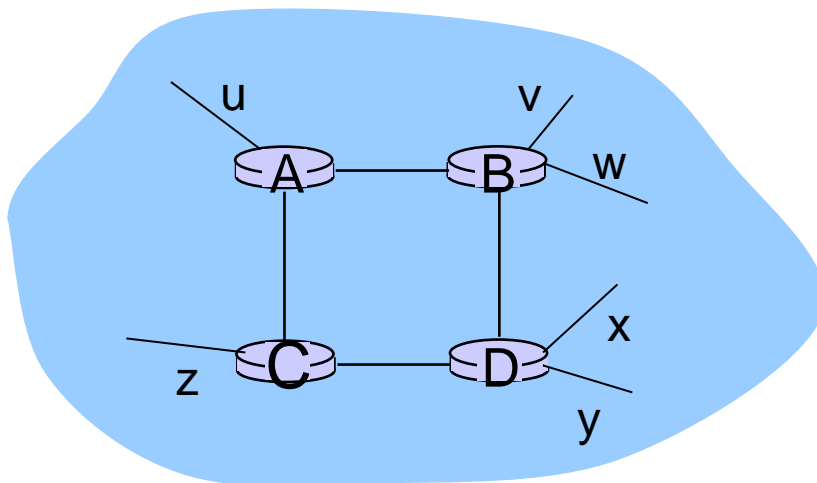
▪ this is also job of inter-AS routing protocol!

# Intra-AS Routing

❖ also known as *interior gateway protocols (IGP)*

❖ most common intra-AS routing protocols:

- RIP: Routing Information Protocol
- OSPF: Open Shortest Path First
- IGRP: Interior Gateway Routing Protocol (Cisco proprietary)
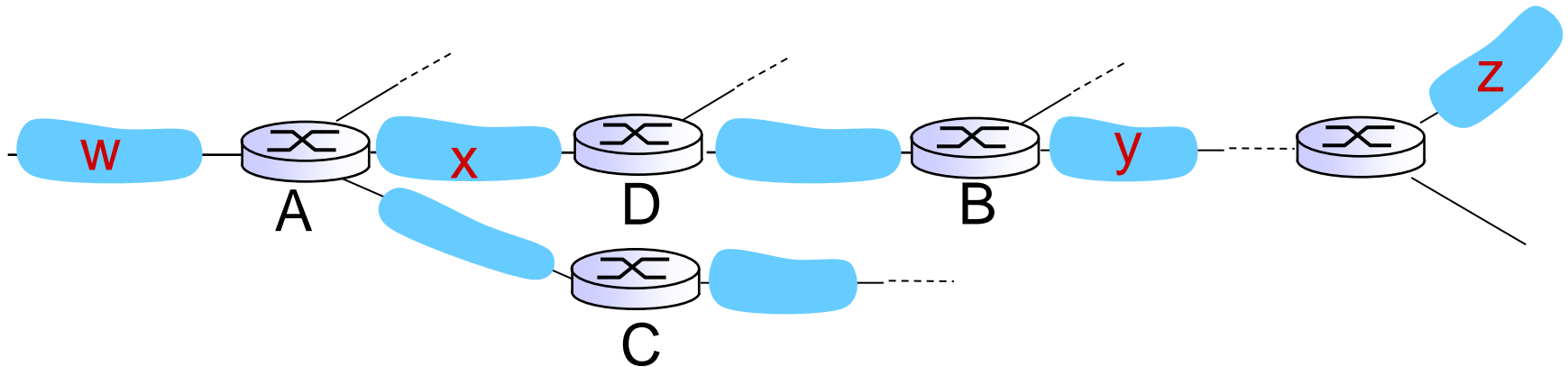
# RIP ( Routing Information Protocol)

❖ included in BSD-UNIX distribution in 1982
❖ distance vector algorithm
 ▪ distance metric: # hops (max = 15 hops), each link has cost 1
 ▪ DVs exchanged with neighbors every 30 sec in response message (aka advertisement)
 ▪ each advertisement: list of up to 25 destination *subnets* (in IP addressing sense)

from router A to destination *subnets:*

| subnet | hops |
|--------|------|
| u | 1 |
| v | 2 |
| w | 2 |
| x | 3 |
| y | 3 |
| z | 2 |

# RIP: example



routing table in router D

| destination subnet | next  router | # hops to dest |
|---|---|---|
| w | A | 2 |
| y | B | 2 |
| z | B | 7 |
| x | -- | 1 |
| …. | …. | .... |

# RIP: example

A-to-D advertisement

| dest | next | hops |
|------|------|------|
| w | - | 1 |
| x | - | 1 |
| z | C | 4 |
| .... | .... | .... |



routing table in router D

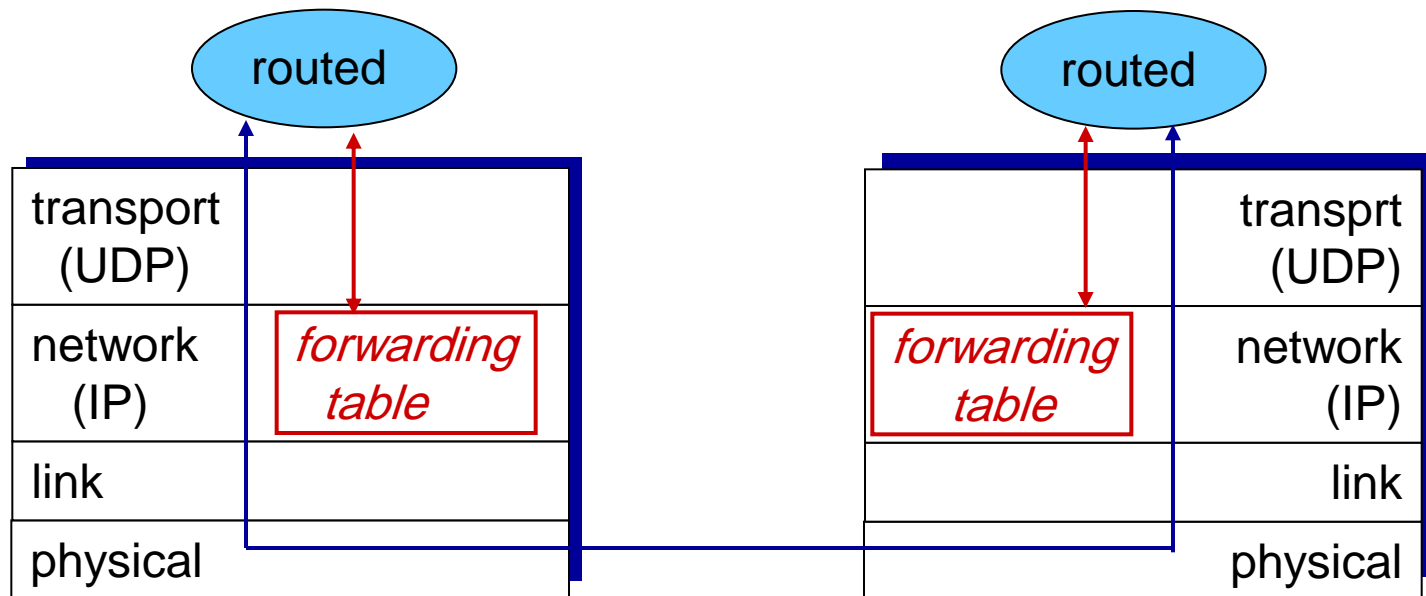| destination subnet | next router | # hops to dest |
|--------------------|-------------|----------------|
| w | A | 2 |
| y | B | 2 |
| z | B  A | 7  5 |
| x | -- | 1 |
| .... | .... | .... |

# RIP: link failure, recovery

if no advertisement heard after 180 sec -->
   neighbor/link declared dead
   - routes via neighbor invalidated
   - new advertisements sent to neighbors
   - neighbors in turn send out new advertisements (if tables changed)
   - link failure info quickly (?) propagates to entire net
   - *poison reverse* used to prevent ping-pong loops (infinite distance = 16 hops)

# RIP table processing

- ❖ RIP routing tables managed by *application-level* process called route-d (daemon)
- ❖ advertisements sent in UDP packets, periodically repeated

# OSPF (Open Shortest Path First)

❖ "open": publicly available
❖ uses link state algorithm
  ▪ LS packet dissemination
  ▪ topology map at each node
  ▪ route computation using Dijkstra's algorithm
❖ OSPF advertisement carries one entry per neighbor
❖ advertisements flooded to *entire* AS
  ▪ carried in OSPF messages directly over IP (rather than TCP or UDP
❖ *IS-IS routing* protocol: nearly identical to OSPF