

### 5.2.2 The Distance-Vector (DV) Routing Algorithm

Whereas the LS algorithm is an algorithm using global information, the **distance-vector (DV)** algorithm is iterative, asynchronous, and distributed. It is *distributed* in that each node receives some information from one or more of its *directly attached* neighbors, performs a calculation, and then distributes the results of its calculation back to its neighbors. It is *iterative* in that this process continues on until no more information is exchanged between neighbors. (Interestingly, the algorithm is also self-terminating—there is no signal that the computation should stop; it just stops.) The algorithm is *asynchronous* in that it does not require all of the nodes to operate in lockstep with each other. We'll see that an asynchronous, iterative, self-terminating, distributed algorithm is much more interesting and fun than a centralized algorithm!

Before we present the DV algorithm, it will prove beneficial to discuss an important relationship that exists among the costs of the least-cost paths. Let  $d_x(y)$  be the cost of the least-cost path from node  $x$  to node  $y$ . Then the least costs are related by the celebrated Bellman-Ford equation, namely,

$$d_x(y) = \min_v \{ c(x, v) + d_v(y) \}, \quad (5.1)$$

where the  $\min_v$  in the equation is taken over all of  $x$ 's neighbors. The Bellman-Ford equation is rather intuitive. Indeed, after traveling from  $x$  to  $v$ , if we then take the least-cost path from  $v$  to  $y$ , the path cost will be  $c(x, v) + d_v(y)$ . Since we must begin by traveling to some neighbor  $v$ , the least cost from  $x$  to  $y$  is the minimum of  $c(x, v) + d_v(y)$  taken over all neighbors  $v$ .

But for those who might be skeptical about the validity of the equation, let's check it for source node  $u$  and destination node  $z$  in Figure 5.3. The source node  $u$  has three neighbors: nodes  $v$ ,  $x$ , and  $w$ . By walking along various paths in the graph, it is easy to see that  $d_v(z) = 5$ ,  $d_x(z) = 3$ , and  $d_w(z) = 3$ . Plugging these values into Equation 5.1, along with the costs  $c(u, v) = 2$ ,  $c(u, x) = 1$ , and  $c(u, w) = 5$ , gives  $d_u(z) = \min \{ 2 + 5, 5 + 3, 1 + 3 \} = 4$ , which is obviously true and which is exactly what the Dijkstra algorithm gave us for the same network. This quick verification should help relieve any skepticism you may have.

The Bellman-Ford equation is not just an intellectual curiosity. It actually has significant practical importance: the solution to the Bellman-Ford equation provides the entries in node  $x$ 's forwarding table. To see this, let  $v^*$  be any neighboring node that achieves the minimum in Equation 5.1. Then, if node  $x$  wants to send a packet to node  $y$  along a least-cost path, it should first forward the packet to node  $v^*$ . Thus, node  $x$ 's forwarding table would specify node  $v^*$  as the next-hop router for the ultimate destination  $y$ . Another important practical contribution of the Bellman-Ford equation is that it suggests the form of the neighbor-to-neighbor communication that will take place in the DV algorithm.

The basic idea is as follows. Each node  $x$  begins with  $D_x(y)$ , an estimate of the cost of the least-cost path from itself to node  $y$ , for all nodes,  $y$ , in  $N$ . Let  $\mathbf{D}_x = [D_x(y): y \text{ in } N]$  be node  $x$ 's distance vector, which is the vector of cost estimates from  $x$  to all other nodes,  $y$ , in  $N$ . With the DV algorithm, each node  $x$  maintains the following routing information:

- For each neighbor  $v$ , the cost  $c(x, v)$  from  $x$  to directly attached neighbor,  $v$
- Node  $x$ 's distance vector, that is,  $\mathbf{D}_x = [D_x(y): y \text{ in } N]$ , containing  $x$ 's estimate of its cost to all destinations,  $y$ , in  $N$
- The distance vectors of each of its neighbors, that is,  $\mathbf{D}_v = [D_v(y): y \text{ in } N]$  for each neighbor  $v$  of  $x$

In the distributed, asynchronous algorithm, from time to time, each node sends a copy of its distance vector to each of its neighbors. When a node  $x$  receives a new distance vector from any of its neighbors  $w$ , it saves  $w$ 's distance vector, and then uses the Bellman-Ford equation to update its own distance vector as follows:

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

If node  $x$ 's distance vector has changed as a result of this update step, node  $x$  will then send its updated distance vector to each of its neighbors, which can in turn update their own distance vectors. Miraculously enough, as long as all the nodes continue to exchange their distance vectors in an asynchronous fashion, each cost estimate  $D_x(y)$  converges to  $d_x(y)$ , the actual cost of the least-cost path from node  $x$  to node  $y$  [Bertsekas 1991]!

### Distance-Vector (DV) Algorithm

At each node,  $x$ :

```

1  Initialization:
2      for all destinations  $y$  in  $N$ :
3           $D_x(y) = c(x, y)$  /* if  $y$  is not a neighbor then  $c(x, y) = \infty$  */
4      for each neighbor  $w$ 
5           $D_w(y) = ?$  for all destinations  $y$  in  $N$ 
6      for each neighbor  $w$ 
7          send distance vector  $\mathbf{D}_x = [D_x(y): y \text{ in } N]$  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          $D_x(y) = \min_v \{ c(x, v) + D_v(y) \}$ 
15
16     if  $D_x(y)$  changed for any destination  $y$ 
17         send distance vector  $\mathbf{D}_x = [D_x(y): y \text{ in } N]$  to all neighbors
18
19 forever
```

In the DV algorithm, a node  $x$  updates its distance-vector estimate when it either sees a cost change in one of its directly attached links or receives a distance-vector update from some neighbor. But to update its own forwarding table for a given destination  $y$ , what node  $x$  really needs to know is not the shortest-path distance to  $y$  but instead the neighboring node  $v^*(y)$  that is the next-hop router along the shortest path to  $y$ . As you might expect, the next-hop router  $v^*(y)$  is the neighbor  $v$  that achieves the minimum in Line 14 of the DV algorithm. (If there are multiple neighbors  $v$  that achieve the minimum, then  $v^*(y)$  can be any of the minimizing neighbors.) Thus, in Lines 13–14, for each destination  $y$ , node  $x$  also determines  $v^*(y)$  and updates its forwarding table for destination  $y$ .

Recall that the LS algorithm is a centralized algorithm in the sense that it requires each node to first obtain a complete map of the network before running the Dijkstra algorithm. The DV algorithm is *decentralized* and does not use such global information. Indeed, the only information a node will have is the costs of the links to its directly attached neighbors and information it receives from these neighbors. Each node waits for an update from any neighbor (Lines 10–11), calculates its new distance vector when receiving an update (Line 14), and distributes its new distance vector to its neighbors (Lines 16–17). DV-like algorithms are used in many routing protocols in practice, including the Internet’s RIP and BGP, ISO IDRP, Novell IPX, and the original ARPAnet.

Figure 5.6 illustrates the operation of the DV algorithm for the simple three-node network shown at the top of the figure. The operation of the algorithm is illustrated in a synchronous manner, where all nodes simultaneously receive distance vectors from their neighbors, compute their new distance vectors, and inform their neighbors if their distance vectors have changed. After studying this example, you should convince yourself that the algorithm operates correctly in an asynchronous manner as well, with node computations and update generation/reception occurring at any time.

The leftmost column of the figure displays three initial **routing tables** for each of the three nodes. For example, the table in the upper-left corner is node  $x$ ’s initial routing table. Within a specific routing table, each row is a distance vector—specifically, each node’s routing table includes its own distance vector and that of each of its neighbors. Thus, the first row in node  $x$ ’s initial routing table is  $\mathbf{D}_x = [D_x(x), D_x(y), D_x(z)] = [0, 2, 7]$ . The second and third rows in this table are the most recently received distance vectors from nodes  $y$  and  $z$ , respectively. Because at initialization node  $x$  has not received anything from node  $y$  or  $z$ , the entries in the second and third rows are initialized to infinity.

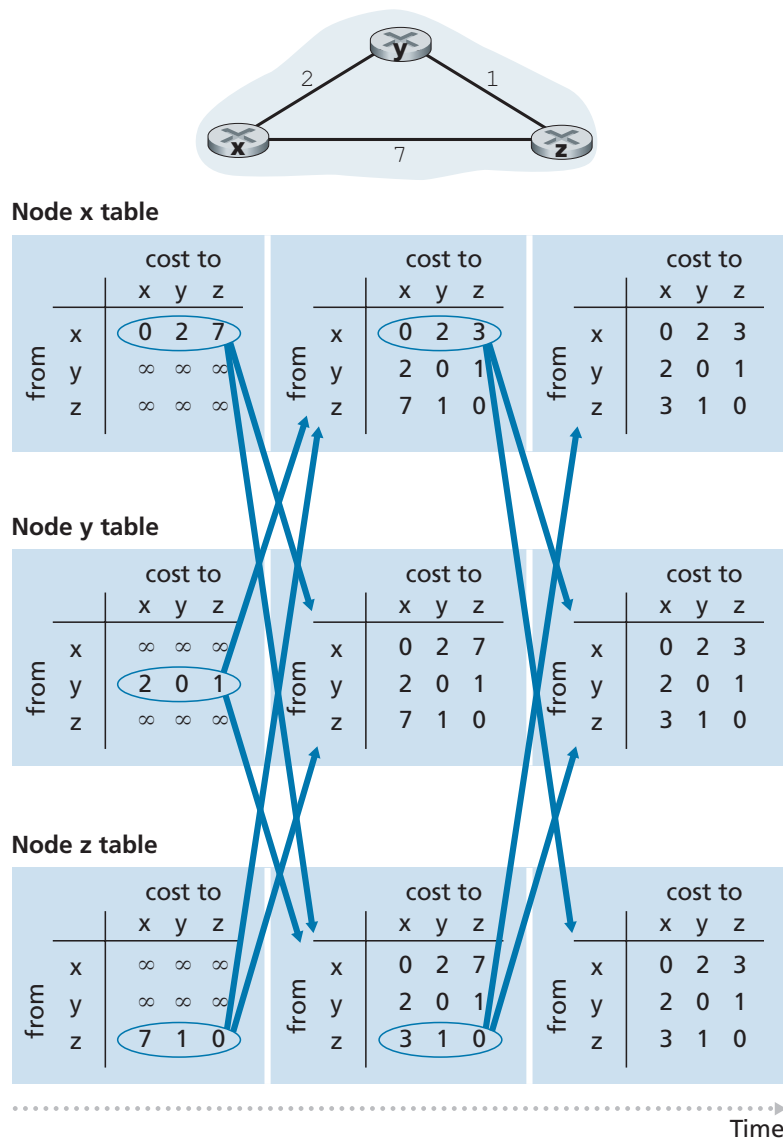
After initialization, each node sends its distance vector to each of its two neighbors. This is illustrated in Figure 5.6 by the arrows from the first column of tables to the second column of tables. For example, node  $x$  sends its distance vector  $\mathbf{D}_x = [0, 2, 7]$  to both nodes  $y$  and  $z$ . After receiving the updates, each node recomputes its own distance vector. For example, node  $x$  computes

$$D_x(x) = 0$$

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

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

The second column therefore displays, for each node, the node's new distance vector along with distance vectors just received from its neighbors. Note, for example, that



**Figure 5.6** ♦ Distance-vector (DV) algorithm in operation

node  $x$ 's estimate for the least cost to node  $z$ ,  $D_x(z)$ , has changed from 7 to 3. Also note that for node  $x$ , neighboring node  $y$  achieves the minimum in line 14 of the DV algorithm; thus, at this stage of the algorithm, we have at node  $x$  that  $v^*(y) = y$  and  $v^*(z) = y$ .

After the nodes recompute their distance vectors, they again send their updated distance vectors to their neighbors (if there has been a change). This is illustrated in Figure 5.6 by the arrows from the second column of tables to the third column of tables. Note that only nodes  $x$  and  $z$  send updates: node  $y$ 's distance vector didn't change so node  $y$  doesn't send an update. After receiving the updates, the nodes then recompute their distance vectors and update their routing tables, which are shown in the third column.

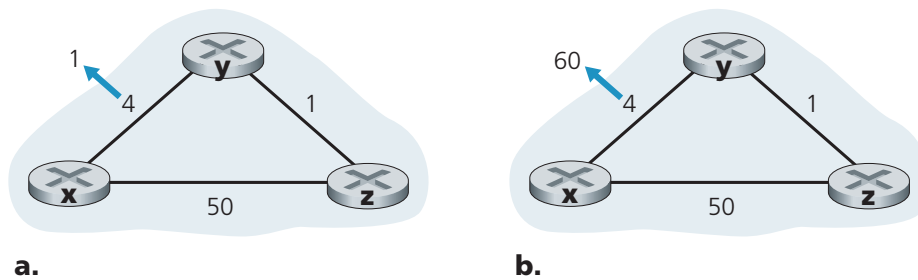
The process of receiving updated distance vectors from neighbors, recomputing routing table entries, and informing neighbors of changed costs of the least-cost path to a destination continues until no update messages are sent. At this point, since no update messages are sent, no further routing table calculations will occur and the algorithm will enter a quiescent state; that is, all nodes will be performing the wait in Lines 10–11 of the DV algorithm. The algorithm remains in the quiescent state until a link cost changes, as discussed next.

### Distance-Vector Algorithm: Link-Cost Changes and Link Failure

When a node running the DV algorithm detects a change in the link cost from itself to a neighbor (Lines 10–11), it updates its distance vector (Lines 13–14) and, if there's a change in the cost of the least-cost path, informs its neighbors (Lines 16–17) of its new distance vector. Figure 5.7(a) illustrates a scenario where the link cost from  $y$  to  $x$  changes from 4 to 1. We focus here only on  $y$ ' and  $z$ 's distance table entries to destination  $x$ . The DV algorithm causes the following sequence of events to occur:

- At time  $t_0$ ,  $y$  detects the link-cost change (the cost has changed from 4 to 1), updates its distance vector, and informs its neighbors of this change since its distance vector has changed.
- At time  $t_1$ ,  $z$  receives the update from  $y$  and updates its table. It computes a new least cost to  $x$  (it has decreased from a cost of 5 to a cost of 2) and sends its new distance vector to its neighbors.
- At time  $t_2$ ,  $y$  receives  $z$ 's update and updates its distance table.  $y$ 's least costs do not change and hence  $y$  does not send any message to  $z$ . The algorithm comes to a quiescent state.

Thus, only two iterations are required for the DV algorithm to reach a quiescent state. The good news about the decreased cost between  $x$  and  $y$  has propagated quickly through the network.



**Figure 5.7** ♦ Changes in link cost

Let's now consider what can happen when a link cost *increases*. Suppose that the link cost between  $x$  and  $y$  increases from 4 to 60, as shown in Figure 5.7(b).

1. Before the link cost changes,  $D_y(x) = 4$ ,  $D_y(z) = 1$ ,  $D_z(y) = 1$ , and  $D_z(x) = 5$ . At time  $t_0$ ,  $y$  detects the link-cost change (the cost has changed from 4 to 60).  $y$  computes its new minimum-cost path to  $x$  to have a cost of

$$D_y(x) = \min \{ c(y,x) + D_x(x), c(y,z) + D_z(x) \} = \min \{ 60 + 0, 1 + 5 \} = 6$$

Of course, with our global view of the network, we can see that this new cost via  $z$  is *wrong*. But the only information node  $y$  has is that its direct cost to  $x$  is 60 and that  $z$  has last told  $y$  that  $z$  could get to  $x$  with a cost of 5. So in order to get to  $x$ ,  $y$  would now route through  $z$ , fully expecting that  $z$  will be able to get to  $x$  with a cost of 5. As of  $t_1$  we have a **routing loop**—in order to get to  $x$ ,  $y$  routes through  $z$ , and  $z$  routes through  $y$ . A routing loop is like a black hole—a packet destined for  $x$  arriving at  $y$  or  $z$  as of  $t_1$  will bounce back and forth between these two nodes forever (or until the forwarding tables are changed).

2. Since node  $y$  has computed a new minimum cost to  $x$ , it informs  $z$  of its new distance vector at time  $t_1$ .
3. Sometime after  $t_1$ ,  $z$  receives  $y$ 's new distance vector, which indicates that  $y$ 's minimum cost to  $x$  is 6.  $z$  knows it can get to  $y$  with a cost of 1 and hence computes a new least cost to  $x$  of  $D_z(x) = \min \{ 50 + 0, 1 + 6 \} = 7$ . Since  $z$ 's least cost to  $x$  has increased, it then informs  $y$  of its new distance vector at  $t_2$ .
4. In a similar manner, after receiving  $z$ 's new distance vector,  $y$  determines  $D_y(x) = 8$  and sends  $z$  its distance vector.  $z$  then determines  $D_z(x) = 9$  and sends  $y$  its distance vector, and so on.

How long will the process continue? You should convince yourself that the loop will persist for 44 iterations (message exchanges between  $y$  and  $z$ )—until  $z$  eventually computes the cost of its path via  $y$  to be greater than 50. At this point,  $z$  will (finally!) determine that its least-cost path to  $x$  is via its direct connection to  $x$ .  $y$  will then

route to  $x$  via  $z$ . The result of the bad news about the increase in link cost has indeed traveled slowly! What would have happened if the link cost  $c(y, x)$  had changed from 4 to 10,000 and the cost  $c(z, x)$  had been 9,999? Because of such scenarios, the problem we have seen is sometimes referred to as the count-to-infinity problem.

### Distance-Vector Algorithm: Adding Poisoned Reverse

The specific looping scenario just described can be avoided using a technique known as *poisoned reverse*. The idea is simple—if  $z$  routes through  $y$  to get to destination  $x$ , then  $z$  will advertise to  $y$  that its distance to  $x$  is infinity, that is,  $z$  will advertise to  $y$  that  $D_z(x) = \infty$  (even though  $z$  knows  $D_z(x) = 5$  in truth).  $z$  will continue telling this little white lie to  $y$  as long as it routes to  $x$  via  $y$ . Since  $y$  believes that  $z$  has no path to  $x$ ,  $y$  will never attempt to route to  $x$  via  $z$ , as long as  $z$  continues to route to  $x$  via  $y$  (and lies about doing so).

Let's now see how poisoned reverse solves the particular looping problem we encountered before in Figure 5.5(b). As a result of the poisoned reverse,  $y$ 's distance table indicates  $D_z(x) = \infty$ . When the cost of the  $(x, y)$  link changes from 4 to 60 at time  $t_0$ ,  $y$  updates its table and continues to route directly to  $x$ , albeit at a higher cost of 60, and informs  $z$  of its new cost to  $x$ , that is,  $D_y(x) = 60$ . After receiving the update at  $t_1$ ,  $z$  immediately shifts its route to  $x$  to be via the direct  $(z, x)$  link at a cost of 50. Since this is a new least-cost path to  $x$ , and since the path no longer passes through  $y$ ,  $z$  now informs  $y$  that  $D_z(x) = 50$  at  $t_2$ . After receiving the update from  $z$ ,  $y$  updates its distance table with  $D_y(x) = 51$ . Also, since  $z$  is now on  $y$ 's least-cost path to  $x$ ,  $y$  poisons the reverse path from  $z$  to  $x$  by informing  $z$  at time  $t_3$  that  $D_y(x) = \infty$  (even though  $y$  knows that  $D_y(x) = 51$  in truth).

Does poisoned reverse solve the general count-to-infinity problem? It does not. You should convince yourself that loops involving three or more nodes (rather than simply two immediately neighboring nodes) will not be detected by the poisoned reverse technique.

### A Comparison of LS and DV Routing Algorithms

The DV and LS algorithms take complementary approaches toward computing routing. In the DV algorithm, each node talks to *only* its directly connected neighbors, but it provides its neighbors with least-cost estimates from itself to *all* the nodes (that it knows about) in the network. The LS algorithm requires global information. Consequently, when implemented in each and every router, for example, as in Figures 4.2 and 5.1, each node would need to communicate with *all* other nodes (via broadcast), but it tells them *only* the costs of its directly connected links. Let's conclude our study of LS and DV algorithms with a quick comparison of some of their attributes. Recall that  $N$  is the set of nodes (routers) and  $E$  is the set of edges (links).

- *Message complexity.* We have seen that LS requires each node to know the cost of each link in the network. This requires  $O(|N| |E|)$  messages to be sent.



Also, whenever a link cost changes, the new link cost must be sent to all nodes. The DV algorithm requires message exchanges between directly connected neighbors at each iteration. We have seen that the time needed for the algorithm to converge can depend on many factors. When link costs change, the DV algorithm will propagate the results of the changed link cost only if the new link cost results in a changed least-cost path for one of the nodes attached to that link.

- *Speed of convergence.* We have seen that our implementation of LS is an  $O(|N|^2)$  algorithm requiring  $O(|N| |E|)$  messages. The DV algorithm can converge slowly and can have routing loops while the algorithm is converging. DV also suffers from the count-to-infinity problem.
- *Robustness.* What can happen if a router fails, misbehaves, or is sabotaged? Under LS, a router could broadcast an incorrect cost for one of its attached links (but no others). A node could also corrupt or drop any packets it received as part of an LS broadcast. But an LS node is computing only its own forwarding tables; other nodes are performing similar calculations for themselves. This means route calculations are somewhat separated under LS, providing a degree of robustness. Under DV, a node can advertise incorrect least-cost paths to any or all destinations. (Indeed, in 1997, a malfunctioning router in a small ISP provided national backbone routers with erroneous routing information. This caused other routers to flood the malfunctioning router with traffic and caused large portions of the Internet to become disconnected for up to several hours [Neumann 1997].) More generally, we note that, at each iteration, a node's calculation in DV is passed on to its neighbor and then indirectly to its neighbor's neighbor on the next iteration. In this sense, an incorrect node calculation can be diffused through the entire network under DV.

In the end, neither algorithm is an obvious winner over the other; indeed, both algorithms are used in the Internet.

## 5.3 Intra-AS Routing in the Internet: OSPF

In our study of routing algorithms so far, we've viewed the network simply as a collection of interconnected routers. One router was indistinguishable from another in the sense that all routers executed the same routing algorithm to compute routing paths through the entire network. In practice, this model and its view of a homogeneous set of routers all executing the same routing algorithm is simplistic for two important reasons:

- *Scale.* As the number of routers becomes large, the overhead involved in communicating, computing, and storing routing information becomes prohibitive. Today's