

Assume that if the sets containing  $x_i$  and  $x_j$  have the same size, then the operation  $\text{UNION}(x_i, x_j)$  appends  $x_j$ 's list onto  $x_i$ 's list.

### 21.2-3

Adapt the aggregate proof of Theorem 21.1 to obtain amortized time bounds of  $O(1)$  for MAKE-SET and FIND-SET and  $O(\lg n)$  for UNION using the linked-list representation and the weighted-union heuristic.

### 21.2-4

Give a tight asymptotic bound on the running time of the sequence of operations in Figure 21.3 assuming the linked-list representation and the weighted-union heuristic.

### 21.2-5

Professor Gompers suspects that it might be possible to keep just one pointer in each set object, rather than two (*head* and *tail*), while keeping the number of pointers in each list element at two. Show that the professor's suspicion is well founded by describing how to represent each set by a linked list such that each operation has the same running time as the operations described in this section. Describe also how the operations work. Your scheme should allow for the weighted-union heuristic, with the same effect as described in this section. (*Hint*: Use the tail of a linked list as its set's representative.)

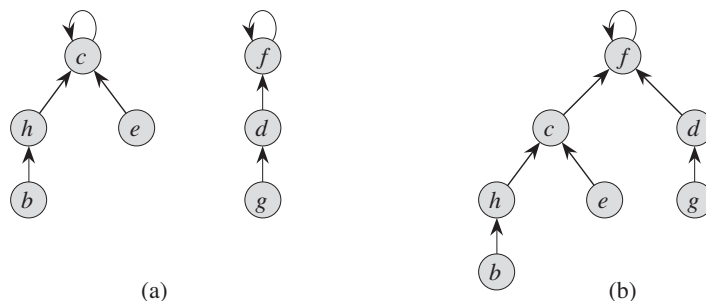
### 21.2-6

Suggest a simple change to the UNION procedure for the linked-list representation that removes the need to keep the *tail* pointer to the last object in each list. Whether or not the weighted-union heuristic is used, your change should not change the asymptotic running time of the UNION procedure. (*Hint*: Rather than appending one list to another, splice them together.)

---

## 21.3 Disjoint-set forests

In a faster implementation of disjoint sets, we represent sets by rooted trees, with each node containing one member and each tree representing one set. In a **disjoint-set forest**, illustrated in Figure 21.4(a), each member points only to its parent. The root of each tree contains the representative and is its own parent. As we shall see, although the straightforward algorithms that use this representation are no faster than ones that use the linked-list representation, by introducing two heuristics—"union by rank" and "path compression"—we can achieve an asymptotically optimal disjoint-set data structure.



**Figure 21.4** A disjoint-set forest. **(a)** Two trees representing the two sets of Figure 21.2. The tree on the left represents the set  $\{b, c, e, h\}$ , with  $c$  as the representative, and the tree on the right represents the set  $\{d, f, g\}$ , with  $f$  as the representative. **(b)** The result of  $\text{UNION}(e, g)$ .

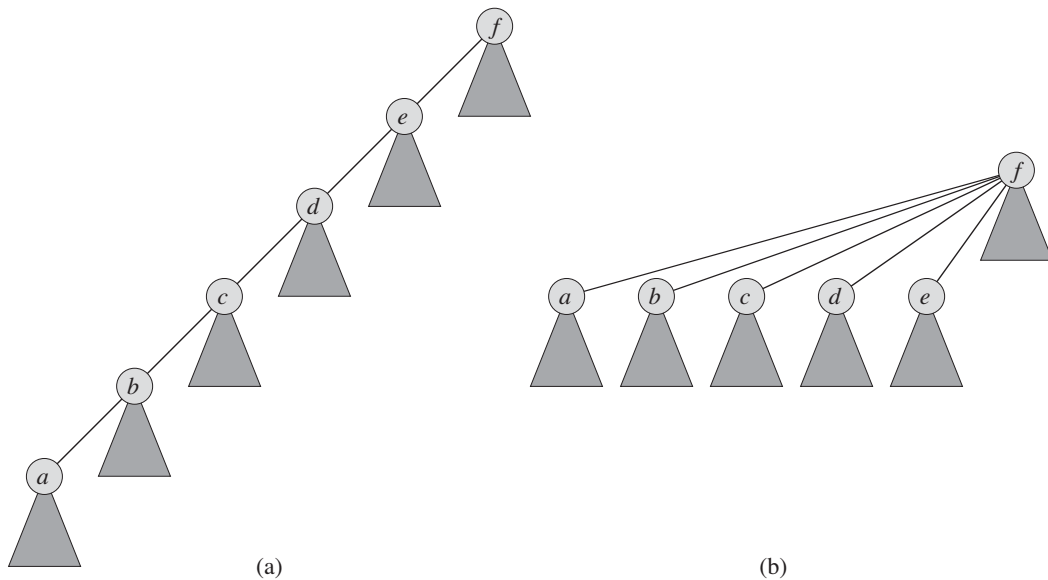
We perform the three disjoint-set operations as follows. A **MAKE-SET** operation simply creates a tree with just one node. We perform a **FIND-SET** operation by following parent pointers until we find the root of the tree. The nodes visited on this simple path toward the root constitute the *find path*. A **UNION** operation, shown in Figure 21.4(b), causes the root of one tree to point to the root of the other.

### Heuristics to improve the running time

So far, we have not improved on the linked-list implementation. A sequence of  $n - 1$  **UNION** operations may create a tree that is just a linear chain of  $n$  nodes. By using two heuristics, however, we can achieve a running time that is almost linear in the total number of operations  $m$ .

The first heuristic, *union by rank*, is similar to the weighted-union heuristic we used with the linked-list representation. The obvious approach would be to make the root of the tree with fewer nodes point to the root of the tree with more nodes. Rather than explicitly keeping track of the size of the subtree rooted at each node, we shall use an approach that eases the analysis. For each node, we maintain a *rank*, which is an upper bound on the height of the node. In union by rank, we make the root with smaller rank point to the root with larger rank during a **UNION** operation.

The second heuristic, *path compression*, is also quite simple and highly effective. As shown in Figure 21.5, we use it during **FIND-SET** operations to make each node on the find path point directly to the root. Path compression does not change any ranks.



**Figure 21.5** Path compression during the operation FIND-SET. Arrows and self-loops at roots are omitted. **(a)** A tree representing a set prior to executing FIND-SET(*a*). Triangles represent subtrees whose roots are the nodes shown. Each node has a pointer to its parent. **(b)** The same set after executing FIND-SET(*a*). Each node on the find path now points directly to the root.

### Pseudocode for disjoint-set forests

To implement a disjoint-set forest with the union-by-rank heuristic, we must keep track of ranks. With each node  $x$ , we maintain the integer value  $x.rank$ , which is an upper bound on the height of  $x$  (the number of edges in the longest simple path between  $x$  and a descendant leaf). When MAKE-SET creates a singleton set, the single node in the corresponding tree has an initial rank of 0. Each FIND-SET operation leaves all ranks unchanged. The UNION operation has two cases, depending on whether the roots of the trees have equal rank. If the roots have unequal rank, we make the root with higher rank the parent of the root with lower rank, but the ranks themselves remain unchanged. If, instead, the roots have equal ranks, we arbitrarily choose one of the roots as the parent and increment its rank.

Let us put this method into pseudocode. We designate the parent of node  $x$  by  $x.p$ . The LINK procedure, a subroutine called by UNION, takes pointers to two roots as inputs.

MAKE-SET( $x$ )

```
1   $x.p = x$ 
2   $x.rank = 0$ 
```

UNION( $x, y$ )

```
1  LINK(FIND-SET( $x$ ), FIND-SET( $y$ ))
```

LINK( $x, y$ )

```
1  if  $x.rank > y.rank$ 
2       $y.p = x$ 
3  else  $x.p = y$ 
4      if  $x.rank == y.rank$ 
5           $y.rank = y.rank + 1$ 
```

The FIND-SET procedure with path compression is quite simple:

FIND-SET( $x$ )

```
1  if  $x \neq x.p$ 
2       $x.p = \text{FIND-SET}(x.p)$ 
3  return  $x.p$ 
```

The FIND-SET procedure is a *two-pass method*: as it recurses, it makes one pass up the find path to find the root, and as the recursion unwinds, it makes a second pass back down the find path to update each node to point directly to the root. Each call of FIND-SET( $x$ ) returns  $x.p$  in line 3. If  $x$  is the root, then FIND-SET skips line 2 and instead returns  $x.p$ , which is  $x$ ; this is the case in which the recursion bottoms out. Otherwise, line 2 executes, and the recursive call with parameter  $x.p$  returns a pointer to the root. Line 2 updates node  $x$  to point directly to the root, and line 3 returns this pointer.

### Effect of the heuristics on the running time

Separately, either union by rank or path compression improves the running time of the operations on disjoint-set forests, and the improvement is even greater when we use the two heuristics together. Alone, union by rank yields a running time of  $O(m \lg n)$  (see Exercise 21.4-4), and this bound is tight (see Exercise 21.3-3). Although we shall not prove it here, for a sequence of  $n$  MAKE-SET operations (and hence at most  $n - 1$  UNION operations) and  $f$  FIND-SET operations, the path-compression heuristic alone gives a worst-case running time of  $\Theta(n + f \cdot (1 + \log_{2+f/n} n))$ .

When we use both union by rank and path compression, the worst-case running time is  $O(m \alpha(n))$ , where  $\alpha(n)$  is a *very* slowly growing function, which we define in Section 21.4. In any conceivable application of a disjoint-set data structure,  $\alpha(n) \leq 4$ ; thus, we can view the running time as linear in  $m$  in all practical situations. Strictly speaking, however, it is superlinear. In Section 21.4, we prove this upper bound.

### Exercises

#### 21.3-1

Redo Exercise 21.2-2 using a disjoint-set forest with union by rank and path compression.

#### 21.3-2

Write a nonrecursive version of FIND-SET with path compression.

#### 21.3-3

Give a sequence of  $m$  MAKE-SET, UNION, and FIND-SET operations,  $n$  of which are MAKE-SET operations, that takes  $\Omega(m \lg n)$  time when we use union by rank only.

#### 21.3-4

Suppose that we wish to add the operation PRINT-SET( $x$ ), which is given a node  $x$  and prints all the members of  $x$ 's set, in any order. Show how we can add just a single attribute to each node in a disjoint-set forest so that PRINT-SET( $x$ ) takes time linear in the number of members of  $x$ 's set and the asymptotic running times of the other operations are unchanged. Assume that we can print each member of the set in  $O(1)$  time.

#### 21.3-5 ★

Show that any sequence of  $m$  MAKE-SET, FIND-SET, and LINK operations, where all the LINK operations appear before any of the FIND-SET operations, takes only  $O(m)$  time if we use both path compression and union by rank. What happens in the same situation if we use only the path-compression heuristic?

## ★ 21.4 Analysis of union by rank with path compression

As noted in Section 21.3, the combined union-by-rank and path-compression heuristic runs in time  $O(m \alpha(n))$  for  $m$  disjoint-set operations on  $n$  elements. In this section, we shall examine the function  $\alpha$  to see just how slowly it grows. Then we prove this running time using the potential method of amortized analysis.

### A very quickly growing function and its very slowly growing inverse

For integers  $k \geq 0$  and  $j \geq 1$ , we define the function  $A_k(j)$  as

$$A_k(j) = \begin{cases} j + 1 & \text{if } k = 0, \\ A_{k-1}^{(j+1)}(j) & \text{if } k \geq 1, \end{cases}$$

where the expression  $A_{k-1}^{(j+1)}(j)$  uses the functional-iteration notation given in Section 3.2. Specifically,  $A_{k-1}^{(0)}(j) = j$  and  $A_{k-1}^{(i)}(j) = A_{k-1}(A_{k-1}^{(i-1)}(j))$  for  $i \geq 1$ . We will refer to the parameter  $k$  as the **level** of the function  $A$ .

The function  $A_k(j)$  strictly increases with both  $j$  and  $k$ . To see just how quickly this function grows, we first obtain closed-form expressions for  $A_1(j)$  and  $A_2(j)$ .

#### Lemma 21.2

For any integer  $j \geq 1$ , we have  $A_1(j) = 2j + 1$ .

**Proof** We first use induction on  $i$  to show that  $A_0^{(i)}(j) = j + i$ . For the base case, we have  $A_0^{(0)}(j) = j = j + 0$ . For the inductive step, assume that  $A_0^{(i-1)}(j) = j + (i - 1)$ . Then  $A_0^{(i)}(j) = A_0(A_0^{(i-1)}(j)) = (j + (i - 1)) + 1 = j + i$ . Finally, we note that  $A_1(j) = A_0^{(j+1)}(j) = j + (j + 1) = 2j + 1$ . ■

#### Lemma 21.3

For any integer  $j \geq 1$ , we have  $A_2(j) = 2^{j+1}(j + 1) - 1$ .

**Proof** We first use induction on  $i$  to show that  $A_1^{(i)}(j) = 2^i(j + 1) - 1$ . For the base case, we have  $A_1^{(0)}(j) = j = 2^0(j + 1) - 1$ . For the inductive step, assume that  $A_1^{(i-1)}(j) = 2^{i-1}(j + 1) - 1$ . Then  $A_1^{(i)}(j) = A_1(A_1^{(i-1)}(j)) = A_1(2^{i-1}(j + 1) - 1) = 2 \cdot (2^{i-1}(j + 1) - 1) + 1 = 2^i(j + 1) - 2 + 1 = 2^i(j + 1) - 1$ . Finally, we note that  $A_2(j) = A_1^{(j+1)}(j) = 2^{j+1}(j + 1) - 1$ . ■

Now we can see how quickly  $A_k(j)$  grows by simply examining  $A_k(1)$  for levels  $k = 0, 1, 2, 3, 4$ . From the definition of  $A_0(k)$  and the above lemmas, we have  $A_0(1) = 1 + 1 = 2$ ,  $A_1(1) = 2 \cdot 1 + 1 = 3$ , and  $A_2(1) = 2^{1+1} \cdot (1 + 1) - 1 = 7$ .

We also have

$$\begin{aligned}
 A_3(1) &= A_2^{(2)}(1) \\
 &= A_2(A_2(1)) \\
 &= A_2(7) \\
 &= 2^8 \cdot 8 - 1 \\
 &= 2^{11} - 1 \\
 &= 2047
 \end{aligned}$$

and

$$\begin{aligned}
 A_4(1) &= A_3^{(2)}(1) \\
 &= A_3(A_3(1)) \\
 &= A_3(2047) \\
 &= A_2^{(2048)}(2047) \\
 &\gg A_2(2047) \\
 &= 2^{2048} \cdot 2048 - 1 \\
 &> 2^{2048} \\
 &= (2^4)^{512} \\
 &= 16^{512} \\
 &\gg 10^{80} ,
 \end{aligned}$$

which is the estimated number of atoms in the observable universe. (The symbol “ $\gg$ ” denotes the “much-greater-than” relation.)

We define the inverse of the function  $A_k(n)$ , for integer  $n \geq 0$ , by

$$\alpha(n) = \min \{k : A_k(1) \geq n\} .$$

In words,  $\alpha(n)$  is the lowest level  $k$  for which  $A_k(1)$  is at least  $n$ . From the above values of  $A_k(1)$ , we see that

$$\alpha(n) = \begin{cases} 0 & \text{for } 0 \leq n \leq 2 , \\ 1 & \text{for } n = 3 , \\ 2 & \text{for } 4 \leq n \leq 7 , \\ 3 & \text{for } 8 \leq n \leq 2047 , \\ 4 & \text{for } 2048 \leq n \leq A_4(1) . \end{cases}$$

It is only for values of  $n$  so large that the term “astronomical” understates them (greater than  $A_4(1)$ , a huge number) that  $\alpha(n) > 4$ , and so  $\alpha(n) \leq 4$  for all practical purposes.

### Properties of ranks

In the remainder of this section, we prove an  $O(m\alpha(n))$  bound on the running time of the disjoint-set operations with union by rank and path compression. In order to prove this bound, we first prove some simple properties of ranks.

#### **Lemma 21.4**

For all nodes  $x$ , we have  $x.rank \leq x.p.rank$ , with strict inequality if  $x \neq x.p$ . The value of  $x.rank$  is initially 0 and increases through time until  $x \neq x.p$ ; from then on,  $x.rank$  does not change. The value of  $x.p.rank$  monotonically increases over time.

**Proof** The proof is a straightforward induction on the number of operations, using the implementations of MAKE-SET, UNION, and FIND-SET that appear in Section 21.3. We leave it as Exercise 21.4-1. ■

#### **Corollary 21.5**

As we follow the simple path from any node toward a root, the node ranks strictly increase. ■

#### **Lemma 21.6**

Every node has rank at most  $n - 1$ .

**Proof** Each node's rank starts at 0, and it increases only upon LINK operations. Because there are at most  $n - 1$  UNION operations, there are also at most  $n - 1$  LINK operations. Because each LINK operation either leaves all ranks alone or increases some node's rank by 1, all ranks are at most  $n - 1$ . ■

Lemma 21.6 provides a weak bound on ranks. In fact, every node has rank at most  $\lceil \lg n \rceil$  (see Exercise 21.4-2). The looser bound of Lemma 21.6 will suffice for our purposes, however.

### Proving the time bound

We shall use the potential method of amortized analysis (see Section 17.3) to prove the  $O(m\alpha(n))$  time bound. In performing the amortized analysis, we will find it convenient to assume that we invoke the LINK operation rather than the UNION operation. That is, since the parameters of the LINK procedure are pointers to two roots, we act as though we perform the appropriate FIND-SET operations separately. The following lemma shows that even if we count the extra FIND-SET operations induced by UNION calls, the asymptotic running time remains unchanged.



**Lemma 21.7**

Suppose we convert a sequence  $S'$  of  $m'$  MAKE-SET, UNION, and FIND-SET operations into a sequence  $S$  of  $m$  MAKE-SET, LINK, and FIND-SET operations by turning each UNION into two FIND-SET operations followed by a LINK. Then, if sequence  $S$  runs in  $O(m \alpha(n))$  time, sequence  $S'$  runs in  $O(m' \alpha(n))$  time.

**Proof** Since each UNION operation in sequence  $S'$  is converted into three operations in  $S$ , we have  $m' \leq m \leq 3m'$ . Since  $m = O(m')$ , an  $O(m \alpha(n))$  time bound for the converted sequence  $S$  implies an  $O(m' \alpha(n))$  time bound for the original sequence  $S'$ . ■

In the remainder of this section, we shall assume that the initial sequence of  $m'$  MAKE-SET, UNION, and FIND-SET operations has been converted to a sequence of  $m$  MAKE-SET, LINK, and FIND-SET operations. We now prove an  $O(m \alpha(n))$  time bound for the converted sequence and appeal to Lemma 21.7 to prove the  $O(m' \alpha(n))$  running time of the original sequence of  $m'$  operations.

**Potential function**

The potential function we use assigns a potential  $\phi_q(x)$  to each node  $x$  in the disjoint-set forest after  $q$  operations. We sum the node potentials for the potential of the entire forest:  $\Phi_q = \sum_x \phi_q(x)$ , where  $\Phi_q$  denotes the potential of the forest after  $q$  operations. The forest is empty prior to the first operation, and we arbitrarily set  $\Phi_0 = 0$ . No potential  $\Phi_q$  will ever be negative.

The value of  $\phi_q(x)$  depends on whether  $x$  is a tree root after the  $q$ th operation. If it is, or if  $x.rank = 0$ , then  $\phi_q(x) = \alpha(n) \cdot x.rank$ .

Now suppose that after the  $q$ th operation,  $x$  is not a root and that  $x.rank \geq 1$ . We need to define two auxiliary functions on  $x$  before we can define  $\phi_q(x)$ . First we define

$$\text{level}(x) = \max \{k : x.p.rank \geq A_k(x.rank)\} .$$

That is,  $\text{level}(x)$  is the greatest level  $k$  for which  $A_k$ , applied to  $x$ 's rank, is no greater than  $x$ 's parent's rank.

We claim that

$$0 \leq \text{level}(x) < \alpha(n) , \tag{21.1}$$

which we see as follows. We have

$$\begin{aligned} x.p.rank &\geq x.rank + 1 \quad (\text{by Lemma 21.4}) \\ &= A_0(x.rank) \quad (\text{by definition of } A_0(j)) , \end{aligned}$$

which implies that  $\text{level}(x) \geq 0$ , and we have

$$\begin{aligned}
A_{\alpha(n)}(x.rank) &\geq A_{\alpha(n)}(1) && \text{(because } A_k(j) \text{ is strictly increasing)} \\
&\geq n && \text{(by the definition of } \alpha(n)) \\
&> x.p.rank && \text{(by Lemma 21.6) ,}
\end{aligned}$$

which implies that  $\text{level}(x) < \alpha(n)$ . Note that because  $x.p.rank$  monotonically increases over time, so does  $\text{level}(x)$ .

The second auxiliary function applies when  $x.rank \geq 1$ :

$$\text{iter}(x) = \max \{i : x.p.rank \geq A_{\text{level}(x)}^{(i)}(x.rank)\} .$$

That is,  $\text{iter}(x)$  is the largest number of times we can iteratively apply  $A_{\text{level}(x)}$ , applied initially to  $x$ 's rank, before we get a value greater than  $x$ 's parent's rank.

We claim that when  $x.rank \geq 1$ , we have

$$1 \leq \text{iter}(x) \leq x.rank , \tag{21.2}$$

which we see as follows. We have

$$\begin{aligned}
x.p.rank &\geq A_{\text{level}(x)}(x.rank) && \text{(by definition of } \text{level}(x)) \\
&= A_{\text{level}(x)}^{(1)}(x.rank) && \text{(by definition of functional iteration) ,}
\end{aligned}$$

which implies that  $\text{iter}(x) \geq 1$ , and we have

$$\begin{aligned}
A_{\text{level}(x)}^{(x.rank+1)}(x.rank) &= A_{\text{level}(x)+1}(x.rank) && \text{(by definition of } A_k(j)) \\
&> x.p.rank && \text{(by definition of } \text{level}(x)) ,
\end{aligned}$$

which implies that  $\text{iter}(x) \leq x.rank$ . Note that because  $x.p.rank$  monotonically increases over time, in order for  $\text{iter}(x)$  to decrease,  $\text{level}(x)$  must increase. As long as  $\text{level}(x)$  remains unchanged,  $\text{iter}(x)$  must either increase or remain unchanged.

With these auxiliary functions in place, we are ready to define the potential of node  $x$  after  $q$  operations:

$$\phi_q(x) = \begin{cases} \alpha(n) \cdot x.rank & \text{if } x \text{ is a root or } x.rank = 0 , \\ (\alpha(n) - \text{level}(x)) \cdot x.rank - \text{iter}(x) & \text{if } x \text{ is not a root and } x.rank \geq 1 . \end{cases}$$

We next investigate some useful properties of node potentials.

### **Lemma 21.8**

For every node  $x$ , and for all operation counts  $q$ , we have

$$0 \leq \phi_q(x) \leq \alpha(n) \cdot x.rank .$$

**Proof** If  $x$  is a root or  $x.rank = 0$ , then  $\phi_q(x) = \alpha(n) \cdot x.rank$  by definition. Now suppose that  $x$  is not a root and that  $x.rank \geq 1$ . We obtain a lower bound on  $\phi_q(x)$  by maximizing  $level(x)$  and  $iter(x)$ . By the bound (21.1),  $level(x) \leq \alpha(n) - 1$ , and by the bound (21.2),  $iter(x) \leq x.rank$ . Thus,

$$\begin{aligned} \phi_q(x) &= (\alpha(n) - level(x)) \cdot x.rank - iter(x) \\ &\geq (\alpha(n) - (\alpha(n) - 1)) \cdot x.rank - x.rank \\ &= x.rank - x.rank \\ &= 0. \end{aligned}$$

Similarly, we obtain an upper bound on  $\phi_q(x)$  by minimizing  $level(x)$  and  $iter(x)$ . By the bound (21.1),  $level(x) \geq 0$ , and by the bound (21.2),  $iter(x) \geq 1$ . Thus,

$$\begin{aligned} \phi_q(x) &\leq (\alpha(n) - 0) \cdot x.rank - 1 \\ &= \alpha(n) \cdot x.rank - 1 \\ &< \alpha(n) \cdot x.rank. \end{aligned} \quad \blacksquare$$

### Corollary 21.9

If node  $x$  is not a root and  $x.rank > 0$ , then  $\phi_q(x) < \alpha(n) \cdot x.rank$ .  $\blacksquare$

## Potential changes and amortized costs of operations

We are now ready to examine how the disjoint-set operations affect node potentials. With an understanding of the change in potential due to each operation, we can determine each operation's amortized cost.

### Lemma 21.10

Let  $x$  be a node that is not a root, and suppose that the  $q$ th operation is either a LINK or FIND-SET. Then after the  $q$ th operation,  $\phi_q(x) \leq \phi_{q-1}(x)$ . Moreover, if  $x.rank \geq 1$  and either  $level(x)$  or  $iter(x)$  changes due to the  $q$ th operation, then  $\phi_q(x) \leq \phi_{q-1}(x) - 1$ . That is,  $x$ 's potential cannot increase, and if it has positive rank and either  $level(x)$  or  $iter(x)$  changes, then  $x$ 's potential drops by at least 1.

**Proof** Because  $x$  is not a root, the  $q$ th operation does not change  $x.rank$ , and because  $n$  does not change after the initial  $n$  MAKE-SET operations,  $\alpha(n)$  remains unchanged as well. Hence, these components of the formula for  $x$ 's potential remain the same after the  $q$ th operation. If  $x.rank = 0$ , then  $\phi_q(x) = \phi_{q-1}(x) = 0$ . Now assume that  $x.rank \geq 1$ .

Recall that  $level(x)$  monotonically increases over time. If the  $q$ th operation leaves  $level(x)$  unchanged, then  $iter(x)$  either increases or remains unchanged. If both  $level(x)$  and  $iter(x)$  are unchanged, then  $\phi_q(x) = \phi_{q-1}(x)$ . If  $level(x)$

is unchanged and  $\text{iter}(x)$  increases, then it increases by at least 1, and so  $\phi_q(x) \leq \phi_{q-1}(x) - 1$ .

Finally, if the  $q$ th operation increases  $\text{level}(x)$ , it increases by at least 1, so that the value of the term  $(\alpha(n) - \text{level}(x)) \cdot x.\text{rank}$  drops by at least  $x.\text{rank}$ . Because  $\text{level}(x)$  increased, the value of  $\text{iter}(x)$  might drop, but according to the bound (21.2), the drop is by at most  $x.\text{rank} - 1$ . Thus, the increase in potential due to the change in  $\text{iter}(x)$  is less than the decrease in potential due to the change in  $\text{level}(x)$ , and we conclude that  $\phi_q(x) \leq \phi_{q-1}(x) - 1$ . ■

Our final three lemmas show that the amortized cost of each MAKE-SET, LINK, and FIND-SET operation is  $O(\alpha(n))$ . Recall from equation (17.2) that the amortized cost of each operation is its actual cost plus the increase in potential due to the operation.

**Lemma 21.11**

The amortized cost of each MAKE-SET operation is  $O(1)$ .

**Proof** Suppose that the  $q$ th operation is MAKE-SET( $x$ ). This operation creates node  $x$  with rank 0, so that  $\phi_q(x) = 0$ . No other ranks or potentials change, and so  $\Phi_q = \Phi_{q-1}$ . Noting that the actual cost of the MAKE-SET operation is  $O(1)$  completes the proof. ■

**Lemma 21.12**

The amortized cost of each LINK operation is  $O(\alpha(n))$ .

**Proof** Suppose that the  $q$ th operation is LINK( $x, y$ ). The actual cost of the LINK operation is  $O(1)$ . Without loss of generality, suppose that the LINK makes  $y$  the parent of  $x$ .

To determine the change in potential due to the LINK, we note that the only nodes whose potentials may change are  $x$ ,  $y$ , and the children of  $y$  just prior to the operation. We shall show that the only node whose potential can increase due to the LINK is  $y$ , and that its increase is at most  $\alpha(n)$ :

- By Lemma 21.10, any node that is  $y$ 's child just before the LINK cannot have its potential increase due to the LINK.
- From the definition of  $\phi_q(x)$ , we see that, since  $x$  was a root just before the  $q$ th operation,  $\phi_{q-1}(x) = \alpha(n) \cdot x.\text{rank}$ . If  $x.\text{rank} = 0$ , then  $\phi_q(x) = \phi_{q-1}(x) = 0$ . Otherwise,

$$\begin{aligned} \phi_q(x) &< \alpha(n) \cdot x.\text{rank} \quad (\text{by Corollary 21.9}) \\ &= \phi_{q-1}(x), \end{aligned}$$

and so  $x$ 's potential decreases.

- Because  $y$  is a root prior to the LINK,  $\phi_{q-1}(y) = \alpha(n) \cdot y.rank$ . The LINK operation leaves  $y$  as a root, and it either leaves  $y$ 's rank alone or it increases  $y$ 's rank by 1. Therefore, either  $\phi_q(y) = \phi_{q-1}(y)$  or  $\phi_q(y) = \phi_{q-1}(y) + \alpha(n)$ .

The increase in potential due to the LINK operation, therefore, is at most  $\alpha(n)$ . The amortized cost of the LINK operation is  $O(1) + \alpha(n) = O(\alpha(n))$ . ■

### Lemma 21.13

The amortized cost of each FIND-SET operation is  $O(\alpha(n))$ .

**Proof** Suppose that the  $q$ th operation is a FIND-SET and that the find path contains  $s$  nodes. The actual cost of the FIND-SET operation is  $O(s)$ . We shall show that no node's potential increases due to the FIND-SET and that at least  $\max(0, s - (\alpha(n) + 2))$  nodes on the find path have their potential decrease by at least 1.

To see that no node's potential increases, we first appeal to Lemma 21.10 for all nodes other than the root. If  $x$  is the root, then its potential is  $\alpha(n) \cdot x.rank$ , which does not change.

Now we show that at least  $\max(0, s - (\alpha(n) + 2))$  nodes have their potential decrease by at least 1. Let  $x$  be a node on the find path such that  $x.rank > 0$  and  $x$  is followed somewhere on the find path by another node  $y$  that is not a root, where  $\text{level}(y) = \text{level}(x)$  just before the FIND-SET operation. (Node  $y$  need not *immediately* follow  $x$  on the find path.) All but at most  $\alpha(n) + 2$  nodes on the find path satisfy these constraints on  $x$ . Those that do not satisfy them are the first node on the find path (if it has rank 0), the last node on the path (i.e., the root), and the last node  $w$  on the path for which  $\text{level}(w) = k$ , for each  $k = 0, 1, 2, \dots, \alpha(n) - 1$ .

Let us fix such a node  $x$ , and we shall show that  $x$ 's potential decreases by at least 1. Let  $k = \text{level}(x) = \text{level}(y)$ . Just prior to the path compression caused by the FIND-SET, we have

$$\begin{aligned} x.p.rank &\geq A_k^{(\text{iter}(x))}(x.rank) && \text{(by definition of iter}(x)) \text{ ,} \\ y.p.rank &\geq A_k(y.rank) && \text{(by definition of level}(y)) \text{ ,} \\ y.rank &\geq x.p.rank && \text{(by Corollary 21.5 and because} \\ &&& \text{y follows x on the find path) .} \end{aligned}$$

Putting these inequalities together and letting  $i$  be the value of  $\text{iter}(x)$  before path compression, we have

$$\begin{aligned} y.p.rank &\geq A_k(y.rank) \\ &\geq A_k(x.p.rank) && \text{(because } A_k(j) \text{ is strictly increasing)} \\ &\geq A_k(A_k^{(\text{iter}(x))}(x.rank)) \\ &= A_k^{(i+1)}(x.rank) . \end{aligned}$$

Because path compression will make  $x$  and  $y$  have the same parent, we know that after path compression,  $x.p.rank = y.p.rank$  and that the path compression does not decrease  $y.p.rank$ . Since  $x.rank$  does not change, after path compression we have that  $x.p.rank \geq A_k^{(i+1)}(x.rank)$ . Thus, path compression will cause either  $\text{iter}(x)$  to increase (to at least  $i + 1$ ) or  $\text{level}(x)$  to increase (which occurs if  $\text{iter}(x)$  increases to at least  $x.rank + 1$ ). In either case, by Lemma 21.10, we have  $\phi_q(x) \leq \phi_{q-1}(x) - 1$ . Hence,  $x$ 's potential decreases by at least 1.

The amortized cost of the FIND-SET operation is the actual cost plus the change in potential. The actual cost is  $O(s)$ , and we have shown that the total potential decreases by at least  $\max(0, s - (\alpha(n) + 2))$ . The amortized cost, therefore, is at most  $O(s) - (s - (\alpha(n) + 2)) = O(s) - s + O(\alpha(n)) = O(\alpha(n))$ , since we can scale up the units of potential to dominate the constant hidden in  $O(s)$ . ■

Putting the preceding lemmas together yields the following theorem.

**Theorem 21.14**

A sequence of  $m$  MAKE-SET, UNION, and FIND-SET operations,  $n$  of which are MAKE-SET operations, can be performed on a disjoint-set forest with union by rank and path compression in worst-case time  $O(m \alpha(n))$ .

**Proof** Immediate from Lemmas 21.7, 21.11, 21.12, and 21.13. ■

**Exercises**

**21.4-1**

Prove Lemma 21.4.

**21.4-2**

Prove that every node has rank at most  $\lfloor \lg n \rfloor$ .

**21.4-3**

In light of Exercise 21.4-2, how many bits are necessary to store  $x.rank$  for each node  $x$ ?

**21.4-4**

Using Exercise 21.4-2, give a simple proof that operations on a disjoint-set forest with union by rank but without path compression run in  $O(m \lg n)$  time.

**21.4-5**

Professor Dante reasons that because node ranks increase strictly along a simple path to the root, node levels must monotonically increase along the path. In other

words, if  $x.rank > 0$  and  $x.p$  is not a root, then  $level(x) \leq level(x.p)$ . Is the professor correct?

#### 21.4-6 ★

Consider the function  $\alpha'(n) = \min \{k : A_k(1) \geq \lg(n+1)\}$ . Show that  $\alpha'(n) \leq 3$  for all practical values of  $n$  and, using Exercise 21.4-2, show how to modify the potential-function argument to prove that we can perform a sequence of  $m$  MAKE-SET, UNION, and FIND-SET operations,  $n$  of which are MAKE-SET operations, on a disjoint-set forest with union by rank and path compression in worst-case time  $O(m \alpha'(n))$ .

## Problems

### 21-1 Off-line minimum

The *off-line minimum problem* asks us to maintain a dynamic set  $T$  of elements from the domain  $\{1, 2, \dots, n\}$  under the operations INSERT and EXTRACT-MIN. We are given a sequence  $S$  of  $n$  INSERT and  $m$  EXTRACT-MIN calls, where each key in  $\{1, 2, \dots, n\}$  is inserted exactly once. We wish to determine which key is returned by each EXTRACT-MIN call. Specifically, we wish to fill in an array *extracted*[1.. $m$ ], where for  $i = 1, 2, \dots, m$ , *extracted*[ $i$ ] is the key returned by the  $i$ th EXTRACT-MIN call. The problem is “off-line” in the sense that we are allowed to process the entire sequence  $S$  before determining any of the returned keys.

- a.* In the following instance of the off-line minimum problem, each operation INSERT( $i$ ) is represented by the value of  $i$  and each EXTRACT-MIN is represented by the letter E:

4, 8, E, 3, E, 9, 2, 6, E, E, E, 1, 7, E, 5 .

Fill in the correct values in the *extracted* array.

To develop an algorithm for this problem, we break the sequence  $S$  into homogeneous subsequences. That is, we represent  $S$  by

$I_1, E, I_2, E, I_3, \dots, I_m, E, I_{m+1}$  ,

where each E represents a single EXTRACT-MIN call and each  $I_j$  represents a (possibly empty) sequence of INSERT calls. For each subsequence  $I_j$ , we initially place the keys inserted by these operations into a set  $K_j$ , which is empty if  $I_j$  is empty. We then do the following:

OFF-LINE-MINIMUM( $m, n$ )

```

1  for  $i = 1$  to  $n$ 
2      determine  $j$  such that  $i \in K_j$ 
3      if  $j \neq m + 1$ 
4           $extracted[j] = i$ 
5          let  $l$  be the smallest value greater than  $j$ 
              for which set  $K_l$  exists
6           $K_l = K_j \cup K_l$ , destroying  $K_j$ 
7  return  $extracted$ 

```

- b.* Argue that the array *extracted* returned by OFF-LINE-MINIMUM is correct.
- c.* Describe how to implement OFF-LINE-MINIMUM efficiently with a disjoint-set data structure. Give a tight bound on the worst-case running time of your implementation.

### 21-2 Depth determination

In the *depth-determination problem*, we maintain a forest  $\mathcal{F} = \{T_i\}$  of rooted trees under three operations:

MAKE-TREE( $v$ ) creates a tree whose only node is  $v$ .

FIND-DEPTH( $v$ ) returns the depth of node  $v$  within its tree.

GRAFT( $r, v$ ) makes node  $r$ , which is assumed to be the root of a tree, become the child of node  $v$ , which is assumed to be in a different tree than  $r$  but may or may not itself be a root.

- a.* Suppose that we use a tree representation similar to a disjoint-set forest:  $v.p$  is the parent of node  $v$ , except that  $v.p = v$  if  $v$  is a root. Suppose further that we implement GRAFT( $r, v$ ) by setting  $r.p = v$  and FIND-DEPTH( $v$ ) by following the find path up to the root, returning a count of all nodes other than  $v$  encountered. Show that the worst-case running time of a sequence of  $m$  MAKE-TREE, FIND-DEPTH, and GRAFT operations is  $\Theta(m^2)$ .

By using the union-by-rank and path-compression heuristics, we can reduce the worst-case running time. We use the disjoint-set forest  $\mathcal{S} = \{S_i\}$ , where each set  $S_i$  (which is itself a tree) corresponds to a tree  $T_i$  in the forest  $\mathcal{F}$ . The tree structure within a set  $S_i$ , however, does not necessarily correspond to that of  $T_i$ . In fact, the implementation of  $S_i$  does not record the exact parent-child relationships but nevertheless allows us to determine any node's depth in  $T_i$ .

The key idea is to maintain in each node  $v$  a “pseudodistance”  $v.d$ , which is defined so that the sum of the pseudodistances along the simple path from  $v$  to the



root of its set  $S_i$  equals the depth of  $v$  in  $T_i$ . That is, if the simple path from  $v$  to its root in  $S_i$  is  $v_0, v_1, \dots, v_k$ , where  $v_0 = v$  and  $v_k$  is  $S_i$ 's root, then the depth of  $v$  in  $T_i$  is  $\sum_{j=0}^k v_j.d$ .

- b. Give an implementation of MAKE-TREE.
- c. Show how to modify FIND-SET to implement FIND-DEPTH. Your implementation should perform path compression, and its running time should be linear in the length of the find path. Make sure that your implementation updates pseudodistances correctly.
- d. Show how to implement GRAFT( $r, v$ ), which combines the sets containing  $r$  and  $v$ , by modifying the UNION and LINK procedures. Make sure that your implementation updates pseudodistances correctly. Note that the root of a set  $S_i$  is not necessarily the root of the corresponding tree  $T_i$ .
- e. Give a tight bound on the worst-case running time of a sequence of  $m$  MAKE-TREE, FIND-DEPTH, and GRAFT operations,  $n$  of which are MAKE-TREE operations.

### 21-3 Tarjan's off-line least-common-ancestors algorithm

The **least common ancestor** of two nodes  $u$  and  $v$  in a rooted tree  $T$  is the node  $w$  that is an ancestor of both  $u$  and  $v$  and that has the greatest depth in  $T$ . In the **off-line least-common-ancestors problem**, we are given a rooted tree  $T$  and an arbitrary set  $P = \{\{u, v\}\}$  of unordered pairs of nodes in  $T$ , and we wish to determine the least common ancestor of each pair in  $P$ .

To solve the off-line least-common-ancestors problem, the following procedure performs a tree walk of  $T$  with the initial call  $\text{LCA}(T.\text{root})$ . We assume that each node is colored WHITE prior to the walk.

LCA( $u$ )

```

1  MAKE-SET( $u$ )
2  FIND-SET( $u$ ).ancestor =  $u$ 
3  for each child  $v$  of  $u$  in  $T$ 
4      LCA( $v$ )
5      UNION( $u, v$ )
6      FIND-SET( $u$ ).ancestor =  $u$ 
7   $u.color = \text{BLACK}$ 
8  for each node  $v$  such that  $\{u, v\} \in P$ 
9      if  $v.color == \text{BLACK}$ 
10         print "The least common ancestor of"
               $u$  "and"  $v$  "is" FIND-SET( $v$ ).ancestor
```

- a.* Argue that line 10 executes exactly once for each pair  $\{u, v\} \in P$ .
- b.* Argue that at the time of the call  $\text{LCA}(u)$ , the number of sets in the disjoint-set data structure equals the depth of  $u$  in  $T$ .
- c.* Prove that  $\text{LCA}$  correctly prints the least common ancestor of  $u$  and  $v$  for each pair  $\{u, v\} \in P$ .
- d.* Analyze the running time of  $\text{LCA}$ , assuming that we use the implementation of the disjoint-set data structure in Section 21.3.

---

## Chapter notes

Many of the important results for disjoint-set data structures are due at least in part to R. E. Tarjan. Using aggregate analysis, Tarjan [328, 330] gave the first tight upper bound in terms of the very slowly growing inverse  $\hat{\alpha}(m, n)$  of Ackermann's function. (The function  $A_k(j)$  given in Section 21.4 is similar to Ackermann's function, and the function  $\alpha(n)$  is similar to the inverse. Both  $\alpha(n)$  and  $\hat{\alpha}(m, n)$  are at most 4 for all conceivable values of  $m$  and  $n$ .) An  $O(m \lg^* n)$  upper bound was proven earlier by Hopcroft and Ullman [5, 179]. The treatment in Section 21.4 is adapted from a later analysis by Tarjan [332], which is in turn based on an analysis by Kozen [220]. Harfst and Reingold [161] give a potential-based version of Tarjan's earlier bound.

Tarjan and van Leeuwen [333] discuss variants on the path-compression heuristic, including "one-pass methods," which sometimes offer better constant factors in their performance than do two-pass methods. As with Tarjan's earlier analyses of the basic path-compression heuristic, the analyses by Tarjan and van Leeuwen are aggregate. Harfst and Reingold [161] later showed how to make a small change to the potential function to adapt their path-compression analysis to these one-pass variants. Gabow and Tarjan [121] show that in certain applications, the disjoint-set operations can be made to run in  $O(m)$  time.

Tarjan [329] showed that a lower bound of  $\Omega(m \hat{\alpha}(m, n))$  time is required for operations on any disjoint-set data structure satisfying certain technical conditions. This lower bound was later generalized by Fredman and Saks [113], who showed that in the worst case,  $\Omega(m \hat{\alpha}(m, n))$  ( $\lg n$ )-bit words of memory must be accessed.

---

## ***VI Graph Algorithms***

---

## Introduction

Graph problems pervade computer science, and algorithms for working with them are fundamental to the field. Hundreds of interesting computational problems are couched in terms of graphs. In this part, we touch on a few of the more significant ones.

Chapter 22 shows how we can represent a graph in a computer and then discusses algorithms based on searching a graph using either breadth-first search or depth-first search. The chapter gives two applications of depth-first search: topologically sorting a directed acyclic graph and decomposing a directed graph into its strongly connected components.

Chapter 23 describes how to compute a minimum-weight spanning tree of a graph: the least-weight way of connecting all of the vertices together when each edge has an associated weight. The algorithms for computing minimum spanning trees serve as good examples of greedy algorithms (see Chapter 16).

Chapters 24 and 25 consider how to compute shortest paths between vertices when each edge has an associated length or “weight.” Chapter 24 shows how to find shortest paths from a given source vertex to all other vertices, and Chapter 25 examines methods to compute shortest paths between every pair of vertices.

Finally, Chapter 26 shows how to compute a maximum flow of material in a flow network, which is a directed graph having a specified source vertex of material, a specified sink vertex, and specified capacities for the amount of material that can traverse each directed edge. This general problem arises in many forms, and a good algorithm for computing maximum flows can help solve a variety of related problems efficiently.

When we characterize the running time of a graph algorithm on a given graph  $G = (V, E)$ , we usually measure the size of the input in terms of the number of vertices  $|V|$  and the number of edges  $|E|$  of the graph. That is, we describe the size of the input with two parameters, not just one. We adopt a common notational convention for these parameters. Inside asymptotic notation (such as  $O$ -notation or  $\Theta$ -notation), and *only* inside such notation, the symbol  $V$  denotes  $|V|$  and the symbol  $E$  denotes  $|E|$ . For example, we might say, “the algorithm runs in time  $O(VE)$ ,” meaning that the algorithm runs in time  $O(|V| |E|)$ . This convention makes the running-time formulas easier to read, without risk of ambiguity.

Another convention we adopt appears in pseudocode. We denote the vertex set of a graph  $G$  by  $G.V$  and its edge set by  $G.E$ . That is, the pseudocode views vertex and edge sets as attributes of a graph.

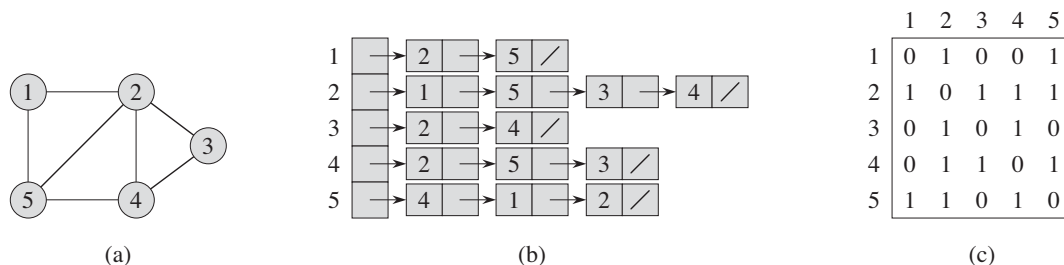
This chapter presents methods for representing a graph and for searching a graph. Searching a graph means systematically following the edges of the graph so as to visit the vertices of the graph. A graph-searching algorithm can discover much about the structure of a graph. Many algorithms begin by searching their input graph to obtain this structural information. Several other graph algorithms elaborate on basic graph searching. Techniques for searching a graph lie at the heart of the field of graph algorithms.

Section 22.1 discusses the two most common computational representations of graphs: as adjacency lists and as adjacency matrices. Section 22.2 presents a simple graph-searching algorithm called breadth-first search and shows how to create a breadth-first tree. Section 22.3 presents depth-first search and proves some standard results about the order in which depth-first search visits vertices. Section 22.4 provides our first real application of depth-first search: topologically sorting a directed acyclic graph. A second application of depth-first search, finding the strongly connected components of a directed graph, is the topic of Section 22.5.

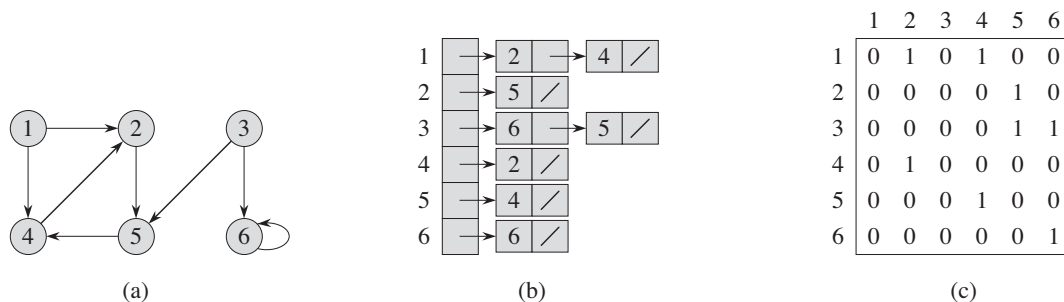
---

### 22.1 Representations of graphs

We can choose between two standard ways to represent a graph  $G = (V, E)$ : as a collection of adjacency lists or as an adjacency matrix. Either way applies to both directed and undirected graphs. Because the adjacency-list representation provides a compact way to represent *sparse* graphs—those for which  $|E|$  is much less than  $|V|^2$ —it is usually the method of choice. Most of the graph algorithms presented in this book assume that an input graph is represented in adjacency-list form. We may prefer an adjacency-matrix representation, however, when the graph is *dense*— $|E|$  is close to  $|V|^2$ —or when we need to be able to tell quickly if there is an edge connecting two given vertices. For example, two of the all-pairs



**Figure 22.1** Two representations of an undirected graph. (a) An undirected graph  $G$  with 5 vertices and 7 edges. (b) An adjacency-list representation of  $G$ . (c) The adjacency-matrix representation of  $G$ .



**Figure 22.2** Two representations of a directed graph. (a) A directed graph  $G$  with 6 vertices and 8 edges. (b) An adjacency-list representation of  $G$ . (c) The adjacency-matrix representation of  $G$ .

shortest-paths algorithms presented in Chapter 25 assume that their input graphs are represented by adjacency matrices.

The **adjacency-list representation** of a graph  $G = (V, E)$  consists of an array  $Adj$  of  $|V|$  lists, one for each vertex in  $V$ . For each  $u \in V$ , the adjacency list  $Adj[u]$  contains all the vertices  $v$  such that there is an edge  $(u, v) \in E$ . That is,  $Adj[u]$  consists of all the vertices adjacent to  $u$  in  $G$ . (Alternatively, it may contain pointers to these vertices.) Since the adjacency lists represent the edges of a graph, in pseudocode we treat the array  $Adj$  as an attribute of the graph, just as we treat the edge set  $E$ . In pseudocode, therefore, we will see notation such as  $G.Adj[u]$ . Figure 22.1(b) is an adjacency-list representation of the undirected graph in Figure 22.1(a). Similarly, Figure 22.2(b) is an adjacency-list representation of the directed graph in Figure 22.2(a).

If  $G$  is a directed graph, the sum of the lengths of all the adjacency lists is  $|E|$ , since an edge of the form  $(u, v)$  is represented by having  $v$  appear in  $Adj[u]$ . If  $G$  is

an undirected graph, the sum of the lengths of all the adjacency lists is  $2|E|$ , since if  $(u, v)$  is an undirected edge, then  $u$  appears in  $v$ 's adjacency list and vice versa. For both directed and undirected graphs, the adjacency-list representation has the desirable property that the amount of memory it requires is  $\Theta(V + E)$ .

We can readily adapt adjacency lists to represent **weighted graphs**, that is, graphs for which each edge has an associated **weight**, typically given by a **weight function**  $w : E \rightarrow \mathbb{R}$ . For example, let  $G = (V, E)$  be a weighted graph with weight function  $w$ . We simply store the weight  $w(u, v)$  of the edge  $(u, v) \in E$  with vertex  $v$  in  $u$ 's adjacency list. The adjacency-list representation is quite robust in that we can modify it to support many other graph variants.

A potential disadvantage of the adjacency-list representation is that it provides no quicker way to determine whether a given edge  $(u, v)$  is present in the graph than to search for  $v$  in the adjacency list  $Adj[u]$ . An adjacency-matrix representation of the graph remedies this disadvantage, but at the cost of using asymptotically more memory. (See Exercise 22.1-8 for suggestions of variations on adjacency lists that permit faster edge lookup.)

For the **adjacency-matrix representation** of a graph  $G = (V, E)$ , we assume that the vertices are numbered  $1, 2, \dots, |V|$  in some arbitrary manner. Then the adjacency-matrix representation of a graph  $G$  consists of a  $|V| \times |V|$  matrix  $A = (a_{ij})$  such that

$$a_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

Figures 22.1(c) and 22.2(c) are the adjacency matrices of the undirected and directed graphs in Figures 22.1(a) and 22.2(a), respectively. The adjacency matrix of a graph requires  $\Theta(V^2)$  memory, independent of the number of edges in the graph.

Observe the symmetry along the main diagonal of the adjacency matrix in Figure 22.1(c). Since in an undirected graph,  $(u, v)$  and  $(v, u)$  represent the same edge, the adjacency matrix  $A$  of an undirected graph is its own transpose:  $A = A^T$ . In some applications, it pays to store only the entries on and above the diagonal of the adjacency matrix, thereby cutting the memory needed to store the graph almost in half.

Like the adjacency-list representation of a graph, an adjacency matrix can represent a weighted graph. For example, if  $G = (V, E)$  is a weighted graph with edge-weight function  $w$ , we can simply store the weight  $w(u, v)$  of the edge  $(u, v) \in E$  as the entry in row  $u$  and column  $v$  of the adjacency matrix. If an edge does not exist, we can store a NIL value as its corresponding matrix entry, though for many problems it is convenient to use a value such as 0 or  $\infty$ .

Although the adjacency-list representation is asymptotically at least as space-efficient as the adjacency-matrix representation, adjacency matrices are simpler, and so we may prefer them when graphs are reasonably small. Moreover, adja-



acency matrices carry a further advantage for unweighted graphs: they require only one bit per entry.

### Representing attributes

Most algorithms that operate on graphs need to maintain attributes for vertices and/or edges. We indicate these attributes using our usual notation, such as  $v.d$  for an attribute  $d$  of a vertex  $v$ . When we indicate edges as pairs of vertices, we use the same style of notation. For example, if edges have an attribute  $f$ , then we denote this attribute for edge  $(u, v)$  by  $(u, v).f$ . For the purpose of presenting and understanding algorithms, our attribute notation suffices.

Implementing vertex and edge attributes in real programs can be another story entirely. There is no one best way to store and access vertex and edge attributes. For a given situation, your decision will likely depend on the programming language you are using, the algorithm you are implementing, and how the rest of your program uses the graph. If you represent a graph using adjacency lists, one design represents vertex attributes in additional arrays, such as an array  $d[1..|V|]$  that parallels the  $Adj$  array. If the vertices adjacent to  $u$  are in  $Adj[u]$ , then what we call the attribute  $u.d$  would actually be stored in the array entry  $d[u]$ . Many other ways of implementing attributes are possible. For example, in an object-oriented programming language, vertex attributes might be represented as instance variables within a subclass of a `Vertex` class.

### Exercises

#### 22.1-1

Given an adjacency-list representation of a directed graph, how long does it take to compute the out-degree of every vertex? How long does it take to compute the in-degrees?

#### 22.1-2

Give an adjacency-list representation for a complete binary tree on 7 vertices. Give an equivalent adjacency-matrix representation. Assume that vertices are numbered from 1 to 7 as in a binary heap.

#### 22.1-3

The **transpose** of a directed graph  $G = (V, E)$  is the graph  $G^T = (V, E^T)$ , where  $E^T = \{(v, u) \in V \times V : (u, v) \in E\}$ . Thus,  $G^T$  is  $G$  with all its edges reversed. Describe efficient algorithms for computing  $G^T$  from  $G$ , for both the adjacency-list and adjacency-matrix representations of  $G$ . Analyze the running times of your algorithms.

**22.1-4**

Given an adjacency-list representation of a multigraph  $G = (V, E)$ , describe an  $O(V + E)$ -time algorithm to compute the adjacency-list representation of the “equivalent” undirected graph  $G' = (V, E')$ , where  $E'$  consists of the edges in  $E$  with all multiple edges between two vertices replaced by a single edge and with all self-loops removed.

**22.1-5**

The **square** of a directed graph  $G = (V, E)$  is the graph  $G^2 = (V, E^2)$  such that  $(u, v) \in E^2$  if and only if  $G$  contains a path with at most two edges between  $u$  and  $v$ . Describe efficient algorithms for computing  $G^2$  from  $G$  for both the adjacency-list and adjacency-matrix representations of  $G$ . Analyze the running times of your algorithms.

**22.1-6**

Most graph algorithms that take an adjacency-matrix representation as input require time  $\Omega(V^2)$ , but there are some exceptions. Show how to determine whether a directed graph  $G$  contains a **universal sink**—a vertex with in-degree  $|V| - 1$  and out-degree 0—in time  $O(V)$ , given an adjacency matrix for  $G$ .

**22.1-7**

The **incidence matrix** of a directed graph  $G = (V, E)$  with no self-loops is a  $|V| \times |E|$  matrix  $B = (b_{ij})$  such that

$$b_{ij} = \begin{cases} -1 & \text{if edge } j \text{ leaves vertex } i, \\ 1 & \text{if edge } j \text{ enters vertex } i, \\ 0 & \text{otherwise.} \end{cases}$$

Describe what the entries of the matrix product  $BB^T$  represent, where  $B^T$  is the transpose of  $B$ .

**22.1-8**

Suppose that instead of a linked list, each array entry  $Adj[u]$  is a hash table containing the vertices  $v$  for which  $(u, v) \in E$ . If all edge lookups are equally likely, what is the expected time to determine whether an edge is in the graph? What disadvantages does this scheme have? Suggest an alternate data structure for each edge list that solves these problems. Does your alternative have disadvantages compared to the hash table?

---

## 22.2 Breadth-first search

**Breadth-first search** is one of the simplest algorithms for searching a graph and the archetype for many important graph algorithms. Prim’s minimum-spanning-tree algorithm (Section 23.2) and Dijkstra’s single-source shortest-paths algorithm (Section 24.3) use ideas similar to those in breadth-first search.

Given a graph  $G = (V, E)$  and a distinguished **source** vertex  $s$ , breadth-first search systematically explores the edges of  $G$  to “discover” every vertex that is reachable from  $s$ . It computes the distance (smallest number of edges) from  $s$  to each reachable vertex. It also produces a “breadth-first tree” with root  $s$  that contains all reachable vertices. For any vertex  $v$  reachable from  $s$ , the simple path in the breadth-first tree from  $s$  to  $v$  corresponds to a “shortest path” from  $s$  to  $v$  in  $G$ , that is, a path containing the smallest number of edges. The algorithm works on both directed and undirected graphs.

Breadth-first search is so named because it expands the frontier between discovered and undiscovered vertices uniformly across the breadth of the frontier. That is, the algorithm discovers all vertices at distance  $k$  from  $s$  before discovering any vertices at distance  $k + 1$ .

To keep track of progress, breadth-first search colors each vertex white, gray, or black. All vertices start out white and may later become gray and then black. A vertex is **discovered** the first time it is encountered during the search, at which time it becomes nonwhite. Gray and black vertices, therefore, have been discovered, but breadth-first search distinguishes between them to ensure that the search proceeds in a breadth-first manner.<sup>1</sup> If  $(u, v) \in E$  and vertex  $u$  is black, then vertex  $v$  is either gray or black; that is, all vertices adjacent to black vertices have been discovered. Gray vertices may have some adjacent white vertices; they represent the frontier between discovered and undiscovered vertices.

Breadth-first search constructs a breadth-first tree, initially containing only its root, which is the source vertex  $s$ . Whenever the search discovers a white vertex  $v$  in the course of scanning the adjacency list of an already discovered vertex  $u$ , the vertex  $v$  and the edge  $(u, v)$  are added to the tree. We say that  $u$  is the **predecessor** or **parent** of  $v$  in the breadth-first tree. Since a vertex is discovered at most once, it has at most one parent. Ancestor and descendant relationships in the breadth-first tree are defined relative to the root  $s$  as usual: if  $u$  is on the simple path in the tree from the root  $s$  to vertex  $v$ , then  $u$  is an ancestor of  $v$  and  $v$  is a descendant of  $u$ .

---

<sup>1</sup>We distinguish between gray and black vertices to help us understand how breadth-first search operates. In fact, as Exercise 22.2-3 shows, we would get the same result even if we did not distinguish between gray and black vertices.

The breadth-first-search procedure BFS below assumes that the input graph  $G = (V, E)$  is represented using adjacency lists. It attaches several additional attributes to each vertex in the graph. We store the color of each vertex  $u \in V$  in the attribute  $u.color$  and the predecessor of  $u$  in the attribute  $u.\pi$ . If  $u$  has no predecessor (for example, if  $u = s$  or  $u$  has not been discovered), then  $u.\pi = \text{NIL}$ . The attribute  $u.d$  holds the distance from the source  $s$  to vertex  $u$  computed by the algorithm. The algorithm also uses a first-in, first-out queue  $Q$  (see Section 10.1) to manage the set of gray vertices.

BFS( $G, s$ )

```

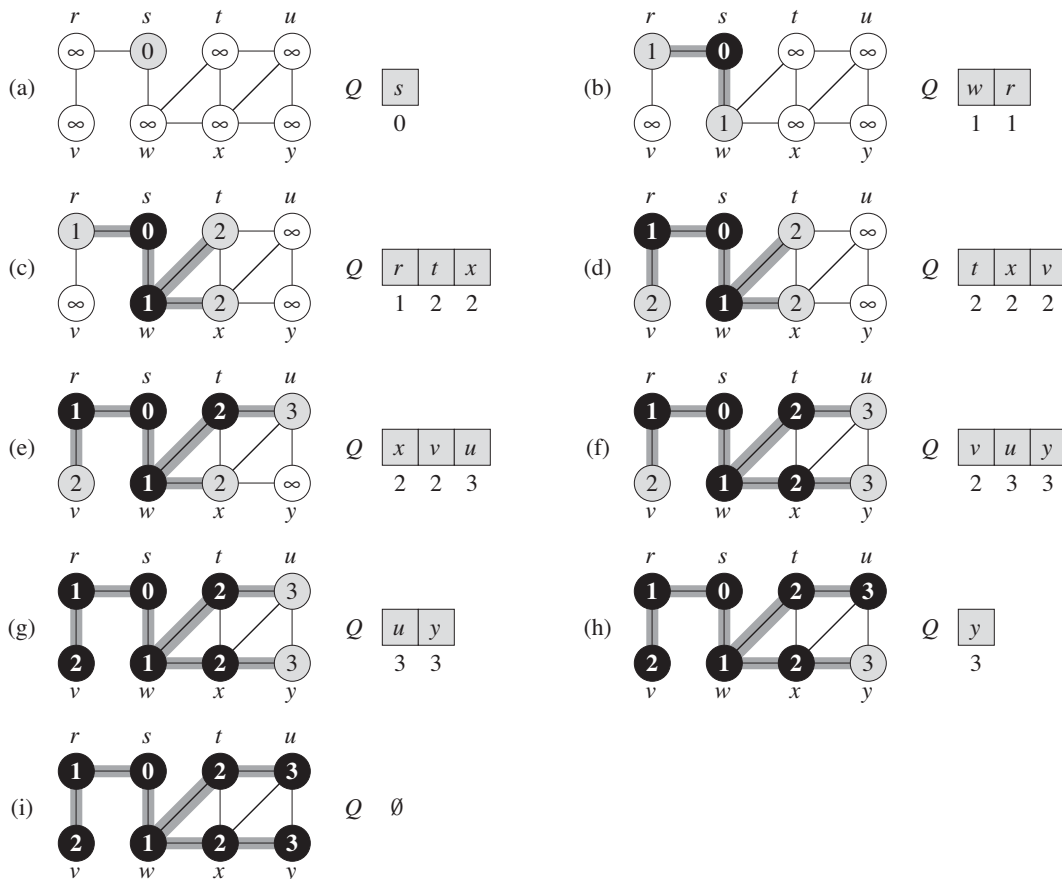
1  for each vertex  $u \in G.V - \{s\}$ 
2       $u.color = \text{WHITE}$ 
3       $u.d = \infty$ 
4       $u.\pi = \text{NIL}$ 
5   $s.color = \text{GRAY}$ 
6   $s.d = 0$ 
7   $s.\pi = \text{NIL}$ 
8   $Q = \emptyset$ 
9  ENQUEUE( $Q, s$ )
10 while  $Q \neq \emptyset$ 
11      $u = \text{DEQUEUE}(Q)$ 
12     for each  $v \in G.Adj[u]$ 
13         if  $v.color == \text{WHITE}$ 
14              $v.color = \text{GRAY}$ 
15              $v.d = u.d + 1$ 
16              $v.\pi = u$ 
17             ENQUEUE( $Q, v$ )
18      $u.color = \text{BLACK}$ 
```

Figure 22.3 illustrates the progress of BFS on a sample graph.

The procedure BFS works as follows. With the exception of the source vertex  $s$ , lines 1–4 paint every vertex white, set  $u.d$  to be infinity for each vertex  $u$ , and set the parent of every vertex to be NIL. Line 5 paints  $s$  gray, since we consider it to be discovered as the procedure begins. Line 6 initializes  $s.d$  to 0, and line 7 sets the predecessor of the source to be NIL. Lines 8–9 initialize  $Q$  to the queue containing just the vertex  $s$ .

The **while** loop of lines 10–18 iterates as long as there remain gray vertices, which are discovered vertices that have not yet had their adjacency lists fully examined. This **while** loop maintains the following invariant:

At the test in line 10, the queue  $Q$  consists of the set of gray vertices.



**Figure 22.3** The operation of BFS on an undirected graph. Tree edges are shown shaded as they are produced by BFS. The value of  $u.d$  appears within each vertex  $u$ . The queue  $Q$  is shown at the beginning of each iteration of the **while** loop of lines 10–18. Vertex distances appear below vertices in the queue.

Although we won't use this loop invariant to prove correctness, it is easy to see that it holds prior to the first iteration and that each iteration of the loop maintains the invariant. Prior to the first iteration, the only gray vertex, and the only vertex in  $Q$ , is the source vertex  $s$ . Line 11 determines the gray vertex  $u$  at the head of the queue  $Q$  and removes it from  $Q$ . The **for** loop of lines 12–17 considers each vertex  $v$  in the adjacency list of  $u$ . If  $v$  is white, then it has not yet been discovered, and the procedure discovers it by executing lines 14–17. The procedure paints vertex  $v$  gray, sets its distance  $v.d$  to  $u.d + 1$ , records  $u$  as its parent  $v.\pi$ , and places it at the tail of the queue  $Q$ . Once the procedure has examined all the vertices on  $u$ 's

adjacency list, it blackens  $u$  in line 18. The loop invariant is maintained because whenever a vertex is painted gray (in line 14) it is also enqueued (in line 17), and whenever a vertex is dequeued (in line 11) it is also painted black (in line 18).

The results of breadth-first search may depend upon the order in which the neighbors of a given vertex are visited in line 12: the breadth-first tree may vary, but the distances  $d$  computed by the algorithm will not. (See Exercise 22.2-5.)

### Analysis

Before proving the various properties of breadth-first search, we take on the somewhat easier job of analyzing its running time on an input graph  $G = (V, E)$ . We use aggregate analysis, as we saw in Section 17.1. After initialization, breadth-first search never whitens a vertex, and thus the test in line 13 ensures that each vertex is enqueued at most once, and hence dequeued at most once. The operations of enqueueing and dequeuing take  $O(1)$  time, and so the total time devoted to queue operations is  $O(V)$ . Because the procedure scans the adjacency list of each vertex only when the vertex is dequeued, it scans each adjacency list at most once. Since the sum of the lengths of all the adjacency lists is  $\Theta(E)$ , the total time spent in scanning adjacency lists is  $O(E)$ . The overhead for initialization is  $O(V)$ , and thus the total running time of the BFS procedure is  $O(V + E)$ . Thus, breadth-first search runs in time linear in the size of the adjacency-list representation of  $G$ .

### Shortest paths

At the beginning of this section, we claimed that breadth-first search finds the distance to each reachable vertex in a graph  $G = (V, E)$  from a given source vertex  $s \in V$ . Define the **shortest-path distance**  $\delta(s, v)$  from  $s$  to  $v$  as the minimum number of edges in any path from vertex  $s$  to vertex  $v$ ; if there is no path from  $s$  to  $v$ , then  $\delta(s, v) = \infty$ . We call a path of length  $\delta(s, v)$  from  $s$  to  $v$  a **shortest path**<sup>2</sup> from  $s$  to  $v$ . Before showing that breadth-first search correctly computes shortest-path distances, we investigate an important property of shortest-path distances.

---

<sup>2</sup>In Chapters 24 and 25, we shall generalize our study of shortest paths to weighted graphs, in which every edge has a real-valued weight and the weight of a path is the sum of the weights of its constituent edges. The graphs considered in the present chapter are unweighted or, equivalently, all edges have unit weight.

**Lemma 22.1**

Let  $G = (V, E)$  be a directed or undirected graph, and let  $s \in V$  be an arbitrary vertex. Then, for any edge  $(u, v) \in E$ ,

$$\delta(s, v) \leq \delta(s, u) + 1 .$$

**Proof** If  $u$  is reachable from  $s$ , then so is  $v$ . In this case, the shortest path from  $s$  to  $v$  cannot be longer than the shortest path from  $s$  to  $u$  followed by the edge  $(u, v)$ , and thus the inequality holds. If  $u$  is not reachable from  $s$ , then  $\delta(s, u) = \infty$ , and the inequality holds. ■

We want to show that BFS properly computes  $v.d = \delta(s, v)$  for each vertex  $v \in V$ . We first show that  $v.d$  bounds  $\delta(s, v)$  from above.

**Lemma 22.2**

Let  $G = (V, E)$  be a directed or undirected graph, and suppose that BFS is run on  $G$  from a given source vertex  $s \in V$ . Then upon termination, for each vertex  $v \in V$ , the value  $v.d$  computed by BFS satisfies  $v.d \geq \delta(s, v)$ .

**Proof** We use induction on the number of ENQUEUE operations. Our inductive hypothesis is that  $v.d \geq \delta(s, v)$  for all  $v \in V$ .

The basis of the induction is the situation immediately after enqueueing  $s$  in line 9 of BFS. The inductive hypothesis holds here, because  $s.d = 0 = \delta(s, s)$  and  $v.d = \infty \geq \delta(s, v)$  for all  $v \in V - \{s\}$ .

For the inductive step, consider a white vertex  $v$  that is discovered during the search from a vertex  $u$ . The inductive hypothesis implies that  $u.d \geq \delta(s, u)$ . From the assignment performed by line 15 and from Lemma 22.1, we obtain

$$\begin{aligned} v.d &= u.d + 1 \\ &\geq \delta(s, u) + 1 \\ &\geq \delta(s, v) . \end{aligned}$$

Vertex  $v$  is then enqueued, and it is never enqueued again because it is also grayed and the **then** clause of lines 14–17 is executed only for white vertices. Thus, the value of  $v.d$  never changes again, and the inductive hypothesis is maintained. ■

To prove that  $v.d = \delta(s, v)$ , we must first show more precisely how the queue  $Q$  operates during the course of BFS. The next lemma shows that at all times, the queue holds at most two distinct  $d$  values.

**Lemma 22.3**

Suppose that during the execution of BFS on a graph  $G = (V, E)$ , the queue  $Q$  contains the vertices  $\langle v_1, v_2, \dots, v_r \rangle$ , where  $v_1$  is the head of  $Q$  and  $v_r$  is the tail. Then,  $v_r.d \leq v_1.d + 1$  and  $v_i.d \leq v_{i+1}.d$  for  $i = 1, 2, \dots, r - 1$ .

**Proof** The proof is by induction on the number of queue operations. Initially, when the queue contains only  $s$ , the lemma certainly holds.

For the inductive step, we must prove that the lemma holds after both dequeuing and enqueueing a vertex. If the head  $v_1$  of the queue is dequeued,  $v_2$  becomes the new head. (If the queue becomes empty, then the lemma holds vacuously.) By the inductive hypothesis,  $v_1.d \leq v_2.d$ . But then we have  $v_r.d \leq v_1.d + 1 \leq v_2.d + 1$ , and the remaining inequalities are unaffected. Thus, the lemma follows with  $v_2$  as the head.

In order to understand what happens upon enqueueing a vertex, we need to examine the code more closely. When we enqueue a vertex  $v$  in line 17 of BFS, it becomes  $v_{r+1}$ . At that time, we have already removed vertex  $u$ , whose adjacency list is currently being scanned, from the queue  $Q$ , and by the inductive hypothesis, the new head  $v_1$  has  $v_1.d \geq u.d$ . Thus,  $v_{r+1}.d = v.d = u.d + 1 \leq v_1.d + 1$ . From the inductive hypothesis, we also have  $v_r.d \leq u.d + 1$ , and so  $v_r.d \leq u.d + 1 = v.d = v_{r+1}.d$ , and the remaining inequalities are unaffected. Thus, the lemma follows when  $v$  is enqueued. ■

The following corollary shows that the  $d$  values at the time that vertices are enqueued are monotonically increasing over time.

**Corollary 22.4**

Suppose that vertices  $v_i$  and  $v_j$  are enqueued during the execution of BFS, and that  $v_i$  is enqueued before  $v_j$ . Then  $v_i.d \leq v_j.d$  at the time that  $v_j$  is enqueued.

**Proof** Immediate from Lemma 22.3 and the property that each vertex receives a finite  $d$  value at most once during the course of BFS. ■

We can now prove that breadth-first search correctly finds shortest-path distances.

**Theorem 22.5 (Correctness of breadth-first search)**

Let  $G = (V, E)$  be a directed or undirected graph, and suppose that BFS is run on  $G$  from a given source vertex  $s \in V$ . Then, during its execution, BFS discovers every vertex  $v \in V$  that is reachable from the source  $s$ , and upon termination,  $v.d = \delta(s, v)$  for all  $v \in V$ . Moreover, for any vertex  $v \neq s$  that is reachable



from  $s$ , one of the shortest paths from  $s$  to  $v$  is a shortest path from  $s$  to  $v.\pi$  followed by the edge  $(v.\pi, v)$ .

**Proof** Assume, for the purpose of contradiction, that some vertex receives a  $d$  value not equal to its shortest-path distance. Let  $v$  be the vertex with minimum  $\delta(s, v)$  that receives such an incorrect  $d$  value; clearly  $v \neq s$ . By Lemma 22.2,  $v.d \geq \delta(s, v)$ , and thus we have that  $v.d > \delta(s, v)$ . Vertex  $v$  must be reachable from  $s$ , for if it is not, then  $\delta(s, v) = \infty \geq v.d$ . Let  $u$  be the vertex immediately preceding  $v$  on a shortest path from  $s$  to  $v$ , so that  $\delta(s, v) = \delta(s, u) + 1$ . Because  $\delta(s, u) < \delta(s, v)$ , and because of how we chose  $v$ , we have  $u.d = \delta(s, u)$ . Putting these properties together, we have

$$v.d > \delta(s, v) = \delta(s, u) + 1 = u.d + 1. \quad (22.1)$$

Now consider the time when BFS chooses to dequeue vertex  $u$  from  $Q$  in line 11. At this time, vertex  $v$  is either white, gray, or black. We shall show that in each of these cases, we derive a contradiction to inequality (22.1). If  $v$  is white, then line 15 sets  $v.d = u.d + 1$ , contradicting inequality (22.1). If  $v$  is black, then it was already removed from the queue and, by Corollary 22.4, we have  $v.d \leq u.d$ , again contradicting inequality (22.1). If  $v$  is gray, then it was painted gray upon dequeuing some vertex  $w$ , which was removed from  $Q$  earlier than  $u$  and for which  $v.d = w.d + 1$ . By Corollary 22.4, however,  $w.d \leq u.d$ , and so we have  $v.d = w.d + 1 \leq u.d + 1$ , once again contradicting inequality (22.1).

Thus we conclude that  $v.d = \delta(s, v)$  for all  $v \in V$ . All vertices  $v$  reachable from  $s$  must be discovered, for otherwise they would have  $\infty = v.d > \delta(s, v)$ . To conclude the proof of the theorem, observe that if  $v.\pi = u$ , then  $v.d = u.d + 1$ . Thus, we can obtain a shortest path from  $s$  to  $v$  by taking a shortest path from  $s$  to  $v.\pi$  and then traversing the edge  $(v.\pi, v)$ . ■

### Breadth-first trees

The procedure BFS builds a breadth-first tree as it searches the graph, as Figure 22.3 illustrates. The tree corresponds to the  $\pi$  attributes. More formally, for a graph  $G = (V, E)$  with source  $s$ , we define the **predecessor subgraph** of  $G$  as  $G_\pi = (V_\pi, E_\pi)$ , where

$$V_\pi = \{v \in V : v.\pi \neq \text{NIL}\} \cup \{s\}$$

and

$$E_\pi = \{(v.\pi, v) : v \in V_\pi - \{s\}\}.$$

The predecessor subgraph  $G_\pi$  is a **breadth-first tree** if  $V_\pi$  consists of the vertices reachable from  $s$  and, for all  $v \in V_\pi$ , the subgraph  $G_\pi$  contains a unique simple

path from  $s$  to  $v$  that is also a shortest path from  $s$  to  $v$  in  $G$ . A breadth-first tree is in fact a tree, since it is connected and  $|E_\pi| = |V_\pi| - 1$  (see Theorem B.2). We call the edges in  $E_\pi$  **tree edges**.

The following lemma shows that the predecessor subgraph produced by the BFS procedure is a breadth-first tree.

**Lemma 22.6**

When applied to a directed or undirected graph  $G = (V, E)$ , procedure BFS constructs  $\pi$  so that the predecessor subgraph  $G_\pi = (V_\pi, E_\pi)$  is a breadth-first tree.

**Proof** Line 16 of BFS sets  $v.\pi = u$  if and only if  $(u, v) \in E$  and  $\delta(s, v) < \infty$ —that is, if  $v$  is reachable from  $s$ —and thus  $V_\pi$  consists of the vertices in  $V$  reachable from  $s$ . Since  $G_\pi$  forms a tree, by Theorem B.2, it contains a unique simple path from  $s$  to each vertex in  $V_\pi$ . By applying Theorem 22.5 inductively, we conclude that every such path is a shortest path in  $G$ . ■

The following procedure prints out the vertices on a shortest path from  $s$  to  $v$ , assuming that BFS has already computed a breadth-first tree:

```

PRINT-PATH( $G, s, v$ )
1  if  $v == s$ 
2    print  $s$ 
3  elseif  $v.\pi == \text{NIL}$ 
4    print “no path from”  $s$  “to”  $v$  “exists”
5  else PRINT-PATH( $G, s, v.\pi$ )
6    print  $v$ 

```

This procedure runs in time linear in the number of vertices in the path printed, since each recursive call is for a path one vertex shorter.

## Exercises

### 22.2-1

Show the  $d$  and  $\pi$  values that result from running breadth-first search on the directed graph of Figure 22.2(a), using vertex 3 as the source.

### 22.2-2

Show the  $d$  and  $\pi$  values that result from running breadth-first search on the undirected graph of Figure 22.3, using vertex  $u$  as the source.

**22.2-3**

Show that using a single bit to store each vertex color suffices by arguing that the BFS procedure would produce the same result if lines 5 and 14 were removed.

**22.2-4**

What is the running time of BFS if we represent its input graph by an adjacency matrix and modify the algorithm to handle this form of input?

**22.2-5**

Argue that in a breadth-first search, the value  $u.d$  assigned to a vertex  $u$  is independent of the order in which the vertices appear in each adjacency list. Using Figure 22.3 as an example, show that the breadth-first tree computed by BFS can depend on the ordering within adjacency lists.

**22.2-6**

Give an example of a directed graph  $G = (V, E)$ , a source vertex  $s \in V$ , and a set of tree edges  $E_\pi \subseteq E$  such that for each vertex  $v \in V$ , the unique simple path in the graph  $(V, E_\pi)$  from  $s$  to  $v$  is a shortest path in  $G$ , yet the set of edges  $E_\pi$  cannot be produced by running BFS on  $G$ , no matter how the vertices are ordered in each adjacency list.

**22.2-7**

There are two types of professional wrestlers: “babyfaces” (“good guys”) and “heels” (“bad guys”). Between any pair of professional wrestlers, there may or may not be a rivalry. Suppose we have  $n$  professional wrestlers and we have a list of  $r$  pairs of wrestlers for which there are rivalries. Give an  $O(n + r)$ -time algorithm that determines whether it is possible to designate some of the wrestlers as babyfaces and the remainder as heels such that each rivalry is between a babyface and a heel. If it is possible to perform such a designation, your algorithm should produce it.

**22.2-8 ★**

The *diameter* of a tree  $T = (V, E)$  is defined as  $\max_{u,v \in V} \delta(u, v)$ , that is, the largest of all shortest-path distances in the tree. Give an efficient algorithm to compute the diameter of a tree, and analyze the running time of your algorithm.

**22.2-9**

Let  $G = (V, E)$  be a connected, undirected graph. Give an  $O(V + E)$ -time algorithm to compute a path in  $G$  that traverses each edge in  $E$  exactly once in each direction. Describe how you can find your way out of a maze if you are given a large supply of pennies.