

---

## *IV Advanced Design and Analysis Techniques*

---

## Introduction

This part covers three important techniques used in designing and analyzing efficient algorithms: dynamic programming (Chapter 15), greedy algorithms (Chapter 16), and amortized analysis (Chapter 17). Earlier parts have presented other widely applicable techniques, such as divide-and-conquer, randomization, and how to solve recurrences. The techniques in this part are somewhat more sophisticated, but they help us to attack many computational problems. The themes introduced in this part will recur later in this book.

Dynamic programming typically applies to optimization problems in which we make a set of choices in order to arrive at an optimal solution. As we make each choice, subproblems of the same form often arise. Dynamic programming is effective when a given subproblem may arise from more than one partial set of choices; the key technique is to store the solution to each such subproblem in case it should reappear. Chapter 15 shows how this simple idea can sometimes transform exponential-time algorithms into polynomial-time algorithms.

Like dynamic-programming algorithms, greedy algorithms typically apply to optimization problems in which we make a set of choices in order to arrive at an optimal solution. The idea of a greedy algorithm is to make each choice in a locally optimal manner. A simple example is coin-changing: to minimize the number of U.S. coins needed to make change for a given amount, we can repeatedly select the largest-denomination coin that is not larger than the amount that remains. A greedy approach provides an optimal solution for many such problems much more quickly than would a dynamic-programming approach. We cannot always easily tell whether a greedy approach will be effective, however. Chapter 16 introduces

matroid theory, which provides a mathematical basis that can help us to show that a greedy algorithm yields an optimal solution.

We use amortized analysis to analyze certain algorithms that perform a sequence of similar operations. Instead of bounding the cost of the sequence of operations by bounding the actual cost of each operation separately, an amortized analysis provides a bound on the actual cost of the entire sequence. One advantage of this approach is that although some operations might be expensive, many others might be cheap. In other words, many of the operations might run in well under the worst-case time. Amortized analysis is not just an analysis tool, however; it is also a way of thinking about the design of algorithms, since the design of an algorithm and the analysis of its running time are often closely intertwined. Chapter 17 introduces three ways to perform an amortized analysis of an algorithm.

---

## 15      Dynamic Programming

Dynamic programming, like the divide-and-conquer method, solves problems by combining the solutions to subproblems. (“Programming” in this context refers to a tabular method, not to writing computer code.) As we saw in Chapters 2 and 4, divide-and-conquer algorithms partition the problem into disjoint subproblems, solve the subproblems recursively, and then combine their solutions to solve the original problem. In contrast, dynamic programming applies when the subproblems overlap—that is, when subproblems share subsubproblems. In this context, a divide-and-conquer algorithm does more work than necessary, repeatedly solving the common subsubproblems. A dynamic-programming algorithm solves each subsubproblem just once and then saves its answer in a table, thereby avoiding the work of recomputing the answer every time it solves each subsubproblem.

We typically apply dynamic programming to *optimization problems*. Such problems can have many possible solutions. Each solution has a value, and we wish to find a solution with the optimal (minimum or maximum) value. We call such a solution *an* optimal solution to the problem, as opposed to *the* optimal solution, since there may be several solutions that achieve the optimal value.

When developing a dynamic-programming algorithm, we follow a sequence of four steps:

1. Characterize the structure of an optimal solution.
2. Recursively define the value of an optimal solution.
3. Compute the value of an optimal solution, typically in a bottom-up fashion.
4. Construct an optimal solution from computed information.

Steps 1–3 form the basis of a dynamic-programming solution to a problem. If we need only the value of an optimal solution, and not the solution itself, then we can omit step 4. When we do perform step 4, we sometimes maintain additional information during step 3 so that we can easily construct an optimal solution.

The sections that follow use the dynamic-programming method to solve some optimization problems. Section 15.1 examines the problem of cutting a rod into

rods of smaller length in way that maximizes their total value. Section 15.2 asks how we can multiply a chain of matrices while performing the fewest total scalar multiplications. Given these examples of dynamic programming, Section 15.3 discusses two key characteristics that a problem must have for dynamic programming to be a viable solution technique. Section 15.4 then shows how to find the longest common subsequence of two sequences via dynamic programming. Finally, Section 15.5 uses dynamic programming to construct binary search trees that are optimal, given a known distribution of keys to be looked up.

---

## 15.1 Rod cutting

Our first example uses dynamic programming to solve a simple problem in deciding where to cut steel rods. Serling Enterprises buys long steel rods and cuts them into shorter rods, which it then sells. Each cut is free. The management of Serling Enterprises wants to know the best way to cut up the rods.

We assume that we know, for  $i = 1, 2, \dots$ , the price  $p_i$  in dollars that Serling Enterprises charges for a rod of length  $i$  inches. Rod lengths are always an integral number of inches. Figure 15.1 gives a sample price table.

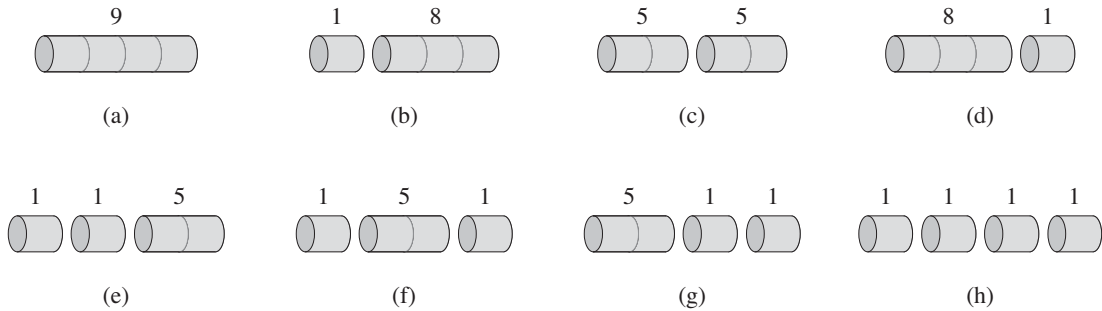
The **rod-cutting problem** is the following. Given a rod of length  $n$  inches and a table of prices  $p_i$  for  $i = 1, 2, \dots, n$ , determine the maximum revenue  $r_n$  obtainable by cutting up the rod and selling the pieces. Note that if the price  $p_n$  for a rod of length  $n$  is large enough, an optimal solution may require no cutting at all.

Consider the case when  $n = 4$ . Figure 15.2 shows all the ways to cut up a rod of 4 inches in length, including the way with no cuts at all. We see that cutting a 4-inch rod into two 2-inch pieces produces revenue  $p_2 + p_2 = 5 + 5 = 10$ , which is optimal.

We can cut up a rod of length  $n$  in  $2^{n-1}$  different ways, since we have an independent option of cutting, or not cutting, at distance  $i$  inches from the left end,

length $i$	1	2	3	4	5	6	7	8	9	10
price $p_i$	1	5	8	9	10	17	17	20	24	30

**Figure 15.1** A sample price table for rods. Each rod of length  $i$  inches earns the company  $p_i$  dollars of revenue.



**Figure 15.2** The 8 possible ways of cutting up a rod of length 4. Above each piece is the value of that piece, according to the sample price chart of Figure 15.1. The optimal strategy is part (c)—cutting the rod into two pieces of length 2—which has total value 10.

for  $i = 1, 2, \dots, n - 1$ .<sup>1</sup> We denote a decomposition into pieces using ordinary additive notation, so that  $7 = 2 + 2 + 3$  indicates that a rod of length 7 is cut into three pieces—two of length 2 and one of length 3. If an optimal solution cuts the rod into  $k$  pieces, for some  $1 \leq k \leq n$ , then an optimal decomposition

$$n = i_1 + i_2 + \dots + i_k$$

of the rod into pieces of lengths  $i_1, i_2, \dots, i_k$  provides maximum corresponding revenue

$$r_n = p_{i_1} + p_{i_2} + \dots + p_{i_k}.$$

For our sample problem, we can determine the optimal revenue figures  $r_i$ , for  $i = 1, 2, \dots, 10$ , by inspection, with the corresponding optimal decompositions

---

<sup>1</sup>If we required the pieces to be cut in order of nondecreasing size, there would be fewer ways to consider. For  $n = 4$ , we would consider only 5 such ways: parts (a), (b), (c), (e), and (h) in Figure 15.2. The number of ways is called the **partition function**; it is approximately equal to  $e^{\pi\sqrt{2n/3}}/4n\sqrt{3}$ . This quantity is less than  $2^{n-1}$ , but still much greater than any polynomial in  $n$ . We shall not pursue this line of inquiry further, however.

$$\begin{aligned}
r_1 &= 1 && \text{from solution } 1 = 1 \quad (\text{no cuts}) , \\
r_2 &= 5 && \text{from solution } 2 = 2 \quad (\text{no cuts}) , \\
r_3 &= 8 && \text{from solution } 3 = 3 \quad (\text{no cuts}) , \\
r_4 &= 10 && \text{from solution } 4 = 2 + 2 , \\
r_5 &= 13 && \text{from solution } 5 = 2 + 3 , \\
r_6 &= 17 && \text{from solution } 6 = 6 \quad (\text{no cuts}) , \\
r_7 &= 18 && \text{from solution } 7 = 1 + 6 \text{ or } 7 = 2 + 2 + 3 , \\
r_8 &= 22 && \text{from solution } 8 = 2 + 6 , \\
r_9 &= 25 && \text{from solution } 9 = 3 + 6 , \\
r_{10} &= 30 && \text{from solution } 10 = 10 \quad (\text{no cuts}) .
\end{aligned}$$

More generally, we can frame the values  $r_n$  for  $n \geq 1$  in terms of optimal revenues from shorter rods:

$$r_n = \max (p_n, r_1 + r_{n-1}, r_2 + r_{n-2}, \dots, r_{n-1} + r_1) . \quad (15.1)$$

The first argument,  $p_n$ , corresponds to making no cuts at all and selling the rod of length  $n$  as is. The other  $n - 1$  arguments to max correspond to the maximum revenue obtained by making an initial cut of the rod into two pieces of size  $i$  and  $n - i$ , for each  $i = 1, 2, \dots, n - 1$ , and then optimally cutting up those pieces further, obtaining revenues  $r_i$  and  $r_{n-i}$  from those two pieces. Since we don't know ahead of time which value of  $i$  optimizes revenue, we have to consider all possible values for  $i$  and pick the one that maximizes revenue. We also have the option of picking no  $i$  at all if we can obtain more revenue by selling the rod uncut.

Note that to solve the original problem of size  $n$ , we solve smaller problems of the same type, but of smaller sizes. Once we make the first cut, we may consider the two pieces as independent instances of the rod-cutting problem. The overall optimal solution incorporates optimal solutions to the two related subproblems, maximizing revenue from each of those two pieces. We say that the rod-cutting problem exhibits **optimal substructure**: optimal solutions to a problem incorporate optimal solutions to related subproblems, which we may solve independently.

In a related, but slightly simpler, way to arrange a recursive structure for the rod-cutting problem, we view a decomposition as consisting of a first piece of length  $i$  cut off the left-hand end, and then a right-hand remainder of length  $n - i$ . Only the remainder, and not the first piece, may be further divided. We may view every decomposition of a length- $n$  rod in this way: as a first piece followed by some decomposition of the remainder. When doing so, we can couch the solution with no cuts at all as saying that the first piece has size  $i = n$  and revenue  $p_n$  and that the remainder has size 0 with corresponding revenue  $r_0 = 0$ . We thus obtain the following simpler version of equation (15.1):

$$r_n = \max_{1 \leq i \leq n} (p_i + r_{n-i}) . \quad (15.2)$$

In this formulation, an optimal solution embodies the solution to only *one* related subproblem—the remainder—rather than two.

### Recursive top-down implementation

The following procedure implements the computation implicit in equation (15.2) in a straightforward, top-down, recursive manner.

```

CUT-ROD( $p, n$ )
1  if  $n == 0$ 
2      return 0
3   $q = -\infty$ 
4  for  $i = 1$  to  $n$ 
5       $q = \max(q, p[i] + \text{CUT-ROD}(p, n - i))$ 
6  return  $q$ 

```

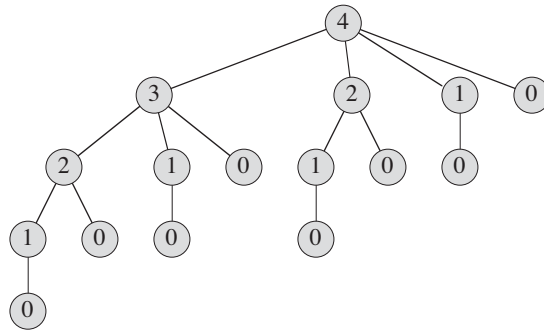
Procedure CUT-ROD takes as input an array  $p[1..n]$  of prices and an integer  $n$ , and it returns the maximum revenue possible for a rod of length  $n$ . If  $n = 0$ , no revenue is possible, and so CUT-ROD returns 0 in line 2. Line 3 initializes the maximum revenue  $q$  to  $-\infty$ , so that the **for** loop in lines 4–5 correctly computes  $q = \max_{1 \leq i \leq n} (p_i + \text{CUT-ROD}(p, n - i))$ ; line 6 then returns this value. A simple induction on  $n$  proves that this answer is equal to the desired answer  $r_n$ , using equation (15.2).

If you were to code up CUT-ROD in your favorite programming language and run it on your computer, you would find that once the input size becomes moderately large, your program would take a long time to run. For  $n = 40$ , you would find that your program takes at least several minutes, and most likely more than an hour. In fact, you would find that each time you increase  $n$  by 1, your program's running time would approximately double.

Why is CUT-ROD so inefficient? The problem is that CUT-ROD calls itself recursively over and over again with the same parameter values; it solves the same subproblems repeatedly. Figure 15.3 illustrates what happens for  $n = 4$ : CUT-ROD( $p, n$ ) calls CUT-ROD( $p, n - i$ ) for  $i = 1, 2, \dots, n$ . Equivalently, CUT-ROD( $p, n$ ) calls CUT-ROD( $p, j$ ) for each  $j = 0, 1, \dots, n - 1$ . When this process unfolds recursively, the amount of work done, as a function of  $n$ , grows explosively.

To analyze the running time of CUT-ROD, let  $T(n)$  denote the total number of calls made to CUT-ROD when called with its second parameter equal to  $n$ . This expression equals the number of nodes in a subtree whose root is labeled  $n$  in the recursion tree. The count includes the initial call at its root. Thus,  $T(0) = 1$  and





**Figure 15.3** The recursion tree showing recursive calls resulting from a call  $\text{CUT-ROD}(p, n)$  for  $n = 4$ . Each node label gives the size  $n$  of the corresponding subproblem, so that an edge from a parent with label  $s$  to a child with label  $t$  corresponds to cutting off an initial piece of size  $s - t$  and leaving a remaining subproblem of size  $t$ . A path from the root to a leaf corresponds to one of the  $2^{n-1}$  ways of cutting up a rod of length  $n$ . In general, this recursion tree has  $2^n$  nodes and  $2^{n-1}$  leaves.

$$T(n) = 1 + \sum_{j=0}^{n-1} T(j) . \quad (15.3)$$

The initial 1 is for the call at the root, and the term  $T(j)$  counts the number of calls (including recursive calls) due to the call  $\text{CUT-ROD}(p, n - i)$ , where  $j = n - i$ . As Exercise 15.1-1 asks you to show,

$$T(n) = 2^n , \quad (15.4)$$

and so the running time of  $\text{CUT-ROD}$  is exponential in  $n$ .

In retrospect, this exponential running time is not so surprising.  $\text{CUT-ROD}$  explicitly considers all the  $2^{n-1}$  possible ways of cutting up a rod of length  $n$ . The tree of recursive calls has  $2^{n-1}$  leaves, one for each possible way of cutting up the rod. The labels on the simple path from the root to a leaf give the sizes of each remaining right-hand piece before making each cut. That is, the labels give the corresponding cut points, measured from the right-hand end of the rod.

### Using dynamic programming for optimal rod cutting

We now show how to convert  $\text{CUT-ROD}$  into an efficient algorithm, using dynamic programming.

The dynamic-programming method works as follows. Having observed that a naive recursive solution is inefficient because it solves the same subproblems repeatedly, we arrange for each subproblem to be solved only *once*, saving its solution. If we need to refer to this subproblem's solution again later, we can just look it

up, rather than recompute it. Dynamic programming thus uses additional memory to save computation time; it serves an example of a *time-memory trade-off*. The savings may be dramatic: an exponential-time solution may be transformed into a polynomial-time solution. A dynamic-programming approach runs in polynomial time when the number of *distinct* subproblems involved is polynomial in the input size and we can solve each such subproblem in polynomial time.

There are usually two equivalent ways to implement a dynamic-programming approach. We shall illustrate both of them with our rod-cutting example.

The first approach is *top-down with memoization*.<sup>2</sup> In this approach, we write the procedure recursively in a natural manner, but modified to save the result of each subproblem (usually in an array or hash table). The procedure now first checks to see whether it has previously solved this subproblem. If so, it returns the saved value, saving further computation at this level; if not, the procedure computes the value in the usual manner. We say that the recursive procedure has been *memoized*; it “remembers” what results it has computed previously.

The second approach is the *bottom-up method*. This approach typically depends on some natural notion of the “size” of a subproblem, such that solving any particular subproblem depends only on solving “smaller” subproblems. We sort the subproblems by size and solve them in size order, smallest first. When solving a particular subproblem, we have already solved all of the smaller subproblems its solution depends upon, and we have saved their solutions. We solve each subproblem only once, and when we first see it, we have already solved all of its prerequisite subproblems.

These two approaches yield algorithms with the same asymptotic running time, except in unusual circumstances where the top-down approach does not actually recurse to examine all possible subproblems. The bottom-up approach often has much better constant factors, since it has less overhead for procedure calls.

Here is the the pseudocode for the top-down CUT-ROD procedure, with memoization added:

```
MEMOIZED-CUT-ROD( $p, n$ )
1  let  $r[0 \dots n]$  be a new array
2  for  $i = 0$  to  $n$ 
3       $r[i] = -\infty$ 
4  return MEMOIZED-CUT-ROD-AUX( $p, n, r$ )
```

---

<sup>2</sup>This is not a misspelling. The word really is *memoization*, not *memorization*. *Memoization* comes from *memo*, since the technique consists of recording a value so that we can look it up later.

MEMOIZED-CUT-ROD-AUX( $p, n, r$ )

```

1  if  $r[n] \geq 0$ 
2      return  $r[n]$ 
3  if  $n == 0$ 
4       $q = 0$ 
5  else  $q = -\infty$ 
6      for  $i = 1$  to  $n$ 
7           $q = \max(q, p[i] + \text{MEMOIZED-CUT-ROD-AUX}(p, n - i, r))$ 
8   $r[n] = q$ 
9  return  $q$ 

```

Here, the main procedure MEMOIZED-CUT-ROD initializes a new auxiliary array  $r[0..n]$  with the value  $-\infty$ , a convenient choice with which to denote “unknown.” (Known revenue values are always nonnegative.) It then calls its helper routine, MEMOIZED-CUT-ROD-AUX.

The procedure MEMOIZED-CUT-ROD-AUX is just the memoized version of our previous procedure, CUT-ROD. It first checks in line 1 to see whether the desired value is already known and, if it is, then line 2 returns it. Otherwise, lines 3–7 compute the desired value  $q$  in the usual manner, line 8 saves it in  $r[n]$ , and line 9 returns it.

The bottom-up version is even simpler:

BOTTOM-UP-CUT-ROD( $p, n$ )

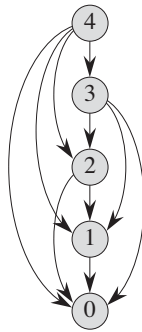
```

1  let  $r[0..n]$  be a new array
2   $r[0] = 0$ 
3  for  $j = 1$  to  $n$ 
4       $q = -\infty$ 
5      for  $i = 1$  to  $j$ 
6           $q = \max(q, p[i] + r[j - i])$ 
7       $r[j] = q$ 
8  return  $r[n]$ 

```

For the bottom-up dynamic-programming approach, BOTTOM-UP-CUT-ROD uses the natural ordering of the subproblems: a problem of size  $i$  is “smaller” than a subproblem of size  $j$  if  $i < j$ . Thus, the procedure solves subproblems of sizes  $j = 0, 1, \dots, n$ , in that order.

Line 1 of procedure BOTTOM-UP-CUT-ROD creates a new array  $r[0..n]$  in which to save the results of the subproblems, and line 2 initializes  $r[0]$  to 0, since a rod of length 0 earns no revenue. Lines 3–6 solve each subproblem of size  $j$ , for  $j = 1, 2, \dots, n$ , in order of increasing size. The approach used to solve a problem of a particular size  $j$  is the same as that used by CUT-ROD, except that line 6 now



**Figure 15.4** The subproblem graph for the rod-cutting problem with  $n = 4$ . The vertex labels give the sizes of the corresponding subproblems. A directed edge  $(x, y)$  indicates that we need a solution to subproblem  $y$  when solving subproblem  $x$ . This graph is a reduced version of the tree of Figure 15.3, in which all nodes with the same label are collapsed into a single vertex and all edges go from parent to child.

directly references array entry  $r[j - i]$  instead of making a recursive call to solve the subproblem of size  $j - i$ . Line 7 saves in  $r[j]$  the solution to the subproblem of size  $j$ . Finally, line 8 returns  $r[n]$ , which equals the optimal value  $r_n$ .

The bottom-up and top-down versions have the same asymptotic running time. The running time of procedure BOTTOM-UP-CUT-ROD is  $\Theta(n^2)$ , due to its doubly-nested loop structure. The number of iterations of its inner **for** loop, in lines 5–6, forms an arithmetic series. The running time of its top-down counterpart, MEMOIZED-CUT-ROD, is also  $\Theta(n^2)$ , although this running time may be a little harder to see. Because a recursive call to solve a previously solved subproblem returns immediately, MEMOIZED-CUT-ROD solves each subproblem just once. It solves subproblems for sizes  $0, 1, \dots, n$ . To solve a subproblem of size  $n$ , the **for** loop of lines 6–7 iterates  $n$  times. Thus, the total number of iterations of this **for** loop, over all recursive calls of MEMOIZED-CUT-ROD, forms an arithmetic series, giving a total of  $\Theta(n^2)$  iterations, just like the inner **for** loop of BOTTOM-UP-CUT-ROD. (We actually are using a form of aggregate analysis here. We shall see aggregate analysis in detail in Section 17.1.)

### Subproblem graphs

When we think about a dynamic-programming problem, we should understand the set of subproblems involved and how subproblems depend on one another.

The **subproblem graph** for the problem embodies exactly this information. Figure 15.4 shows the subproblem graph for the rod-cutting problem with  $n = 4$ . It is a directed graph, containing one vertex for each distinct subproblem. The sub-

problem graph has a directed edge from the vertex for subproblem  $x$  to the vertex for subproblem  $y$  if determining an optimal solution for subproblem  $x$  involves directly considering an optimal solution for subproblem  $y$ . For example, the subproblem graph contains an edge from  $x$  to  $y$  if a top-down recursive procedure for solving  $x$  directly calls itself to solve  $y$ . We can think of the subproblem graph as a “reduced” or “collapsed” version of the recursion tree for the top-down recursive method, in which we coalesce all nodes for the same subproblem into a single vertex and direct all edges from parent to child.

The bottom-up method for dynamic programming considers the vertices of the subproblem graph in such an order that we solve the subproblems  $y$  adjacent to a given subproblem  $x$  before we solve subproblem  $x$ . (Recall from Section B.4 that the adjacency relation is not necessarily symmetric.) Using the terminology from Chapter 22, in a bottom-up dynamic-programming algorithm, we consider the vertices of the subproblem graph in an order that is a “reverse topological sort,” or a “topological sort of the transpose” (see Section 22.4) of the subproblem graph. In other words, no subproblem is considered until all of the subproblems it depends upon have been solved. Similarly, using notions from the same chapter, we can view the top-down method (with memoization) for dynamic programming as a “depth-first search” of the subproblem graph (see Section 22.3).

The size of the subproblem graph  $G = (V, E)$  can help us determine the running time of the dynamic programming algorithm. Since we solve each subproblem just once, the running time is the sum of the times needed to solve each subproblem. Typically, the time to compute the solution to a subproblem is proportional to the degree (number of outgoing edges) of the corresponding vertex in the subproblem graph, and the number of subproblems is equal to the number of vertices in the subproblem graph. In this common case, the running time of dynamic programming is linear in the number of vertices and edges.

### Reconstructing a solution

Our dynamic-programming solutions to the rod-cutting problem return the value of an optimal solution, but they do not return an actual solution: a list of piece sizes. We can extend the dynamic-programming approach to record not only the optimal *value* computed for each subproblem, but also a *choice* that led to the optimal value. With this information, we can readily print an optimal solution.

Here is an extended version of BOTTOM-UP-CUT-ROD that computes, for each rod size  $j$ , not only the maximum revenue  $r_j$ , but also  $s_j$ , the optimal size of the first piece to cut off:

EXTENDED-BOTTOM-UP-CUT-ROD( $p, n$ )

```

1  let  $r[0..n]$  and  $s[0..n]$  be new arrays
2   $r[0] = 0$ 
3  for  $j = 1$  to  $n$ 
4       $q = -\infty$ 
5      for  $i = 1$  to  $j$ 
6          if  $q < p[i] + r[j - i]$ 
7               $q = p[i] + r[j - i]$ 
8               $s[j] = i$ 
9       $r[j] = q$ 
10 return  $r$  and  $s$ 

```

This procedure is similar to BOTTOM-UP-CUT-ROD, except that it creates the array  $s$  in line 1, and it updates  $s[j]$  in line 8 to hold the optimal size  $i$  of the first piece to cut off when solving a subproblem of size  $j$ .

The following procedure takes a price table  $p$  and a rod size  $n$ , and it calls EXTENDED-BOTTOM-UP-CUT-ROD to compute the array  $s[1..n]$  of optimal first-piece sizes and then prints out the complete list of piece sizes in an optimal decomposition of a rod of length  $n$ :

PRINT-CUT-ROD-SOLUTION( $p, n$ )

```

1   $(r, s) = \text{EXTENDED-BOTTOM-UP-CUT-ROD}(p, n)$ 
2  while  $n > 0$ 
3      print  $s[n]$ 
4       $n = n - s[n]$ 

```

In our rod-cutting example, the call EXTENDED-BOTTOM-UP-CUT-ROD( $p, 10$ ) would return the following arrays:

$i$	0	1	2	3	4	5	6	7	8	9	10
$r[i]$	0	1	5	8	10	13	17	18	22	25	30
$s[i]$	0	1	2	3	2	2	6	1	2	3	10

A call to PRINT-CUT-ROD-SOLUTION( $p, 10$ ) would print just 10, but a call with  $n = 7$  would print the cuts 1 and 6, corresponding to the first optimal decomposition for  $r_7$  given earlier.

## Exercises

### 15.1-1

Show that equation (15.4) follows from equation (15.3) and the initial condition  $T(0) = 1$ .

**15.1-2**

Show, by means of a counterexample, that the following “greedy” strategy does not always determine an optimal way to cut rods. Define the *density* of a rod of length  $i$  to be  $p_i/i$ , that is, its value per inch. The greedy strategy for a rod of length  $n$  cuts off a first piece of length  $i$ , where  $1 \leq i \leq n$ , having maximum density. It then continues by applying the greedy strategy to the remaining piece of length  $n - i$ .

**15.1-3**

Consider a modification of the rod-cutting problem in which, in addition to a price  $p_i$  for each rod, each cut incurs a fixed cost of  $c$ . The revenue associated with a solution is now the sum of the prices of the pieces minus the costs of making the cuts. Give a dynamic-programming algorithm to solve this modified problem.

**15.1-4**

Modify MEMOIZED-CUT-ROD to return not only the value but the actual solution, too.

**15.1-5**

The Fibonacci numbers are defined by recurrence (3.22). Give an  $O(n)$ -time dynamic-programming algorithm to compute the  $n$ th Fibonacci number. Draw the subproblem graph. How many vertices and edges are in the graph?

---

## 15.2 Matrix-chain multiplication

Our next example of dynamic programming is an algorithm that solves the problem of matrix-chain multiplication. We are given a sequence (chain)  $\langle A_1, A_2, \dots, A_n \rangle$  of  $n$  matrices to be multiplied, and we wish to compute the product

$$A_1 A_2 \cdots A_n. \quad (15.5)$$

We can evaluate the expression (15.5) using the standard algorithm for multiplying pairs of matrices as a subroutine once we have parenthesized it to resolve all ambiguities in how the matrices are multiplied together. Matrix multiplication is associative, and so all parenthesizations yield the same product. A product of matrices is *fully parenthesized* if it is either a single matrix or the product of two fully parenthesized matrix products, surrounded by parentheses. For example, if the chain of matrices is  $\langle A_1, A_2, A_3, A_4 \rangle$ , then we can fully parenthesize the product  $A_1 A_2 A_3 A_4$  in five distinct ways:

$(A_1(A_2(A_3A_4)))$  ,  
 $(A_1((A_2A_3)A_4))$  ,  
 $((A_1A_2)(A_3A_4))$  ,  
 $((A_1(A_2A_3))A_4)$  ,  
 $((A_1A_2)A_3)A_4)$  .

How we parenthesize a chain of matrices can have a dramatic impact on the cost of evaluating the product. Consider first the cost of multiplying two matrices. The standard algorithm is given by the following pseudocode, which generalizes the SQUARE-MATRIX-MULTIPLY procedure from Section 4.2. The attributes *rows* and *columns* are the numbers of rows and columns in a matrix.

MATRIX-MULTIPLY( $A, B$ )

```

1  if  $A.columns \neq B.rows$ 
2      error “incompatible dimensions”
3  else let  $C$  be a new  $A.rows \times B.columns$  matrix
4      for  $i = 1$  to  $A.rows$ 
5          for  $j = 1$  to  $B.columns$ 
6               $c_{ij} = 0$ 
7              for  $k = 1$  to  $A.columns$ 
8                   $c_{ij} = c_{ij} + a_{ik} \cdot b_{kj}$ 
9  return  $C$ 
  
```

We can multiply two matrices  $A$  and  $B$  only if they are **compatible**: the number of columns of  $A$  must equal the number of rows of  $B$ . If  $A$  is a  $p \times q$  matrix and  $B$  is a  $q \times r$  matrix, the resulting matrix  $C$  is a  $p \times r$  matrix. The time to compute  $C$  is dominated by the number of scalar multiplications in line 8, which is  $pqr$ . In what follows, we shall express costs in terms of the number of scalar multiplications.

To illustrate the different costs incurred by different parenthesizations of a matrix product, consider the problem of a chain  $\langle A_1, A_2, A_3 \rangle$  of three matrices. Suppose that the dimensions of the matrices are  $10 \times 100$ ,  $100 \times 5$ , and  $5 \times 50$ , respectively. If we multiply according to the parenthesization  $((A_1A_2)A_3)$ , we perform  $10 \cdot 100 \cdot 5 = 5000$  scalar multiplications to compute the  $10 \times 5$  matrix product  $A_1A_2$ , plus another  $10 \cdot 5 \cdot 50 = 2500$  scalar multiplications to multiply this matrix by  $A_3$ , for a total of 7500 scalar multiplications. If instead we multiply according to the parenthesization  $(A_1(A_2A_3))$ , we perform  $100 \cdot 5 \cdot 50 = 25,000$  scalar multiplications to compute the  $100 \times 50$  matrix product  $A_2A_3$ , plus another  $10 \cdot 100 \cdot 50 = 50,000$  scalar multiplications to multiply  $A_1$  by this matrix, for a total of 75,000 scalar multiplications. Thus, computing the product according to the first parenthesization is 10 times faster.

We state the **matrix-chain multiplication problem** as follows: given a chain  $\langle A_1, A_2, \dots, A_n \rangle$  of  $n$  matrices, where for  $i = 1, 2, \dots, n$ , matrix  $A_i$  has dimension



$p_{i-1} \times p_i$ , fully parenthesize the product  $A_1 A_2 \cdots A_n$  in a way that minimizes the number of scalar multiplications.

Note that in the matrix-chain multiplication problem, we are not actually multiplying matrices. Our goal is only to determine an order for multiplying matrices that has the lowest cost. Typically, the time invested in determining this optimal order is more than paid for by the time saved later on when actually performing the matrix multiplications (such as performing only 7500 scalar multiplications instead of 75,000).

### Counting the number of parenthesizations

Before solving the matrix-chain multiplication problem by dynamic programming, let us convince ourselves that exhaustively checking all possible parenthesizations does not yield an efficient algorithm. Denote the number of alternative parenthesizations of a sequence of  $n$  matrices by  $P(n)$ . When  $n = 1$ , we have just one matrix and therefore only one way to fully parenthesize the matrix product. When  $n \geq 2$ , a fully parenthesized matrix product is the product of two fully parenthesized matrix subproducts, and the split between the two subproducts may occur between the  $k$ th and  $(k + 1)$ st matrices for any  $k = 1, 2, \dots, n - 1$ . Thus, we obtain the recurrence

$$P(n) = \begin{cases} 1 & \text{if } n = 1, \\ \sum_{k=1}^{n-1} P(k)P(n-k) & \text{if } n \geq 2. \end{cases} \quad (15.6)$$

Problem 12-4 asked you to show that the solution to a similar recurrence is the sequence of **Catalan numbers**, which grows as  $\Omega(4^n/n^{3/2})$ . A simpler exercise (see Exercise 15.2-3) is to show that the solution to the recurrence (15.6) is  $\Omega(2^n)$ . The number of solutions is thus exponential in  $n$ , and the brute-force method of exhaustive search makes for a poor strategy when determining how to optimally parenthesize a matrix chain.

### Applying dynamic programming

We shall use the dynamic-programming method to determine how to optimally parenthesize a matrix chain. In so doing, we shall follow the four-step sequence that we stated at the beginning of this chapter:

1. Characterize the structure of an optimal solution.
2. Recursively define the value of an optimal solution.
3. Compute the value of an optimal solution.

4. Construct an optimal solution from computed information.

We shall go through these steps in order, demonstrating clearly how we apply each step to the problem.

**Step 1: The structure of an optimal parenthesization**

For our first step in the dynamic-programming paradigm, we find the optimal substructure and then use it to construct an optimal solution to the problem from optimal solutions to subproblems. In the matrix-chain multiplication problem, we can perform this step as follows. For convenience, let us adopt the notation  $A_{i..j}$ , where  $i \leq j$ , for the matrix that results from evaluating the product  $A_i A_{i+1} \cdots A_j$ . Observe that if the problem is nontrivial, i.e.,  $i < j$ , then to parenthesize the product  $A_i A_{i+1} \cdots A_j$ , we must split the product between  $A_k$  and  $A_{k+1}$  for some integer  $k$  in the range  $i \leq k < j$ . That is, for some value of  $k$ , we first compute the matrices  $A_{i..k}$  and  $A_{k+1..j}$  and then multiply them together to produce the final product  $A_{i..j}$ . The cost of parenthesizing this way is the cost of computing the matrix  $A_{i..k}$ , plus the cost of computing  $A_{k+1..j}$ , plus the cost of multiplying them together.

The optimal substructure of this problem is as follows. Suppose that to optimally parenthesize  $A_i A_{i+1} \cdots A_j$ , we split the product between  $A_k$  and  $A_{k+1}$ . Then the way we parenthesize the “prefix” subchain  $A_i A_{i+1} \cdots A_k$  within this optimal parenthesization of  $A_i A_{i+1} \cdots A_j$  must be an optimal parenthesization of  $A_i A_{i+1} \cdots A_k$ . Why? If there were a less costly way to parenthesize  $A_i A_{i+1} \cdots A_k$ , then we could substitute that parenthesization in the optimal parenthesization of  $A_i A_{i+1} \cdots A_j$  to produce another way to parenthesize  $A_i A_{i+1} \cdots A_j$  whose cost was lower than the optimum: a contradiction. A similar observation holds for how we parenthesize the subchain  $A_{k+1} A_{k+2} \cdots A_j$  in the optimal parenthesization of  $A_i A_{i+1} \cdots A_j$ : it must be an optimal parenthesization of  $A_{k+1} A_{k+2} \cdots A_j$ .

Now we use our optimal substructure to show that we can construct an optimal solution to the problem from optimal solutions to subproblems. We have seen that any solution to a nontrivial instance of the matrix-chain multiplication problem requires us to split the product, and that any optimal solution contains within it optimal solutions to subproblem instances. Thus, we can build an optimal solution to an instance of the matrix-chain multiplication problem by splitting the problem into two subproblems (optimally parenthesizing  $A_i A_{i+1} \cdots A_k$  and  $A_{k+1} A_{k+2} \cdots A_j$ ), finding optimal solutions to subproblem instances, and then combining these optimal subproblem solutions. We must ensure that when we search for the correct place to split the product, we have considered all possible places, so that we are sure of having examined the optimal one.

### Step 2: A recursive solution

Next, we define the cost of an optimal solution recursively in terms of the optimal solutions to subproblems. For the matrix-chain multiplication problem, we pick as our subproblems the problems of determining the minimum cost of parenthesizing  $A_i A_{i+1} \cdots A_j$  for  $1 \leq i \leq j \leq n$ . Let  $m[i, j]$  be the minimum number of scalar multiplications needed to compute the matrix  $A_{i..j}$ ; for the full problem, the lowest-cost way to compute  $A_{1..n}$  would thus be  $m[1, n]$ .

We can define  $m[i, j]$  recursively as follows. If  $i = j$ , the problem is trivial; the chain consists of just one matrix  $A_{i..i} = A_i$ , so that no scalar multiplications are necessary to compute the product. Thus,  $m[i, i] = 0$  for  $i = 1, 2, \dots, n$ . To compute  $m[i, j]$  when  $i < j$ , we take advantage of the structure of an optimal solution from step 1. Let us assume that to optimally parenthesize, we split the product  $A_i A_{i+1} \cdots A_j$  between  $A_k$  and  $A_{k+1}$ , where  $i \leq k < j$ . Then,  $m[i, j]$  equals the minimum cost for computing the subproducts  $A_{i..k}$  and  $A_{k+1..j}$ , plus the cost of multiplying these two matrices together. Recalling that each matrix  $A_i$  is  $p_{i-1} \times p_i$ , we see that computing the matrix product  $A_{i..k} A_{k+1..j}$  takes  $p_{i-1} p_k p_j$  scalar multiplications. Thus, we obtain

$$m[i, j] = m[i, k] + m[k + 1, j] + p_{i-1} p_k p_j .$$

This recursive equation assumes that we know the value of  $k$ , which we do not. There are only  $j - i$  possible values for  $k$ , however, namely  $k = i, i + 1, \dots, j - 1$ . Since the optimal parenthesization must use one of these values for  $k$ , we need only check them all to find the best. Thus, our recursive definition for the minimum cost of parenthesizing the product  $A_i A_{i+1} \cdots A_j$  becomes

$$m[i, j] = \begin{cases} 0 & \text{if } i = j , \\ \min_{i \leq k < j} \{m[i, k] + m[k + 1, j] + p_{i-1} p_k p_j\} & \text{if } i < j . \end{cases} \quad (15.7)$$

The  $m[i, j]$  values give the costs of optimal solutions to subproblems, but they do not provide all the information we need to construct an optimal solution. To help us do so, we define  $s[i, j]$  to be a value of  $k$  at which we split the product  $A_i A_{i+1} \cdots A_j$  in an optimal parenthesization. That is,  $s[i, j]$  equals a value  $k$  such that  $m[i, j] = m[i, k] + m[k + 1, j] + p_{i-1} p_k p_j$ .

### Step 3: Computing the optimal costs

At this point, we could easily write a recursive algorithm based on recurrence (15.7) to compute the minimum cost  $m[1, n]$  for multiplying  $A_1 A_2 \cdots A_n$ . As we saw for the rod-cutting problem, and as we shall see in Section 15.3, this recursive algorithm takes exponential time, which is no better than the brute-force method of checking each way of parenthesizing the product.

Observe that we have relatively few distinct subproblems: one subproblem for each choice of  $i$  and  $j$  satisfying  $1 \leq i \leq j \leq n$ , or  $\binom{n}{2} + n = \Theta(n^2)$  in all. A recursive algorithm may encounter each subproblem many times in different branches of its recursion tree. This property of overlapping subproblems is the second hallmark of when dynamic programming applies (the first hallmark being optimal substructure).

Instead of computing the solution to recurrence (15.7) recursively, we compute the optimal cost by using a tabular, bottom-up approach. (We present the corresponding top-down approach using memoization in Section 15.3.)

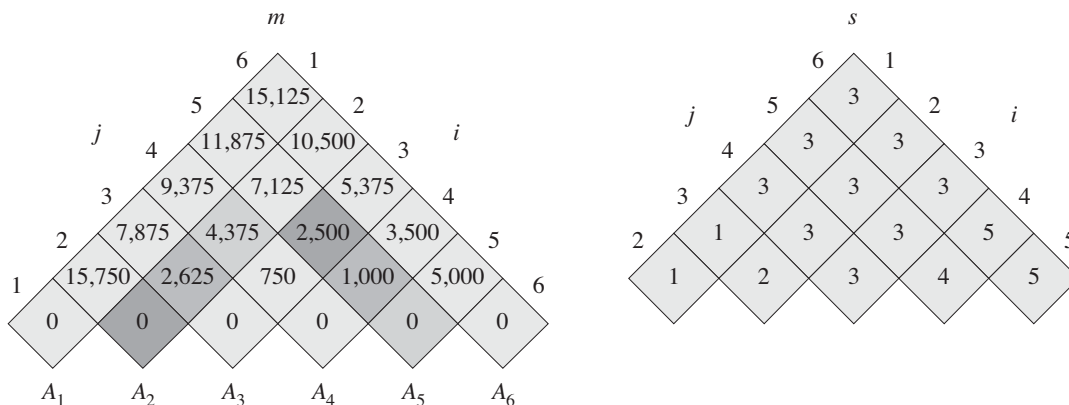
We shall implement the tabular, bottom-up method in the procedure MATRIX-CHAIN-ORDER, which appears below. This procedure assumes that matrix  $A_i$  has dimensions  $p_{i-1} \times p_i$  for  $i = 1, 2, \dots, n$ . Its input is a sequence  $p = \langle p_0, p_1, \dots, p_n \rangle$ , where  $p.length = n + 1$ . The procedure uses an auxiliary table  $m[1..n, 1..n]$  for storing the  $m[i, j]$  costs and another auxiliary table  $s[1..n - 1, 2..n]$  that records which index of  $k$  achieved the optimal cost in computing  $m[i, j]$ . We shall use the table  $s$  to construct an optimal solution.

In order to implement the bottom-up approach, we must determine which entries of the table we refer to when computing  $m[i, j]$ . Equation (15.7) shows that the cost  $m[i, j]$  of computing a matrix-chain product of  $j - i + 1$  matrices depends only on the costs of computing matrix-chain products of fewer than  $j - i + 1$  matrices. That is, for  $k = i, i + 1, \dots, j - 1$ , the matrix  $A_{i..k}$  is a product of  $k - i + 1 < j - i + 1$  matrices and the matrix  $A_{k+1..j}$  is a product of  $j - k < j - i + 1$  matrices. Thus, the algorithm should fill in the table  $m$  in a manner that corresponds to solving the parenthesization problem on matrix chains of increasing length. For the subproblem of optimally parenthesizing the chain  $A_i A_{i+1} \cdots A_j$ , we consider the subproblem size to be the length  $j - i + 1$  of the chain.

MATRIX-CHAIN-ORDER( $p$ )

```

1   $n = p.length - 1$ 
2  let  $m[1..n, 1..n]$  and  $s[1..n - 1, 2..n]$  be new tables
3  for  $i = 1$  to  $n$ 
4       $m[i, i] = 0$ 
5  for  $l = 2$  to  $n$            //  $l$  is the chain length
6      for  $i = 1$  to  $n - l + 1$ 
7           $j = i + l - 1$ 
8           $m[i, j] = \infty$ 
9          for  $k = i$  to  $j - 1$ 
10              $q = m[i, k] + m[k + 1, j] + p_{i-1}p_kp_j$ 
11             if  $q < m[i, j]$ 
12                  $m[i, j] = q$ 
13                  $s[i, j] = k$ 
14  return  $m$  and  $s$ 
```



**Figure 15.5** The  $m$  and  $s$  tables computed by MATRIX-CHAIN-ORDER for  $n = 6$  and the following matrix dimensions:

matrix	$A_1$	$A_2$	$A_3$	$A_4$	$A_5$	$A_6$
dimension	$30 \times 35$	$35 \times 15$	$15 \times 5$	$5 \times 10$	$10 \times 20$	$20 \times 25$

The tables are rotated so that the main diagonal runs horizontally. The  $m$  table uses only the main diagonal and upper triangle, and the  $s$  table uses only the upper triangle. The minimum number of scalar multiplications to multiply the 6 matrices is  $m[1, 6] = 15,125$ . Of the darker entries, the pairs that have the same shading are taken together in line 10 when computing

$$\begin{aligned}
 m[2, 5] &= \min \begin{cases} m[2, 2] + m[3, 5] + p_1 p_2 p_5 = 0 + 2500 + 35 \cdot 15 \cdot 20 = 13,000, \\ m[2, 3] + m[4, 5] + p_1 p_3 p_5 = 2625 + 1000 + 35 \cdot 5 \cdot 20 = 7125, \\ m[2, 4] + m[5, 5] + p_1 p_4 p_5 = 4375 + 0 + 35 \cdot 10 \cdot 20 = 11,375 \end{cases} \\
 &= 7125.
 \end{aligned}$$

The algorithm first computes  $m[i, i] = 0$  for  $i = 1, 2, \dots, n$  (the minimum costs for chains of length 1) in lines 3–4. It then uses recurrence (15.7) to compute  $m[i, i + 1]$  for  $i = 1, 2, \dots, n - 1$  (the minimum costs for chains of length  $l = 2$ ) during the first execution of the **for** loop in lines 5–13. The second time through the loop, it computes  $m[i, i + 2]$  for  $i = 1, 2, \dots, n - 2$  (the minimum costs for chains of length  $l = 3$ ), and so forth. At each step, the  $m[i, j]$  cost computed in lines 10–13 depends only on table entries  $m[i, k]$  and  $m[k + 1, j]$  already computed.

Figure 15.5 illustrates this procedure on a chain of  $n = 6$  matrices. Since we have defined  $m[i, j]$  only for  $i \leq j$ , only the portion of the table  $m$  strictly above the main diagonal is used. The figure shows the table rotated to make the main diagonal run horizontally. The matrix chain is listed along the bottom. Using this layout, we can find the minimum cost  $m[i, j]$  for multiplying a subchain  $A_i A_{i+1} \cdots A_j$  of matrices at the intersection of lines running northeast from  $A_i$  and

northwest from  $A_j$ . Each horizontal row in the table contains the entries for matrix chains of the same length. MATRIX-CHAIN-ORDER computes the rows from bottom to top and from left to right within each row. It computes each entry  $m[i, j]$  using the products  $p_{i-1}p_kp_j$  for  $k = i, i + 1, \dots, j - 1$  and all entries southwest and southeast from  $m[i, j]$ .

A simple inspection of the nested loop structure of MATRIX-CHAIN-ORDER yields a running time of  $O(n^3)$  for the algorithm. The loops are nested three deep, and each loop index ( $l$ ,  $i$ , and  $k$ ) takes on at most  $n - 1$  values. Exercise 15.2-5 asks you to show that the running time of this algorithm is in fact also  $\Omega(n^3)$ . The algorithm requires  $\Theta(n^2)$  space to store the  $m$  and  $s$  tables. Thus, MATRIX-CHAIN-ORDER is much more efficient than the exponential-time method of enumerating all possible parenthesizations and checking each one.

#### Step 4: Constructing an optimal solution

Although MATRIX-CHAIN-ORDER determines the optimal number of scalar multiplications needed to compute a matrix-chain product, it does not directly show how to multiply the matrices. The table  $s[1..n - 1, 2..n]$  gives us the information we need to do so. Each entry  $s[i, j]$  records a value of  $k$  such that an optimal parenthesization of  $A_iA_{i+1}\cdots A_j$  splits the product between  $A_k$  and  $A_{k+1}$ . Thus, we know that the final matrix multiplication in computing  $A_{1..n}$  optimally is  $A_{1..s[1,n]}A_{s[1,n]+1..n}$ . We can determine the earlier matrix multiplications recursively, since  $s[1, s[1, n]]$  determines the last matrix multiplication when computing  $A_{1..s[1,n]}$  and  $s[s[1, n] + 1, n]$  determines the last matrix multiplication when computing  $A_{s[1,n]+1..n}$ . The following recursive procedure prints an optimal parenthesization of  $\langle A_i, A_{i+1}, \dots, A_j \rangle$ , given the  $s$  table computed by MATRIX-CHAIN-ORDER and the indices  $i$  and  $j$ . The initial call PRINT-OPTIMAL-PARENS( $s, 1, n$ ) prints an optimal parenthesization of  $\langle A_1, A_2, \dots, A_n \rangle$ .

```

PRINT-OPTIMAL-PARENS( $s, i, j$ )
1  if  $i == j$ 
2      print " $A$ " $i$ 
3  else print "("
4      PRINT-OPTIMAL-PARENS( $s, i, s[i, j]$ )
5      PRINT-OPTIMAL-PARENS( $s, s[i, j] + 1, j$ )
6      print ")"

```

In the example of Figure 15.5, the call PRINT-OPTIMAL-PARENS( $s, 1, 6$ ) prints the parenthesization  $((A_1(A_2A_3))((A_4A_5)A_6))$ .

**Exercises****15.2-1**

Find an optimal parenthesization of a matrix-chain product whose sequence of dimensions is  $\langle 5, 10, 3, 12, 5, 50, 6 \rangle$ .

**15.2-2**

Give a recursive algorithm `MATRIX-CHAIN-MULTIPLY( $A, s, i, j$ )` that actually performs the optimal matrix-chain multiplication, given the sequence of matrices  $\langle A_1, A_2, \dots, A_n \rangle$ , the  $s$  table computed by `MATRIX-CHAIN-ORDER`, and the indices  $i$  and  $j$ . (The initial call would be `MATRIX-CHAIN-MULTIPLY( $A, s, 1, n$ )`.)

**15.2-3**

Use the substitution method to show that the solution to the recurrence (15.6) is  $\Omega(2^n)$ .

**15.2-4**

Describe the subproblem graph for matrix-chain multiplication with an input chain of length  $n$ . How many vertices does it have? How many edges does it have, and which edges are they?

**15.2-5**

Let  $R(i, j)$  be the number of times that table entry  $m[i, j]$  is referenced while computing other table entries in a call of `MATRIX-CHAIN-ORDER`. Show that the total number of references for the entire table is

$$\sum_{i=1}^n \sum_{j=i}^n R(i, j) = \frac{n^3 - n}{3}.$$

(*Hint:* You may find equation (A.3) useful.)

**15.2-6**

Show that a full parenthesization of an  $n$ -element expression has exactly  $n - 1$  pairs of parentheses.

---

**15.3 Elements of dynamic programming**

Although we have just worked through two examples of the dynamic-programming method, you might still be wondering just when the method applies. From an engineering perspective, when should we look for a dynamic-programming solution to a problem? In this section, we examine the two key ingredients that an opti-

mization problem must have in order for dynamic programming to apply: optimal substructure and overlapping subproblems. We also revisit and discuss more fully how memoization might help us take advantage of the overlapping-subproblems property in a top-down recursive approach.

### Optimal substructure

The first step in solving an optimization problem by dynamic programming is to characterize the structure of an optimal solution. Recall that a problem exhibits **optimal substructure** if an optimal solution to the problem contains within it optimal solutions to subproblems. Whenever a problem exhibits optimal substructure, we have a good clue that dynamic programming might apply. (As Chapter 16 discusses, it also might mean that a greedy strategy applies, however.) In dynamic programming, we build an optimal solution to the problem from optimal solutions to subproblems. Consequently, we must take care to ensure that the range of subproblems we consider includes those used in an optimal solution.

We discovered optimal substructure in both of the problems we have examined in this chapter so far. In Section 15.1, we observed that the optimal way of cutting up a rod of length  $n$  (if we make any cuts at all) involves optimally cutting up the two pieces resulting from the first cut. In Section 15.2, we observed that an optimal parenthesization of  $A_i A_{i+1} \cdots A_j$  that splits the product between  $A_k$  and  $A_{k+1}$  contains within it optimal solutions to the problems of parenthesizing  $A_i A_{i+1} \cdots A_k$  and  $A_{k+1} A_{k+2} \cdots A_j$ .

You will find yourself following a common pattern in discovering optimal substructure:

1. You show that a solution to the problem consists of making a choice, such as choosing an initial cut in a rod or choosing an index at which to split the matrix chain. Making this choice leaves one or more subproblems to be solved.
2. You suppose that for a given problem, you are given the choice that leads to an optimal solution. You do not concern yourself yet with how to determine this choice. You just assume that it has been given to you.
3. Given this choice, you determine which subproblems ensue and how to best characterize the resulting space of subproblems.
4. You show that the solutions to the subproblems used within an optimal solution to the problem must themselves be optimal by using a “cut-and-paste” technique. You do so by supposing that each of the subproblem solutions is not optimal and then deriving a contradiction. In particular, by “cutting out” the nonoptimal solution to each subproblem and “pasting in” the optimal one, you show that you can get a better solution to the original problem, thus contradicting your supposition that you already had an optimal solution. If an optimal



solution gives rise to more than one subproblem, they are typically so similar that you can modify the cut-and-paste argument for one to apply to the others with little effort.

To characterize the space of subproblems, a good rule of thumb says to try to keep the space as simple as possible and then expand it as necessary. For example, the space of subproblems that we considered for the rod-cutting problem contained the problems of optimally cutting up a rod of length  $i$  for each size  $i$ . This subproblem space worked well, and we had no need to try a more general space of subproblems.

Conversely, suppose that we had tried to constrain our subproblem space for matrix-chain multiplication to matrix products of the form  $A_1 A_2 \cdots A_j$ . As before, an optimal parenthesization must split this product between  $A_k$  and  $A_{k+1}$  for some  $1 \leq k < j$ . Unless we could guarantee that  $k$  always equals  $j - 1$ , we would find that we had subproblems of the form  $A_1 A_2 \cdots A_k$  and  $A_{k+1} A_{k+2} \cdots A_j$ , and that the latter subproblem is not of the form  $A_1 A_2 \cdots A_j$ . For this problem, we needed to allow our subproblems to vary at “both ends,” that is, to allow both  $i$  and  $j$  to vary in the subproblem  $A_i A_{i+1} \cdots A_j$ .

Optimal substructure varies across problem domains in two ways:

1. how many subproblems an optimal solution to the original problem uses, and
2. how many choices we have in determining which subproblem(s) to use in an optimal solution.

In the rod-cutting problem, an optimal solution for cutting up a rod of size  $n$  uses just one subproblem (of size  $n - i$ ), but we must consider  $n$  choices for  $i$  in order to determine which one yields an optimal solution. Matrix-chain multiplication for the subchain  $A_i A_{i+1} \cdots A_j$  serves as an example with two subproblems and  $j - i$  choices. For a given matrix  $A_k$  at which we split the product, we have two subproblems—parenthesizing  $A_i A_{i+1} \cdots A_k$  and parenthesizing  $A_{k+1} A_{k+2} \cdots A_j$ —and we must solve *both* of them optimally. Once we determine the optimal solutions to subproblems, we choose from among  $j - i$  candidates for the index  $k$ .

Informally, the running time of a dynamic-programming algorithm depends on the product of two factors: the number of subproblems overall and how many choices we look at for each subproblem. In rod cutting, we had  $\Theta(n)$  subproblems overall, and at most  $n$  choices to examine for each, yielding an  $O(n^2)$  running time. Matrix-chain multiplication had  $\Theta(n^2)$  subproblems overall, and in each we had at most  $n - 1$  choices, giving an  $O(n^3)$  running time (actually, a  $\Theta(n^3)$  running time, by Exercise 15.2-5).

Usually, the subproblem graph gives an alternative way to perform the same analysis. Each vertex corresponds to a subproblem, and the choices for a sub-

problem are the edges incident to that subproblem. Recall that in rod cutting, the subproblem graph had  $n$  vertices and at most  $n$  edges per vertex, yielding an  $O(n^2)$  running time. For matrix-chain multiplication, if we were to draw the subproblem graph, it would have  $\Theta(n^2)$  vertices and each vertex would have degree at most  $n - 1$ , giving a total of  $O(n^3)$  vertices and edges.

Dynamic programming often uses optimal substructure in a bottom-up fashion. That is, we first find optimal solutions to subproblems and, having solved the subproblems, we find an optimal solution to the problem. Finding an optimal solution to the problem entails making a choice among subproblems as to which we will use in solving the problem. The cost of the problem solution is usually the subproblem costs plus a cost that is directly attributable to the choice itself. In rod cutting, for example, first we solved the subproblems of determining optimal ways to cut up rods of length  $i$  for  $i = 0, 1, \dots, n - 1$ , and then we determined which such subproblem yielded an optimal solution for a rod of length  $n$ , using equation (15.2). The cost attributable to the choice itself is the term  $p_i$  in equation (15.2). In matrix-chain multiplication, we determined optimal parenthesizations of subchains of  $A_i A_{i+1} \cdots A_j$ , and then we chose the matrix  $A_k$  at which to split the product. The cost attributable to the choice itself is the term  $p_{i-1} p_k p_j$ .

In Chapter 16, we shall examine “greedy algorithms,” which have many similarities to dynamic programming. In particular, problems to which greedy algorithms apply have optimal substructure. One major difference between greedy algorithms and dynamic programming is that instead of first finding optimal solutions to subproblems and then making an informed choice, greedy algorithms first make a “greedy” choice—the choice that looks best at the time—and then solve a resulting subproblem, without bothering to solve all possible related smaller subproblems. Surprisingly, in some cases this strategy works!

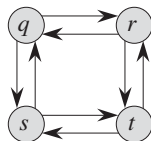
### *Subtleties*

You should be careful not to assume that optimal substructure applies when it does not. Consider the following two problems in which we are given a directed graph  $G = (V, E)$  and vertices  $u, v \in V$ .

**Unweighted shortest path:**<sup>3</sup> Find a path from  $u$  to  $v$  consisting of the fewest edges. Such a path must be simple, since removing a cycle from a path produces a path with fewer edges.

---

<sup>3</sup>We use the term “unweighted” to distinguish this problem from that of finding shortest paths with weighted edges, which we shall see in Chapters 24 and 25. We can use the breadth-first search technique of Chapter 22 to solve the unweighted problem.



**Figure 15.6** A directed graph showing that the problem of finding a longest simple path in an unweighted directed graph does not have optimal substructure. The path  $q \rightarrow r \rightarrow t$  is a longest simple path from  $q$  to  $t$ , but the subpath  $q \rightarrow r$  is not a longest simple path from  $q$  to  $r$ , nor is the subpath  $r \rightarrow t$  a longest simple path from  $r$  to  $t$ .

**Unweighted longest simple path:** Find a simple path from  $u$  to  $v$  consisting of the most edges. We need to include the requirement of simplicity because otherwise we can traverse a cycle as many times as we like to create paths with an arbitrarily large number of edges.

The unweighted shortest-path problem exhibits optimal substructure, as follows. Suppose that  $u \neq v$ , so that the problem is nontrivial. Then, any path  $p$  from  $u$  to  $v$  must contain an intermediate vertex, say  $w$ . (Note that  $w$  may be  $u$  or  $v$ .) Thus, we can decompose the path  $u \xrightarrow{p} v$  into subpaths  $u \xrightarrow{p_1} w \xrightarrow{p_2} v$ . Clearly, the number of edges in  $p$  equals the number of edges in  $p_1$  plus the number of edges in  $p_2$ . We claim that if  $p$  is an optimal (i.e., shortest) path from  $u$  to  $v$ , then  $p_1$  must be a shortest path from  $u$  to  $w$ . Why? We use a “cut-and-paste” argument: if there were another path, say  $p'_1$ , from  $u$  to  $w$  with fewer edges than  $p_1$ , then we could cut out  $p_1$  and paste in  $p'_1$  to produce a path  $u \xrightarrow{p'_1} w \xrightarrow{p_2} v$  with fewer edges than  $p$ , thus contradicting  $p$ ’s optimality. Symmetrically,  $p_2$  must be a shortest path from  $w$  to  $v$ . Thus, we can find a shortest path from  $u$  to  $v$  by considering all intermediate vertices  $w$ , finding a shortest path from  $u$  to  $w$  and a shortest path from  $w$  to  $v$ , and choosing an intermediate vertex  $w$  that yields the overall shortest path. In Section 25.2, we use a variant of this observation of optimal substructure to find a shortest path between every pair of vertices on a weighted, directed graph.

You might be tempted to assume that the problem of finding an unweighted longest simple path exhibits optimal substructure as well. After all, if we decompose a longest simple path  $u \xrightarrow{p} v$  into subpaths  $u \xrightarrow{p_1} w \xrightarrow{p_2} v$ , then mustn’t  $p_1$  be a longest simple path from  $u$  to  $w$ , and mustn’t  $p_2$  be a longest simple path from  $w$  to  $v$ ? The answer is no! Figure 15.6 supplies an example. Consider the path  $q \rightarrow r \rightarrow t$ , which is a longest simple path from  $q$  to  $t$ . Is  $q \rightarrow r$  a longest simple path from  $q$  to  $r$ ? No, for the path  $q \rightarrow s \rightarrow t \rightarrow r$  is a simple path that is longer. Is  $r \rightarrow t$  a longest simple path from  $r$  to  $t$ ? No again, for the path  $r \rightarrow q \rightarrow s \rightarrow t$  is a simple path that is longer.

This example shows that for longest simple paths, not only does the problem lack optimal substructure, but we cannot necessarily assemble a “legal” solution to the problem from solutions to subproblems. If we combine the longest simple paths  $q \rightarrow s \rightarrow t \rightarrow r$  and  $r \rightarrow q \rightarrow s \rightarrow t$ , we get the path  $q \rightarrow s \rightarrow t \rightarrow r \rightarrow q \rightarrow s \rightarrow t$ , which is not simple. Indeed, the problem of finding an unweighted longest simple path does not appear to have any sort of optimal substructure. No efficient dynamic-programming algorithm for this problem has ever been found. In fact, this problem is NP-complete, which—as we shall see in Chapter 34—means that we are unlikely to find a way to solve it in polynomial time.

Why is the substructure of a longest simple path so different from that of a shortest path? Although a solution to a problem for both longest and shortest paths uses two subproblems, the subproblems in finding the longest simple path are not *independent*, whereas for shortest paths they are. What do we mean by subproblems being independent? We mean that the solution to one subproblem does not affect the solution to another subproblem of the same problem. For the example of Figure 15.6, we have the problem of finding a longest simple path from  $q$  to  $t$  with two subproblems: finding longest simple paths from  $q$  to  $r$  and from  $r$  to  $t$ . For the first of these subproblems, we choose the path  $q \rightarrow s \rightarrow t \rightarrow r$ , and so we have also used the vertices  $s$  and  $t$ . We can no longer use these vertices in the second subproblem, since the combination of the two solutions to subproblems would yield a path that is not simple. If we cannot use vertex  $t$  in the second problem, then we cannot solve it at all, since  $t$  is required to be on the path that we find, and it is not the vertex at which we are “splicing” together the subproblem solutions (that vertex being  $r$ ). Because we use vertices  $s$  and  $t$  in one subproblem solution, we cannot use them in the other subproblem solution. We must use at least one of them to solve the other subproblem, however, and we must use both of them to solve it optimally. Thus, we say that these subproblems are not independent. Looked at another way, using resources in solving one subproblem (those resources being vertices) renders them unavailable for the other subproblem.

Why, then, are the subproblems independent for finding a shortest path? The answer is that by nature, the subproblems do not share resources. We claim that if a vertex  $w$  is on a shortest path  $p$  from  $u$  to  $v$ , then we can splice together *any* shortest path  $u \xrightarrow{p_1} w$  and *any* shortest path  $w \xrightarrow{p_2} v$  to produce a shortest path from  $u$  to  $v$ . We are assured that, other than  $w$ , no vertex can appear in both paths  $p_1$  and  $p_2$ . Why? Suppose that some vertex  $x \neq w$  appears in both  $p_1$  and  $p_2$ , so that we can decompose  $p_1$  as  $u \xrightarrow{p_{ux}} x \rightsquigarrow w$  and  $p_2$  as  $w \rightsquigarrow x \xrightarrow{p_{xv}} v$ . By the optimal substructure of this problem, path  $p$  has as many edges as  $p_1$  and  $p_2$  together; let’s say that  $p$  has  $e$  edges. Now let us construct a path  $p' = u \xrightarrow{p_{ux}} x \xrightarrow{p_{xv}} v$  from  $u$  to  $v$ . Because we have excised the paths from  $x$  to  $w$  and from  $w$  to  $x$ , each of which contains at least one edge, path  $p'$  contains at most  $e - 2$  edges, which contradicts

the assumption that  $p$  is a shortest path. Thus, we are assured that the subproblems for the shortest-path problem are independent.

Both problems examined in Sections 15.1 and 15.2 have independent subproblems. In matrix-chain multiplication, the subproblems are multiplying subchains  $A_i A_{i+1} \cdots A_k$  and  $A_{k+1} A_{k+2} \cdots A_j$ . These subchains are disjoint, so that no matrix could possibly be included in both of them. In rod cutting, to determine the best way to cut up a rod of length  $n$ , we look at the best ways of cutting up rods of length  $i$  for  $i = 0, 1, \dots, n-1$ . Because an optimal solution to the length- $n$  problem includes just one of these subproblem solutions (after we have cut off the first piece), independence of subproblems is not an issue.

### Overlapping subproblems

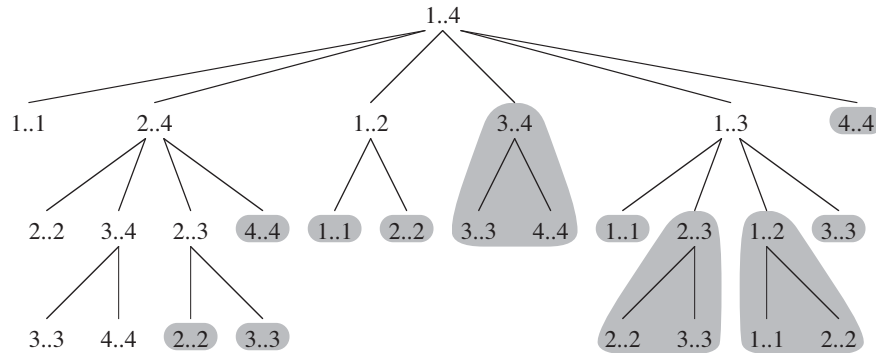
The second ingredient that an optimization problem must have for dynamic programming to apply is that the space of subproblems must be “small” in the sense that a recursive algorithm for the problem solves the same subproblems over and over, rather than always generating new subproblems. Typically, the total number of distinct subproblems is a polynomial in the input size. When a recursive algorithm revisits the same problem repeatedly, we say that the optimization problem has *overlapping subproblems*.<sup>4</sup> In contrast, a problem for which a divide-and-conquer approach is suitable usually generates brand-new problems at each step of the recursion. Dynamic-programming algorithms typically take advantage of overlapping subproblems by solving each subproblem once and then storing the solution in a table where it can be looked up when needed, using constant time per lookup.

In Section 15.1, we briefly examined how a recursive solution to rod cutting makes exponentially many calls to find solutions of smaller subproblems. Our dynamic-programming solution takes an exponential-time recursive algorithm down to quadratic time.

To illustrate the overlapping-subproblems property in greater detail, let us re-examine the matrix-chain multiplication problem. Referring back to Figure 15.5, observe that MATRIX-CHAIN-ORDER repeatedly looks up the solution to subproblems in lower rows when solving subproblems in higher rows. For example, it references entry  $m[3, 4]$  four times: during the computations of  $m[2, 4]$ ,  $m[1, 4]$ ,

---

<sup>4</sup>It may seem strange that dynamic programming relies on subproblems being both independent and overlapping. Although these requirements may sound contradictory, they describe two different notions, rather than two points on the same axis. Two subproblems of the same problem are independent if they do not share resources. Two subproblems are overlapping if they are really the same subproblem that occurs as a subproblem of different problems.



**Figure 15.7** The recursion tree for the computation of `RECURSIVE-MATRIX-CHAIN( $p, 1, 4$ )`. Each node contains the parameters  $i$  and  $j$ . The computations performed in a shaded subtree are replaced by a single table lookup in `MEMOIZED-MATRIX-CHAIN`.

$m[3, 5]$ , and  $m[3, 6]$ . If we were to recompute  $m[3, 4]$  each time, rather than just looking it up, the running time would increase dramatically. To see how, consider the following (inefficient) recursive procedure that determines  $m[i, j]$ , the minimum number of scalar multiplications needed to compute the matrix-chain product  $A_{i..j} = A_i A_{i+1} \cdots A_j$ . The procedure is based directly on the recurrence (15.7).

`RECURSIVE-MATRIX-CHAIN( $p, i, j$ )`

```

1  if  $i == j$ 
2      return 0
3   $m[i, j] = \infty$ 
4  for  $k = i$  to  $j - 1$ 
5       $q = \text{RECURSIVE-MATRIX-CHAIN}(p, i, k)$ 
          +  $\text{RECURSIVE-MATRIX-CHAIN}(p, k + 1, j)$ 
          +  $p_{i-1} p_k p_j$ 
6      if  $q < m[i, j]$ 
7           $m[i, j] = q$ 
8  return  $m[i, j]$ 
```

Figure 15.7 shows the recursion tree produced by the call `RECURSIVE-MATRIX-CHAIN( $p, 1, 4$ )`. Each node is labeled by the values of the parameters  $i$  and  $j$ . Observe that some pairs of values occur many times.

In fact, we can show that the time to compute  $m[1, n]$  by this recursive procedure is at least exponential in  $n$ . Let  $T(n)$  denote the time taken by `RECURSIVE-MATRIX-CHAIN` to compute an optimal parenthesization of a chain of  $n$  matrices. Because the execution of lines 1–2 and of lines 6–7 each take at least unit time, as

does the multiplication in line 5, inspection of the procedure yields the recurrence

$$\begin{aligned} T(1) &\geq 1, \\ T(n) &\geq 1 + \sum_{k=1}^{n-1} (T(k) + T(n-k) + 1) \quad \text{for } n > 1. \end{aligned}$$

Noting that for  $i = 1, 2, \dots, n-1$ , each term  $T(i)$  appears once as  $T(k)$  and once as  $T(n-k)$ , and collecting the  $n-1$  1s in the summation together with the 1 out front, we can rewrite the recurrence as

$$T(n) \geq 2 \sum_{i=1}^{n-1} T(i) + n. \quad (15.8)$$

We shall prove that  $T(n) = \Omega(2^n)$  using the substitution method. Specifically, we shall show that  $T(n) \geq 2^{n-1}$  for all  $n \geq 1$ . The basis is easy, since  $T(1) \geq 1 = 2^0$ . Inductively, for  $n \geq 2$  we have

$$\begin{aligned} T(n) &\geq 2 \sum_{i=1}^{n-1} 2^{i-1} + n \\ &= 2 \sum_{i=0}^{n-2} 2^i + n \\ &= 2(2^{n-1} - 1) + n \quad (\text{by equation (A.5)}) \\ &= 2^n - 2 + n \\ &\geq 2^{n-1}, \end{aligned}$$

which completes the proof. Thus, the total amount of work performed by the call `RECURSIVE-MATRIX-CHAIN( $p, 1, n$ )` is at least exponential in  $n$ .

Compare this top-down, recursive algorithm (without memoization) with the bottom-up dynamic-programming algorithm. The latter is more efficient because it takes advantage of the overlapping-subproblems property. Matrix-chain multiplication has only  $\Theta(n^2)$  distinct subproblems, and the dynamic-programming algorithm solves each exactly once. The recursive algorithm, on the other hand, must again solve each subproblem every time it reappears in the recursion tree. Whenever a recursion tree for the natural recursive solution to a problem contains the same subproblem repeatedly, and the total number of distinct subproblems is small, dynamic programming can improve efficiency, sometimes dramatically.

### Reconstructing an optimal solution

As a practical matter, we often store which choice we made in each subproblem in a table so that we do not have to reconstruct this information from the costs that we stored.

For matrix-chain multiplication, the table  $s[i, j]$  saves us a significant amount of work when reconstructing an optimal solution. Suppose that we did not maintain the  $s[i, j]$  table, having filled in only the table  $m[i, j]$  containing optimal subproblem costs. We choose from among  $j - i$  possibilities when we determine which subproblems to use in an optimal solution to parenthesizing  $A_i A_{i+1} \cdots A_j$ , and  $j - i$  is not a constant. Therefore, it would take  $\Theta(j - i) = \omega(1)$  time to reconstruct which subproblems we chose for a solution to a given problem. By storing in  $s[i, j]$  the index of the matrix at which we split the product  $A_i A_{i+1} \cdots A_j$ , we can reconstruct each choice in  $O(1)$  time.

### Memoization

As we saw for the rod-cutting problem, there is an alternative approach to dynamic programming that often offers the efficiency of the bottom-up dynamic-programming approach while maintaining a top-down strategy. The idea is to **memoize** the natural, but inefficient, recursive algorithm. As in the bottom-up approach, we maintain a table with subproblem solutions, but the control structure for filling in the table is more like the recursive algorithm.

A memoized recursive algorithm maintains an entry in a table for the solution to each subproblem. Each table entry initially contains a special value to indicate that the entry has yet to be filled in. When the subproblem is first encountered as the recursive algorithm unfolds, its solution is computed and then stored in the table. Each subsequent time that we encounter this subproblem, we simply look up the value stored in the table and return it.<sup>5</sup>

Here is a memoized version of **RECURSIVE-MATRIX-CHAIN**. Note where it resembles the memoized top-down method for the rod-cutting problem.

---

<sup>5</sup>This approach presupposes that we know the set of all possible subproblem parameters and that we have established the relationship between table positions and subproblems. Another, more general, approach is to memoize by using hashing with the subproblem parameters as keys.



MEMOIZED-MATRIX-CHAIN( $p$ )

```

1   $n = p.length - 1$ 
2  let  $m[1..n, 1..n]$  be a new table
3  for  $i = 1$  to  $n$ 
4      for  $j = i$  to  $n$ 
5           $m[i, j] = \infty$ 
6  return LOOKUP-CHAIN( $m, p, 1, n$ )

```

LOOKUP-CHAIN( $m, p, i, j$ )

```

1  if  $m[i, j] < \infty$ 
2      return  $m[i, j]$ 
3  if  $i == j$ 
4       $m[i, j] = 0$ 
5  else for  $k = i$  to  $j - 1$ 
6       $q = \text{LOOKUP-CHAIN}(m, p, i, k)$ 
         $+ \text{LOOKUP-CHAIN}(m, p, k + 1, j) + p_{i-1}p_kp_j$ 
7      if  $q < m[i, j]$ 
8           $m[i, j] = q$ 
9  return  $m[i, j]$ 

```

The MEMOIZED-MATRIX-CHAIN procedure, like MATRIX-CHAIN-ORDER, maintains a table  $m[1..n, 1..n]$  of computed values of  $m[i, j]$ , the minimum number of scalar multiplications needed to compute the matrix  $A_{i..j}$ . Each table entry initially contains the value  $\infty$  to indicate that the entry has yet to be filled in. Upon calling LOOKUP-CHAIN( $m, p, i, j$ ), if line 1 finds that  $m[i, j] < \infty$ , then the procedure simply returns the previously computed cost  $m[i, j]$  in line 2. Otherwise, the cost is computed as in RECURSIVE-MATRIX-CHAIN, stored in  $m[i, j]$ , and returned. Thus, LOOKUP-CHAIN( $m, p, i, j$ ) always returns the value of  $m[i, j]$ , but it computes it only upon the first call of LOOKUP-CHAIN with these specific values of  $i$  and  $j$ .

Figure 15.7 illustrates how MEMOIZED-MATRIX-CHAIN saves time compared with RECURSIVE-MATRIX-CHAIN. Shaded subtrees represent values that it looks up rather than recomputes.

Like the bottom-up dynamic-programming algorithm MATRIX-CHAIN-ORDER, the procedure MEMOIZED-MATRIX-CHAIN runs in  $O(n^3)$  time. Line 5 of MEMOIZED-MATRIX-CHAIN executes  $\Theta(n^2)$  times. We can categorize the calls of LOOKUP-CHAIN into two types:

1. calls in which  $m[i, j] = \infty$ , so that lines 3–9 execute, and
2. calls in which  $m[i, j] < \infty$ , so that LOOKUP-CHAIN simply returns in line 2.

There are  $\Theta(n^2)$  calls of the first type, one per table entry. All calls of the second type are made as recursive calls by calls of the first type. Whenever a given call of LOOKUP-CHAIN makes recursive calls, it makes  $O(n)$  of them. Therefore, there are  $O(n^3)$  calls of the second type in all. Each call of the second type takes  $O(1)$  time, and each call of the first type takes  $O(n)$  time plus the time spent in its recursive calls. The total time, therefore, is  $O(n^3)$ . Memoization thus turns an  $\Omega(2^n)$ -time algorithm into an  $O(n^3)$ -time algorithm.

In summary, we can solve the matrix-chain multiplication problem by either a top-down, memoized dynamic-programming algorithm or a bottom-up dynamic-programming algorithm in  $O(n^3)$  time. Both methods take advantage of the overlapping-subproblems property. There are only  $\Theta(n^2)$  distinct subproblems in total, and either of these methods computes the solution to each subproblem only once. Without memoization, the natural recursive algorithm runs in exponential time, since solved subproblems are repeatedly solved.

In general practice, if all subproblems must be solved at least once, a bottom-up dynamic-programming algorithm usually outperforms the corresponding top-down memoized algorithm by a constant factor, because the bottom-up algorithm has no overhead for recursion and less overhead for maintaining the table. Moreover, for some problems we can exploit the regular pattern of table accesses in the dynamic-programming algorithm to reduce time or space requirements even further. Alternatively, if some subproblems in the subproblem space need not be solved at all, the memoized solution has the advantage of solving only those subproblems that are definitely required.

## Exercises

### 15.3-1

Which is a more efficient way to determine the optimal number of multiplications in a matrix-chain multiplication problem: enumerating all the ways of parenthesizing the product and computing the number of multiplications for each, or running RECURSIVE-MATRIX-CHAIN? Justify your answer.

### 15.3-2

Draw the recursion tree for the MERGE-SORT procedure from Section 2.3.1 on an array of 16 elements. Explain why memoization fails to speed up a good divide-and-conquer algorithm such as MERGE-SORT.

### 15.3-3

Consider a variant of the matrix-chain multiplication problem in which the goal is to parenthesize the sequence of matrices so as to maximize, rather than minimize,

the number of scalar multiplications. Does this problem exhibit optimal substructure?

### 15.3-4

As stated, in dynamic programming we first solve the subproblems and then choose which of them to use in an optimal solution to the problem. Professor Capulet claims that we do not always need to solve all the subproblems in order to find an optimal solution. She suggests that we can find an optimal solution to the matrix-chain multiplication problem by always choosing the matrix  $A_k$  at which to split the subproduct  $A_i A_{i+1} \cdots A_j$  (by selecting  $k$  to minimize the quantity  $p_{i-1} p_k p_j$ ) *before* solving the subproblems. Find an instance of the matrix-chain multiplication problem for which this greedy approach yields a suboptimal solution.

### 15.3-5

Suppose that in the rod-cutting problem of Section 15.1, we also had limit  $l_i$  on the number of pieces of length  $i$  that we are allowed to produce, for  $i = 1, 2, \dots, n$ . Show that the optimal-substructure property described in Section 15.1 no longer holds.

### 15.3-6

Imagine that you wish to exchange one currency for another. You realize that instead of directly exchanging one currency for another, you might be better off making a series of trades through other currencies, winding up with the currency you want. Suppose that you can trade  $n$  different currencies, numbered  $1, 2, \dots, n$ , where you start with currency 1 and wish to wind up with currency  $n$ . You are given, for each pair of currencies  $i$  and  $j$ , an exchange rate  $r_{ij}$ , meaning that if you start with  $d$  units of currency  $i$ , you can trade for  $dr_{ij}$  units of currency  $j$ . A sequence of trades may entail a commission, which depends on the number of trades you make. Let  $c_k$  be the commission that you are charged when you make  $k$  trades. Show that, if  $c_k = 0$  for all  $k = 1, 2, \dots, n$ , then the problem of finding the best sequence of exchanges from currency 1 to currency  $n$  exhibits optimal substructure. Then show that if commissions  $c_k$  are arbitrary values, then the problem of finding the best sequence of exchanges from currency 1 to currency  $n$  does not necessarily exhibit optimal substructure.

---

## 15.4 Longest common subsequence

Biological applications often need to compare the DNA of two (or more) different organisms. A strand of DNA consists of a string of molecules called

**bases**, where the possible bases are adenine, guanine, cytosine, and thymine. Representing each of these bases by its initial letter, we can express a strand of DNA as a string over the finite set  $\{A, C, G, T\}$ . (See Appendix C for the definition of a string.) For example, the DNA of one organism may be  $S_1 = \text{ACCGGTCGAGTGC GCGGAAGCCGGCCGAA}$ , and the DNA of another organism may be  $S_2 = \text{GTCGTT CGGAATGCCGTTGCTCTGTAAA}$ . One reason to compare two strands of DNA is to determine how “similar” the two strands are, as some measure of how closely related the two organisms are. We can, and do, define similarity in many different ways. For example, we can say that two DNA strands are similar if one is a substring of the other. (Chapter 32 explores algorithms to solve this problem.) In our example, neither  $S_1$  nor  $S_2$  is a substring of the other. Alternatively, we could say that two strands are similar if the number of changes needed to turn one into the other is small. (Problem 15-5 looks at this notion.) Yet another way to measure the similarity of strands  $S_1$  and  $S_2$  is by finding a third strand  $S_3$  in which the bases in  $S_3$  appear in each of  $S_1$  and  $S_2$ ; these bases must appear in the same order, but not necessarily consecutively. The longer the strand  $S_3$  we can find, the more similar  $S_1$  and  $S_2$  are. In our example, the longest strand  $S_3$  is  $\text{GTCGTCGGAAGCCGGCCGAA}$ .

We formalize this last notion of similarity as the longest-common-subsequence problem. A subsequence of a given sequence is just the given sequence with zero or more elements left out. Formally, given a sequence  $X = \langle x_1, x_2, \dots, x_m \rangle$ , another sequence  $Z = \langle z_1, z_2, \dots, z_k \rangle$  is a **subsequence** of  $X$  if there exists a strictly increasing sequence  $\langle i_1, i_2, \dots, i_k \rangle$  of indices of  $X$  such that for all  $j = 1, 2, \dots, k$ , we have  $x_{i_j} = z_j$ . For example,  $Z = \langle B, C, D, B \rangle$  is a subsequence of  $X = \langle A, B, C, B, D, A, B \rangle$  with corresponding index sequence  $\langle 2, 3, 5, 7 \rangle$ .

Given two sequences  $X$  and  $Y$ , we say that a sequence  $Z$  is a **common subsequence** of  $X$  and  $Y$  if  $Z$  is a subsequence of both  $X$  and  $Y$ . For example, if  $X = \langle A, B, C, B, D, A, B \rangle$  and  $Y = \langle B, D, C, A, B, A \rangle$ , the sequence  $\langle B, C, A \rangle$  is a common subsequence of both  $X$  and  $Y$ . The sequence  $\langle B, C, A \rangle$  is not a **longest common subsequence (LCS)** of  $X$  and  $Y$ , however, since it has length 3 and the sequence  $\langle B, C, B, A \rangle$ , which is also common to both  $X$  and  $Y$ , has length 4. The sequence  $\langle B, C, B, A \rangle$  is an LCS of  $X$  and  $Y$ , as is the sequence  $\langle B, D, A, B \rangle$ , since  $X$  and  $Y$  have no common subsequence of length 5 or greater.

In the **longest-common-subsequence problem**, we are given two sequences  $X = \langle x_1, x_2, \dots, x_m \rangle$  and  $Y = \langle y_1, y_2, \dots, y_n \rangle$  and wish to find a maximum-length common subsequence of  $X$  and  $Y$ . This section shows how to efficiently solve the LCS problem using dynamic programming.

### Step 1: Characterizing a longest common subsequence

In a brute-force approach to solving the LCS problem, we would enumerate all subsequences of  $X$  and check each subsequence to see whether it is also a subsequence of  $Y$ , keeping track of the longest subsequence we find. Each subsequence of  $X$  corresponds to a subset of the indices  $\{1, 2, \dots, m\}$  of  $X$ . Because  $X$  has  $2^m$  subsequences, this approach requires exponential time, making it impractical for long sequences.

The LCS problem has an optimal-substructure property, however, as the following theorem shows. As we shall see, the natural classes of subproblems correspond to pairs of “prefixes” of the two input sequences. To be precise, given a sequence  $X = \langle x_1, x_2, \dots, x_m \rangle$ , we define the  $i$ th **prefix** of  $X$ , for  $i = 0, 1, \dots, m$ , as  $X_i = \langle x_1, x_2, \dots, x_i \rangle$ . For example, if  $X = \langle A, B, C, B, D, A, B \rangle$ , then  $X_4 = \langle A, B, C, B \rangle$  and  $X_0$  is the empty sequence.

#### **Theorem 15.1 (Optimal substructure of an LCS)**

Let  $X = \langle x_1, x_2, \dots, x_m \rangle$  and  $Y = \langle y_1, y_2, \dots, y_n \rangle$  be sequences, and let  $Z = \langle z_1, z_2, \dots, z_k \rangle$  be any LCS of  $X$  and  $Y$ .

1. If  $x_m = y_n$ , then  $z_k = x_m = y_n$  and  $Z_{k-1}$  is an LCS of  $X_{m-1}$  and  $Y_{n-1}$ .
2. If  $x_m \neq y_n$ , then  $z_k \neq x_m$  implies that  $Z$  is an LCS of  $X_{m-1}$  and  $Y$ .
3. If  $x_m \neq y_n$ , then  $z_k \neq y_n$  implies that  $Z$  is an LCS of  $X$  and  $Y_{n-1}$ .

**Proof** (1) If  $z_k \neq x_m$ , then we could append  $x_m = y_n$  to  $Z$  to obtain a common subsequence of  $X$  and  $Y$  of length  $k + 1$ , contradicting the supposition that  $Z$  is a *longest* common subsequence of  $X$  and  $Y$ . Thus, we must have  $z_k = x_m = y_n$ . Now, the prefix  $Z_{k-1}$  is a length- $(k - 1)$  common subsequence of  $X_{m-1}$  and  $Y_{n-1}$ . We wish to show that it is an LCS. Suppose for the purpose of contradiction that there exists a common subsequence  $W$  of  $X_{m-1}$  and  $Y_{n-1}$  with length greater than  $k - 1$ . Then, appending  $x_m = y_n$  to  $W$  produces a common subsequence of  $X$  and  $Y$  whose length is greater than  $k$ , which is a contradiction.

(2) If  $z_k \neq x_m$ , then  $Z$  is a common subsequence of  $X_{m-1}$  and  $Y$ . If there were a common subsequence  $W$  of  $X_{m-1}$  and  $Y$  with length greater than  $k$ , then  $W$  would also be a common subsequence of  $X_m$  and  $Y$ , contradicting the assumption that  $Z$  is an LCS of  $X$  and  $Y$ .

(3) The proof is symmetric to (2). ■

The way that Theorem 15.1 characterizes longest common subsequences tells us that an LCS of two sequences contains within it an LCS of prefixes of the two sequences. Thus, the LCS problem has an optimal-substructure property. A recur-

sive solution also has the overlapping-subproblems property, as we shall see in a moment.

### Step 2: A recursive solution

Theorem 15.1 implies that we should examine either one or two subproblems when finding an LCS of  $X = \langle x_1, x_2, \dots, x_m \rangle$  and  $Y = \langle y_1, y_2, \dots, y_n \rangle$ . If  $x_m = y_n$ , we must find an LCS of  $X_{m-1}$  and  $Y_{n-1}$ . Appending  $x_m = y_n$  to this LCS yields an LCS of  $X$  and  $Y$ . If  $x_m \neq y_n$ , then we must solve two subproblems: finding an LCS of  $X_{m-1}$  and  $Y$  and finding an LCS of  $X$  and  $Y_{n-1}$ . Whichever of these two LCSs is longer is an LCS of  $X$  and  $Y$ . Because these cases exhaust all possibilities, we know that one of the optimal subproblem solutions must appear within an LCS of  $X$  and  $Y$ .

We can readily see the overlapping-subproblems property in the LCS problem. To find an LCS of  $X$  and  $Y$ , we may need to find the LCSs of  $X$  and  $Y_{n-1}$  and of  $X_{m-1}$  and  $Y$ . But each of these subproblems has the subsubproblem of finding an LCS of  $X_{m-1}$  and  $Y_{n-1}$ . Many other subproblems share subsubproblems.

As in the matrix-chain multiplication problem, our recursive solution to the LCS problem involves establishing a recurrence for the value of an optimal solution. Let us define  $c[i, j]$  to be the length of an LCS of the sequences  $X_i$  and  $Y_j$ . If either  $i = 0$  or  $j = 0$ , one of the sequences has length 0, and so the LCS has length 0. The optimal substructure of the LCS problem gives the recursive formula

$$c[i, j] = \begin{cases} 0 & \text{if } i = 0 \text{ or } j = 0, \\ c[i-1, j-1] + 1 & \text{if } i, j > 0 \text{ and } x_i = y_j, \\ \max(c[i, j-1], c[i-1, j]) & \text{if } i, j > 0 \text{ and } x_i \neq y_j. \end{cases} \quad (15.9)$$

Observe that in this recursive formulation, a condition in the problem restricts which subproblems we may consider. When  $x_i = y_j$ , we can and should consider the subproblem of finding an LCS of  $X_{i-1}$  and  $Y_{j-1}$ . Otherwise, we instead consider the two subproblems of finding an LCS of  $X_i$  and  $Y_{j-1}$  and of  $X_{i-1}$  and  $Y_j$ . In the previous dynamic-programming algorithms we have examined—for rod cutting and matrix-chain multiplication—we ruled out no subproblems due to conditions in the problem. Finding an LCS is not the only dynamic-programming algorithm that rules out subproblems based on conditions in the problem. For example, the edit-distance problem (see Problem 15-5) has this characteristic.

### Step 3: Computing the length of an LCS

Based on equation (15.9), we could easily write an exponential-time recursive algorithm to compute the length of an LCS of two sequences. Since the LCS problem

has only  $\Theta(mn)$  distinct subproblems, however, we can use dynamic programming to compute the solutions bottom up.

Procedure **LCS-LENGTH** takes two sequences  $X = \langle x_1, x_2, \dots, x_m \rangle$  and  $Y = \langle y_1, y_2, \dots, y_n \rangle$  as inputs. It stores the  $c[i, j]$  values in a table  $c[0..m, 0..n]$ , and it computes the entries in **row-major** order. (That is, the procedure fills in the first row of  $c$  from left to right, then the second row, and so on.) The procedure also maintains the table  $b[1..m, 1..n]$  to help us construct an optimal solution. Intuitively,  $b[i, j]$  points to the table entry corresponding to the optimal subproblem solution chosen when computing  $c[i, j]$ . The procedure returns the  $b$  and  $c$  tables;  $c[m, n]$  contains the length of an LCS of  $X$  and  $Y$ .

```

LCS-LENGTH( $X, Y$ )
1   $m = X.length$ 
2   $n = Y.length$ 
3  let  $b[1..m, 1..n]$  and  $c[0..m, 0..n]$  be new tables
4  for  $i = 1$  to  $m$ 
5       $c[i, 0] = 0$ 
6  for  $j = 0$  to  $n$ 
7       $c[0, j] = 0$ 
8  for  $i = 1$  to  $m$ 
9      for  $j = 1$  to  $n$ 
10         if  $x_i == y_j$ 
11              $c[i, j] = c[i - 1, j - 1] + 1$ 
12              $b[i, j] = \nwarrow$ 
13         elseif  $c[i - 1, j] \geq c[i, j - 1]$ 
14              $c[i, j] = c[i - 1, j]$ 
15              $b[i, j] = \uparrow$ 
16         else  $c[i, j] = c[i, j - 1]$ 
17              $b[i, j] = \leftarrow$ 
18  return  $c$  and  $b$ 

```

Figure 15.8 shows the tables produced by **LCS-LENGTH** on the sequences  $X = \langle A, B, C, B, D, A, B \rangle$  and  $Y = \langle B, D, C, A, B, A \rangle$ . The running time of the procedure is  $\Theta(mn)$ , since each table entry takes  $\Theta(1)$  time to compute.

#### Step 4: Constructing an LCS

The  $b$  table returned by **LCS-LENGTH** enables us to quickly construct an LCS of  $X = \langle x_1, x_2, \dots, x_m \rangle$  and  $Y = \langle y_1, y_2, \dots, y_n \rangle$ . We simply begin at  $b[m, n]$  and trace through the table by following the arrows. Whenever we encounter a “ $\nwarrow$ ” in entry  $b[i, j]$ , it implies that  $x_i = y_j$  is an element of the LCS that **LCS-LENGTH**

		$j$	0	1	2	3	4	5	6
$i$	$x_i$	$y_j$		<b>B</b>	D	<b>C</b>	A	<b>B</b>	<b>A</b>
0			0	0	0	0	0	0	0
1	A		0	↑	↑	↑	↖	←	↖
2	<b>B</b>		0	↖	1	←	1	↖	←
3	<b>C</b>		0	↑	↑	↑	↖	2	↑
4	<b>B</b>		0	↖	1	↑	2	↖	←
5	D		0	↑	↖	2	↑	2	↑
6	<b>A</b>		0	↑	↑	↑	↖	3	↖
7	B		0	↖	↑	↑	↑	4	↑

**Figure 15.8** The  $c$  and  $b$  tables computed by LCS-LENGTH on the sequences  $X = \langle A, B, C, B, D, A, B \rangle$  and  $Y = \langle B, D, C, A, B, A \rangle$ . The square in row  $i$  and column  $j$  contains the value of  $c[i, j]$  and the appropriate arrow for the value of  $b[i, j]$ . The entry 4 in  $c[7, 6]$ —the lower right-hand corner of the table—is the length of an LCS  $\langle B, C, B, A \rangle$  of  $X$  and  $Y$ . For  $i, j > 0$ , entry  $c[i, j]$  depends only on whether  $x_i = y_j$  and the values in entries  $c[i - 1, j]$ ,  $c[i, j - 1]$ , and  $c[i - 1, j - 1]$ , which are computed before  $c[i, j]$ . To reconstruct the elements of an LCS, follow the  $b[i, j]$  arrows from the lower right-hand corner; the sequence is shaded. Each “↖” on the shaded sequence corresponds to an entry (highlighted) for which  $x_i = y_j$  is a member of an LCS.

found. With this method, we encounter the elements of this LCS in reverse order. The following recursive procedure prints out an LCS of  $X$  and  $Y$  in the proper, forward order. The initial call is PRINT-LCS( $b, X, X.length, Y.length$ ).

```

PRINT-LCS( $b, X, i, j$ )
1  if  $i == 0$  or  $j == 0$ 
2      return
3  if  $b[i, j] == \text{“}\nwarrow\text{”}$ 
4      PRINT-LCS( $b, X, i - 1, j - 1$ )
5      print  $x_i$ 
6  elseif  $b[i, j] == \text{“}\uparrow\text{”}$ 
7      PRINT-LCS( $b, X, i - 1, j$ )
8  else PRINT-LCS( $b, X, i, j - 1$ )

```

For the  $b$  table in Figure 15.8, this procedure prints  $BCBA$ . The procedure takes time  $O(m + n)$ , since it decrements at least one of  $i$  and  $j$  in each recursive call.



### Improving the code

Once you have developed an algorithm, you will often find that you can improve on the time or space it uses. Some changes can simplify the code and improve constant factors but otherwise yield no asymptotic improvement in performance. Others can yield substantial asymptotic savings in time and space.

In the LCS algorithm, for example, we can eliminate the  $b$  table altogether. Each  $c[i, j]$  entry depends on only three other  $c$  table entries:  $c[i - 1, j - 1]$ ,  $c[i - 1, j]$ , and  $c[i, j - 1]$ . Given the value of  $c[i, j]$ , we can determine in  $O(1)$  time which of these three values was used to compute  $c[i, j]$ , without inspecting table  $b$ . Thus, we can reconstruct an LCS in  $O(m + n)$  time using a procedure similar to PRINT-LCS. (Exercise 15.4-2 asks you to give the pseudocode.) Although we save  $\Theta(mn)$  space by this method, the auxiliary space requirement for computing an LCS does not asymptotically decrease, since we need  $\Theta(mn)$  space for the  $c$  table anyway.

We can, however, reduce the asymptotic space requirements for LCS-LENGTH, since it needs only two rows of table  $c$  at a time: the row being computed and the previous row. (In fact, as Exercise 15.4-4 asks you to show, we can use only slightly more than the space for one row of  $c$  to compute the length of an LCS.) This improvement works if we need only the length of an LCS; if we need to reconstruct the elements of an LCS, the smaller table does not keep enough information to retrace our steps in  $O(m + n)$  time.

### Exercises

#### 15.4-1

Determine an LCS of  $\langle 1, 0, 0, 1, 0, 1, 0, 1 \rangle$  and  $\langle 0, 1, 0, 1, 1, 0, 1, 1, 0 \rangle$ .

#### 15.4-2

Give pseudocode to reconstruct an LCS from the completed  $c$  table and the original sequences  $X = \langle x_1, x_2, \dots, x_m \rangle$  and  $Y = \langle y_1, y_2, \dots, y_n \rangle$  in  $O(m + n)$  time, without using the  $b$  table.

#### 15.4-3

Give a memoized version of LCS-LENGTH that runs in  $O(mn)$  time.

#### 15.4-4

Show how to compute the length of an LCS using only  $2 \cdot \min(m, n)$  entries in the  $c$  table plus  $O(1)$  additional space. Then show how to do the same thing, but using  $\min(m, n)$  entries plus  $O(1)$  additional space.

**15.4-5**

Give an  $O(n^2)$ -time algorithm to find the longest monotonically increasing subsequence of a sequence of  $n$  numbers.

**15.4-6 ★**

Give an  $O(n \lg n)$ -time algorithm to find the longest monotonically increasing subsequence of a sequence of  $n$  numbers. (*Hint:* Observe that the last element of a candidate subsequence of length  $i$  is at least as large as the last element of a candidate subsequence of length  $i - 1$ . Maintain candidate subsequences by linking them through the input sequence.)

---

**15.5 Optimal binary search trees**

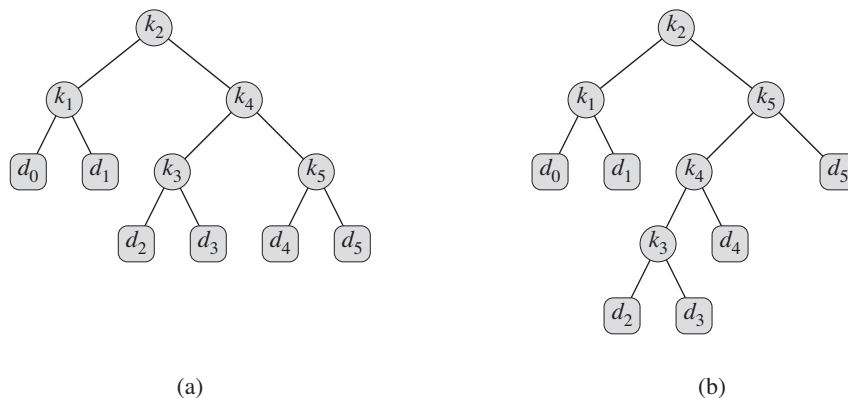
Suppose that we are designing a program to translate text from English to French. For each occurrence of each English word in the text, we need to look up its French equivalent. We could perform these lookup operations by building a binary search tree with  $n$  English words as keys and their French equivalents as satellite data. Because we will search the tree for each individual word in the text, we want the total time spent searching to be as low as possible. We could ensure an  $O(\lg n)$  search time per occurrence by using a red-black tree or any other balanced binary search tree. Words appear with different frequencies, however, and a frequently used word such as *the* may appear far from the root while a rarely used word such as *machicolation* appears near the root. Such an organization would slow down the translation, since the number of nodes visited when searching for a key in a binary search tree equals one plus the depth of the node containing the key. We want words that occur frequently in the text to be placed nearer the root.<sup>6</sup> Moreover, some words in the text might have no French translation,<sup>7</sup> and such words would not appear in the binary search tree at all. How do we organize a binary search tree so as to minimize the number of nodes visited in all searches, given that we know how often each word occurs?

What we need is known as an **optimal binary search tree**. Formally, we are given a sequence  $K = \langle k_1, k_2, \dots, k_n \rangle$  of  $n$  distinct keys in sorted order (so that  $k_1 < k_2 < \dots < k_n$ ), and we wish to build a binary search tree from these keys. For each key  $k_i$ , we have a probability  $p_i$  that a search will be for  $k_i$ . Some searches may be for values not in  $K$ , and so we also have  $n + 1$  “dummy keys”

---

<sup>6</sup>If the subject of the text is castle architecture, we might want *machicolation* to appear near the root.

<sup>7</sup>Yes, *machicolation* has a French counterpart: *mâchicoulis*.



**Figure 15.9** Two binary search trees for a set of  $n = 5$  keys with the following probabilities:

$i$	0	1	2	3	4	5
$p_i$		0.15	0.10	0.05	0.10	0.20
$q_i$	0.05	0.10	0.05	0.05	0.05	0.10

(a) A binary search tree with expected search cost 2.80. (b) A binary search tree with expected search cost 2.75. This tree is optimal.

$d_0, d_1, d_2, \dots, d_n$  representing values not in  $K$ . In particular,  $d_0$  represents all values less than  $k_1$ ,  $d_n$  represents all values greater than  $k_n$ , and for  $i = 1, 2, \dots, n-1$ , the dummy key  $d_i$  represents all values between  $k_i$  and  $k_{i+1}$ . For each dummy key  $d_i$ , we have a probability  $q_i$  that a search will correspond to  $d_i$ . Figure 15.9 shows two binary search trees for a set of  $n = 5$  keys. Each key  $k_i$  is an internal node, and each dummy key  $d_i$  is a leaf. Every search is either successful (finding some key  $k_i$ ) or unsuccessful (finding some dummy key  $d_i$ ), and so we have

$$\sum_{i=1}^n p_i + \sum_{i=0}^n q_i = 1. \quad (15.10)$$

Because we have probabilities of searches for each key and each dummy key, we can determine the expected cost of a search in a given binary search tree  $T$ . Let us assume that the actual cost of a search equals the number of nodes examined, i.e., the depth of the node found by the search in  $T$ , plus 1. Then the expected cost of a search in  $T$  is

$$\begin{aligned} E[\text{search cost in } T] &= \sum_{i=1}^n (\text{depth}_T(k_i) + 1) \cdot p_i + \sum_{i=0}^n (\text{depth}_T(d_i) + 1) \cdot q_i \\ &= 1 + \sum_{i=1}^n \text{depth}_T(k_i) \cdot p_i + \sum_{i=0}^n \text{depth}_T(d_i) \cdot q_i, \end{aligned} \quad (15.11)$$

where  $\text{depth}_T$  denotes a node's depth in the tree  $T$ . The last equality follows from equation (15.10). In Figure 15.9(a), we can calculate the expected search cost node by node:

node	depth	probability	contribution
$k_1$	1	0.15	0.30
$k_2$	0	0.10	0.10
$k_3$	2	0.05	0.15
$k_4$	1	0.10	0.20
$k_5$	2	0.20	0.60
$d_0$	2	0.05	0.15
$d_1$	2	0.10	0.30
$d_2$	3	0.05	0.20
$d_3$	3	0.05	0.20
$d_4$	3	0.05	0.20
$d_5$	3	0.10	0.40
Total			2.80

For a given set of probabilities, we wish to construct a binary search tree whose expected search cost is smallest. We call such a tree an **optimal binary search tree**. Figure 15.9(b) shows an optimal binary search tree for the probabilities given in the figure caption; its expected cost is 2.75. This example shows that an optimal binary search tree is not necessarily a tree whose overall height is smallest. Nor can we necessarily construct an optimal binary search tree by always putting the key with the greatest probability at the root. Here, key  $k_5$  has the greatest search probability of any key, yet the root of the optimal binary search tree shown is  $k_2$ . (The lowest expected cost of any binary search tree with  $k_5$  at the root is 2.85.)

As with matrix-chain multiplication, exhaustive checking of all possibilities fails to yield an efficient algorithm. We can label the nodes of any  $n$ -node binary tree with the keys  $k_1, k_2, \dots, k_n$  to construct a binary search tree, and then add in the dummy keys as leaves. In Problem 12-4, we saw that the number of binary trees with  $n$  nodes is  $\Omega(4^n/n^{3/2})$ , and so we would have to examine an exponential number of binary search trees in an exhaustive search. Not surprisingly, we shall solve this problem with dynamic programming.

### Step 1: The structure of an optimal binary search tree

To characterize the optimal substructure of optimal binary search trees, we start with an observation about subtrees. Consider any subtree of a binary search tree. It must contain keys in a contiguous range  $k_i, \dots, k_j$ , for some  $1 \leq i \leq j \leq n$ . In addition, a subtree that contains keys  $k_i, \dots, k_j$  must also have as its leaves the dummy keys  $d_{i-1}, \dots, d_j$ .

Now we can state the optimal substructure: if an optimal binary search tree  $T$  has a subtree  $T'$  containing keys  $k_i, \dots, k_j$ , then this subtree  $T'$  must be optimal as

well for the subproblem with keys  $k_i, \dots, k_j$  and dummy keys  $d_{i-1}, \dots, d_j$ . The usual cut-and-paste argument applies. If there were a subtree  $T''$  whose expected cost is lower than that of  $T'$ , then we could cut  $T'$  out of  $T$  and paste in  $T''$ , resulting in a binary search tree of lower expected cost than  $T$ , thus contradicting the optimality of  $T$ .

We need to use the optimal substructure to show that we can construct an optimal solution to the problem from optimal solutions to subproblems. Given keys  $k_i, \dots, k_j$ , one of these keys, say  $k_r$  ( $i \leq r \leq j$ ), is the root of an optimal subtree containing these keys. The left subtree of the root  $k_r$  contains the keys  $k_i, \dots, k_{r-1}$  (and dummy keys  $d_{i-1}, \dots, d_{r-1}$ ), and the right subtree contains the keys  $k_{r+1}, \dots, k_j$  (and dummy keys  $d_r, \dots, d_j$ ). As long as we examine all candidate roots  $k_r$ , where  $i \leq r \leq j$ , and we determine all optimal binary search trees containing  $k_i, \dots, k_{r-1}$  and those containing  $k_{r+1}, \dots, k_j$ , we are guaranteed that we will find an optimal binary search tree.

There is one detail worth noting about “empty” subtrees. Suppose that in a subtree with keys  $k_i, \dots, k_j$ , we select  $k_i$  as the root. By the above argument,  $k_i$ ’s left subtree contains the keys  $k_i, \dots, k_{i-1}$ . We interpret this sequence as containing no keys. Bear in mind, however, that subtrees also contain dummy keys. We adopt the convention that a subtree containing keys  $k_i, \dots, k_{i-1}$  has no actual keys but does contain the single dummy key  $d_{i-1}$ . Symmetrically, if we select  $k_j$  as the root, then  $k_j$ ’s right subtree contains the keys  $k_{j+1}, \dots, k_j$ ; this right subtree contains no actual keys, but it does contain the dummy key  $d_j$ .

## Step 2: A recursive solution

We are ready to define the value of an optimal solution recursively. We pick our subproblem domain as finding an optimal binary search tree containing the keys  $k_i, \dots, k_j$ , where  $i \geq 1$ ,  $j \leq n$ , and  $j \geq i - 1$ . (When  $j = i - 1$ , there are no actual keys; we have just the dummy key  $d_{i-1}$ .) Let us define  $e[i, j]$  as the expected cost of searching an optimal binary search tree containing the keys  $k_i, \dots, k_j$ . Ultimately, we wish to compute  $e[1, n]$ .

The easy case occurs when  $j = i - 1$ . Then we have just the dummy key  $d_{i-1}$ . The expected search cost is  $e[i, i - 1] = q_{i-1}$ .

When  $j \geq i$ , we need to select a root  $k_r$  from among  $k_i, \dots, k_j$  and then make an optimal binary search tree with keys  $k_i, \dots, k_{r-1}$  as its left subtree and an optimal binary search tree with keys  $k_{r+1}, \dots, k_j$  as its right subtree. What happens to the expected search cost of a subtree when it becomes a subtree of a node? The depth of each node in the subtree increases by 1. By equation (15.11), the expected search cost of this subtree increases by the sum of all the probabilities in the subtree. For a subtree with keys  $k_i, \dots, k_j$ , let us denote this sum of probabilities as

$$w(i, j) = \sum_{l=i}^j p_l + \sum_{l=i-1}^j q_l . \quad (15.12)$$

Thus, if  $k_r$  is the root of an optimal subtree containing keys  $k_i, \dots, k_j$ , we have

$$e[i, j] = p_r + (e[i, r-1] + w(i, r-1)) + (e[r+1, j] + w(r+1, j)) .$$

Noting that

$$w(i, j) = w(i, r-1) + p_r + w(r+1, j) ,$$

we rewrite  $e[i, j]$  as

$$e[i, j] = e[i, r-1] + e[r+1, j] + w(i, j) . \quad (15.13)$$

The recursive equation (15.13) assumes that we know which node  $k_r$  to use as the root. We choose the root that gives the lowest expected search cost, giving us our final recursive formulation:

$$e[i, j] = \begin{cases} q_{i-1} & \text{if } j = i - 1 , \\ \min_{i \leq r \leq j} \{e[i, r-1] + e[r+1, j] + w(i, j)\} & \text{if } i \leq j . \end{cases} \quad (15.14)$$

The  $e[i, j]$  values give the expected search costs in optimal binary search trees. To help us keep track of the structure of optimal binary search trees, we define  $root[i, j]$ , for  $1 \leq i \leq j \leq n$ , to be the index  $r$  for which  $k_r$  is the root of an optimal binary search tree containing keys  $k_i, \dots, k_j$ . Although we will see how to compute the values of  $root[i, j]$ , we leave the construction of an optimal binary search tree from these values as Exercise 15.5-1.

### Step 3: Computing the expected search cost of an optimal binary search tree

At this point, you may have noticed some similarities between our characterizations of optimal binary search trees and matrix-chain multiplication. For both problem domains, our subproblems consist of contiguous index subranges. A direct, recursive implementation of equation (15.14) would be as inefficient as a direct, recursive matrix-chain multiplication algorithm. Instead, we store the  $e[i, j]$  values in a table  $e[1 \dots n+1, 0 \dots n]$ . The first index needs to run to  $n+1$  rather than  $n$  because in order to have a subtree containing only the dummy key  $d_n$ , we need to compute and store  $e[n+1, n]$ . The second index needs to start from 0 because in order to have a subtree containing only the dummy key  $d_0$ , we need to compute and store  $e[1, 0]$ . We use only the entries  $e[i, j]$  for which  $j \geq i-1$ . We also use a table  $root[i, j]$ , for recording the root of the subtree containing keys  $k_i, \dots, k_j$ . This table uses only the entries for which  $1 \leq i \leq j \leq n$ .

We will need one other table for efficiency. Rather than compute the value of  $w(i, j)$  from scratch every time we are computing  $e[i, j]$ —which would take

$\Theta(j - i)$  additions—we store these values in a table  $w[1 \dots n + 1, 0 \dots n]$ . For the base case, we compute  $w[i, i - 1] = q_{i-1}$  for  $1 \leq i \leq n + 1$ . For  $j \geq i$ , we compute

$$w[i, j] = w[i, j - 1] + p_j + q_j. \quad (15.15)$$

Thus, we can compute the  $\Theta(n^2)$  values of  $w[i, j]$  in  $\Theta(1)$  time each.

The pseudocode that follows takes as inputs the probabilities  $p_1, \dots, p_n$  and  $q_0, \dots, q_n$  and the size  $n$ , and it returns the tables  $e$  and  $root$ .

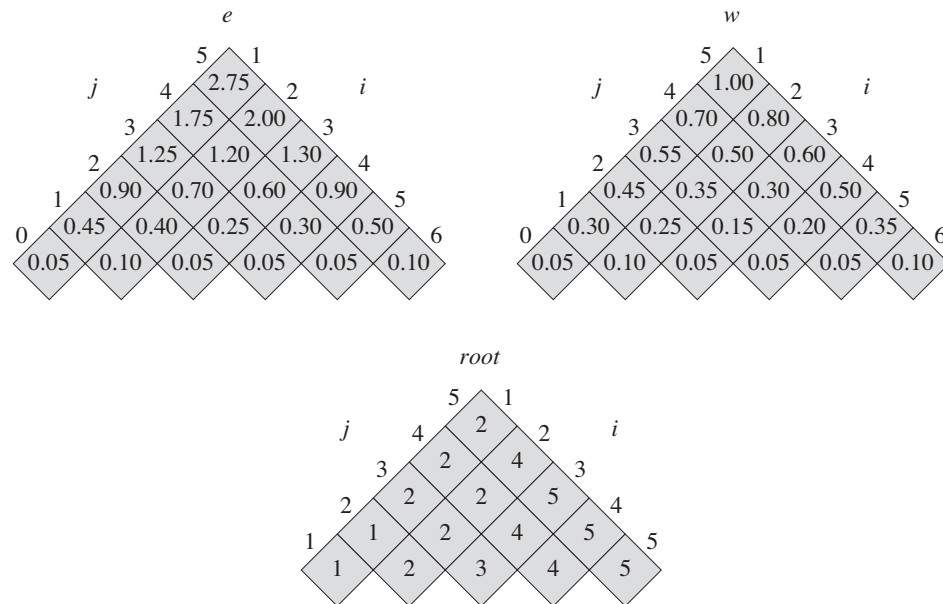
OPTIMAL-BST( $p, q, n$ )

```

1  let  $e[1 \dots n + 1, 0 \dots n]$ ,  $w[1 \dots n + 1, 0 \dots n]$ ,
   and  $root[1 \dots n, 1 \dots n]$  be new tables
2  for  $i = 1$  to  $n + 1$ 
3       $e[i, i - 1] = q_{i-1}$ 
4       $w[i, i - 1] = q_{i-1}$ 
5  for  $l = 1$  to  $n$ 
6      for  $i = 1$  to  $n - l + 1$ 
7           $j = i + l - 1$ 
8           $e[i, j] = \infty$ 
9           $w[i, j] = w[i, j - 1] + p_j + q_j$ 
10         for  $r = i$  to  $j$ 
11              $t = e[i, r - 1] + e[r + 1, j] + w[i, j]$ 
12             if  $t < e[i, j]$ 
13                  $e[i, j] = t$ 
14                  $root[i, j] = r$ 
15  return  $e$  and  $root$ 
```

From the description above and the similarity to the MATRIX-CHAIN-ORDER procedure in Section 15.2, you should find the operation of this procedure to be fairly straightforward. The **for** loop of lines 2–4 initializes the values of  $e[i, i - 1]$  and  $w[i, i - 1]$ . The **for** loop of lines 5–14 then uses the recurrences (15.14) and (15.15) to compute  $e[i, j]$  and  $w[i, j]$  for all  $1 \leq i \leq j \leq n$ . In the first iteration, when  $l = 1$ , the loop computes  $e[i, i]$  and  $w[i, i]$  for  $i = 1, 2, \dots, n$ . The second iteration, with  $l = 2$ , computes  $e[i, i + 1]$  and  $w[i, i + 1]$  for  $i = 1, 2, \dots, n - 1$ , and so forth. The innermost **for** loop, in lines 10–14, tries each candidate index  $r$  to determine which key  $k_r$  to use as the root of an optimal binary search tree containing keys  $k_i, \dots, k_j$ . This **for** loop saves the current value of the index  $r$  in  $root[i, j]$  whenever it finds a better key to use as the root.

Figure 15.10 shows the tables  $e[i, j]$ ,  $w[i, j]$ , and  $root[i, j]$  computed by the procedure OPTIMAL-BST on the key distribution shown in Figure 15.9. As in the matrix-chain multiplication example of Figure 15.5, the tables are rotated to make



**Figure 15.10** The tables  $e[i, j]$ ,  $w[i, j]$ , and  $root[i, j]$  computed by OPTIMAL-BST on the key distribution shown in Figure 15.9. The tables are rotated so that the diagonals run horizontally.

the diagonals run horizontally. OPTIMAL-BST computes the rows from bottom to top and from left to right within each row.

The OPTIMAL-BST procedure takes  $\Theta(n^3)$  time, just like MATRIX-CHAIN-ORDER. We can easily see that its running time is  $O(n^3)$ , since its **for** loops are nested three deep and each loop index takes on at most  $n$  values. The loop indices in OPTIMAL-BST do not have exactly the same bounds as those in MATRIX-CHAIN-ORDER, but they are within at most 1 in all directions. Thus, like MATRIX-CHAIN-ORDER, the OPTIMAL-BST procedure takes  $\Omega(n^3)$  time.

## Exercises

### 15.5-1

Write pseudocode for the procedure CONSTRUCT-OPTIMAL-BST( $root$ ) which, given the table  $root$ , outputs the structure of an optimal binary search tree. For the example in Figure 15.10, your procedure should print out the structure



$k_2$  is the root  
 $k_1$  is the left child of  $k_2$   
 $d_0$  is the left child of  $k_1$   
 $d_1$  is the right child of  $k_1$   
 $k_5$  is the right child of  $k_2$   
 $k_4$  is the left child of  $k_5$   
 $k_3$  is the left child of  $k_4$   
 $d_2$  is the left child of  $k_3$   
 $d_3$  is the right child of  $k_3$   
 $d_4$  is the right child of  $k_4$   
 $d_5$  is the right child of  $k_5$

corresponding to the optimal binary search tree shown in Figure 15.9(b).

### 15.5-2

Determine the cost and structure of an optimal binary search tree for a set of  $n = 7$  keys with the following probabilities:

$i$	0	1	2	3	4	5	6	7
$p_i$		0.04	0.06	0.08	0.02	0.10	0.12	0.14
$q_i$	0.06	0.06	0.06	0.06	0.05	0.05	0.05	0.05

### 15.5-3

Suppose that instead of maintaining the table  $w[i, j]$ , we computed the value of  $w(i, j)$  directly from equation (15.12) in line 9 of OPTIMAL-BST and used this computed value in line 11. How would this change affect the asymptotic running time of OPTIMAL-BST?

### 15.5-4 ★

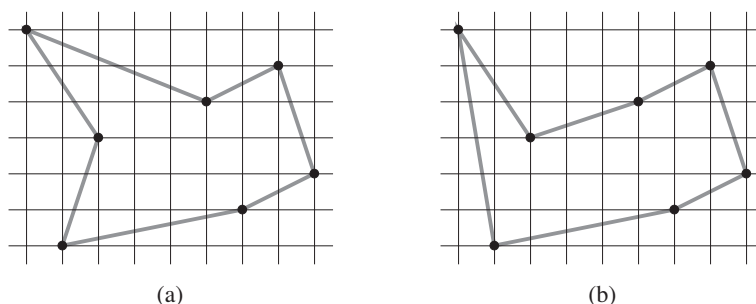
Knuth [212] has shown that there are always roots of optimal subtrees such that  $root[i, j - 1] \leq root[i, j] \leq root[i + 1, j]$  for all  $1 \leq i < j \leq n$ . Use this fact to modify the OPTIMAL-BST procedure to run in  $\Theta(n^2)$  time.

---

## Problems

### 15-1 Longest simple path in a directed acyclic graph

Suppose that we are given a directed acyclic graph  $G = (V, E)$  with real-valued edge weights and two distinguished vertices  $s$  and  $t$ . Describe a dynamic-programming approach for finding a longest weighted simple path from  $s$  to  $t$ . What does the subproblem graph look like? What is the efficiency of your algorithm?



**Figure 15.11** Seven points in the plane, shown on a unit grid. **(a)** The shortest closed tour, with length approximately 24.89. This tour is not bitonic. **(b)** The shortest bitonic tour for the same set of points. Its length is approximately 25.58.

### 15-2 Longest palindrome subsequence

A *palindrome* is a nonempty string over some alphabet that reads the same forward and backward. Examples of palindromes are all strings of length 1, `civic`, `racecar`, and `aibohphobia` (fear of palindromes).

Give an efficient algorithm to find the longest palindrome that is a subsequence of a given input string. For example, given the input `character`, your algorithm should return `carac`. What is the running time of your algorithm?

### 15-3 Bitonic euclidean traveling-salesman problem

In the *euclidean traveling-salesman problem*, we are given a set of  $n$  points in the plane, and we wish to find the shortest closed tour that connects all  $n$  points. Figure 15.11(a) shows the solution to a 7-point problem. The general problem is NP-hard, and its solution is therefore believed to require more than polynomial time (see Chapter 34).

J. L. Bentley has suggested that we simplify the problem by restricting our attention to *bitonic tours*, that is, tours that start at the leftmost point, go strictly rightward to the rightmost point, and then go strictly leftward back to the starting point. Figure 15.11(b) shows the shortest bitonic tour of the same 7 points. In this case, a polynomial-time algorithm is possible.

Describe an  $O(n^2)$ -time algorithm for determining an optimal bitonic tour. You may assume that no two points have the same  $x$ -coordinate and that all operations on real numbers take unit time. (*Hint*: Scan left to right, maintaining optimal possibilities for the two parts of the tour.)

### 15-4 Printing neatly

Consider the problem of neatly printing a paragraph with a monospaced font (all characters having the same width) on a printer. The input text is a sequence of  $n$

words of lengths  $l_1, l_2, \dots, l_n$ , measured in characters. We want to print this paragraph neatly on a number of lines that hold a maximum of  $M$  characters each. Our criterion of “neatness” is as follows. If a given line contains words  $i$  through  $j$ , where  $i \leq j$ , and we leave exactly one space between words, the number of extra space characters at the end of the line is  $M - j + i - \sum_{k=i}^j l_k$ , which must be nonnegative so that the words fit on the line. We wish to minimize the sum, over all lines except the last, of the cubes of the numbers of extra space characters at the ends of lines. Give a dynamic-programming algorithm to print a paragraph of  $n$  words neatly on a printer. Analyze the running time and space requirements of your algorithm.

### 15-5 Edit distance

In order to transform one source string of text  $x[1..m]$  to a target string  $y[1..n]$ , we can perform various transformation operations. Our goal is, given  $x$  and  $y$ , to produce a series of transformations that change  $x$  to  $y$ . We use an array  $z$ —assumed to be large enough to hold all the characters it will need—to hold the intermediate results. Initially,  $z$  is empty, and at termination, we should have  $z[j] = y[j]$  for  $j = 1, 2, \dots, n$ . We maintain current indices  $i$  into  $x$  and  $j$  into  $z$ , and the operations are allowed to alter  $z$  and these indices. Initially,  $i = j = 1$ . We are required to examine every character in  $x$  during the transformation, which means that at the end of the sequence of transformation operations, we must have  $i = m + 1$ .

We may choose from among six transformation operations:

**Copy** a character from  $x$  to  $z$  by setting  $z[j] = x[i]$  and then incrementing both  $i$  and  $j$ . This operation examines  $x[i]$ .

**Replace** a character from  $x$  by another character  $c$ , by setting  $z[j] = c$ , and then incrementing both  $i$  and  $j$ . This operation examines  $x[i]$ .

**Delete** a character from  $x$  by incrementing  $i$  but leaving  $j$  alone. This operation examines  $x[i]$ .

**Insert** the character  $c$  into  $z$  by setting  $z[j] = c$  and then incrementing  $j$ , but leaving  $i$  alone. This operation examines no characters of  $x$ .

**Twiddle** (i.e., exchange) the next two characters by copying them from  $x$  to  $z$  but in the opposite order; we do so by setting  $z[j] = x[i + 1]$  and  $z[j + 1] = x[i]$  and then setting  $i = i + 2$  and  $j = j + 2$ . This operation examines  $x[i]$  and  $x[i + 1]$ .

**Kill** the remainder of  $x$  by setting  $i = m + 1$ . This operation examines all characters in  $x$  that have not yet been examined. This operation, if performed, must be the final operation.

As an example, one way to transform the source string `algorithm` to the target string `altruistic` is to use the following sequence of operations, where the underlined characters are  $x[i]$  and  $z[j]$  after the operation:

Operation	$x$	$z$
<i>initial strings</i>	<u>a</u> lgorithm	—
copy	a <u>l</u> gorithm	a_
copy	al <u>g</u> orithm	al_
replace by t	alg <u>o</u> rithm	alt_
delete	algor <u>i</u> thm	alt_
copy	algor <u>i</u> thm	altr_
insert u	algori <u>h</u> m	altru_
insert i	algori <u>t</u> h	altrui_
insert s	algori <u>t</u> h	altruis_
twiddle	algorith <u>m</u>	altruisti_
insert c	algorith <u>m</u>	altruistic_
kill	algorithm_	altruistic_

Note that there are several other sequences of transformation operations that transform `algorithm` to `altruistic`.

Each of the transformation operations has an associated cost. The cost of an operation depends on the specific application, but we assume that each operation's cost is a constant that is known to us. We also assume that the individual costs of the copy and replace operations are less than the combined costs of the delete and insert operations; otherwise, the copy and replace operations would not be used. The cost of a given sequence of transformation operations is the sum of the costs of the individual operations in the sequence. For the sequence above, the cost of transforming `algorithm` to `altruistic` is

$$(3 \cdot \text{cost}(\text{copy})) + \text{cost}(\text{replace}) + \text{cost}(\text{delete}) + (4 \cdot \text{cost}(\text{insert})) \\ + \text{cost}(\text{twiddle}) + \text{cost}(\text{kill}) .$$

- a. Given two sequences  $x[1..m]$  and  $y[1..n]$  and set of transformation-operation costs, the **edit distance** from  $x$  to  $y$  is the cost of the least expensive operation sequence that transforms  $x$  to  $y$ . Describe a dynamic-programming algorithm that finds the edit distance from  $x[1..m]$  to  $y[1..n]$  and prints an optimal operation sequence. Analyze the running time and space requirements of your algorithm.

The edit-distance problem generalizes the problem of aligning two DNA sequences (see, for example, Setubal and Meidanis [310, Section 3.2]). There are several methods for measuring the similarity of two DNA sequences by aligning them. One such method to align two sequences  $x$  and  $y$  consists of inserting spaces at

arbitrary locations in the two sequences (including at either end) so that the resulting sequences  $x'$  and  $y'$  have the same length but do not have a space in the same position (i.e., for no position  $j$  are both  $x'[j]$  and  $y'[j]$  a space). Then we assign a “score” to each position. Position  $j$  receives a score as follows:

- $+1$  if  $x'[j] = y'[j]$  and neither is a space,
- $-1$  if  $x'[j] \neq y'[j]$  and neither is a space,
- $-2$  if either  $x'[j]$  or  $y'[j]$  is a space.

The score for the alignment is the sum of the scores of the individual positions. For example, given the sequences  $x = \text{GATCGGCAT}$  and  $y = \text{CAATGTGAATC}$ , one alignment is

```
G  ATCG  GCAT
CAAT GTGAATC
-*****-+++
```

A  $+$  under a position indicates a score of  $+1$  for that position, a  $-$  indicates a score of  $-1$ , and a  $*$  indicates a score of  $-2$ , so that this alignment has a total score of  $6 \cdot 1 - 2 \cdot 1 - 4 \cdot 2 = -4$ .

- b.* Explain how to cast the problem of finding an optimal alignment as an edit distance problem using a subset of the transformation operations copy, replace, delete, insert, twiddle, and kill.

### 15-6 Planning a company party

Professor Stewart is consulting for the president of a corporation that is planning a company party. The company has a hierarchical structure; that is, the supervisor relation forms a tree rooted at the president. The personnel office has ranked each employee with a conviviality rating, which is a real number. In order to make the party fun for all attendees, the president does not want both an employee and his or her immediate supervisor to attend.

Professor Stewart is given the tree that describes the structure of the corporation, using the left-child, right-sibling representation described in Section 10.4. Each node of the tree holds, in addition to the pointers, the name of an employee and that employee’s conviviality ranking. Describe an algorithm to make up a guest list that maximizes the sum of the conviviality ratings of the guests. Analyze the running time of your algorithm.

### 15-7 Viterbi algorithm

We can use dynamic programming on a directed graph  $G = (V, E)$  for speech recognition. Each edge  $(u, v) \in E$  is labeled with a sound  $\sigma(u, v)$  from a finite set  $\Sigma$  of sounds. The labeled graph is a formal model of a person speaking

a restricted language. Each path in the graph starting from a distinguished vertex  $v_0 \in V$  corresponds to a possible sequence of sounds produced by the model. We define the label of a directed path to be the concatenation of the labels of the edges on that path.

- a. Describe an efficient algorithm that, given an edge-labeled graph  $G$  with distinguished vertex  $v_0$  and a sequence  $s = \langle \sigma_1, \sigma_2, \dots, \sigma_k \rangle$  of sounds from  $\Sigma$ , returns a path in  $G$  that begins at  $v_0$  and has  $s$  as its label, if any such path exists. Otherwise, the algorithm should return NO-SUCH-PATH. Analyze the running time of your algorithm. (*Hint:* You may find concepts from Chapter 22 useful.)

Now, suppose that every edge  $(u, v) \in E$  has an associated nonnegative probability  $p(u, v)$  of traversing the edge  $(u, v)$  from vertex  $u$  and thus producing the corresponding sound. The sum of the probabilities of the edges leaving any vertex equals 1. The probability of a path is defined to be the product of the probabilities of its edges. We can view the probability of a path beginning at  $v_0$  as the probability that a “random walk” beginning at  $v_0$  will follow the specified path, where we randomly choose which edge to take leaving a vertex  $u$  according to the probabilities of the available edges leaving  $u$ .

- b. Extend your answer to part (a) so that if a path is returned, it is a *most probable path* starting at  $v_0$  and having label  $s$ . Analyze the running time of your algorithm.

### 15-8 Image compression by seam carving

We are given a color picture consisting of an  $m \times n$  array  $A[1..m, 1..n]$  of pixels, where each pixel specifies a triple of red, green, and blue (RGB) intensities. Suppose that we wish to compress this picture slightly. Specifically, we wish to remove one pixel from each of the  $m$  rows, so that the whole picture becomes one pixel narrower. To avoid disturbing visual effects, however, we require that the pixels removed in two adjacent rows be in the same or adjacent columns; the pixels removed form a “seam” from the top row to the bottom row where successive pixels in the seam are adjacent vertically or diagonally.

- a. Show that the number of such possible seams grows at least exponentially in  $m$ , assuming that  $n > 1$ .
- b. Suppose now that along with each pixel  $A[i, j]$ , we have calculated a real-valued disruption measure  $d[i, j]$ , indicating how disruptive it would be to remove pixel  $A[i, j]$ . Intuitively, the lower a pixel’s disruption measure, the more similar the pixel is to its neighbors. Suppose further that we define the disruption measure of a seam to be the sum of the disruption measures of its pixels.

Give an algorithm to find a seam with the lowest disruption measure. How efficient is your algorithm?

### 15-9 *Breaking a string*

A certain string-processing language allows a programmer to break a string into two pieces. Because this operation copies the string, it costs  $n$  time units to break a string of  $n$  characters into two pieces. Suppose a programmer wants to break a string into many pieces. The order in which the breaks occur can affect the total amount of time used. For example, suppose that the programmer wants to break a 20-character string after characters 2, 8, and 10 (numbering the characters in ascending order from the left-hand end, starting from 1). If she programs the breaks to occur in left-to-right order, then the first break costs 20 time units, the second break costs 18 time units (breaking the string from characters 3 to 20 at character 8), and the third break costs 12 time units, totaling 50 time units. If she programs the breaks to occur in right-to-left order, however, then the first break costs 20 time units, the second break costs 10 time units, and the third break costs 8 time units, totaling 38 time units. In yet another order, she could break first at 8 (costing 20), then break the left piece at 2 (costing 8), and finally the right piece at 10 (costing 12), for a total cost of 40.

Design an algorithm that, given the numbers of characters after which to break, determines a least-cost way to sequence those breaks. More formally, given a string  $S$  with  $n$  characters and an array  $L[1..m]$  containing the break points, compute the lowest cost for a sequence of breaks, along with a sequence of breaks that achieves this cost.

### 15-10 *Planning an investment strategy*

Your knowledge of algorithms helps you obtain an exciting job with the Acme Computer Company, along with a \$10,000 signing bonus. You decide to invest this money with the goal of maximizing your return at the end of 10 years. You decide to use the Amalgamated Investment Company to manage your investments. Amalgamated Investments requires you to observe the following rules. It offers  $n$  different investments, numbered 1 through  $n$ . In each year  $j$ , investment  $i$  provides a return rate of  $r_{ij}$ . In other words, if you invest  $d$  dollars in investment  $i$  in year  $j$ , then at the end of year  $j$ , you have  $dr_{ij}$  dollars. The return rates are guaranteed, that is, you are given all the return rates for the next 10 years for each investment. You make investment decisions only once per year. At the end of each year, you can leave the money made in the previous year in the same investments, or you can shift money to other investments, by either shifting money between existing investments or moving money to a new investment. If you do not move your money between two consecutive years, you pay a fee of  $f_1$  dollars, whereas if you switch your money, you pay a fee of  $f_2$  dollars, where  $f_2 > f_1$ .

- a. The problem, as stated, allows you to invest your money in multiple investments in each year. Prove that there exists an optimal investment strategy that, in each year, puts all the money into a single investment. (Recall that an optimal investment strategy maximizes the amount of money after 10 years and is not concerned with any other objectives, such as minimizing risk.)
- b. Prove that the problem of planning your optimal investment strategy exhibits optimal substructure.
- c. Design an algorithm that plans your optimal investment strategy. What is the running time of your algorithm?
- d. Suppose that Amalgamated Investments imposed the additional restriction that, at any point, you can have no more than \$15,000 in any one investment. Show that the problem of maximizing your income at the end of 10 years no longer exhibits optimal substructure.

### 15-11 Inventory planning

The Rinky Dink Company makes machines that resurface ice rinks. The demand for such products varies from month to month, and so the company needs to develop a strategy to plan its manufacturing given the fluctuating, but predictable, demand. The company wishes to design a plan for the next  $n$  months. For each month  $i$ , the company knows the demand  $d_i$ , that is, the number of machines that it will sell. Let  $D = \sum_{i=1}^n d_i$  be the total demand over the next  $n$  months. The company keeps a full-time staff who provide labor to manufacture up to  $m$  machines per month. If the company needs to make more than  $m$  machines in a given month, it can hire additional, part-time labor, at a cost that works out to  $c$  dollars per machine. Furthermore, if, at the end of a month, the company is holding any unsold machines, it must pay inventory costs. The cost for holding  $j$  machines is given as a function  $h(j)$  for  $j = 1, 2, \dots, D$ , where  $h(j) \geq 0$  for  $1 \leq j \leq D$  and  $h(j) \leq h(j+1)$  for  $1 \leq j \leq D-1$ .

Give an algorithm that calculates a plan for the company that minimizes its costs while fulfilling all the demand. The running time should be polyomial in  $n$  and  $D$ .

### 15-12 Signing free-agent baseball players

Suppose that you are the general manager for a major-league baseball team. During the off-season, you need to sign some free-agent players for your team. The team owner has given you a budget of  $\$X$  to spend on free agents. You are allowed to spend less than  $\$X$  altogether, but the owner will fire you if you spend any more than  $\$X$ .



You are considering  $N$  different positions, and for each position,  $P$  free-agent players who play that position are available.<sup>8</sup> Because you do not want to overload your roster with too many players at any position, for each position you may sign at most one free agent who plays that position. (If you do not sign any players at a particular position, then you plan to stick with the players you already have at that position.)

To determine how valuable a player is going to be, you decide to use a sabermetric statistic<sup>9</sup> known as “VORP,” or “value over replacement player.” A player with a higher VORP is more valuable than a player with a lower VORP. A player with a higher VORP is not necessarily more expensive to sign than a player with a lower VORP, because factors other than a player’s value determine how much it costs to sign him.

For each available free-agent player, you have three pieces of information:

- the player’s position,
- the amount of money it will cost to sign the player, and
- the player’s VORP.

Devise an algorithm that maximizes the total VORP of the players you sign while spending no more than  $\$X$  altogether. You may assume that each player signs for a multiple of \$100,000. Your algorithm should output the total VORP of the players you sign, the total amount of money you spend, and a list of which players you sign. Analyze the running time and space requirement of your algorithm.

---

## Chapter notes

R. Bellman began the systematic study of dynamic programming in 1955. The word “programming,” both here and in linear programming, refers to using a tabular solution method. Although optimization techniques incorporating elements of dynamic programming were known earlier, Bellman provided the area with a solid mathematical basis [37].

---

<sup>8</sup>Although there are nine positions on a baseball team,  $N$  is not necessarily equal to 9 because some general managers have particular ways of thinking about positions. For example, a general manager might consider right-handed pitchers and left-handed pitchers to be separate “positions,” as well as starting pitchers, long relief pitchers (relief pitchers who can pitch several innings), and short relief pitchers (relief pitchers who normally pitch at most only one inning).

<sup>9</sup>*Sabermetrics* is the application of statistical analysis to baseball records. It provides several ways to compare the relative values of individual players.

Galil and Park [125] classify dynamic-programming algorithms according to the size of the table and the number of other table entries each entry depends on. They call a dynamic-programming algorithm  $tD/eD$  if its table size is  $O(n^t)$  and each entry depends on  $O(n^e)$  other entries. For example, the matrix-chain multiplication algorithm in Section 15.2 would be  $2D/1D$ , and the longest-common-subsequence algorithm in Section 15.4 would be  $2D/0D$ .

Hu and Shing [182, 183] give an  $O(n \lg n)$ -time algorithm for the matrix-chain multiplication problem.

The  $O(mn)$ -time algorithm for the longest-common-subsequence problem appears to be a folk algorithm. Knuth [70] posed the question of whether subquadratic algorithms for the LCS problem exist. Masek and Paterson [244] answered this question in the affirmative by giving an algorithm that runs in  $O(mn/\lg n)$  time, where  $n \leq m$  and the sequences are drawn from a set of bounded size. For the special case in which no element appears more than once in an input sequence, Szymanski [326] shows how to solve the problem in  $O((n + m) \lg(n + m))$  time. Many of these results extend to the problem of computing string edit distances (Problem 15-5).

An early paper on variable-length binary encodings by Gilbert and Moore [133] had applications to constructing optimal binary search trees for the case in which all probabilities  $p_i$  are 0; this paper contains an  $O(n^3)$ -time algorithm. Aho, Hopcroft, and Ullman [5] present the algorithm from Section 15.5. Exercise 15.5-4 is due to Knuth [212]. Hu and Tucker [184] devised an algorithm for the case in which all probabilities  $p_i$  are 0 that uses  $O(n^2)$  time and  $O(n)$  space; subsequently, Knuth [211] reduced the time to  $O(n \lg n)$ .

Problem 15-8 is due to Avidan and Shamir [27], who have posted on the Web a wonderful video illustrating this image-compression technique.

---

## 16 Greedy Algorithms

Algorithms for optimization problems typically go through a sequence of steps, with a set of choices at each step. For many optimization problems, using dynamic programming to determine the best choices is overkill; simpler, more efficient algorithms will do. A *greedy algorithm* always makes the choice that looks best at the moment. That is, it makes a locally optimal choice in the hope that this choice will lead to a globally optimal solution. This chapter explores optimization problems for which greedy algorithms provide optimal solutions. Before reading this chapter, you should read about dynamic programming in Chapter 15, particularly Section 15.3.

Greedy algorithms do not always yield optimal solutions, but for many problems they do. We shall first examine, in Section 16.1, a simple but nontrivial problem, the activity-selection problem, for which a greedy algorithm efficiently computes an optimal solution. We shall arrive at the greedy algorithm by first considering a dynamic-programming approach and then showing that we can always make greedy choices to arrive at an optimal solution. Section 16.2 reviews the basic elements of the greedy approach, giving a direct approach for proving greedy algorithms correct. Section 16.3 presents an important application of greedy techniques: designing data-compression (Huffman) codes. In Section 16.4, we investigate some of the theory underlying combinatorial structures called “matroids,” for which a greedy algorithm always produces an optimal solution. Finally, Section 16.5 applies matroids to solve a problem of scheduling unit-time tasks with deadlines and penalties.

The greedy method is quite powerful and works well for a wide range of problems. Later chapters will present many algorithms that we can view as applications of the greedy method, including minimum-spanning-tree algorithms (Chapter 23), Dijkstra’s algorithm for shortest paths from a single source (Chapter 24), and Chvátal’s greedy set-covering heuristic (Chapter 35). Minimum-spanning-tree algorithms furnish a classic example of the greedy method. Although you can read

this chapter and Chapter 23 independently of each other, you might find it useful to read them together.

## 16.1 An activity-selection problem

Our first example is the problem of scheduling several competing activities that require exclusive use of a common resource, with a goal of selecting a maximum-size set of mutually compatible activities. Suppose we have a set  $S = \{a_1, a_2, \dots, a_n\}$  of  $n$  proposed **activities** that wish to use a resource, such as a lecture hall, which can serve only one activity at a time. Each activity  $a_i$  has a **start time**  $s_i$  and a **finish time**  $f_i$ , where  $0 \leq s_i < f_i < \infty$ . If selected, activity  $a_i$  takes place during the half-open time interval  $[s_i, f_i)$ . Activities  $a_i$  and  $a_j$  are **compatible** if the intervals  $[s_i, f_i)$  and  $[s_j, f_j)$  do not overlap. That is,  $a_i$  and  $a_j$  are compatible if  $s_i \geq f_j$  or  $s_j \geq f_i$ . In the **activity-selection problem**, we wish to select a maximum-size subset of mutually compatible activities. We assume that the activities are sorted in monotonically increasing order of finish time:

$$f_1 \leq f_2 \leq f_3 \leq \dots \leq f_{n-1} \leq f_n. \quad (16.1)$$

(We shall see later the advantage that this assumption provides.) For example, consider the following set  $S$  of activities:

$i$	1	2	3	4	5	6	7	8	9	10	11
$s_i$	1	3	0	5	3	5	6	8	8	2	12
$f_i$	4	5	6	7	9	9	10	11	12	14	16

For this example, the subset  $\{a_3, a_9, a_{11}\}$  consists of mutually compatible activities. It is not a maximum subset, however, since the subset  $\{a_1, a_4, a_8, a_{11}\}$  is larger. In fact,  $\{a_1, a_4, a_8, a_{11}\}$  is a largest subset of mutually compatible activities; another largest subset is  $\{a_2, a_4, a_9, a_{11}\}$ .

We shall solve this problem in several steps. We start by thinking about a dynamic-programming solution, in which we consider several choices when determining which subproblems to use in an optimal solution. We shall then observe that we need to consider only one choice—the greedy choice—and that when we make the greedy choice, only one subproblem remains. Based on these observations, we shall develop a recursive greedy algorithm to solve the activity-scheduling problem. We shall complete the process of developing a greedy solution by converting the recursive algorithm to an iterative one. Although the steps we shall go through in this section are slightly more involved than is typical when developing a greedy algorithm, they illustrate the relationship between greedy algorithms and dynamic programming.

### The optimal substructure of the activity-selection problem

We can easily verify that the activity-selection problem exhibits optimal substructure. Let us denote by  $S_{ij}$  the set of activities that start after activity  $a_i$  finishes and that finish before activity  $a_j$  starts. Suppose that we wish to find a maximum set of mutually compatible activities in  $S_{ij}$ , and suppose further that such a maximum set is  $A_{ij}$ , which includes some activity  $a_k$ . By including  $a_k$  in an optimal solution, we are left with two subproblems: finding mutually compatible activities in the set  $S_{ik}$  (activities that start after activity  $a_i$  finishes and that finish before activity  $a_k$  starts) and finding mutually compatible activities in the set  $S_{kj}$  (activities that start after activity  $a_k$  finishes and that finish before activity  $a_j$  starts). Let  $A_{ik} = A_{ij} \cap S_{ik}$  and  $A_{kj} = A_{ij} \cap S_{kj}$ , so that  $A_{ik}$  contains the activities in  $A_{ij}$  that finish before  $a_k$  starts and  $A_{kj}$  contains the activities in  $A_{ij}$  that start after  $a_k$  finishes. Thus, we have  $A_{ij} = A_{ik} \cup \{a_k\} \cup A_{kj}$ , and so the maximum-size set  $A_{ij}$  of mutually compatible activities in  $S_{ij}$  consists of  $|A_{ij}| = |A_{ik}| + |A_{kj}| + 1$  activities.

The usual cut-and-paste argument shows that the optimal solution  $A_{ij}$  must also include optimal solutions to the two subproblems for  $S_{ik}$  and  $S_{kj}$ . If we could find a set  $A'_{kj}$  of mutually compatible activities in  $S_{kj}$  where  $|A'_{kj}| > |A_{kj}|$ , then we could use  $A'_{kj}$ , rather than  $A_{kj}$ , in a solution to the subproblem for  $S_{ij}$ . We would have constructed a set of  $|A_{ik}| + |A'_{kj}| + 1 > |A_{ik}| + |A_{kj}| + 1 = |A_{ij}|$  mutually compatible activities, which contradicts the assumption that  $A_{ij}$  is an optimal solution. A symmetric argument applies to the activities in  $S_{ik}$ .

This way of characterizing optimal substructure suggests that we might solve the activity-selection problem by dynamic programming. If we denote the size of an optimal solution for the set  $S_{ij}$  by  $c[i, j]$ , then we would have the recurrence

$$c[i, j] = c[i, k] + c[k, j] + 1.$$

Of course, if we did not know that an optimal solution for the set  $S_{ij}$  includes activity  $a_k$ , we would have to examine all activities in  $S_{ij}$  to find which one to choose, so that

$$c[i, j] = \begin{cases} 0 & \text{if } S_{ij} = \emptyset, \\ \max_{a_k \in S_{ij}} \{c[i, k] + c[k, j] + 1\} & \text{if } S_{ij} \neq \emptyset. \end{cases} \quad (16.2)$$

We could then develop a recursive algorithm and memoize it, or we could work bottom-up and fill in table entries as we go along. But we would be overlooking another important characteristic of the activity-selection problem that we can use to great advantage.

### Making the greedy choice

What if we could choose an activity to add to our optimal solution without having to first solve all the subproblems? That could save us from having to consider all the choices inherent in recurrence (16.2). In fact, for the activity-selection problem, we need consider only one choice: the greedy choice.

What do we mean by the greedy choice for the activity-selection problem? Intuition suggests that we should choose an activity that leaves the resource available for as many other activities as possible. Now, of the activities we end up choosing, one of them must be the first one to finish. Our intuition tells us, therefore, to choose the activity in  $S$  with the earliest finish time, since that would leave the resource available for as many of the activities that follow it as possible. (If more than one activity in  $S$  has the earliest finish time, then we can choose any such activity.) In other words, since the activities are sorted in monotonically increasing order by finish time, the greedy choice is activity  $a_1$ . Choosing the first activity to finish is not the only way to think of making a greedy choice for this problem; Exercise 16.1-3 asks you to explore other possibilities.

If we make the greedy choice, we have only one remaining subproblem to solve: finding activities that start after  $a_1$  finishes. Why don't we have to consider activities that finish before  $a_1$  starts? We have that  $s_1 < f_1$ , and  $f_1$  is the earliest finish time of any activity, and therefore no activity can have a finish time less than or equal to  $s_1$ . Thus, all activities that are compatible with activity  $a_1$  must start after  $a_1$  finishes.

Furthermore, we have already established that the activity-selection problem exhibits optimal substructure. Let  $S_k = \{a_i \in S : s_i \geq f_k\}$  be the set of activities that start after activity  $a_k$  finishes. If we make the greedy choice of activity  $a_1$ , then  $S_1$  remains as the only subproblem to solve.<sup>1</sup> Optimal substructure tells us that if  $a_1$  is in the optimal solution, then an optimal solution to the original problem consists of activity  $a_1$  and all the activities in an optimal solution to the subproblem  $S_1$ .

One big question remains: is our intuition correct? Is the greedy choice—in which we choose the first activity to finish—always part of some optimal solution? The following theorem shows that it is.

---

<sup>1</sup>We sometimes refer to the sets  $S_k$  as subproblems rather than as just sets of activities. It will always be clear from the context whether we are referring to  $S_k$  as a set of activities or as a subproblem whose input is that set.

**Theorem 16.1**

Consider any nonempty subproblem  $S_k$ , and let  $a_m$  be an activity in  $S_k$  with the earliest finish time. Then  $a_m$  is included in some maximum-size subset of mutually compatible activities of  $S_k$ .

**Proof** Let  $A_k$  be a maximum-size subset of mutually compatible activities in  $S_k$ , and let  $a_j$  be the activity in  $A_k$  with the earliest finish time. If  $a_j = a_m$ , we are done, since we have shown that  $a_m$  is in some maximum-size subset of mutually compatible activities of  $S_k$ . If  $a_j \neq a_m$ , let the set  $A'_k = A_k - \{a_j\} \cup \{a_m\}$  be  $A_k$  but substituting  $a_m$  for  $a_j$ . The activities in  $A'_k$  are disjoint, which follows because the activities in  $A_k$  are disjoint,  $a_j$  is the first activity in  $A_k$  to finish, and  $f_m \leq f_j$ . Since  $|A'_k| = |A_k|$ , we conclude that  $A'_k$  is a maximum-size subset of mutually compatible activities of  $S_k$ , and it includes  $a_m$ . ■

Thus, we see that although we might be able to solve the activity-selection problem with dynamic programming, we don't need to. (Besides, we have not yet examined whether the activity-selection problem even has overlapping subproblems.) Instead, we can repeatedly choose the activity that finishes first, keep only the activities compatible with this activity, and repeat until no activities remain. Moreover, because we always choose the activity with the earliest finish time, the finish times of the activities we choose must strictly increase. We can consider each activity just once overall, in monotonically increasing order of finish times.

An algorithm to solve the activity-selection problem does not need to work bottom-up, like a table-based dynamic-programming algorithm. Instead, it can work top-down, choosing an activity to put into the optimal solution and then solving the subproblem of choosing activities from those that are compatible with those already chosen. Greedy algorithms typically have this top-down design: make a choice and then solve a subproblem, rather than the bottom-up technique of solving subproblems before making a choice.

**A recursive greedy algorithm**

Now that we have seen how to bypass the dynamic-programming approach and instead use a top-down, greedy algorithm, we can write a straightforward, recursive procedure to solve the activity-selection problem. The procedure `RECURSIVE-ACTIVITY-SELECTOR` takes the start and finish times of the activities, represented as arrays  $s$  and  $f$ ,<sup>2</sup> the index  $k$  that defines the subproblem  $S_k$  it is to solve, and

---

<sup>2</sup>Because the pseudocode takes  $s$  and  $f$  as arrays, it indexes into them with square brackets rather than subscripts.

the size  $n$  of the original problem. It returns a maximum-size set of mutually compatible activities in  $S_k$ . We assume that the  $n$  input activities are already ordered by monotonically increasing finish time, according to equation (16.1). If not, we can sort them into this order in  $O(n \lg n)$  time, breaking ties arbitrarily. In order to start, we add the fictitious activity  $a_0$  with  $f_0 = 0$ , so that subproblem  $S_0$  is the entire set of activities  $S$ . The initial call, which solves the entire problem, is `RECURSIVE-ACTIVITY-SELECTOR( $s, f, 0, n$ )`.

`RECURSIVE-ACTIVITY-SELECTOR( $s, f, k, n$ )`

```

1   $m = k + 1$ 
2  while  $m \leq n$  and  $s[m] < f[k]$       // find the first activity in  $S_k$  to finish
3       $m = m + 1$ 
4  if  $m \leq n$ 
5      return  $\{a_m\} \cup \text{RECURSIVE-ACTIVITY-SELECTOR}(s, f, m, n)$ 
6  else return  $\emptyset$ 
```

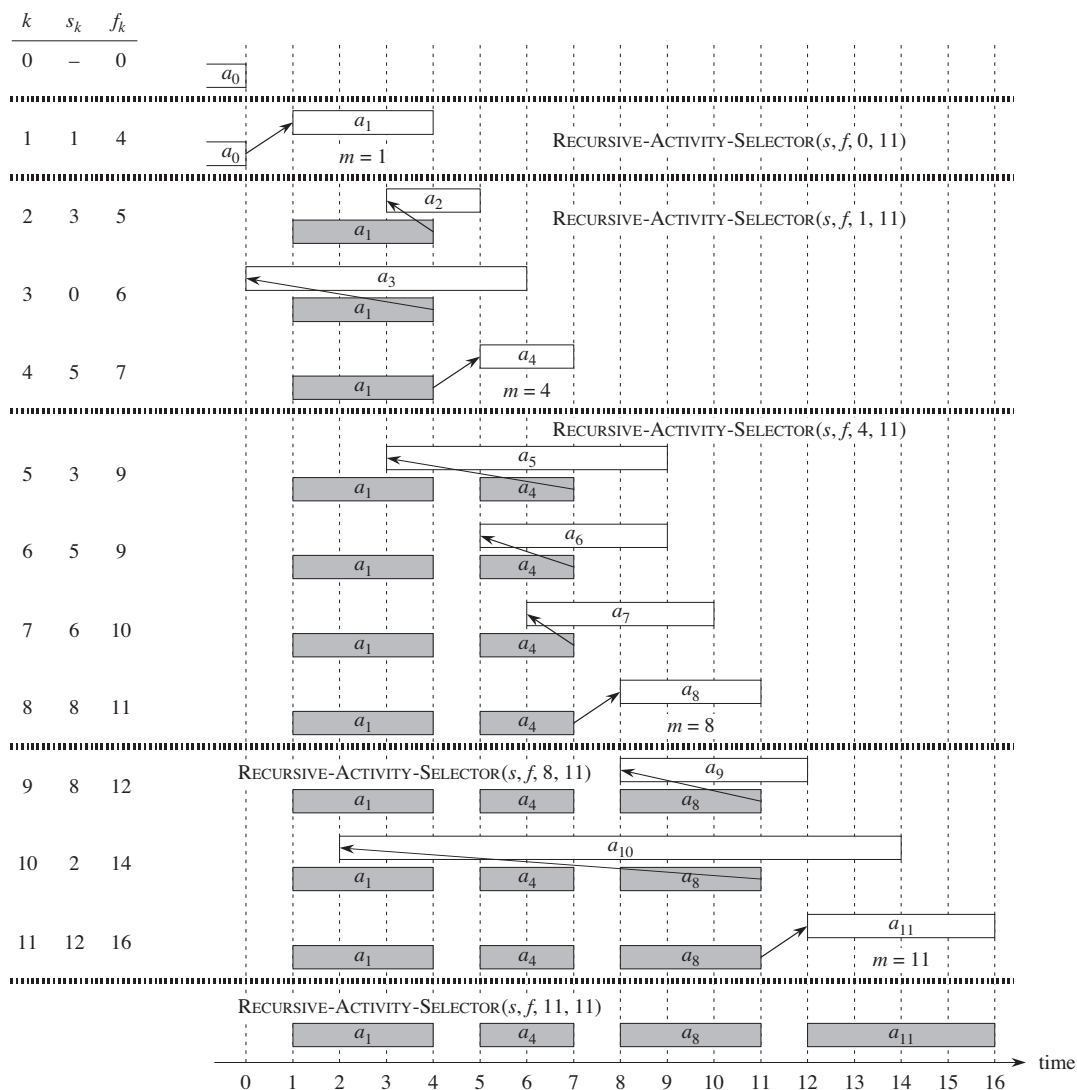
Figure 16.1 shows the operation of the algorithm. In a given recursive call `RECURSIVE-ACTIVITY-SELECTOR( $s, f, k, n$ )`, the **while** loop of lines 2–3 looks for the first activity in  $S_k$  to finish. The loop examines  $a_{k+1}, a_{k+2}, \dots, a_n$ , until it finds the first activity  $a_m$  that is compatible with  $a_k$ ; such an activity has  $s_m \geq f_k$ . If the loop terminates because it finds such an activity, line 5 returns the union of  $\{a_m\}$  and the maximum-size subset of  $S_m$  returned by the recursive call `RECURSIVE-ACTIVITY-SELECTOR( $s, f, m, n$ )`. Alternatively, the loop may terminate because  $m > n$ , in which case we have examined all activities in  $S_k$  without finding one that is compatible with  $a_k$ . In this case,  $S_k = \emptyset$ , and so the procedure returns  $\emptyset$  in line 6.

Assuming that the activities have already been sorted by finish times, the running time of the call `RECURSIVE-ACTIVITY-SELECTOR( $s, f, 0, n$ )` is  $\Theta(n)$ , which we can see as follows. Over all recursive calls, each activity is examined exactly once in the **while** loop test of line 2. In particular, activity  $a_i$  is examined in the last call made in which  $k < i$ .

### An iterative greedy algorithm

We easily can convert our recursive procedure to an iterative one. The procedure `RECURSIVE-ACTIVITY-SELECTOR` is almost “tail recursive” (see Problem 7-4): it ends with a recursive call to itself followed by a union operation. It is usually a straightforward task to transform a tail-recursive procedure to an iterative form; in fact, some compilers for certain programming languages perform this task automatically. As written, `RECURSIVE-ACTIVITY-SELECTOR` works for subproblems  $S_k$ , i.e., subproblems that consist of the last activities to finish.





**Figure 16.1** The operation of RECURSIVE-ACTIVITY-SELECTOR on the 11 activities given earlier. Activities considered in each recursive call appear between horizontal lines. The fictitious activity  $a_0$  finishes at time 0, and the initial call RECURSIVE-ACTIVITY-SELECTOR( $s, f, 0, 11$ ), selects activity  $a_1$ . In each recursive call, the activities that have already been selected are shaded, and the activity shown in white is being considered. If the starting time of an activity occurs before the finish time of the most recently added activity (the arrow between them points left), it is rejected. Otherwise (the arrow points directly up or to the right), it is selected. The last recursive call, RECURSIVE-ACTIVITY-SELECTOR( $s, f, 11, 11$ ), returns  $\emptyset$ . The resulting set of selected activities is  $\{a_1, a_4, a_8, a_{11}\}$ .

The procedure GREEDY-ACTIVITY-SELECTOR is an iterative version of the procedure RECURSIVE-ACTIVITY-SELECTOR. It also assumes that the input activities are ordered by monotonically increasing finish time. It collects selected activities into a set  $A$  and returns this set when it is done.

GREEDY-ACTIVITY-SELECTOR( $s, f$ )

```

1   $n = s.length$ 
2   $A = \{a_1\}$ 
3   $k = 1$ 
4  for  $m = 2$  to  $n$ 
5      if  $s[m] \geq f[k]$ 
6           $A = A \cup \{a_m\}$ 
7           $k = m$ 
8  return  $A$ 
```

The procedure works as follows. The variable  $k$  indexes the most recent addition to  $A$ , corresponding to the activity  $a_k$  in the recursive version. Since we consider the activities in order of monotonically increasing finish time,  $f_k$  is always the maximum finish time of any activity in  $A$ . That is,

$$f_k = \max \{f_i : a_i \in A\} . \quad (16.3)$$

Lines 2–3 select activity  $a_1$ , initialize  $A$  to contain just this activity, and initialize  $k$  to index this activity. The **for** loop of lines 4–7 finds the earliest activity in  $S_k$  to finish. The loop considers each activity  $a_m$  in turn and adds  $a_m$  to  $A$  if it is compatible with all previously selected activities; such an activity is the earliest in  $S_k$  to finish. To see whether activity  $a_m$  is compatible with every activity currently in  $A$ , it suffices by equation (16.3) to check (in line 5) that its start time  $s_m$  is not earlier than the finish time  $f_k$  of the activity most recently added to  $A$ . If activity  $a_m$  is compatible, then lines 6–7 add activity  $a_m$  to  $A$  and set  $k$  to  $m$ . The set  $A$  returned by the call GREEDY-ACTIVITY-SELECTOR( $s, f$ ) is precisely the set returned by the call RECURSIVE-ACTIVITY-SELECTOR( $s, f, 0, n$ ).

Like the recursive version, GREEDY-ACTIVITY-SELECTOR schedules a set of  $n$  activities in  $\Theta(n)$  time, assuming that the activities were already sorted initially by their finish times.

## Exercises

### 16.1-1

Give a dynamic-programming algorithm for the activity-selection problem, based on recurrence (16.2). Have your algorithm compute the sizes  $c[i, j]$  as defined above and also produce the maximum-size subset of mutually compatible activities.

Assume that the inputs have been sorted as in equation (16.1). Compare the running time of your solution to the running time of GREEDY-ACTIVITY-SELECTOR.

**16.1-2**

Suppose that instead of always selecting the first activity to finish, we instead select the last activity to start that is compatible with all previously selected activities. Describe how this approach is a greedy algorithm, and prove that it yields an optimal solution.

**16.1-3**

Not just any greedy approach to the activity-selection problem produces a maximum-size set of mutually compatible activities. Give an example to show that the approach of selecting the activity of least duration from among those that are compatible with previously selected activities does not work. Do the same for the approaches of always selecting the compatible activity that overlaps the fewest other remaining activities and always selecting the compatible remaining activity with the earliest start time.

**16.1-4**

Suppose that we have a set of activities to schedule among a large number of lecture halls, where any activity can take place in any lecture hall. We wish to schedule all the activities using as few lecture halls as possible. Give an efficient greedy algorithm to determine which activity should use which lecture hall.

(This problem is also known as the *interval-graph coloring problem*. We can create an interval graph whose vertices are the given activities and whose edges connect incompatible activities. The smallest number of colors required to color every vertex so that no two adjacent vertices have the same color corresponds to finding the fewest lecture halls needed to schedule all of the given activities.)

**16.1-5**

Consider a modification to the activity-selection problem in which each activity  $a_i$  has, in addition to a start and finish time, a value  $v_i$ . The objective is no longer to maximize the number of activities scheduled, but instead to maximize the total value of the activities scheduled. That is, we wish to choose a set  $A$  of compatible activities such that  $\sum_{a_k \in A} v_k$  is maximized. Give a polynomial-time algorithm for this problem.

---

## 16.2 Elements of the greedy strategy

A greedy algorithm obtains an optimal solution to a problem by making a sequence of choices. At each decision point, the algorithm makes choice that seems best at the moment. This heuristic strategy does not always produce an optimal solution, but as we saw in the activity-selection problem, sometimes it does. This section discusses some of the general properties of greedy methods.

The process that we followed in Section 16.1 to develop a greedy algorithm was a bit more involved than is typical. We went through the following steps:

1. Determine the optimal substructure of the problem.
2. Develop a recursive solution. (For the activity-selection problem, we formulated recurrence (16.2), but we bypassed developing a recursive algorithm based on this recurrence.)
3. Show that if we make the greedy choice, then only one subproblem remains.
4. Prove that it is always safe to make the greedy choice. (Steps 3 and 4 can occur in either order.)
5. Develop a recursive algorithm that implements the greedy strategy.
6. Convert the recursive algorithm to an iterative algorithm.

In going through these steps, we saw in great detail the dynamic-programming underpinnings of a greedy algorithm. For example, in the activity-selection problem, we first defined the subproblems  $S_{ij}$ , where both  $i$  and  $j$  varied. We then found that if we always made the greedy choice, we could restrict the subproblems to be of the form  $S_k$ .

Alternatively, we could have fashioned our optimal substructure with a greedy choice in mind, so that the choice leaves just one subproblem to solve. In the activity-selection problem, we could have started by dropping the second subscript and defining subproblems of the form  $S_k$ . Then, we could have proven that a greedy choice (the first activity  $a_m$  to finish in  $S_k$ ), combined with an optimal solution to the remaining set  $S_m$  of compatible activities, yields an optimal solution to  $S_k$ . More generally, we design greedy algorithms according to the following sequence of steps:

1. Cast the optimization problem as one in which we make a choice and are left with one subproblem to solve.
2. Prove that there is always an optimal solution to the original problem that makes the greedy choice, so that the greedy choice is always safe.

3. Demonstrate optimal substructure by showing that, having made the greedy choice, what remains is a subproblem with the property that if we combine an optimal solution to the subproblem with the greedy choice we have made, we arrive at an optimal solution to the original problem.

We shall use this more direct process in later sections of this chapter. Nevertheless, beneath every greedy algorithm, there is almost always a more cumbersome dynamic-programming solution.

How can we tell whether a greedy algorithm will solve a particular optimization problem? No way works all the time, but the greedy-choice property and optimal substructure are the two key ingredients. If we can demonstrate that the problem has these properties, then we are well on the way to developing a greedy algorithm for it.

### Greedy-choice property

The first key ingredient is the *greedy-choice property*: we can assemble a globally optimal solution by making locally optimal (greedy) choices. In other words, when we are considering which choice to make, we make the choice that looks best in the current problem, without considering results from subproblems.

Here is where greedy algorithms differ from dynamic programming. In dynamic programming, we make a choice at each step, but the choice usually depends on the solutions to subproblems. Consequently, we typically solve dynamic-programming problems in a bottom-up manner, progressing from smaller subproblems to larger subproblems. (Alternatively, we can solve them top down, but memoizing. Of course, even though the code works top down, we still must solve the subproblems before making a choice.) In a greedy algorithm, we make whatever choice seems best at the moment and then solve the subproblem that remains. The choice made by a greedy algorithm may depend on choices so far, but it cannot depend on any future choices or on the solutions to subproblems. Thus, unlike dynamic programming, which solves the subproblems before making the first choice, a greedy algorithm makes its first choice before solving any subproblems. A dynamic-programming algorithm proceeds bottom up, whereas a greedy strategy usually progresses in a top-down fashion, making one greedy choice after another, reducing each given problem instance to a smaller one.

Of course, we must prove that a greedy choice at each step yields a globally optimal solution. Typically, as in the case of Theorem 16.1, the proof examines a globally optimal solution to some subproblem. It then shows how to modify the solution to substitute the greedy choice for some other choice, resulting in one similar, but smaller, subproblem.

We can usually make the greedy choice more efficiently than when we have to consider a wider set of choices. For example, in the activity-selection problem, as-

suming that we had already sorted the activities in monotonically increasing order of finish times, we needed to examine each activity just once. By preprocessing the input or by using an appropriate data structure (often a priority queue), we often can make greedy choices quickly, thus yielding an efficient algorithm.

### Optimal substructure

A problem exhibits *optimal substructure* if an optimal solution to the problem contains within it optimal solutions to subproblems. This property is a key ingredient of assessing the applicability of dynamic programming as well as greedy algorithms. As an example of optimal substructure, recall how we demonstrated in Section 16.1 that if an optimal solution to subproblem  $S_{ij}$  includes an activity  $a_k$ , then it must also contain optimal solutions to the subproblems  $S_{ik}$  and  $S_{kj}$ . Given this optimal substructure, we argued that if we knew which activity to use as  $a_k$ , we could construct an optimal solution to  $S_{ij}$  by selecting  $a_k$  along with all activities in optimal solutions to the subproblems  $S_{ik}$  and  $S_{kj}$ . Based on this observation of optimal substructure, we were able to devise the recurrence (16.2) that described the value of an optimal solution.

We usually use a more direct approach regarding optimal substructure when applying it to greedy algorithms. As mentioned above, we have the luxury of assuming that we arrived at a subproblem by having made the greedy choice in the original problem. All we really need to do is argue that an optimal solution to the subproblem, combined with the greedy choice already made, yields an optimal solution to the original problem. This scheme implicitly uses induction on the subproblems to prove that making the greedy choice at every step produces an optimal solution.

### Greedy versus dynamic programming

Because both the greedy and dynamic-programming strategies exploit optimal substructure, you might be tempted to generate a dynamic-programming solution to a problem when a greedy solution suffices or, conversely, you might mistakenly think that a greedy solution works when in fact a dynamic-programming solution is required. To illustrate the subtleties between the two techniques, let us investigate two variants of a classical optimization problem.

The **0-1 knapsack problem** is the following. A thief robbing a store finds  $n$  items. The  $i$ th item is worth  $v_i$  dollars and weighs  $w_i$  pounds, where  $v_i$  and  $w_i$  are integers. The thief wants to take as valuable a load as possible, but he can carry at most  $W$  pounds in his knapsack, for some integer  $W$ . Which items should he take? (We call this the 0-1 knapsack problem because for each item, the thief must either

take it or leave it behind; he cannot take a fractional amount of an item or take an item more than once.)

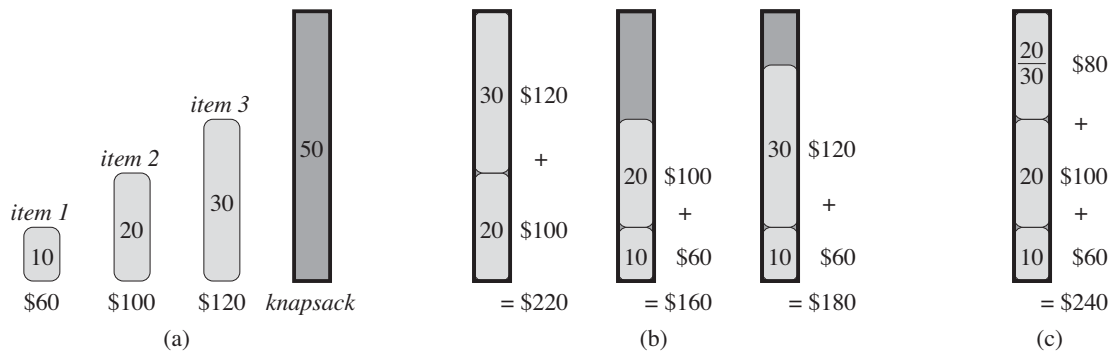
In the *fractional knapsack problem*, the setup is the same, but the thief can take fractions of items, rather than having to make a binary (0-1) choice for each item. You can think of an item in the 0-1 knapsack problem as being like a gold ingot and an item in the fractional knapsack problem as more like gold dust.

Both knapsack problems exhibit the optimal-substructure property. For the 0-1 problem, consider the most valuable load that weighs at most  $W$  pounds. If we remove item  $j$  from this load, the remaining load must be the most valuable load weighing at most  $W - w_j$  that the thief can take from the  $n - 1$  original items excluding  $j$ . For the comparable fractional problem, consider that if we remove a weight  $w$  of one item  $j$  from the optimal load, the remaining load must be the most valuable load weighing at most  $W - w$  that the thief can take from the  $n - 1$  original items plus  $w_j - w$  pounds of item  $j$ .

Although the problems are similar, we can solve the fractional knapsack problem by a greedy strategy, but we cannot solve the 0-1 problem by such a strategy. To solve the fractional problem, we first compute the value per pound  $v_i/w_i$  for each item. Obeying a greedy strategy, the thief begins by taking as much as possible of the item with the greatest value per pound. If the supply of that item is exhausted and he can still carry more, he takes as much as possible of the item with the next greatest value per pound, and so forth, until he reaches his weight limit  $W$ . Thus, by sorting the items by value per pound, the greedy algorithm runs in  $O(n \lg n)$  time. We leave the proof that the fractional knapsack problem has the greedy-choice property as Exercise 16.2-1.

To see that this greedy strategy does not work for the 0-1 knapsack problem, consider the problem instance illustrated in Figure 16.2(a). This example has 3 items and a knapsack that can hold 50 pounds. Item 1 weighs 10 pounds and is worth 60 dollars. Item 2 weighs 20 pounds and is worth 100 dollars. Item 3 weighs 30 pounds and is worth 120 dollars. Thus, the value per pound of item 1 is 6 dollars per pound, which is greater than the value per pound of either item 2 (5 dollars per pound) or item 3 (4 dollars per pound). The greedy strategy, therefore, would take item 1 first. As you can see from the case analysis in Figure 16.2(b), however, the optimal solution takes items 2 and 3, leaving item 1 behind. The two possible solutions that take item 1 are both suboptimal.

For the comparable fractional problem, however, the greedy strategy, which takes item 1 first, does yield an optimal solution, as shown in Figure 16.2(c). Taking item 1 doesn't work in the 0-1 problem because the thief is unable to fill his knapsack to capacity, and the empty space lowers the effective value per pound of his load. In the 0-1 problem, when we consider whether to include an item in the knapsack, we must compare the solution to the subproblem that includes the item with the solution to the subproblem that excludes the item before we can make the



**Figure 16.2** An example showing that the greedy strategy does not work for the 0-1 knapsack problem. (a) The thief must select a subset of the three items shown whose weight must not exceed 50 pounds. (b) The optimal subset includes items 2 and 3. Any solution with item 1 is suboptimal, even though item 1 has the greatest value per pound. (c) For the fractional knapsack problem, taking the items in order of greatest value per pound yields an optimal solution.

choice. The problem formulated in this way gives rise to many overlapping sub-problems—a hallmark of dynamic programming, and indeed, as Exercise 16.2-2 asks you to show, we can use dynamic programming to solve the 0-1 problem.

## Exercises

### 16.2-1

Prove that the fractional knapsack problem has the greedy-choice property.

### 16.2-2

Give a dynamic-programming solution to the 0-1 knapsack problem that runs in  $O(nW)$  time, where  $n$  is the number of items and  $W$  is the maximum weight of items that the thief can put in his knapsack.

### 16.2-3

Suppose that in a 0-1 knapsack problem, the order of the items when sorted by increasing weight is the same as their order when sorted by decreasing value. Give an efficient algorithm to find an optimal solution to this variant of the knapsack problem, and argue that your algorithm is correct.

### 16.2-4

Professor Gekko has always dreamed of inline skating across North Dakota. He plans to cross the state on highway U.S. 2, which runs from Grand Forks, on the eastern border with Minnesota, to Williston, near the western border with Montana.



The professor can carry two liters of water, and he can skate  $m$  miles before running out of water. (Because North Dakota is relatively flat, the professor does not have to worry about drinking water at a greater rate on uphill sections than on flat or downhill sections.) The professor will start in Grand Forks with two full liters of water. His official North Dakota state map shows all the places along U.S. 2 at which he can refill his water and the distances between these locations.

The professor's goal is to minimize the number of water stops along his route across the state. Give an efficient method by which he can determine which water stops he should make. Prove that your strategy yields an optimal solution, and give its running time.

### 16.2-5

Describe an efficient algorithm that, given a set  $\{x_1, x_2, \dots, x_n\}$  of points on the real line, determines the smallest set of unit-length closed intervals that contains all of the given points. Argue that your algorithm is correct.

### 16.2-6 ★

Show how to solve the fractional knapsack problem in  $O(n)$  time.

### 16.2-7

Suppose you are given two sets  $A$  and  $B$ , each containing  $n$  positive integers. You can choose to reorder each set however you like. After reordering, let  $a_i$  be the  $i$ th element of set  $A$ , and let  $b_i$  be the  $i$ th element of set  $B$ . You then receive a payoff of  $\prod_{i=1}^n a_i^{b_i}$ . Give an algorithm that will maximize your payoff. Prove that your algorithm maximizes the payoff, and state its running time.

---

## 16.3 Huffman codes

Huffman codes compress data very effectively: savings of 20% to 90% are typical, depending on the characteristics of the data being compressed. We consider the data to be a sequence of characters. Huffman's greedy algorithm uses a table giving how often each character occurs (i.e., its frequency) to build up an optimal way of representing each character as a binary string.

Suppose we have a 100,000-character data file that we wish to store compactly. We observe that the characters in the file occur with the frequencies given by Figure 16.3. That is, only 6 different characters appear, and the character **a** occurs 45,000 times.

We have many options for how to represent such a file of information. Here, we consider the problem of designing a *binary character code* (or *code* for short)

	a	b	c	d	e	f
Frequency (in thousands)	45	13	12	16	9	5
Fixed-length codeword	000	001	010	011	100	101
Variable-length codeword	0	101	100	111	1101	1100

**Figure 16.3** A character-coding problem. A data file of 100,000 characters contains only the characters a–f, with the frequencies indicated. If we assign each character a 3-bit codeword, we can encode the file in 300,000 bits. Using the variable-length code shown, we can encode the file in only 224,000 bits.

in which each character is represented by a unique binary string, which we call a **codeword**. If we use a **fixed-length code**, we need 3 bits to represent 6 characters:  $a = 000$ ,  $b = 001$ , ...,  $f = 101$ . This method requires 300,000 bits to code the entire file. Can we do better?

A **variable-length code** can do considerably better than a fixed-length code, by giving frequent characters short codewords and infrequent characters long codewords. Figure 16.3 shows such a code; here the 1-bit string 0 represents  $a$ , and the 4-bit string 1100 represents  $f$ . This code requires

$$(45 \cdot 1 + 13 \cdot 3 + 12 \cdot 3 + 16 \cdot 3 + 9 \cdot 4 + 5 \cdot 4) \cdot 1,000 = 224,000 \text{ bits}$$

to represent the file, a savings of approximately 25%. In fact, this is an optimal character code for this file, as we shall see.

### Prefix codes

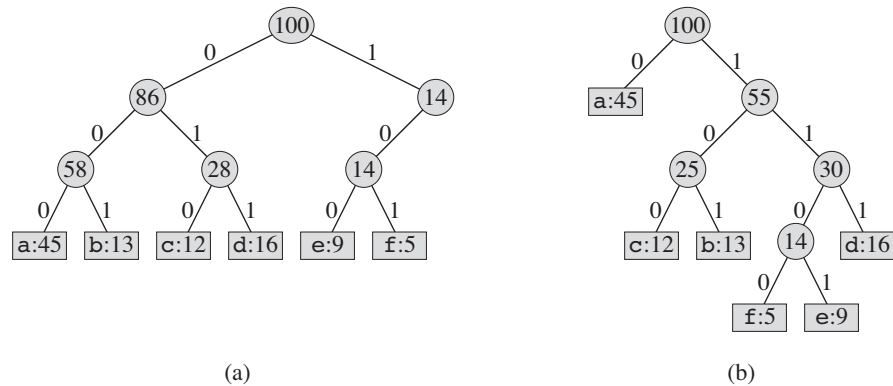
We consider here only codes in which no codeword is also a prefix of some other codeword. Such codes are called **prefix codes**.<sup>3</sup> Although we won't prove it here, a prefix code can always achieve the optimal data compression among any character code, and so we suffer no loss of generality by restricting our attention to prefix codes.

Encoding is always simple for any binary character code; we just concatenate the codewords representing each character of the file. For example, with the variable-length prefix code of Figure 16.3, we code the 3-character file  $abc$  as  $0 \cdot 101 \cdot 100 = 0101100$ , where “ $\cdot$ ” denotes concatenation.

Prefix codes are desirable because they simplify decoding. Since no codeword is a prefix of any other, the codeword that begins an encoded file is unambiguous. We can simply identify the initial codeword, translate it back to the original char-

---

<sup>3</sup>Perhaps “prefix-free codes” would be a better name, but the term “prefix codes” is standard in the literature.



**Figure 16.4** Trees corresponding to the coding schemes in Figure 16.3. Each leaf is labeled with a character and its frequency of occurrence. Each internal node is labeled with the sum of the frequencies of the leaves in its subtree. (a) The tree corresponding to the fixed-length code  $a = 000, \dots, f = 101$ . (b) The tree corresponding to the optimal prefix code  $a = 0, b = 101, \dots, f = 1100$ .

acter, and repeat the decoding process on the remainder of the encoded file. In our example, the string 001011101 parses uniquely as  $0 \cdot 0 \cdot 101 \cdot 1101$ , which decodes to **aabe**.

The decoding process needs a convenient representation for the prefix code so that we can easily pick off the initial codeword. A binary tree whose leaves are the given characters provides one such representation. We interpret the binary codeword for a character as the simple path from the root to that character, where 0 means “go to the left child” and 1 means “go to the right child.” Figure 16.4 shows the trees for the two codes of our example. Note that these are not binary search trees, since the leaves need not appear in sorted order and internal nodes do not contain character keys.

An optimal code for a file is always represented by a *full* binary tree, in which every nonleaf node has two children (see Exercise 16.3-2). The fixed-length code in our example is not optimal since its tree, shown in Figure 16.4(a), is not a full binary tree: it contains codewords beginning 10..., but none beginning 11.... Since we can now restrict our attention to full binary trees, we can say that if  $C$  is the alphabet from which the characters are drawn and all character frequencies are positive, then the tree for an optimal prefix code has exactly  $|C|$  leaves, one for each letter of the alphabet, and exactly  $|C| - 1$  internal nodes (see Exercise B.5-3).

Given a tree  $T$  corresponding to a prefix code, we can easily compute the number of bits required to encode a file. For each character  $c$  in the alphabet  $C$ , let the attribute  $c.freq$  denote the frequency of  $c$  in the file and let  $d_T(c)$  denote the depth

of  $c$ 's leaf in the tree. Note that  $d_T(c)$  is also the length of the codeword for character  $c$ . The number of bits required to encode a file is thus

$$B(T) = \sum_{c \in C} c.freq \cdot d_T(c) , \quad (16.4)$$

which we define as the *cost* of the tree  $T$ .

### Constructing a Huffman code

Huffman invented a greedy algorithm that constructs an optimal prefix code called a **Huffman code**. In line with our observations in Section 16.2, its proof of correctness relies on the greedy-choice property and optimal substructure. Rather than demonstrating that these properties hold and then developing pseudocode, we present the pseudocode first. Doing so will help clarify how the algorithm makes greedy choices.

In the pseudocode that follows, we assume that  $C$  is a set of  $n$  characters and that each character  $c \in C$  is an object with an attribute  $c.freq$  giving its frequency. The algorithm builds the tree  $T$  corresponding to the optimal code in a bottom-up manner. It begins with a set of  $|C|$  leaves and performs a sequence of  $|C| - 1$  “merging” operations to create the final tree. The algorithm uses a min-priority queue  $Q$ , keyed on the *freq* attribute, to identify the two least-frequent objects to merge together. When we merge two objects, the result is a new object whose frequency is the sum of the frequencies of the two objects that were merged.

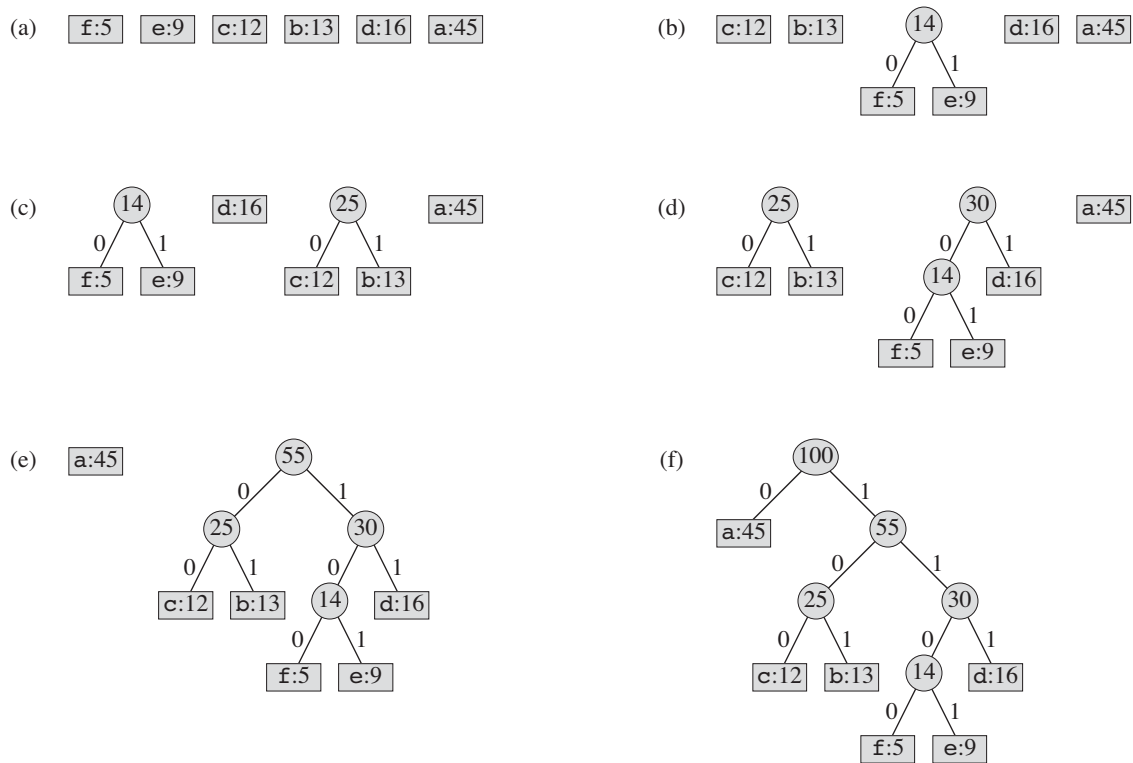
HUFFMAN( $C$ )

```

1   $n = |C|$ 
2   $Q = C$ 
3  for  $i = 1$  to  $n - 1$ 
4      allocate a new node  $z$ 
5       $z.left = x = \text{EXTRACT-MIN}(Q)$ 
6       $z.right = y = \text{EXTRACT-MIN}(Q)$ 
7       $z.freq = x.freq + y.freq$ 
8       $\text{INSERT}(Q, z)$ 
9  return  $\text{EXTRACT-MIN}(Q)$     // return the root of the tree
```

For our example, Huffman's algorithm proceeds as shown in Figure 16.5. Since the alphabet contains 6 letters, the initial queue size is  $n = 6$ , and 5 merge steps build the tree. The final tree represents the optimal prefix code. The codeword for a letter is the sequence of edge labels on the simple path from the root to the letter.

Line 2 initializes the min-priority queue  $Q$  with the characters in  $C$ . The **for** loop in lines 3–8 repeatedly extracts the two nodes  $x$  and  $y$  of lowest frequency



**Figure 16.5** The steps of Huffman's algorithm for the frequencies given in Figure 16.3. Each part shows the contents of the queue sorted into increasing order by frequency. At each step, the two trees with lowest frequencies are merged. Leaves are shown as rectangles containing a character and its frequency. Internal nodes are shown as circles containing the sum of the frequencies of their children. An edge connecting an internal node with its children is labeled 0 if it is an edge to a left child and 1 if it is an edge to a right child. The codeword for a letter is the sequence of labels on the edges connecting the root to the leaf for that letter. (a) The initial set of  $n = 6$  nodes, one for each letter. (b)–(e) Intermediate stages. (f) The final tree.

from the queue, replacing them in the queue with a new node  $z$  representing their merger. The frequency of  $z$  is computed as the sum of the frequencies of  $x$  and  $y$  in line 7. The node  $z$  has  $x$  as its left child and  $y$  as its right child. (This order is arbitrary; switching the left and right child of any node yields a different code of the same cost.) After  $n - 1$  mergers, line 9 returns the one node left in the queue, which is the root of the code tree.

Although the algorithm would produce the same result if we were to excise the variables  $x$  and  $y$ —assigning directly to  $z.left$  and  $z.right$  in lines 5 and 6, and changing line 7 to  $z.freq = z.left.freq + z.right.freq$ —we shall use the node

names  $x$  and  $y$  in the proof of correctness. Therefore, we find it convenient to leave them in.

To analyze the running time of Huffman's algorithm, we assume that  $Q$  is implemented as a binary min-heap (see Chapter 6). For a set  $C$  of  $n$  characters, we can initialize  $Q$  in line 2 in  $O(n)$  time using the BUILD-MIN-HEAP procedure discussed in Section 6.3. The **for** loop in lines 3–8 executes exactly  $n - 1$  times, and since each heap operation requires time  $O(\lg n)$ , the loop contributes  $O(n \lg n)$  to the running time. Thus, the total running time of HUFFMAN on a set of  $n$  characters is  $O(n \lg n)$ . We can reduce the running time to  $O(n \lg \lg n)$  by replacing the binary min-heap with a van Emde Boas tree (see Chapter 20).

### Correctness of Huffman's algorithm

To prove that the greedy algorithm HUFFMAN is correct, we show that the problem of determining an optimal prefix code exhibits the greedy-choice and optimal-substructure properties. The next lemma shows that the greedy-choice property holds.

#### Lemma 16.2

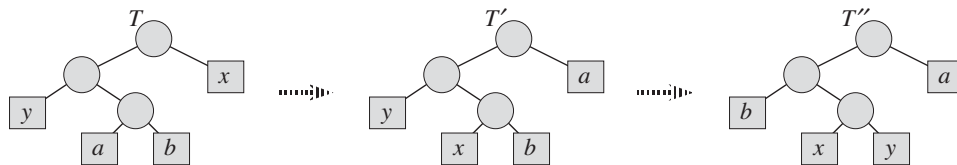
Let  $C$  be an alphabet in which each character  $c \in C$  has frequency  $c.freq$ . Let  $x$  and  $y$  be two characters in  $C$  having the lowest frequencies. Then there exists an optimal prefix code for  $C$  in which the codewords for  $x$  and  $y$  have the same length and differ only in the last bit.

**Proof** The idea of the proof is to take the tree  $T$  representing an arbitrary optimal prefix code and modify it to make a tree representing another optimal prefix code such that the characters  $x$  and  $y$  appear as sibling leaves of maximum depth in the new tree. If we can construct such a tree, then the codewords for  $x$  and  $y$  will have the same length and differ only in the last bit.

Let  $a$  and  $b$  be two characters that are sibling leaves of maximum depth in  $T$ . Without loss of generality, we assume that  $a.freq \leq b.freq$  and  $x.freq \leq y.freq$ . Since  $x.freq$  and  $y.freq$  are the two lowest leaf frequencies, in order, and  $a.freq$  and  $b.freq$  are two arbitrary frequencies, in order, we have  $x.freq \leq a.freq$  and  $y.freq \leq b.freq$ .

In the remainder of the proof, it is possible that we could have  $x.freq = a.freq$  or  $y.freq = b.freq$ . However, if we had  $x.freq = b.freq$ , then we would also have  $a.freq = b.freq = x.freq = y.freq$  (see Exercise 16.3-1), and the lemma would be trivially true. Thus, we will assume that  $x.freq \neq b.freq$ , which means that  $x \neq b$ .

As Figure 16.6 shows, we exchange the positions in  $T$  of  $a$  and  $x$  to produce a tree  $T'$ , and then we exchange the positions in  $T'$  of  $b$  and  $y$  to produce a tree  $T''$



**Figure 16.6** An illustration of the key step in the proof of Lemma 16.2. In the optimal tree  $T$ , leaves  $a$  and  $b$  are two siblings of maximum depth. Leaves  $x$  and  $y$  are the two characters with the lowest frequencies; they appear in arbitrary positions in  $T$ . Assuming that  $x \neq b$ , swapping leaves  $a$  and  $x$  produces tree  $T'$ , and then swapping leaves  $b$  and  $y$  produces tree  $T''$ . Since each swap does not increase the cost, the resulting tree  $T''$  is also an optimal tree.

in which  $x$  and  $y$  are sibling leaves of maximum depth. (Note that if  $x = b$  but  $y \neq a$ , then tree  $T''$  does not have  $x$  and  $y$  as sibling leaves of maximum depth. Because we assume that  $x \neq b$ , this situation cannot occur.) By equation (16.4), the difference in cost between  $T$  and  $T'$  is

$$\begin{aligned}
 B(T) - B(T') &= \sum_{c \in C} c.\text{freq} \cdot d_T(c) - \sum_{c \in C} c.\text{freq} \cdot d_{T'}(c) \\
 &= x.\text{freq} \cdot d_T(x) + a.\text{freq} \cdot d_T(a) - x.\text{freq} \cdot d_{T'}(x) - a.\text{freq} \cdot d_{T'}(a) \\
 &= x.\text{freq} \cdot d_T(x) + a.\text{freq} \cdot d_T(a) - x.\text{freq} \cdot d_T(a) - a.\text{freq} \cdot d_T(x) \\
 &= (a.\text{freq} - x.\text{freq})(d_T(a) - d_T(x)) \\
 &\geq 0,
 \end{aligned}$$

because both  $a.\text{freq} - x.\text{freq}$  and  $d_T(a) - d_T(x)$  are nonnegative. More specifically,  $a.\text{freq} - x.\text{freq}$  is nonnegative because  $x$  is a minimum-frequency leaf, and  $d_T(a) - d_T(x)$  is nonnegative because  $a$  is a leaf of maximum depth in  $T$ . Similarly, exchanging  $y$  and  $b$  does not increase the cost, and so  $B(T') - B(T'')$  is nonnegative. Therefore,  $B(T'') \leq B(T)$ , and since  $T$  is optimal, we have  $B(T) \leq B(T'')$ , which implies  $B(T'') = B(T)$ . Thus,  $T''$  is an optimal tree in which  $x$  and  $y$  appear as sibling leaves of maximum depth, from which the lemma follows. ■

Lemma 16.2 implies that the process of building up an optimal tree by mergers can, without loss of generality, begin with the greedy choice of merging together those two characters of lowest frequency. Why is this a greedy choice? We can view the cost of a single merger as being the sum of the frequencies of the two items being merged. Exercise 16.3-4 shows that the total cost of the tree constructed equals the sum of the costs of its mergers. Of all possible mergers at each step, HUFFMAN chooses the one that incurs the least cost.

The next lemma shows that the problem of constructing optimal prefix codes has the optimal-substructure property.

**Lemma 16.3**

Let  $C$  be a given alphabet with frequency  $c.freq$  defined for each character  $c \in C$ . Let  $x$  and  $y$  be two characters in  $C$  with minimum frequency. Let  $C'$  be the alphabet  $C$  with the characters  $x$  and  $y$  removed and a new character  $z$  added, so that  $C' = C - \{x, y\} \cup \{z\}$ . Define  $f$  for  $C'$  as for  $C$ , except that  $z.freq = x.freq + y.freq$ . Let  $T'$  be any tree representing an optimal prefix code for the alphabet  $C'$ . Then the tree  $T$ , obtained from  $T'$  by replacing the leaf node for  $z$  with an internal node having  $x$  and  $y$  as children, represents an optimal prefix code for the alphabet  $C$ .

**Proof** We first show how to express the cost  $B(T)$  of tree  $T$  in terms of the cost  $B(T')$  of tree  $T'$ , by considering the component costs in equation (16.4). For each character  $c \in C - \{x, y\}$ , we have that  $d_T(c) = d_{T'}(c)$ , and hence  $c.freq \cdot d_T(c) = c.freq \cdot d_{T'}(c)$ . Since  $d_T(x) = d_T(y) = d_{T'}(z) + 1$ , we have

$$\begin{aligned} x.freq \cdot d_T(x) + y.freq \cdot d_T(y) &= (x.freq + y.freq)(d_{T'}(z) + 1) \\ &= z.freq \cdot d_{T'}(z) + (x.freq + y.freq), \end{aligned}$$

from which we conclude that

$$B(T) = B(T') + x.freq + y.freq$$

or, equivalently,

$$B(T') = B(T) - x.freq - y.freq.$$

We now prove the lemma by contradiction. Suppose that  $T$  does not represent an optimal prefix code for  $C$ . Then there exists an optimal tree  $T''$  such that  $B(T'') < B(T)$ . Without loss of generality (by Lemma 16.2),  $T''$  has  $x$  and  $y$  as siblings. Let  $T'''$  be the tree  $T''$  with the common parent of  $x$  and  $y$  replaced by a leaf  $z$  with frequency  $z.freq = x.freq + y.freq$ . Then

$$\begin{aligned} B(T''') &= B(T'') - x.freq - y.freq \\ &< B(T) - x.freq - y.freq \\ &= B(T'), \end{aligned}$$

yielding a contradiction to the assumption that  $T'$  represents an optimal prefix code for  $C'$ . Thus,  $T$  must represent an optimal prefix code for the alphabet  $C$ . ■

**Theorem 16.4**

Procedure HUFFMAN produces an optimal prefix code.

**Proof** Immediate from Lemmas 16.2 and 16.3. ■



**Exercises****16.3-1**

Explain why, in the proof of Lemma 16.2, if  $x.freq = b.freq$ , then we must have  $a.freq = b.freq = x.freq = y.freq$ .

**16.3-2**

Prove that a binary tree that is not full cannot correspond to an optimal prefix code.

**16.3-3**

What is an optimal Huffman code for the following set of frequencies, based on the first 8 Fibonacci numbers?

a:1 b:1 c:2 d:3 e:5 f:8 g:13 h:21

Can you generalize your answer to find the optimal code when the frequencies are the first  $n$  Fibonacci numbers?

**16.3-4**

Prove that we can also express the total cost of a tree for a code as the sum, over all internal nodes, of the combined frequencies of the two children of the node.

**16.3-5**

Prove that if we order the characters in an alphabet so that their frequencies are monotonically decreasing, then there exists an optimal code whose codeword lengths are monotonically increasing.

**16.3-6**

Suppose we have an optimal prefix code on a set  $C = \{0, 1, \dots, n-1\}$  of characters and we wish to transmit this code using as few bits as possible. Show how to represent any optimal prefix code on  $C$  using only  $2n - 1 + n \lceil \lg n \rceil$  bits. (*Hint:* Use  $2n - 1$  bits to specify the structure of the tree, as discovered by a walk of the tree.)

**16.3-7**

Generalize Huffman's algorithm to ternary codewords (i.e., codewords using the symbols 0, 1, and 2), and prove that it yields optimal ternary codes.

**16.3-8**

Suppose that a data file contains a sequence of 8-bit characters such that all 256 characters are about equally common: the maximum character frequency is less than twice the minimum character frequency. Prove that Huffman coding in this case is no more efficient than using an ordinary 8-bit fixed-length code.

**16.3-9**

Show that no compression scheme can expect to compress a file of randomly chosen 8-bit characters by even a single bit. (*Hint:* Compare the number of possible files with the number of possible encoded files.)

---

**★ 16.4 Matroids and greedy methods**

In this section, we sketch a beautiful theory about greedy algorithms. This theory describes many situations in which the greedy method yields optimal solutions. It involves combinatorial structures known as “matroids.” Although this theory does not cover all cases for which a greedy method applies (for example, it does not cover the activity-selection problem of Section 16.1 or the Huffman-coding problem of Section 16.3), it does cover many cases of practical interest. Furthermore, this theory has been extended to cover many applications; see the notes at the end of this chapter for references.

**Matroids**

A **matroid** is an ordered pair  $M = (S, \mathcal{I})$  satisfying the following conditions.

1.  $S$  is a finite set.
2.  $\mathcal{I}$  is a nonempty family of subsets of  $S$ , called the **independent** subsets of  $S$ , such that if  $B \in \mathcal{I}$  and  $A \subseteq B$ , then  $A \in \mathcal{I}$ . We say that  $\mathcal{I}$  is **hereditary** if it satisfies this property. Note that the empty set  $\emptyset$  is necessarily a member of  $\mathcal{I}$ .
3. If  $A \in \mathcal{I}$ ,  $B \in \mathcal{I}$ , and  $|A| < |B|$ , then there exists some element  $x \in B - A$  such that  $A \cup \{x\} \in \mathcal{I}$ . We say that  $M$  satisfies the **exchange property**.

The word “matroid” is due to Hassler Whitney. He was studying **matric matroids**, in which the elements of  $S$  are the rows of a given matrix and a set of rows is independent if they are linearly independent in the usual sense. As Exercise 16.4-2 asks you to show, this structure defines a matroid.

As another example of matroids, consider the **graphic matroid**  $M_G = (S_G, \mathcal{I}_G)$  defined in terms of a given undirected graph  $G = (V, E)$  as follows:

- The set  $S_G$  is defined to be  $E$ , the set of edges of  $G$ .
- If  $A$  is a subset of  $E$ , then  $A \in \mathcal{I}_G$  if and only if  $A$  is acyclic. That is, a set of edges  $A$  is independent if and only if the subgraph  $G_A = (V, A)$  forms a forest.

The graphic matroid  $M_G$  is closely related to the minimum-spanning-tree problem, which Chapter 23 covers in detail.

**Theorem 16.5**

If  $G = (V, E)$  is an undirected graph, then  $M_G = (S_G, \mathcal{I}_G)$  is a matroid.

**Proof** Clearly,  $S_G = E$  is a finite set. Furthermore,  $\mathcal{I}_G$  is hereditary, since a subset of a forest is a forest. Putting it another way, removing edges from an acyclic set of edges cannot create cycles.

Thus, it remains to show that  $M_G$  satisfies the exchange property. Suppose that  $G_A = (V, A)$  and  $G_B = (V, B)$  are forests of  $G$  and that  $|B| > |A|$ . That is,  $A$  and  $B$  are acyclic sets of edges, and  $B$  contains more edges than  $A$  does.

We claim that a forest  $F = (V_F, E_F)$  contains exactly  $|V_F| - |E_F|$  trees. To see why, suppose that  $F$  consists of  $t$  trees, where the  $i$ th tree contains  $v_i$  vertices and  $e_i$  edges. Then, we have

$$\begin{aligned} |E_F| &= \sum_{i=1}^t e_i \\ &= \sum_{i=1}^t (v_i - 1) \quad (\text{by Theorem B.2}) \\ &= \sum_{i=1}^t v_i - t \\ &= |V_F| - t, \end{aligned}$$

which implies that  $t = |V_F| - |E_F|$ . Thus, forest  $G_A$  contains  $|V| - |A|$  trees, and forest  $G_B$  contains  $|V| - |B|$  trees.

Since forest  $G_B$  has fewer trees than forest  $G_A$  does, forest  $G_B$  must contain some tree  $T$  whose vertices are in two different trees in forest  $G_A$ . Moreover, since  $T$  is connected, it must contain an edge  $(u, v)$  such that vertices  $u$  and  $v$  are in different trees in forest  $G_A$ . Since the edge  $(u, v)$  connects vertices in two different trees in forest  $G_A$ , we can add the edge  $(u, v)$  to forest  $G_A$  without creating a cycle. Therefore,  $M_G$  satisfies the exchange property, completing the proof that  $M_G$  is a matroid. ■

Given a matroid  $M = (S, \mathcal{I})$ , we call an element  $x \notin A$  an *extension* of  $A \in \mathcal{I}$  if we can add  $x$  to  $A$  while preserving independence; that is,  $x$  is an extension of  $A$  if  $A \cup \{x\} \in \mathcal{I}$ . As an example, consider a graphic matroid  $M_G$ . If  $A$  is an independent set of edges, then edge  $e$  is an extension of  $A$  if and only if  $e$  is not in  $A$  and the addition of  $e$  to  $A$  does not create a cycle.

If  $A$  is an independent subset in a matroid  $M$ , we say that  $A$  is *maximal* if it has no extensions. That is,  $A$  is maximal if it is not contained in any larger independent subset of  $M$ . The following property is often useful.

**Theorem 16.6**

All maximal independent subsets in a matroid have the same size.

**Proof** Suppose to the contrary that  $A$  is a maximal independent subset of  $M$  and there exists another larger maximal independent subset  $B$  of  $M$ . Then, the exchange property implies that for some  $x \in B - A$ , we can extend  $A$  to a larger independent set  $A \cup \{x\}$ , contradicting the assumption that  $A$  is maximal. ■

As an illustration of this theorem, consider a graphic matroid  $M_G$  for a connected, undirected graph  $G$ . Every maximal independent subset of  $M_G$  must be a free tree with exactly  $|V| - 1$  edges that connects all the vertices of  $G$ . Such a tree is called a *spanning tree* of  $G$ .

We say that a matroid  $M = (S, \mathcal{I})$  is **weighted** if it is associated with a weight function  $w$  that assigns a strictly positive weight  $w(x)$  to each element  $x \in S$ . The weight function  $w$  extends to subsets of  $S$  by summation:

$$w(A) = \sum_{x \in A} w(x)$$

for any  $A \subseteq S$ . For example, if we let  $w(e)$  denote the weight of an edge  $e$  in a graphic matroid  $M_G$ , then  $w(A)$  is the total weight of the edges in edge set  $A$ .

**Greedy algorithms on a weighted matroid**

Many problems for which a greedy approach provides optimal solutions can be formulated in terms of finding a maximum-weight independent subset in a weighted matroid. That is, we are given a weighted matroid  $M = (S, \mathcal{I})$ , and we wish to find an independent set  $A \in \mathcal{I}$  such that  $w(A)$  is maximized. We call such a subset that is independent and has maximum possible weight an **optimal** subset of the matroid. Because the weight  $w(x)$  of any element  $x \in S$  is positive, an optimal subset is always a maximal independent subset—it always helps to make  $A$  as large as possible.

For example, in the **minimum-spanning-tree problem**, we are given a connected undirected graph  $G = (V, E)$  and a length function  $w$  such that  $w(e)$  is the (positive) length of edge  $e$ . (We use the term “length” here to refer to the original edge weights for the graph, reserving the term “weight” to refer to the weights in the associated matroid.) We wish to find a subset of the edges that connects all of the vertices together and has minimum total length. To view this as a problem of finding an optimal subset of a matroid, consider the weighted matroid  $M_G$  with weight function  $w'$ , where  $w'(e) = w_0 - w(e)$  and  $w_0$  is larger than the maximum length of any edge. In this weighted matroid, all weights are positive and an optimal subset is a spanning tree of minimum total length in the original graph. More specifically, each maximal independent subset  $A$  corresponds to a spanning tree

with  $|V| - 1$  edges, and since

$$\begin{aligned}
 w'(A) &= \sum_{e \in A} w'(e) \\
 &= \sum_{e \in A} (w_0 - w(e)) \\
 &= (|V| - 1)w_0 - \sum_{e \in A} w(e) \\
 &= (|V| - 1)w_0 - w(A)
 \end{aligned}$$

for any maximal independent subset  $A$ , an independent subset that maximizes the quantity  $w'(A)$  must minimize  $w(A)$ . Thus, any algorithm that can find an optimal subset  $A$  in an arbitrary matroid can solve the minimum-spanning-tree problem.

Chapter 23 gives algorithms for the minimum-spanning-tree problem, but here we give a greedy algorithm that works for any weighted matroid. The algorithm takes as input a weighted matroid  $M = (S, \mathcal{I})$  with an associated positive weight function  $w$ , and it returns an optimal subset  $A$ . In our pseudocode, we denote the components of  $M$  by  $M.S$  and  $M.\mathcal{I}$  and the weight function by  $w$ . The algorithm is greedy because it considers in turn each element  $x \in S$ , in order of monotonically decreasing weight, and immediately adds it to the set  $A$  being accumulated if  $A \cup \{x\}$  is independent.

GREEDY( $M, w$ )

```

1   $A = \emptyset$ 
2  sort  $M.S$  into monotonically decreasing order by weight  $w$ 
3  for each  $x \in M.S$ , taken in monotonically decreasing order by weight  $w(x)$ 
4      if  $A \cup \{x\} \in M.\mathcal{I}$ 
5           $A = A \cup \{x\}$ 
6  return  $A$ 

```

Line 4 checks whether adding each element  $x$  to  $A$  would maintain  $A$  as an independent set. If  $A$  would remain independent, then line 5 adds  $x$  to  $A$ . Otherwise,  $x$  is discarded. Since the empty set is independent, and since each iteration of the **for** loop maintains  $A$ 's independence, the subset  $A$  is always independent, by induction. Therefore, GREEDY always returns an independent subset  $A$ . We shall see in a moment that  $A$  is a subset of maximum possible weight, so that  $A$  is an optimal subset.

The running time of GREEDY is easy to analyze. Let  $n$  denote  $|S|$ . The sorting phase of GREEDY takes time  $O(n \lg n)$ . Line 4 executes exactly  $n$  times, once for each element of  $S$ . Each execution of line 4 requires a check on whether or not the set  $A \cup \{x\}$  is independent. If each such check takes time  $O(f(n))$ , the entire algorithm runs in time  $O(n \lg n + nf(n))$ .

We now prove that GREEDY returns an optimal subset.

**Lemma 16.7 (Matroids exhibit the greedy-choice property)**

Suppose that  $M = (S, \mathcal{I})$  is a weighted matroid with weight function  $w$  and that  $S$  is sorted into monotonically decreasing order by weight. Let  $x$  be the first element of  $S$  such that  $\{x\}$  is independent, if any such  $x$  exists. If  $x$  exists, then there exists an optimal subset  $A$  of  $S$  that contains  $x$ .

**Proof** If no such  $x$  exists, then the only independent subset is the empty set and the lemma is vacuously true. Otherwise, let  $B$  be any nonempty optimal subset. Assume that  $x \notin B$ ; otherwise, letting  $A = B$  gives an optimal subset of  $S$  that contains  $x$ .

No element of  $B$  has weight greater than  $w(x)$ . To see why, observe that  $y \in B$  implies that  $\{y\}$  is independent, since  $B \in \mathcal{I}$  and  $\mathcal{I}$  is hereditary. Our choice of  $x$  therefore ensures that  $w(x) \geq w(y)$  for any  $y \in B$ .

Construct the set  $A$  as follows. Begin with  $A = \{x\}$ . By the choice of  $x$ , set  $A$  is independent. Using the exchange property, repeatedly find a new element of  $B$  that we can add to  $A$  until  $|A| = |B|$ , while preserving the independence of  $A$ . At that point,  $A$  and  $B$  are the same except that  $A$  has  $x$  and  $B$  has some other element  $y$ . That is,  $A = B - \{y\} \cup \{x\}$  for some  $y \in B$ , and so

$$\begin{aligned} w(A) &= w(B) - w(y) + w(x) \\ &\geq w(B). \end{aligned}$$

Because set  $B$  is optimal, set  $A$ , which contains  $x$ , must also be optimal. ■

We next show that if an element is not an option initially, then it cannot be an option later.

**Lemma 16.8**

Let  $M = (S, \mathcal{I})$  be any matroid. If  $x$  is an element of  $S$  that is an extension of some independent subset  $A$  of  $S$ , then  $x$  is also an extension of  $\emptyset$ .

**Proof** Since  $x$  is an extension of  $A$ , we have that  $A \cup \{x\}$  is independent. Since  $\mathcal{I}$  is hereditary,  $\{x\}$  must be independent. Thus,  $x$  is an extension of  $\emptyset$ . ■

**Corollary 16.9**

Let  $M = (S, \mathcal{I})$  be any matroid. If  $x$  is an element of  $S$  such that  $x$  is not an extension of  $\emptyset$ , then  $x$  is not an extension of any independent subset  $A$  of  $S$ .

**Proof** This corollary is simply the contrapositive of Lemma 16.8. ■

Corollary 16.9 says that any element that cannot be used immediately can never be used. Therefore, GREEDY cannot make an error by passing over any initial elements in  $S$  that are not an extension of  $\emptyset$ , since they can never be used.

**Lemma 16.10 (Matroids exhibit the optimal-substructure property)**

Let  $x$  be the first element of  $S$  chosen by GREEDY for the weighted matroid  $M = (S, \mathcal{I})$ . The remaining problem of finding a maximum-weight independent subset containing  $x$  reduces to finding a maximum-weight independent subset of the weighted matroid  $M' = (S', \mathcal{I}')$ , where

$$\begin{aligned} S' &= \{y \in S : \{x, y\} \in \mathcal{I}\} , \\ \mathcal{I}' &= \{B \subseteq S - \{x\} : B \cup \{x\} \in \mathcal{I}\} , \end{aligned}$$

and the weight function for  $M'$  is the weight function for  $M$ , restricted to  $S'$ . (We call  $M'$  the **contraction** of  $M$  by the element  $x$ .)

**Proof** If  $A$  is any maximum-weight independent subset of  $M$  containing  $x$ , then  $A' = A - \{x\}$  is an independent subset of  $M'$ . Conversely, any independent subset  $A'$  of  $M'$  yields an independent subset  $A = A' \cup \{x\}$  of  $M$ . Since we have in both cases that  $w(A) = w(A') + w(x)$ , a maximum-weight solution in  $M$  containing  $x$  yields a maximum-weight solution in  $M'$ , and vice versa. ■

**Theorem 16.11 (Correctness of the greedy algorithm on matroids)**

If  $M = (S, \mathcal{I})$  is a weighted matroid with weight function  $w$ , then GREEDY( $M, w$ ) returns an optimal subset.

**Proof** By Corollary 16.9, any elements that GREEDY passes over initially because they are not extensions of  $\emptyset$  can be forgotten about, since they can never be useful. Once GREEDY selects the first element  $x$ , Lemma 16.7 implies that the algorithm does not err by adding  $x$  to  $A$ , since there exists an optimal subset containing  $x$ . Finally, Lemma 16.10 implies that the remaining problem is one of finding an optimal subset in the matroid  $M'$  that is the contraction of  $M$  by  $x$ . After the procedure GREEDY sets  $A$  to  $\{x\}$ , we can interpret all of its remaining steps as acting in the matroid  $M' = (S', \mathcal{I}')$ , because  $B$  is independent in  $M'$  if and only if  $B \cup \{x\}$  is independent in  $M$ , for all sets  $B \in \mathcal{I}'$ . Thus, the subsequent operation of GREEDY will find a maximum-weight independent subset for  $M'$ , and the overall operation of GREEDY will find a maximum-weight independent subset for  $M$ . ■

## Exercises

### 16.4-1

Show that  $(S, \mathcal{I}_k)$  is a matroid, where  $S$  is any finite set and  $\mathcal{I}_k$  is the set of all subsets of  $S$  of size at most  $k$ , where  $k \leq |S|$ .

### 16.4-2 ★

Given an  $m \times n$  matrix  $T$  over some field (such as the reals), show that  $(S, \mathcal{I})$  is a matroid, where  $S$  is the set of columns of  $T$  and  $A \in \mathcal{I}$  if and only if the columns in  $A$  are linearly independent.

### 16.4-3 ★

Show that if  $(S, \mathcal{I})$  is a matroid, then  $(S, \mathcal{I}')$  is a matroid, where

$$\mathcal{I}' = \{A' : S - A' \text{ contains some maximal } A \in \mathcal{I}\}.$$

That is, the maximal independent sets of  $(S, \mathcal{I}')$  are just the complements of the maximal independent sets of  $(S, \mathcal{I})$ .

### 16.4-4 ★

Let  $S$  be a finite set and let  $S_1, S_2, \dots, S_k$  be a partition of  $S$  into nonempty disjoint subsets. Define the structure  $(S, \mathcal{I})$  by the condition that  $\mathcal{I} = \{A : |A \cap S_i| \leq 1 \text{ for } i = 1, 2, \dots, k\}$ . Show that  $(S, \mathcal{I})$  is a matroid. That is, the set of all sets  $A$  that contain at most one member of each subset in the partition determines the independent sets of a matroid.

### 16.4-5

Show how to transform the weight function of a weighted matroid problem, where the desired optimal solution is a *minimum-weight* maximal independent subset, to make it a standard weighted-matroid problem. Argue carefully that your transformation is correct.

---

## ★ 16.5 A task-scheduling problem as a matroid

An interesting problem that we can solve using matroids is the problem of optimally scheduling unit-time tasks on a single processor, where each task has a deadline, along with a penalty paid if the task misses its deadline. The problem looks complicated, but we can solve it in a surprisingly simple manner by casting it as a matroid and using a greedy algorithm.

A **unit-time task** is a job, such as a program to be run on a computer, that requires exactly one unit of time to complete. Given a finite set  $S$  of unit-time tasks, a



**schedule** for  $S$  is a permutation of  $S$  specifying the order in which to perform these tasks. The first task in the schedule begins at time 0 and finishes at time 1, the second task begins at time 1 and finishes at time 2, and so on.

The problem of **scheduling unit-time tasks with deadlines and penalties for a single processor** has the following inputs:

- a set  $S = \{a_1, a_2, \dots, a_n\}$  of  $n$  unit-time tasks;
- a set of  $n$  integer **deadlines**  $d_1, d_2, \dots, d_n$ , such that each  $d_i$  satisfies  $1 \leq d_i \leq n$  and task  $a_i$  is supposed to finish by time  $d_i$ ; and
- a set of  $n$  nonnegative weights or **penalties**  $w_1, w_2, \dots, w_n$ , such that we incur a penalty of  $w_i$  if task  $a_i$  is not finished by time  $d_i$ , and we incur no penalty if a task finishes by its deadline.

We wish to find a schedule for  $S$  that minimizes the total penalty incurred for missed deadlines.

Consider a given schedule. We say that a task is **late** in this schedule if it finishes after its deadline. Otherwise, the task is **early** in the schedule. We can always transform an arbitrary schedule into **early-first form**, in which the early tasks precede the late tasks. To see why, note that if some early task  $a_i$  follows some late task  $a_j$ , then we can switch the positions of  $a_i$  and  $a_j$ , and  $a_i$  will still be early and  $a_j$  will still be late.

Furthermore, we claim that we can always transform an arbitrary schedule into **canonical form**, in which the early tasks precede the late tasks and we schedule the early tasks in order of monotonically increasing deadlines. To do so, we put the schedule into early-first form. Then, as long as there exist two early tasks  $a_i$  and  $a_j$  finishing at respective times  $k$  and  $k + 1$  in the schedule such that  $d_j < d_i$ , we swap the positions of  $a_i$  and  $a_j$ . Since  $a_j$  is early before the swap,  $k + 1 \leq d_j$ . Therefore,  $k + 1 < d_i$ , and so  $a_i$  is still early after the swap. Because task  $a_j$  is moved earlier in the schedule, it remains early after the swap.

The search for an optimal schedule thus reduces to finding a set  $A$  of tasks that we assign to be early in the optimal schedule. Having determined  $A$ , we can create the actual schedule by listing the elements of  $A$  in order of monotonically increasing deadlines, then listing the late tasks (i.e.,  $S - A$ ) in any order, producing a canonical ordering of the optimal schedule.

We say that a set  $A$  of tasks is **independent** if there exists a schedule for these tasks such that no tasks are late. Clearly, the set of early tasks for a schedule forms an independent set of tasks. Let  $\mathcal{I}$  denote the set of all independent sets of tasks.

Consider the problem of determining whether a given set  $A$  of tasks is independent. For  $t = 0, 1, 2, \dots, n$ , let  $N_t(A)$  denote the number of tasks in  $A$  whose deadline is  $t$  or earlier. Note that  $N_0(A) = 0$  for any set  $A$ .

**Lemma 16.12**

For any set of tasks  $A$ , the following statements are equivalent.

1. The set  $A$  is independent.
2. For  $t = 0, 1, 2, \dots, n$ , we have  $N_t(A) \leq t$ .
3. If the tasks in  $A$  are scheduled in order of monotonically increasing deadlines, then no task is late.

**Proof** To show that (1) implies (2), we prove the contrapositive: if  $N_t(A) > t$  for some  $t$ , then there is no way to make a schedule with no late tasks for set  $A$ , because more than  $t$  tasks must finish before time  $t$ . Therefore, (1) implies (2). If (2) holds, then (3) must follow: there is no way to “get stuck” when scheduling the tasks in order of monotonically increasing deadlines, since (2) implies that the  $i$ th largest deadline is at least  $i$ . Finally, (3) trivially implies (1). ■

Using property 2 of Lemma 16.12, we can easily compute whether or not a given set of tasks is independent (see Exercise 16.5-2).

The problem of minimizing the sum of the penalties of the late tasks is the same as the problem of maximizing the sum of the penalties of the early tasks. The following theorem thus ensures that we can use the greedy algorithm to find an independent set  $A$  of tasks with the maximum total penalty.

**Theorem 16.13**

If  $S$  is a set of unit-time tasks with deadlines, and  $\mathcal{I}$  is the set of all independent sets of tasks, then the corresponding system  $(S, \mathcal{I})$  is a matroid.

**Proof** Every subset of an independent set of tasks is certainly independent. To prove the exchange property, suppose that  $B$  and  $A$  are independent sets of tasks and that  $|B| > |A|$ . Let  $k$  be the largest  $t$  such that  $N_t(B) \leq N_t(A)$ . (Such a value of  $t$  exists, since  $N_0(A) = N_0(B) = 0$ .) Since  $N_n(B) = |B|$  and  $N_n(A) = |A|$ , but  $|B| > |A|$ , we must have that  $k < n$  and that  $N_j(B) > N_j(A)$  for all  $j$  in the range  $k + 1 \leq j \leq n$ . Therefore,  $B$  contains more tasks with deadline  $k + 1$  than  $A$  does. Let  $a_i$  be a task in  $B - A$  with deadline  $k + 1$ . Let  $A' = A \cup \{a_i\}$ .

We now show that  $A'$  must be independent by using property 2 of Lemma 16.12. For  $0 \leq t \leq k$ , we have  $N_t(A') = N_t(A) \leq t$ , since  $A$  is independent. For  $k < t \leq n$ , we have  $N_t(A') \leq N_t(B) \leq t$ , since  $B$  is independent. Therefore,  $A'$  is independent, completing our proof that  $(S, \mathcal{I})$  is a matroid. ■

By Theorem 16.11, we can use a greedy algorithm to find a maximum-weight independent set of tasks  $A$ . We can then create an optimal schedule having the tasks in  $A$  as its early tasks. This method is an efficient algorithm for scheduling

	Task						
$a_i$	1	2	3	4	5	6	7
$d_i$	4	2	4	3	1	4	6
$w_i$	70	60	50	40	30	20	10

**Figure 16.7** An instance of the problem of scheduling unit-time tasks with deadlines and penalties for a single processor.

unit-time tasks with deadlines and penalties for a single processor. The running time is  $O(n^2)$  using GREEDY, since each of the  $O(n)$  independence checks made by that algorithm takes time  $O(n)$  (see Exercise 16.5-2). Problem 16-4 gives a faster implementation.

Figure 16.7 demonstrates an example of the problem of scheduling unit-time tasks with deadlines and penalties for a single processor. In this example, the greedy algorithm selects, in order, tasks  $a_1$ ,  $a_2$ ,  $a_3$ , and  $a_4$ , then rejects  $a_5$  (because  $N_4(\{a_1, a_2, a_3, a_4, a_5\}) = 5$ ) and  $a_6$  (because  $N_4(\{a_1, a_2, a_3, a_4, a_6\}) = 5$ ), and finally accepts  $a_7$ . The final optimal schedule is

$$\langle a_2, a_4, a_1, a_3, a_7, a_5, a_6 \rangle ,$$

which has a total penalty incurred of  $w_5 + w_6 = 50$ .

## Exercises

### 16.5-1

Solve the instance of the scheduling problem given in Figure 16.7, but with each penalty  $w_i$  replaced by  $80 - w_i$ .

### 16.5-2

Show how to use property 2 of Lemma 16.12 to determine in time  $O(|A|)$  whether or not a given set  $A$  of tasks is independent.

---

## Problems

### 16-1 Coin changing

Consider the problem of making change for  $n$  cents using the fewest number of coins. Assume that each coin's value is an integer.

- Describe a greedy algorithm to make change consisting of quarters, dimes, nickels, and pennies. Prove that your algorithm yields an optimal solution.

- b. Suppose that the available coins are in the denominations that are powers of  $c$ , i.e., the denominations are  $c^0, c^1, \dots, c^k$  for some integers  $c > 1$  and  $k \geq 1$ . Show that the greedy algorithm always yields an optimal solution.
- c. Give a set of coin denominations for which the greedy algorithm does not yield an optimal solution. Your set should include a penny so that there is a solution for every value of  $n$ .
- d. Give an  $O(nk)$ -time algorithm that makes change for any set of  $k$  different coin denominations, assuming that one of the coins is a penny.

### 16-2 Scheduling to minimize average completion time

Suppose you are given a set  $S = \{a_1, a_2, \dots, a_n\}$  of tasks, where task  $a_i$  requires  $p_i$  units of processing time to complete, once it has started. You have one computer on which to run these tasks, and the computer can run only one task at a time. Let  $c_i$  be the **completion time** of task  $a_i$ , that is, the time at which task  $a_i$  completes processing. Your goal is to minimize the average completion time, that is, to minimize  $(1/n) \sum_{i=1}^n c_i$ . For example, suppose there are two tasks,  $a_1$  and  $a_2$ , with  $p_1 = 3$  and  $p_2 = 5$ , and consider the schedule in which  $a_2$  runs first, followed by  $a_1$ . Then  $c_2 = 5$ ,  $c_1 = 8$ , and the average completion time is  $(5 + 8)/2 = 6.5$ . If task  $a_1$  runs first, however, then  $c_1 = 3$ ,  $c_2 = 8$ , and the average completion time is  $(3 + 8)/2 = 5.5$ .

- a. Give an algorithm that schedules the tasks so as to minimize the average completion time. Each task must run non-preemptively, that is, once task  $a_i$  starts, it must run continuously for  $p_i$  units of time. Prove that your algorithm minimizes the average completion time, and state the running time of your algorithm.
- b. Suppose now that the tasks are not all available at once. That is, each task cannot start until its **release time**  $r_i$ . Suppose also that we allow **preemption**, so that a task can be suspended and restarted at a later time. For example, a task  $a_i$  with processing time  $p_i = 6$  and release time  $r_i = 1$  might start running at time 1 and be preempted at time 4. It might then resume at time 10 but be preempted at time 11, and it might finally resume at time 13 and complete at time 15. Task  $a_i$  has run for a total of 6 time units, but its running time has been divided into three pieces. In this scenario,  $a_i$ 's completion time is 15. Give an algorithm that schedules the tasks so as to minimize the average completion time in this new scenario. Prove that your algorithm minimizes the average completion time, and state the running time of your algorithm.

**16-3 Acyclic subgraphs**

- a. The **incidence matrix** for an undirected graph  $G = (V, E)$  is a  $|V| \times |E|$  matrix  $M$  such that  $M_{ve} = 1$  if edge  $e$  is incident on vertex  $v$ , and  $M_{ve} = 0$  otherwise. Argue that a set of columns of  $M$  is linearly independent over the field of integers modulo 2 if and only if the corresponding set of edges is acyclic. Then, use the result of Exercise 16.4-2 to provide an alternate proof that  $(E, \mathcal{I})$  of part (a) is a matroid.
- b. Suppose that we associate a nonnegative weight  $w(e)$  with each edge in an undirected graph  $G = (V, E)$ . Give an efficient algorithm to find an acyclic subset of  $E$  of maximum total weight.
- c. Let  $G(V, E)$  be an arbitrary directed graph, and let  $(E, \mathcal{I})$  be defined so that  $A \in \mathcal{I}$  if and only if  $A$  does not contain any directed cycles. Give an example of a directed graph  $G$  such that the associated system  $(E, \mathcal{I})$  is not a matroid. Specify which defining condition for a matroid fails to hold.
- d. The **incidence matrix** for a directed graph  $G = (V, E)$  with no self-loops is a  $|V| \times |E|$  matrix  $M$  such that  $M_{ve} = -1$  if edge  $e$  leaves vertex  $v$ ,  $M_{ve} = 1$  if edge  $e$  enters vertex  $v$ , and  $M_{ve} = 0$  otherwise. Argue that if a set of columns of  $M$  is linearly independent, then the corresponding set of edges does not contain a directed cycle.
- e. Exercise 16.4-2 tells us that the set of linearly independent sets of columns of any matrix  $M$  forms a matroid. Explain carefully why the results of parts (d) and (e) are not contradictory. How can there fail to be a perfect correspondence between the notion of a set of edges being acyclic and the notion of the associated set of columns of the incidence matrix being linearly independent?

**16-4 Scheduling variations**

Consider the following algorithm for the problem from Section 16.5 of scheduling unit-time tasks with deadlines and penalties. Let all  $n$  time slots be initially empty, where time slot  $i$  is the unit-length slot of time that finishes at time  $i$ . We consider the tasks in order of monotonically decreasing penalty. When considering task  $a_j$ , if there exists a time slot at or before  $a_j$ 's deadline  $d_j$  that is still empty, assign  $a_j$  to the latest such slot, filling it. If there is no such slot, assign task  $a_j$  to the latest of the as yet unfilled slots.

- a. Argue that this algorithm always gives an optimal answer.
- b. Use the fast disjoint-set forest presented in Section 21.3 to implement the algorithm efficiently. Assume that the set of input tasks has already been sorted into

monotonically decreasing order by penalty. Analyze the running time of your implementation.

### 16-5 Off-line caching

Modern computers use a cache to store a small amount of data in a fast memory. Even though a program may access large amounts of data, by storing a small subset of the main memory in the *cache*—a small but faster memory—overall access time can greatly decrease. When a computer program executes, it makes a sequence  $\langle r_1, r_2, \dots, r_n \rangle$  of  $n$  memory requests, where each request is for a particular data element. For example, a program that accesses 4 distinct elements  $\{a, b, c, d\}$  might make the sequence of requests  $\langle d, b, d, b, d, a, c, d, b, a, c, b \rangle$ . Let  $k$  be the size of the cache. When the cache contains  $k$  elements and the program requests the  $(k + 1)$ st element, the system must decide, for this and each subsequent request, which  $k$  elements to keep in the cache. More precisely, for each request  $r_i$ , the cache-management algorithm checks whether element  $r_i$  is already in the cache. If it is, then we have a *cache hit*; otherwise, we have a *cache miss*. Upon a cache miss, the system retrieves  $r_i$  from the main memory, and the cache-management algorithm must decide whether to keep  $r_i$  in the cache. If it decides to keep  $r_i$  and the cache already holds  $k$  elements, then it must evict one element to make room for  $r_i$ . The cache-management algorithm evicts data with the goal of minimizing the number of cache misses over the entire sequence of requests.

Typically, caching is an on-line problem. That is, we have to make decisions about which data to keep in the cache without knowing the future requests. Here, however, we consider the off-line version of this problem, in which we are given in advance the entire sequence of  $n$  requests and the cache size  $k$ , and we wish to minimize the total number of cache misses.

We can solve this off-line problem by a greedy strategy called *furthest-in-future*, which chooses to evict the item in the cache whose next access in the request sequence comes furthest in the future.

- a. Write pseudocode for a cache manager that uses the furthest-in-future strategy. The input should be a sequence  $\langle r_1, r_2, \dots, r_n \rangle$  of requests and a cache size  $k$ , and the output should be a sequence of decisions about which data element (if any) to evict upon each request. What is the running time of your algorithm?
- b. Show that the off-line caching problem exhibits optimal substructure.
- c. Prove that furthest-in-future produces the minimum possible number of cache misses.

---

**Chapter notes**

Much more material on greedy algorithms and matroids can be found in Lawler [224] and Papadimitriou and Steiglitz [271].

The greedy algorithm first appeared in the combinatorial optimization literature in a 1971 article by Edmonds [101], though the theory of matroids dates back to a 1935 article by Whitney [355].

Our proof of the correctness of the greedy algorithm for the activity-selection problem is based on that of Gavril [131]. The task-scheduling problem is studied in Lawler [224]; Horowitz, Sahni, and Rajasekaran [181]; and Brassard and Bratley [54].

Huffman codes were invented in 1952 [185]; Lelewer and Hirschberg [231] surveys data-compression techniques known as of 1987.

An extension of matroid theory to greedoid theory was pioneered by Korte and Lovász [216, 217, 218, 219], who greatly generalize the theory presented here.

In an *amortized analysis*, we average the time required to perform a sequence of data-structure operations over all the operations performed. With amortized analysis, we can show that the average cost of an operation is small, if we average over a sequence of operations, even though a single operation within the sequence might be expensive. Amortized analysis differs from average-case analysis in that probability is not involved; an amortized analysis guarantees the *average performance of each operation in the worst case*.

The first three sections of this chapter cover the three most common techniques used in amortized analysis. Section 17.1 starts with aggregate analysis, in which we determine an upper bound  $T(n)$  on the total cost of a sequence of  $n$  operations. The average cost per operation is then  $T(n)/n$ . We take the average cost as the amortized cost of each operation, so that all operations have the same amortized cost.

Section 17.2 covers the accounting method, in which we determine an amortized cost of each operation. When there is more than one type of operation, each type of operation may have a different amortized cost. The accounting method overcharges some operations early in the sequence, storing the overcharge as “prepaid credit” on specific objects in the data structure. Later in the sequence, the credit pays for operations that are charged less than they actually cost.

Section 17.3 discusses the potential method, which is like the accounting method in that we determine the amortized cost of each operation and may overcharge operations early on to compensate for undercharges later. The potential method maintains the credit as the “potential energy” of the data structure as a whole instead of associating the credit with individual objects within the data structure.

We shall use two examples to examine these three methods. One is a stack with the additional operation `MULTIPOP`, which pops several objects at once. The other is a binary counter that counts up from 0 by means of the single operation `INCREMENT`.



While reading this chapter, bear in mind that the charges assigned during an amortized analysis are for analysis purposes only. They need not—and should not—appear in the code. If, for example, we assign a credit to an object  $x$  when using the accounting method, we have no need to assign an appropriate amount to some attribute, such as  $x.credit$ , in the code.

When we perform an amortized analysis, we often gain insight into a particular data structure, and this insight can help us optimize the design. In Section 17.4, for example, we shall use the potential method to analyze a dynamically expanding and contracting table.

---

## 17.1 Aggregate analysis

In *aggregate analysis*, we show that for all  $n$ , a sequence of  $n$  operations takes *worst-case* time  $T(n)$  in total. In the worst case, the average cost, or *amortized cost*, per operation is therefore  $T(n)/n$ . Note that this amortized cost applies to each operation, even when there are several types of operations in the sequence. The other two methods we shall study in this chapter, the accounting method and the potential method, may assign different amortized costs to different types of operations.

### Stack operations

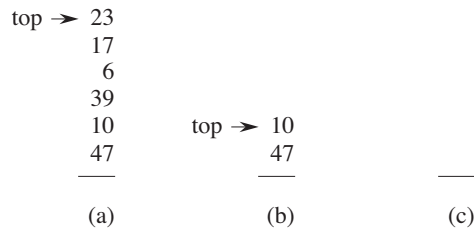
In our first example of aggregate analysis, we analyze stacks that have been augmented with a new operation. Section 10.1 presented the two fundamental stack operations, each of which takes  $O(1)$  time:

PUSH( $S, x$ ) pushes object  $x$  onto stack  $S$ .

POP( $S$ ) pops the top of stack  $S$  and returns the popped object. Calling POP on an empty stack generates an error.

Since each of these operations runs in  $O(1)$  time, let us consider the cost of each to be 1. The total cost of a sequence of  $n$  PUSH and POP operations is therefore  $n$ , and the actual running time for  $n$  operations is therefore  $\Theta(n)$ .

Now we add the stack operation MULTIPOP( $S, k$ ), which removes the  $k$  top objects of stack  $S$ , popping the entire stack if the stack contains fewer than  $k$  objects. Of course, we assume that  $k$  is positive; otherwise the MULTIPOP operation leaves the stack unchanged. In the following pseudocode, the operation STACK-EMPTY returns TRUE if there are no objects currently on the stack, and FALSE otherwise.



**Figure 17.1** The action of MULTIPOP on a stack  $S$ , shown initially in (a). The top 4 objects are popped by MULTIPOP( $S$ , 4), whose result is shown in (b). The next operation is MULTIPOP( $S$ , 7), which empties the stack—shown in (c)—since there were fewer than 7 objects remaining.

MULTIPOP( $S$ ,  $k$ )

```

1  while not STACK-EMPTY( $S$ ) and  $k > 0$ 
2      POP( $S$ )
3       $k = k - 1$ 

```

Figure 17.1 shows an example of MULTIPOP.

What is the running time of MULTIPOP( $S$ ,  $k$ ) on a stack of  $s$  objects? The actual running time is linear in the number of POP operations actually executed, and thus we can analyze MULTIPOP in terms of the abstract costs of 1 each for PUSH and POP. The number of iterations of the **while** loop is the number  $\min(s, k)$  of objects popped off the stack. Each iteration of the loop makes one call to POP in line 2. Thus, the total cost of MULTIPOP is  $\min(s, k)$ , and the actual running time is a linear function of this cost.

Let us analyze a sequence of  $n$  PUSH, POP, and MULTIPOP operations on an initially empty stack. The worst-case cost of a MULTIPOP operation in the sequence is  $O(n)$ , since the stack size is at most  $n$ . The worst-case time of any stack operation is therefore  $O(n)$ , and hence a sequence of  $n$  operations costs  $O(n^2)$ , since we may have  $O(n)$  MULTIPOP operations costing  $O(n)$  each. Although this analysis is correct, the  $O(n^2)$  result, which we obtained by considering the worst-case cost of each operation individually, is not tight.

Using aggregate analysis, we can obtain a better upper bound that considers the entire sequence of  $n$  operations. In fact, although a single MULTIPOP operation can be expensive, any sequence of  $n$  PUSH, POP, and MULTIPOP operations on an initially empty stack can cost at most  $O(n)$ . Why? We can pop each object from the stack at most once for each time we have pushed it onto the stack. Therefore, the number of times that POP can be called on a nonempty stack, including calls within MULTIPOP, is at most the number of PUSH operations, which is at most  $n$ . For any value of  $n$ , any sequence of  $n$  PUSH, POP, and MULTIPOP operations takes a total of  $O(n)$  time. The average cost of an operation is  $O(n)/n = O(1)$ . In aggregate

analysis, we assign the amortized cost of each operation to be the average cost. In this example, therefore, all three stack operations have an amortized cost of  $O(1)$ .

We emphasize again that although we have just shown that the average cost, and hence the running time, of a stack operation is  $O(1)$ , we did not use probabilistic reasoning. We actually showed a *worst-case* bound of  $O(n)$  on a sequence of  $n$  operations. Dividing this total cost by  $n$  yielded the average cost per operation, or the amortized cost.

### Incrementing a binary counter

As another example of aggregate analysis, consider the problem of implementing a  $k$ -bit binary counter that counts upward from 0. We use an array  $A[0 \dots k - 1]$  of bits, where  $A.length = k$ , as the counter. A binary number  $x$  that is stored in the counter has its lowest-order bit in  $A[0]$  and its highest-order bit in  $A[k - 1]$ , so that  $x = \sum_{i=0}^{k-1} A[i] \cdot 2^i$ . Initially,  $x = 0$ , and thus  $A[i] = 0$  for  $i = 0, 1, \dots, k - 1$ . To add 1 (modulo  $2^k$ ) to the value in the counter, we use the following procedure.

```

INCREMENT( $A$ )
1   $i = 0$ 
2  while  $i < A.length$  and  $A[i] == 1$ 
3       $A[i] = 0$ 
4       $i = i + 1$ 
5  if  $i < A.length$ 
6       $A[i] = 1$ 

```

Figure 17.2 shows what happens to a binary counter as we increment it 16 times, starting with the initial value 0 and ending with the value 16. At the start of each iteration of the **while** loop in lines 2–4, we wish to add a 1 into position  $i$ . If  $A[i] = 1$ , then adding 1 flips the bit to 0 in position  $i$  and yields a carry of 1, to be added into position  $i + 1$  on the next iteration of the loop. Otherwise, the loop ends, and then, if  $i < k$ , we know that  $A[i] = 0$ , so that line 6 adds a 1 into position  $i$ , flipping the 0 to a 1. The cost of each INCREMENT operation is linear in the number of bits flipped.

As with the stack example, a cursory analysis yields a bound that is correct but not tight. A single execution of INCREMENT takes time  $\Theta(k)$  in the worst case, in which array  $A$  contains all 1s. Thus, a sequence of  $n$  INCREMENT operations on an initially zero counter takes time  $O(nk)$  in the worst case.

We can tighten our analysis to yield a worst-case cost of  $O(n)$  for a sequence of  $n$  INCREMENT operations by observing that not all bits flip each time INCREMENT is called. As Figure 17.2 shows,  $A[0]$  does flip each time INCREMENT is called. The next bit up,  $A[1]$ , flips only every other time: a sequence of  $n$  INCREMENT

Counter value	A[7]	A[6]	A[5]	A[4]	A[3]	A[2]	A[1]	A[0]	Total cost
0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	1	1
2	0	0	0	0	0	0	1	0	3
3	0	0	0	0	0	0	1	1	4
4	0	0	0	0	0	1	0	0	7
5	0	0	0	0	0	1	0	1	8
6	0	0	0	0	0	1	1	0	10
7	0	0	0	0	0	1	1	1	11
8	0	0	0	0	1	0	0	0	15
9	0	0	0	0	1	0	0	1	16
10	0	0	0	0	1	0	1	0	18
11	0	0	0	0	1	0	1	1	19
12	0	0	0	0	1	1	0	0	22
13	0	0	0	0	1	1	0	1	23
14	0	0	0	0	1	1	1	0	25
15	0	0	0	0	1	1	1	1	26
16	0	0	0	1	0	0	0	0	31

**Figure 17.2** An 8-bit binary counter as its value goes from 0 to 16 by a sequence of 16 INCREMENT operations. Bits that flip to achieve the next value are shaded. The running cost for flipping bits is shown at the right. Notice that the total cost is always less than twice the total number of INCREMENT operations.

operations on an initially zero counter causes  $A[1]$  to flip  $\lfloor n/2 \rfloor$  times. Similarly, bit  $A[2]$  flips only every fourth time, or  $\lfloor n/4 \rfloor$  times in a sequence of  $n$  INCREMENT operations. In general, for  $i = 0, 1, \dots, k-1$ , bit  $A[i]$  flips  $\lfloor n/2^i \rfloor$  times in a sequence of  $n$  INCREMENT operations on an initially zero counter. For  $i \geq k$ , bit  $A[i]$  does not exist, and so it cannot flip. The total number of flips in the sequence is thus

$$\begin{aligned} \sum_{i=0}^{k-1} \left\lfloor \frac{n}{2^i} \right\rfloor &< n \sum_{i=0}^{\infty} \frac{1}{2^i} \\ &= 2n, \end{aligned}$$

by equation (A.6). The worst-case time for a sequence of  $n$  INCREMENT operations on an initially zero counter is therefore  $O(n)$ . The average cost of each operation, and therefore the amortized cost per operation, is  $O(n)/n = O(1)$ .

## Exercises

### 17.1-1

If the set of stack operations included a MULTIPUSH operation, which pushes  $k$  items onto the stack, would the  $O(1)$  bound on the amortized cost of stack operations continue to hold?

### 17.1-2

Show that if a DECREMENT operation were included in the  $k$ -bit counter example,  $n$  operations could cost as much as  $\Theta(nk)$  time.

### 17.1-3

Suppose we perform a sequence of  $n$  operations on a data structure in which the  $i$ th operation costs  $i$  if  $i$  is an exact power of 2, and 1 otherwise. Use aggregate analysis to determine the amortized cost per operation.

---

## 17.2 The accounting method

In the *accounting method* of amortized analysis, we assign differing charges to different operations, with some operations charged more or less than they actually cost. We call the amount we charge an operation its *amortized cost*. When an operation's amortized cost exceeds its actual cost, we assign the difference to specific objects in the data structure as *credit*. Credit can help pay for later operations whose amortized cost is less than their actual cost. Thus, we can view the amortized cost of an operation as being split between its actual cost and credit that is either deposited or used up. Different operations may have different amortized costs. This method differs from aggregate analysis, in which all operations have the same amortized cost.

We must choose the amortized costs of operations carefully. If we want to show that in the worst case the average cost per operation is small by analyzing with amortized costs, we must ensure that the total amortized cost of a sequence of operations provides an upper bound on the total actual cost of the sequence. Moreover, as in aggregate analysis, this relationship must hold for all sequences of operations. If we denote the actual cost of the  $i$ th operation by  $c_i$  and the amortized cost of the  $i$ th operation by  $\hat{c}_i$ , we require

$$\sum_{i=1}^n \hat{c}_i \geq \sum_{i=1}^n c_i \quad (17.1)$$

for all sequences of  $n$  operations. The total credit stored in the data structure is the difference between the total amortized cost and the total actual cost, or

$\sum_{i=1}^n \hat{c}_i - \sum_{i=1}^n c_i$ . By inequality (17.1), the total credit associated with the data structure must be nonnegative at all times. If we ever were to allow the total credit to become negative (the result of undercharging early operations with the promise of repaying the account later on), then the total amortized costs incurred at that time would be below the total actual costs incurred; for the sequence of operations up to that time, the total amortized cost would not be an upper bound on the total actual cost. Thus, we must take care that the total credit in the data structure never becomes negative.

### Stack operations

To illustrate the accounting method of amortized analysis, let us return to the stack example. Recall that the actual costs of the operations were

PUSH            1 ,  
 POP            1 ,  
 MULTIPOP     $\min(k, s)$  ,

where  $k$  is the argument supplied to MULTIPOP and  $s$  is the stack size when it is called. Let us assign the following amortized costs:

PUSH            2 ,  
 POP            0 ,  
 MULTIPOP    0 .

Note that the amortized cost of MULTIPOP is a constant (0), whereas the actual cost is variable. Here, all three amortized costs are constant. In general, the amortized costs of the operations under consideration may differ from each other, and they may even differ asymptotically.

We shall now show that we can pay for any sequence of stack operations by charging the amortized costs. Suppose we use a dollar bill to represent each unit of cost. We start with an empty stack. Recall the analogy of Section 10.1 between the stack data structure and a stack of plates in a cafeteria. When we push a plate on the stack, we use 1 dollar to pay the actual cost of the push and are left with a credit of 1 dollar (out of the 2 dollars charged), which we leave on top of the plate. At any point in time, every plate on the stack has a dollar of credit on it.

The dollar stored on the plate serves as prepayment for the cost of popping it from the stack. When we execute a POP operation, we charge the operation nothing and pay its actual cost using the credit stored in the stack. To pop a plate, we take the dollar of credit off the plate and use it to pay the actual cost of the operation. Thus, by charging the PUSH operation a little bit more, we can charge the POP operation nothing.

Moreover, we can also charge MULTIPOP operations nothing. To pop the first plate, we take the dollar of credit off the plate and use it to pay the actual cost of a POP operation. To pop a second plate, we again have a dollar of credit on the plate to pay for the POP operation, and so on. Thus, we have always charged enough up front to pay for MULTIPOP operations. In other words, since each plate on the stack has 1 dollar of credit on it, and the stack always has a nonnegative number of plates, we have ensured that the amount of credit is always nonnegative. Thus, for *any* sequence of  $n$  PUSH, POP, and MULTIPOP operations, the total amortized cost is an upper bound on the total actual cost. Since the total amortized cost is  $O(n)$ , so is the total actual cost.

### Incrementing a binary counter

As another illustration of the accounting method, we analyze the INCREMENT operation on a binary counter that starts at zero. As we observed earlier, the running time of this operation is proportional to the number of bits flipped, which we shall use as our cost for this example. Let us once again use a dollar bill to represent each unit of cost (the flipping of a bit in this example).

For the amortized analysis, let us charge an amortized cost of 2 dollars to set a bit to 1. When a bit is set, we use 1 dollar (out of the 2 dollars charged) to pay for the actual setting of the bit, and we place the other dollar on the bit as credit to be used later when we flip the bit back to 0. At any point in time, every 1 in the counter has a dollar of credit on it, and thus we can charge nothing to reset a bit to 0; we just pay for the reset with the dollar bill on the bit.

Now we can determine the amortized cost of INCREMENT. The cost of resetting the bits within the **while** loop is paid for by the dollars on the bits that are reset. The INCREMENT procedure sets at most one bit, in line 6, and therefore the amortized cost of an INCREMENT operation is at most 2 dollars. The number of 1s in the counter never becomes negative, and thus the amount of credit stays nonnegative at all times. Thus, for  $n$  INCREMENT operations, the total amortized cost is  $O(n)$ , which bounds the total actual cost.

### Exercises

#### 17.2-1

Suppose we perform a sequence of stack operations on a stack whose size never exceeds  $k$ . After every  $k$  operations, we make a copy of the entire stack for backup purposes. Show that the cost of  $n$  stack operations, including copying the stack, is  $O(n)$  by assigning suitable amortized costs to the various stack operations.

**17.2-2**

Redo Exercise 17.1-3 using an accounting method of analysis.

**17.2-3**

Suppose we wish not only to increment a counter but also to reset it to zero (i.e., make all bits in it 0). Counting the time to examine or modify a bit as  $\Theta(1)$ , show how to implement a counter as an array of bits so that any sequence of  $n$  INCREMENT and RESET operations takes time  $O(n)$  on an initially zero counter. (*Hint*: Keep a pointer to the high-order 1.)

---

**17.3 The potential method**

Instead of representing prepaid work as credit stored with specific objects in the data structure, the *potential method* of amortized analysis represents the prepaid work as “potential energy,” or just “potential,” which can be released to pay for future operations. We associate the potential with the data structure as a whole rather than with specific objects within the data structure.

The potential method works as follows. We will perform  $n$  operations, starting with an initial data structure  $D_0$ . For each  $i = 1, 2, \dots, n$ , we let  $c_i$  be the actual cost of the  $i$ th operation and  $D_i$  be the data structure that results after applying the  $i$ th operation to data structure  $D_{i-1}$ . A *potential function*  $\Phi$  maps each data structure  $D_i$  to a real number  $\Phi(D_i)$ , which is the *potential* associated with data structure  $D_i$ . The *amortized cost*  $\hat{c}_i$  of the  $i$ th operation with respect to potential function  $\Phi$  is defined by

$$\hat{c}_i = c_i + \Phi(D_i) - \Phi(D_{i-1}) . \quad (17.2)$$

The amortized cost of each operation is therefore its actual cost plus the change in potential due to the operation. By equation (17.2), the total amortized cost of the  $n$  operations is

$$\begin{aligned} \sum_{i=1}^n \hat{c}_i &= \sum_{i=1}^n (c_i + \Phi(D_i) - \Phi(D_{i-1})) \\ &= \sum_{i=1}^n c_i + \Phi(D_n) - \Phi(D_0) . \end{aligned} \quad (17.3)$$

The second equality follows from equation (A.9) because the  $\Phi(D_i)$  terms telescope.

If we can define a potential function  $\Phi$  so that  $\Phi(D_n) \geq \Phi(D_0)$ , then the total amortized cost  $\sum_{i=1}^n \hat{c}_i$  gives an upper bound on the total actual cost  $\sum_{i=1}^n c_i$ .



In practice, we do not always know how many operations might be performed. Therefore, if we require that  $\Phi(D_i) \geq \Phi(D_0)$  for all  $i$ , then we guarantee, as in the accounting method, that we pay in advance. We usually just define  $\Phi(D_0)$  to be 0 and then show that  $\Phi(D_i) \geq 0$  for all  $i$ . (See Exercise 17.3-1 for an easy way to handle cases in which  $\Phi(D_0) \neq 0$ .)

Intuitively, if the potential difference  $\Phi(D_i) - \Phi(D_{i-1})$  of the  $i$ th operation is positive, then the amortized cost  $\hat{c}_i$  represents an overcharge to the  $i$ th operation, and the potential of the data structure increases. If the potential difference is negative, then the amortized cost represents an undercharge to the  $i$ th operation, and the decrease in the potential pays for the actual cost of the operation.

The amortized costs defined by equations (17.2) and (17.3) depend on the choice of the potential function  $\Phi$ . Different potential functions may yield different amortized costs yet still be upper bounds on the actual costs. We often find trade-offs that we can make in choosing a potential function; the best potential function to use depends on the desired time bounds.

### Stack operations

To illustrate the potential method, we return once again to the example of the stack operations PUSH, POP, and MULTIPOP. We define the potential function  $\Phi$  on a stack to be the number of objects in the stack. For the empty stack  $D_0$  with which we start, we have  $\Phi(D_0) = 0$ . Since the number of objects in the stack is never negative, the stack  $D_i$  that results after the  $i$ th operation has nonnegative potential, and thus

$$\begin{aligned}\Phi(D_i) &\geq 0 \\ &= \Phi(D_0) .\end{aligned}$$

The total amortized cost of  $n$  operations with respect to  $\Phi$  therefore represents an upper bound on the actual cost.

Let us now compute the amortized costs of the various stack operations. If the  $i$ th operation on a stack containing  $s$  objects is a PUSH operation, then the potential difference is

$$\begin{aligned}\Phi(D_i) - \Phi(D_{i-1}) &= (s + 1) - s \\ &= 1 .\end{aligned}$$

By equation (17.2), the amortized cost of this PUSH operation is

$$\begin{aligned}\hat{c}_i &= c_i + \Phi(D_i) - \Phi(D_{i-1}) \\ &= 1 + 1 \\ &= 2 .\end{aligned}$$

Suppose that the  $i$ th operation on the stack is  $\text{MULTIPOP}(S, k)$ , which causes  $k' = \min(k, s)$  objects to be popped off the stack. The actual cost of the operation is  $k'$ , and the potential difference is

$$\Phi(D_i) - \Phi(D_{i-1}) = -k'.$$

Thus, the amortized cost of the  $\text{MULTIPOP}$  operation is

$$\begin{aligned}\hat{c}_i &= c_i + \Phi(D_i) - \Phi(D_{i-1}) \\ &= k' - k' \\ &= 0.\end{aligned}$$

Similarly, the amortized cost of an ordinary  $\text{POP}$  operation is 0.

The amortized cost of each of the three operations is  $O(1)$ , and thus the total amortized cost of a sequence of  $n$  operations is  $O(n)$ . Since we have already argued that  $\Phi(D_i) \geq \Phi(D_0)$ , the total amortized cost of  $n$  operations is an upper bound on the total actual cost. The worst-case cost of  $n$  operations is therefore  $O(n)$ .

### Incrementing a binary counter

As another example of the potential method, we again look at incrementing a binary counter. This time, we define the potential of the counter after the  $i$ th  $\text{INCREMENT}$  operation to be  $b_i$ , the number of 1s in the counter after the  $i$ th operation.

Let us compute the amortized cost of an  $\text{INCREMENT}$  operation. Suppose that the  $i$ th  $\text{INCREMENT}$  operation resets  $t_i$  bits. The actual cost of the operation is therefore at most  $t_i + 1$ , since in addition to resetting  $t_i$  bits, it sets at most one bit to 1. If  $b_i = 0$ , then the  $i$ th operation resets all  $k$  bits, and so  $b_{i-1} = t_i = k$ . If  $b_i > 0$ , then  $b_i = b_{i-1} - t_i + 1$ . In either case,  $b_i \leq b_{i-1} - t_i + 1$ , and the potential difference is

$$\begin{aligned}\Phi(D_i) - \Phi(D_{i-1}) &\leq (b_{i-1} - t_i + 1) - b_{i-1} \\ &= 1 - t_i.\end{aligned}$$

The amortized cost is therefore

$$\begin{aligned}\hat{c}_i &= c_i + \Phi(D_i) - \Phi(D_{i-1}) \\ &\leq (t_i + 1) + (1 - t_i) \\ &= 2.\end{aligned}$$

If the counter starts at zero, then  $\Phi(D_0) = 0$ . Since  $\Phi(D_i) \geq 0$  for all  $i$ , the total amortized cost of a sequence of  $n$   $\text{INCREMENT}$  operations is an upper bound on the total actual cost, and so the worst-case cost of  $n$   $\text{INCREMENT}$  operations is  $O(n)$ .

The potential method gives us an easy way to analyze the counter even when it does not start at zero. The counter starts with  $b_0$  1s, and after  $n$   $\text{INCREMENT}$

operations it has  $b_n$  1s, where  $0 \leq b_0, b_n \leq k$ . (Recall that  $k$  is the number of bits in the counter.) We can rewrite equation (17.3) as

$$\sum_{i=1}^n c_i = \sum_{i=1}^n \hat{c}_i - \Phi(D_n) + \Phi(D_0). \quad (17.4)$$

We have  $\hat{c}_i \leq 2$  for all  $1 \leq i \leq n$ . Since  $\Phi(D_0) = b_0$  and  $\Phi(D_n) = b_n$ , the total actual cost of  $n$  INCREMENT operations is

$$\begin{aligned} \sum_{i=1}^n c_i &\leq \sum_{i=1}^n 2 - b_n + b_0 \\ &= 2n - b_n + b_0. \end{aligned}$$

Note in particular that since  $b_0 \leq k$ , as long as  $k = O(n)$ , the total actual cost is  $O(n)$ . In other words, if we execute at least  $n = \Omega(k)$  INCREMENT operations, the total actual cost is  $O(n)$ , no matter what initial value the counter contains.

## Exercises

### 17.3-1

Suppose we have a potential function  $\Phi$  such that  $\Phi(D_i) \geq \Phi(D_0)$  for all  $i$ , but  $\Phi(D_0) \neq 0$ . Show that there exists a potential function  $\Phi'$  such that  $\Phi'(D_0) = 0$ ,  $\Phi'(D_i) \geq 0$  for all  $i \geq 1$ , and the amortized costs using  $\Phi'$  are the same as the amortized costs using  $\Phi$ .

### 17.3-2

Redo Exercise 17.1-3 using a potential method of analysis.

### 17.3-3

Consider an ordinary binary min-heap data structure with  $n$  elements supporting the instructions INSERT and EXTRACT-MIN in  $O(\lg n)$  worst-case time. Give a potential function  $\Phi$  such that the amortized cost of INSERT is  $O(\lg n)$  and the amortized cost of EXTRACT-MIN is  $O(1)$ , and show that it works.

### 17.3-4

What is the total cost of executing  $n$  of the stack operations PUSH, POP, and MULTIPOP, assuming that the stack begins with  $s_0$  objects and finishes with  $s_n$  objects?

### 17.3-5

Suppose that a counter begins at a number with  $b$  1s in its binary representation, rather than at 0. Show that the cost of performing  $n$  INCREMENT operations is  $O(n)$  if  $n = \Omega(b)$ . (Do not assume that  $b$  is constant.)

**17.3-6**

Show how to implement a queue with two ordinary stacks (Exercise 10.1-6) so that the amortized cost of each ENQUEUE and each DEQUEUE operation is  $O(1)$ .

**17.3-7**

Design a data structure to support the following two operations for a dynamic multiset  $S$  of integers, which allows duplicate values:

INSERT( $S, x$ ) inserts  $x$  into  $S$ .

DELETE-LARGER-HALF( $S$ ) deletes the largest  $\lceil |S|/2 \rceil$  elements from  $S$ .

Explain how to implement this data structure so that any sequence of  $m$  INSERT and DELETE-LARGER-HALF operations runs in  $O(m)$  time. Your implementation should also include a way to output the elements of  $S$  in  $O(|S|)$  time.

---

**17.4 Dynamic tables**

We do not always know in advance how many objects some applications will store in a table. We might allocate space for a table, only to find out later that it is not enough. We must then reallocate the table with a larger size and copy all objects stored in the original table over into the new, larger table. Similarly, if many objects have been deleted from the table, it may be worthwhile to reallocate the table with a smaller size. In this section, we study this problem of dynamically expanding and contracting a table. Using amortized analysis, we shall show that the amortized cost of insertion and deletion is only  $O(1)$ , even though the actual cost of an operation is large when it triggers an expansion or a contraction. Moreover, we shall see how to guarantee that the unused space in a dynamic table never exceeds a constant fraction of the total space.

We assume that the dynamic table supports the operations TABLE-INSERT and TABLE-DELETE. TABLE-INSERT inserts into the table an item that occupies a single *slot*, that is, a space for one item. Likewise, TABLE-DELETE removes an item from the table, thereby freeing a slot. The details of the data-structuring method used to organize the table are unimportant; we might use a stack (Section 10.1), a heap (Chapter 6), or a hash table (Chapter 11). We might also use an array or collection of arrays to implement object storage, as we did in Section 10.3.

We shall find it convenient to use a concept introduced in our analysis of hashing (Chapter 11). We define the **load factor**  $\alpha(T)$  of a nonempty table  $T$  to be the number of items stored in the table divided by the size (number of slots) of the table. We assign an empty table (one with no items) size 0, and we define its load factor to be 1. If the load factor of a dynamic table is bounded below by a constant,

the unused space in the table is never more than a constant fraction of the total amount of space.

We start by analyzing a dynamic table in which we only insert items. We then consider the more general case in which we both insert and delete items.

### 17.4.1 Table expansion

Let us assume that storage for a table is allocated as an array of slots. A table fills up when all slots have been used or, equivalently, when its load factor is 1.<sup>1</sup> In some software environments, upon attempting to insert an item into a full table, the only alternative is to abort with an error. We shall assume, however, that our software environment, like many modern ones, provides a memory-management system that can allocate and free blocks of storage on request. Thus, upon inserting an item into a full table, we can *expand* the table by allocating a new table with more slots than the old table had. Because we always need the table to reside in contiguous memory, we must allocate a new array for the larger table and then copy items from the old table into the new table.

A common heuristic allocates a new table with twice as many slots as the old one. If the only table operations are insertions, then the load factor of the table is always at least  $1/2$ , and thus the amount of wasted space never exceeds half the total space in the table.

In the following pseudocode, we assume that  $T$  is an object representing the table. The attribute  $T.table$  contains a pointer to the block of storage representing the table,  $T.num$  contains the number of items in the table, and  $T.size$  gives the total number of slots in the table. Initially, the table is empty:  $T.num = T.size = 0$ .

TABLE-INSERT( $T, x$ )

```

1  if  $T.size == 0$ 
2      allocate  $T.table$  with 1 slot
3       $T.size = 1$ 
4  if  $T.num == T.size$ 
5      allocate new-table with  $2 \cdot T.size$  slots
6      insert all items in  $T.table$  into new-table
7      free  $T.table$ 
8       $T.table = \textit{new-table}$ 
9       $T.size = 2 \cdot T.size$ 
10 insert  $x$  into  $T.table$ 
11  $T.num = T.num + 1$ 
```

---

<sup>1</sup>In some situations, such as an open-address hash table, we may wish to consider a table to be full if its load factor equals some constant strictly less than 1. (See Exercise 17.4-1.)

Notice that we have two “insertion” procedures here: the TABLE-INSERT procedure itself and the *elementary insertion* into a table in lines 6 and 10. We can analyze the running time of TABLE-INSERT in terms of the number of elementary insertions by assigning a cost of 1 to each elementary insertion. We assume that the actual running time of TABLE-INSERT is linear in the time to insert individual items, so that the overhead for allocating an initial table in line 2 is constant and the overhead for allocating and freeing storage in lines 5 and 7 is dominated by the cost of transferring items in line 6. We call the event in which lines 5–9 are executed an *expansion*.

Let us analyze a sequence of  $n$  TABLE-INSERT operations on an initially empty table. What is the cost  $c_i$  of the  $i$ th operation? If the current table has room for the new item (or if this is the first operation), then  $c_i = 1$ , since we need only perform the one elementary insertion in line 10. If the current table is full, however, and an expansion occurs, then  $c_i = i$ : the cost is 1 for the elementary insertion in line 10 plus  $i - 1$  for the items that we must copy from the old table to the new table in line 6. If we perform  $n$  operations, the worst-case cost of an operation is  $O(n)$ , which leads to an upper bound of  $O(n^2)$  on the total running time for  $n$  operations.

This bound is not tight, because we rarely expand the table in the course of  $n$  TABLE-INSERT operations. Specifically, the  $i$ th operation causes an expansion only when  $i - 1$  is an exact power of 2. The amortized cost of an operation is in fact  $O(1)$ , as we can show using aggregate analysis. The cost of the  $i$ th operation is

$$c_i = \begin{cases} i & \text{if } i - 1 \text{ is an exact power of } 2, \\ 1 & \text{otherwise.} \end{cases}$$

The total cost of  $n$  TABLE-INSERT operations is therefore

$$\begin{aligned} \sum_{i=1}^n c_i &\leq n + \sum_{j=0}^{\lfloor \lg n \rfloor} 2^j \\ &< n + 2n \\ &= 3n, \end{aligned}$$

because at most  $n$  operations cost 1 and the costs of the remaining operations form a geometric series. Since the total cost of  $n$  TABLE-INSERT operations is bounded by  $3n$ , the amortized cost of a single operation is at most 3.

By using the accounting method, we can gain some feeling for why the amortized cost of a TABLE-INSERT operation should be 3. Intuitively, each item pays for 3 elementary insertions: inserting itself into the current table, moving itself when the table expands, and moving another item that has already been moved once when the table expands. For example, suppose that the size of the table is  $m$  immediately after an expansion. Then the table holds  $m/2$  items, and it contains

no credit. We charge 3 dollars for each insertion. The elementary insertion that occurs immediately costs 1 dollar. We place another dollar as credit on the item inserted. We place the third dollar as credit on one of the  $m/2$  items already in the table. The table will not fill again until we have inserted another  $m/2 - 1$  items, and thus, by the time the table contains  $m$  items and is full, we will have placed a dollar on each item to pay to reinsert it during the expansion.

We can use the potential method to analyze a sequence of  $n$  TABLE-INSERT operations, and we shall use it in Section 17.4.2 to design a TABLE-DELETE operation that has an  $O(1)$  amortized cost as well. We start by defining a potential function  $\Phi$  that is 0 immediately after an expansion but builds to the table size by the time the table is full, so that we can pay for the next expansion by the potential. The function

$$\Phi(T) = 2 \cdot T.num - T.size \quad (17.5)$$

is one possibility. Immediately after an expansion, we have  $T.num = T.size/2$ , and thus  $\Phi(T) = 0$ , as desired. Immediately before an expansion, we have  $T.num = T.size$ , and thus  $\Phi(T) = T.num$ , as desired. The initial value of the potential is 0, and since the table is always at least half full,  $T.num \geq T.size/2$ , which implies that  $\Phi(T)$  is always nonnegative. Thus, the sum of the amortized costs of  $n$  TABLE-INSERT operations gives an upper bound on the sum of the actual costs.

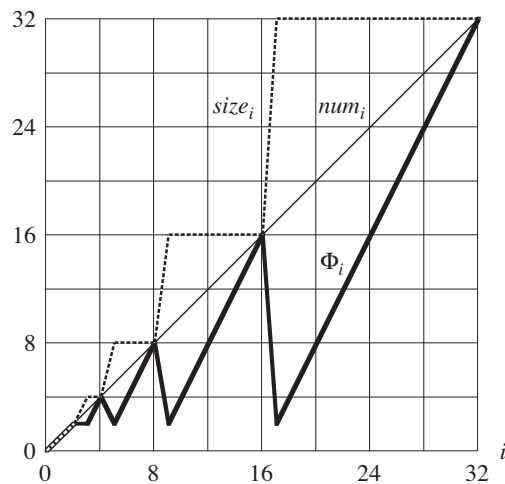
To analyze the amortized cost of the  $i$ th TABLE-INSERT operation, we let  $num_i$  denote the number of items stored in the table after the  $i$ th operation,  $size_i$  denote the total size of the table after the  $i$ th operation, and  $\Phi_i$  denote the potential after the  $i$ th operation. Initially, we have  $num_0 = 0$ ,  $size_0 = 0$ , and  $\Phi_0 = 0$ .

If the  $i$ th TABLE-INSERT operation does not trigger an expansion, then we have  $size_i = size_{i-1}$  and the amortized cost of the operation is

$$\begin{aligned} \hat{c}_i &= c_i + \Phi_i - \Phi_{i-1} \\ &= 1 + (2 \cdot num_i - size_i) - (2 \cdot num_{i-1} - size_{i-1}) \\ &= 1 + (2 \cdot num_i - size_i) - (2(num_i - 1) - size_i) \\ &= 3. \end{aligned}$$

If the  $i$ th operation does trigger an expansion, then we have  $size_i = 2 \cdot size_{i-1}$  and  $size_{i-1} = num_{i-1} = num_i - 1$ , which implies that  $size_i = 2 \cdot (num_i - 1)$ . Thus, the amortized cost of the operation is

$$\begin{aligned} \hat{c}_i &= c_i + \Phi_i - \Phi_{i-1} \\ &= num_i + (2 \cdot num_i - size_i) - (2 \cdot num_{i-1} - size_{i-1}) \\ &= num_i + (2 \cdot num_i - 2 \cdot (num_i - 1)) - (2(num_i - 1) - (num_i - 1)) \\ &= num_i + 2 - (num_i - 1) \\ &= 3. \end{aligned}$$



**Figure 17.3** The effect of a sequence of  $n$  TABLE-INSERT operations on the number  $num_i$  of items in the table, the number  $size_i$  of slots in the table, and the potential  $\Phi_i = 2 \cdot num_i - size_i$ , each being measured after the  $i$ th operation. The thin line shows  $num_i$ , the dashed line shows  $size_i$ , and the thick line shows  $\Phi_i$ . Notice that immediately before an expansion, the potential has built up to the number of items in the table, and therefore it can pay for moving all the items to the new table. Afterwards, the potential drops to 0, but it is immediately increased by 2 upon inserting the item that caused the expansion.

Figure 17.3 plots the values of  $num_i$ ,  $size_i$ , and  $\Phi_i$  against  $i$ . Notice how the potential builds up to pay for expanding the table.

### 17.4.2 Table expansion and contraction

To implement a TABLE-DELETE operation, it is simple enough to remove the specified item from the table. In order to limit the amount of wasted space, however, we might wish to **contract** the table when the load factor becomes too small. Table contraction is analogous to table expansion: when the number of items in the table drops too low, we allocate a new, smaller table and then copy the items from the old table into the new one. We can then free the storage for the old table by returning it to the memory-management system. Ideally, we would like to preserve two properties:

- the load factor of the dynamic table is bounded below by a positive constant, and
- the amortized cost of a table operation is bounded above by a constant.



We assume that we measure the cost in terms of elementary insertions and deletions.

You might think that we should double the table size upon inserting an item into a full table and halve the size when deleting an item would cause the table to become less than half full. This strategy would guarantee that the load factor of the table never drops below  $1/2$ , but unfortunately, it can cause the amortized cost of an operation to be quite large. Consider the following scenario. We perform  $n$  operations on a table  $T$ , where  $n$  is an exact power of 2. The first  $n/2$  operations are insertions, which by our previous analysis cost a total of  $\Theta(n)$ . At the end of this sequence of insertions,  $T.num = T.size = n/2$ . For the second  $n/2$  operations, we perform the following sequence:

insert, delete, delete, insert, insert, delete, delete, insert, insert, . . .

The first insertion causes the table to expand to size  $n$ . The two following deletions cause the table to contract back to size  $n/2$ . Two further insertions cause another expansion, and so forth. The cost of each expansion and contraction is  $\Theta(n)$ , and there are  $\Theta(n)$  of them. Thus, the total cost of the  $n$  operations is  $\Theta(n^2)$ , making the amortized cost of an operation  $\Theta(n)$ .

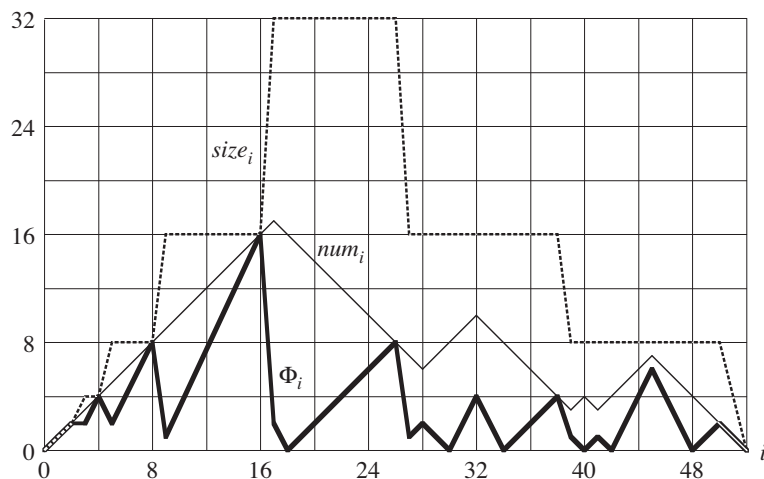
The downside of this strategy is obvious: after expanding the table, we do not delete enough items to pay for a contraction. Likewise, after contracting the table, we do not insert enough items to pay for an expansion.

We can improve upon this strategy by allowing the load factor of the table to drop below  $1/2$ . Specifically, we continue to double the table size upon inserting an item into a full table, but we halve the table size when deleting an item causes the table to become less than  $1/4$  full, rather than  $1/2$  full as before. The load factor of the table is therefore bounded below by the constant  $1/4$ .

Intuitively, we would consider a load factor of  $1/2$  to be ideal, and the table's potential would then be 0. As the load factor deviates from  $1/2$ , the potential increases so that by the time we expand or contract the table, the table has garnered sufficient potential to pay for copying all the items into the newly allocated table. Thus, we will need a potential function that has grown to  $T.num$  by the time that the load factor has either increased to 1 or decreased to  $1/4$ . After either expanding or contracting the table, the load factor goes back to  $1/2$  and the table's potential reduces back to 0.

We omit the code for TABLE-DELETE, since it is analogous to TABLE-INSERT. For our analysis, we shall assume that whenever the number of items in the table drops to 0, we free the storage for the table. That is, if  $T.num = 0$ , then  $T.size = 0$ .

We can now use the potential method to analyze the cost of a sequence of  $n$  TABLE-INSERT and TABLE-DELETE operations. We start by defining a potential function  $\Phi$  that is 0 immediately after an expansion or contraction and builds as the load factor increases to 1 or decreases to  $1/4$ . Let us denote the load fac-



**Figure 17.4** The effect of a sequence of  $n$  TABLE-INSERT and TABLE-DELETE operations on the number  $num_i$  of items in the table, the number  $size_i$  of slots in the table, and the potential

$$\Phi_i = \begin{cases} 2 \cdot num_i - size_i & \text{if } \alpha_i \geq 1/2, \\ size_i/2 - num_i & \text{if } \alpha_i < 1/2, \end{cases}$$

each measured after the  $i$ th operation. The thin line shows  $num_i$ , the dashed line shows  $size_i$ , and the thick line shows  $\Phi_i$ . Notice that immediately before an expansion, the potential has built up to the number of items in the table, and therefore it can pay for moving all the items to the new table. Likewise, immediately before a contraction, the potential has built up to the number of items in the table.

tor of a nonempty table  $T$  by  $\alpha(T) = T.num/T.size$ . Since for an empty table,  $T.num = T.size = 0$  and  $\alpha(T) = 1$ , we always have  $T.num = \alpha(T) \cdot T.size$ , whether the table is empty or not. We shall use as our potential function

$$\Phi(T) = \begin{cases} 2 \cdot T.num - T.size & \text{if } \alpha(T) \geq 1/2, \\ T.size/2 - T.num & \text{if } \alpha(T) < 1/2. \end{cases} \quad (17.6)$$

Observe that the potential of an empty table is 0 and that the potential is never negative. Thus, the total amortized cost of a sequence of operations with respect to  $\Phi$  provides an upper bound on the actual cost of the sequence.

Before proceeding with a precise analysis, we pause to observe some properties of the potential function, as illustrated in Figure 17.4. Notice that when the load factor is  $1/2$ , the potential is 0. When the load factor is 1, we have  $T.size = T.num$ , which implies  $\Phi(T) = T.num$ , and thus the potential can pay for an expansion if an item is inserted. When the load factor is  $1/4$ , we have  $T.size = 4 \cdot T.num$ , which

implies  $\Phi(T) = T.num$ , and thus the potential can pay for a contraction if an item is deleted.

To analyze a sequence of  $n$  TABLE-INSERT and TABLE-DELETE operations, we let  $c_i$  denote the actual cost of the  $i$ th operation,  $\hat{c}_i$  denote its amortized cost with respect to  $\Phi$ ,  $num_i$  denote the number of items stored in the table after the  $i$ th operation,  $size_i$  denote the total size of the table after the  $i$ th operation,  $\alpha_i$  denote the load factor of the table after the  $i$ th operation, and  $\Phi_i$  denote the potential after the  $i$ th operation. Initially,  $num_0 = 0$ ,  $size_0 = 0$ ,  $\alpha_0 = 1$ , and  $\Phi_0 = 0$ .

We start with the case in which the  $i$ th operation is TABLE-INSERT. The analysis is identical to that for table expansion in Section 17.4.1 if  $\alpha_{i-1} \geq 1/2$ . Whether the table expands or not, the amortized cost  $\hat{c}_i$  of the operation is at most 3. If  $\alpha_{i-1} < 1/2$ , the table cannot expand as a result of the operation, since the table expands only when  $\alpha_{i-1} = 1$ . If  $\alpha_i < 1/2$  as well, then the amortized cost of the  $i$ th operation is

$$\begin{aligned}\hat{c}_i &= c_i + \Phi_i - \Phi_{i-1} \\ &= 1 + (size_i/2 - num_i) - (size_{i-1}/2 - num_{i-1}) \\ &= 1 + (size_i/2 - num_i) - (size_i/2 - (num_i - 1)) \\ &= 0.\end{aligned}$$

If  $\alpha_{i-1} < 1/2$  but  $\alpha_i \geq 1/2$ , then

$$\begin{aligned}\hat{c}_i &= c_i + \Phi_i - \Phi_{i-1} \\ &= 1 + (2 \cdot num_i - size_i) - (size_{i-1}/2 - num_{i-1}) \\ &= 1 + (2(num_{i-1} + 1) - size_{i-1}) - (size_{i-1}/2 - num_{i-1}) \\ &= 3 \cdot num_{i-1} - \frac{3}{2}size_{i-1} + 3 \\ &= 3\alpha_{i-1}size_{i-1} - \frac{3}{2}size_{i-1} + 3 \\ &< \frac{3}{2}size_{i-1} - \frac{3}{2}size_{i-1} + 3 \\ &= 3.\end{aligned}$$

Thus, the amortized cost of a TABLE-INSERT operation is at most 3.

We now turn to the case in which the  $i$ th operation is TABLE-DELETE. In this case,  $num_i = num_{i-1} - 1$ . If  $\alpha_{i-1} < 1/2$ , then we must consider whether the operation causes the table to contract. If it does not, then  $size_i = size_{i-1}$  and the amortized cost of the operation is

$$\begin{aligned}\hat{c}_i &= c_i + \Phi_i - \Phi_{i-1} \\ &= 1 + (size_i/2 - num_i) - (size_{i-1}/2 - num_{i-1}) \\ &= 1 + (size_i/2 - num_i) - (size_i/2 - (num_i + 1)) \\ &= 2.\end{aligned}$$

If  $\alpha_{i-1} < 1/2$  and the  $i$ th operation does trigger a contraction, then the actual cost of the operation is  $c_i = \text{num}_i + 1$ , since we delete one item and move  $\text{num}_i$  items. We have  $\text{size}_i/2 = \text{size}_{i-1}/4 = \text{num}_{i-1} = \text{num}_i + 1$ , and the amortized cost of the operation is

$$\begin{aligned}\hat{c}_i &= c_i + \Phi_i - \Phi_{i-1} \\ &= (\text{num}_i + 1) + (\text{size}_i/2 - \text{num}_i) - (\text{size}_{i-1}/2 - \text{num}_{i-1}) \\ &= (\text{num}_i + 1) + ((\text{num}_i + 1) - \text{num}_i) - ((2 \cdot \text{num}_i + 2) - (\text{num}_i + 1)) \\ &= 1.\end{aligned}$$

When the  $i$ th operation is a TABLE-DELETE and  $\alpha_{i-1} \geq 1/2$ , the amortized cost is also bounded above by a constant. We leave the analysis as Exercise 17.4-2.

In summary, since the amortized cost of each operation is bounded above by a constant, the actual time for any sequence of  $n$  operations on a dynamic table is  $O(n)$ .

## Exercises

### 17.4-1

Suppose that we wish to implement a dynamic, open-address hash table. Why might we consider the table to be full when its load factor reaches some value  $\alpha$  that is strictly less than 1? Describe briefly how to make insertion into a dynamic, open-address hash table run in such a way that the expected value of the amortized cost per insertion is  $O(1)$ . Why is the expected value of the actual cost per insertion not necessarily  $O(1)$  for all insertions?

### 17.4-2

Show that if  $\alpha_{i-1} \geq 1/2$  and the  $i$ th operation on a dynamic table is TABLE-DELETE, then the amortized cost of the operation with respect to the potential function (17.6) is bounded above by a constant.

### 17.4-3

Suppose that instead of contracting a table by halving its size when its load factor drops below  $1/4$ , we contract it by multiplying its size by  $2/3$  when its load factor drops below  $1/3$ . Using the potential function

$$\Phi(T) = |2 \cdot T.\text{num} - T.\text{size}|,$$

show that the amortized cost of a TABLE-DELETE that uses this strategy is bounded above by a constant.

---

## Problems

### 17-1 Bit-reversed binary counter

Chapter 30 examines an important algorithm called the fast Fourier transform, or FFT. The first step of the FFT algorithm performs a **bit-reversal permutation** on an input array  $A[0 \dots n-1]$  whose length is  $n = 2^k$  for some nonnegative integer  $k$ . This permutation swaps elements whose indices have binary representations that are the reverse of each other.

We can express each index  $a$  as a  $k$ -bit sequence  $\langle a_{k-1}, a_{k-2}, \dots, a_0 \rangle$ , where  $a = \sum_{i=0}^{k-1} a_i 2^i$ . We define

$$\text{rev}_k(\langle a_{k-1}, a_{k-2}, \dots, a_0 \rangle) = \langle a_0, a_1, \dots, a_{k-1} \rangle;$$

thus,

$$\text{rev}_k(a) = \sum_{i=0}^{k-1} a_{k-i-1} 2^i.$$

For example, if  $n = 16$  (or, equivalently,  $k = 4$ ), then  $\text{rev}_k(3) = 12$ , since the 4-bit representation of 3 is 0011, which when reversed gives 1100, the 4-bit representation of 12.

- a. Given a function  $\text{rev}_k$  that runs in  $\Theta(k)$  time, write an algorithm to perform the bit-reversal permutation on an array of length  $n = 2^k$  in  $O(nk)$  time.

We can use an algorithm based on an amortized analysis to improve the running time of the bit-reversal permutation. We maintain a “bit-reversed counter” and a procedure BIT-REVERSED-INCREMENT that, when given a bit-reversed-counter value  $a$ , produces  $\text{rev}_k(\text{rev}_k(a) + 1)$ . If  $k = 4$ , for example, and the bit-reversed counter starts at 0, then successive calls to BIT-REVERSED-INCREMENT produce the sequence

0000, 1000, 0100, 1100, 0010, 1010,  $\dots = 0, 8, 4, 12, 2, 10, \dots$ .

- b. Assume that the words in your computer store  $k$ -bit values and that in unit time, your computer can manipulate the binary values with operations such as shifting left or right by arbitrary amounts, bitwise-AND, bitwise-OR, etc. Describe an implementation of the BIT-REVERSED-INCREMENT procedure that allows the bit-reversal permutation on an  $n$ -element array to be performed in a total of  $O(n)$  time.
- c. Suppose that you can shift a word left or right by only one bit in unit time. Is it still possible to implement an  $O(n)$ -time bit-reversal permutation?

### 17-2 Making binary search dynamic

Binary search of a sorted array takes logarithmic search time, but the time to insert a new element is linear in the size of the array. We can improve the time for insertion by keeping several sorted arrays.

Specifically, suppose that we wish to support SEARCH and INSERT on a set of  $n$  elements. Let  $k = \lceil \lg(n+1) \rceil$ , and let the binary representation of  $n$  be  $\langle n_{k-1}, n_{k-2}, \dots, n_0 \rangle$ . We have  $k$  sorted arrays  $A_0, A_1, \dots, A_{k-1}$ , where for  $i = 0, 1, \dots, k-1$ , the length of array  $A_i$  is  $2^i$ . Each array is either full or empty, depending on whether  $n_i = 1$  or  $n_i = 0$ , respectively. The total number of elements held in all  $k$  arrays is therefore  $\sum_{i=0}^{k-1} n_i 2^i = n$ . Although each individual array is sorted, elements in different arrays bear no particular relationship to each other.

- a. Describe how to perform the SEARCH operation for this data structure. Analyze its worst-case running time.
- b. Describe how to perform the INSERT operation. Analyze its worst-case and amortized running times.
- c. Discuss how to implement DELETE.

### 17-3 Amortized weight-balanced trees

Consider an ordinary binary search tree augmented by adding to each node  $x$  the attribute  $x.size$  giving the number of keys stored in the subtree rooted at  $x$ . Let  $\alpha$  be a constant in the range  $1/2 \leq \alpha < 1$ . We say that a given node  $x$  is  **$\alpha$ -balanced** if  $x.left.size \leq \alpha \cdot x.size$  and  $x.right.size \leq \alpha \cdot x.size$ . The tree as a whole is  **$\alpha$ -balanced** if every node in the tree is  $\alpha$ -balanced. The following amortized approach to maintaining weight-balanced trees was suggested by G. Varghese.

- a. A  $1/2$ -balanced tree is, in a sense, as balanced as it can be. Given a node  $x$  in an arbitrary binary search tree, show how to rebuild the subtree rooted at  $x$  so that it becomes  $1/2$ -balanced. Your algorithm should run in time  $\Theta(x.size)$ , and it can use  $O(x.size)$  auxiliary storage.
- b. Show that performing a search in an  $n$ -node  $\alpha$ -balanced binary search tree takes  $O(\lg n)$  worst-case time.

For the remainder of this problem, assume that the constant  $\alpha$  is strictly greater than  $1/2$ . Suppose that we implement INSERT and DELETE as usual for an  $n$ -node binary search tree, except that after every such operation, if any node in the tree is no longer  $\alpha$ -balanced, then we “rebuild” the subtree rooted at the highest such node in the tree so that it becomes  $1/2$ -balanced.

We shall analyze this rebuilding scheme using the potential method. For a node  $x$  in a binary search tree  $T$ , we define

$$\Delta(x) = |x.\text{left.size} - x.\text{right.size}| ,$$

and we define the potential of  $T$  as

$$\Phi(T) = c \sum_{x \in T: \Delta(x) \geq 2} \Delta(x) ,$$

where  $c$  is a sufficiently large constant that depends on  $\alpha$ .

- c.* Argue that any binary search tree has nonnegative potential and that a  $1/2$ -balanced tree has potential 0.
- d.* Suppose that  $m$  units of potential can pay for rebuilding an  $m$ -node subtree. How large must  $c$  be in terms of  $\alpha$  in order for it to take  $O(1)$  amortized time to rebuild a subtree that is not  $\alpha$ -balanced?
- e.* Show that inserting a node into or deleting a node from an  $n$ -node  $\alpha$ -balanced tree costs  $O(\lg n)$  amortized time.

#### 17-4 The cost of restructuring red-black trees

There are four basic operations on red-black trees that perform **structural modifications**: node insertions, node deletions, rotations, and color changes. We have seen that RB-INSERT and RB-DELETE use only  $O(1)$  rotations, node insertions, and node deletions to maintain the red-black properties, but they may make many more color changes.

- a.* Describe a legal red-black tree with  $n$  nodes such that calling RB-INSERT to add the  $(n + 1)$ st node causes  $\Omega(\lg n)$  color changes. Then describe a legal red-black tree with  $n$  nodes for which calling RB-DELETE on a particular node causes  $\Omega(\lg n)$  color changes.

Although the worst-case number of color changes per operation can be logarithmic, we shall prove that any sequence of  $m$  RB-INSERT and RB-DELETE operations on an initially empty red-black tree causes  $O(m)$  structural modifications in the worst case. Note that we count each color change as a structural modification.

- b.* Some of the cases handled by the main loop of the code of both RB-INSERT-FIXUP and RB-DELETE-FIXUP are **terminating**: once encountered, they cause the loop to terminate after a constant number of additional operations. For each of the cases of RB-INSERT-FIXUP and RB-DELETE-FIXUP, specify which are terminating and which are not. (*Hint*: Look at Figures 13.5, 13.6, and 13.7.)

We shall first analyze the structural modifications when only insertions are performed. Let  $T$  be a red-black tree, and define  $\Phi(T)$  to be the number of red nodes in  $T$ . Assume that 1 unit of potential can pay for the structural modifications performed by any of the three cases of RB-INSERT-FIXUP.

- c. Let  $T'$  be the result of applying Case 1 of RB-INSERT-FIXUP to  $T$ . Argue that  $\Phi(T') = \Phi(T) - 1$ .
- d. When we insert a node into a red-black tree using RB-INSERT, we can break the operation into three parts. List the structural modifications and potential changes resulting from lines 1–16 of RB-INSERT, from nonterminating cases of RB-INSERT-FIXUP, and from terminating cases of RB-INSERT-FIXUP.
- e. Using part (d), argue that the amortized number of structural modifications performed by any call of RB-INSERT is  $O(1)$ .

We now wish to prove that there are  $O(m)$  structural modifications when there are both insertions and deletions. Let us define, for each node  $x$ ,

$$w(x) = \begin{cases} 0 & \text{if } x \text{ is red ,} \\ 1 & \text{if } x \text{ is black and has no red children ,} \\ 0 & \text{if } x \text{ is black and has one red child ,} \\ 2 & \text{if } x \text{ is black and has two red children .} \end{cases}$$

Now we redefine the potential of a red-black tree  $T$  as

$$\Phi(T) = \sum_{x \in T} w(x) ,$$

and let  $T'$  be the tree that results from applying any nonterminating case of RB-INSERT-FIXUP or RB-DELETE-FIXUP to  $T$ .

- f. Show that  $\Phi(T') \leq \Phi(T) - 1$  for all nonterminating cases of RB-INSERT-FIXUP. Argue that the amortized number of structural modifications performed by any call of RB-INSERT-FIXUP is  $O(1)$ .
- g. Show that  $\Phi(T') \leq \Phi(T) - 1$  for all nonterminating cases of RB-DELETE-FIXUP. Argue that the amortized number of structural modifications performed by any call of RB-DELETE-FIXUP is  $O(1)$ .
- h. Complete the proof that in the worst case, any sequence of  $m$  RB-INSERT and RB-DELETE operations performs  $O(m)$  structural modifications.



### 17-5 Competitive analysis of self-organizing lists with move-to-front

A *self-organizing list* is a linked list of  $n$  elements, in which each element has a unique key. When we search for an element in the list, we are given a key, and we want to find an element with that key.

A self-organizing list has two important properties:

1. To find an element in the list, given its key, we must traverse the list from the beginning until we encounter the element with the given key. If that element is the  $k$ th element from the start of the list, then the cost to find the element is  $k$ .
2. We may reorder the list elements after any operation, according to a given rule with a given cost. We may choose any heuristic we like to decide how to reorder the list.

Assume that we start with a given list of  $n$  elements, and we are given an access sequence  $\sigma = \langle \sigma_1, \sigma_2, \dots, \sigma_m \rangle$  of keys to find, in order. The cost of the sequence is the sum of the costs of the individual accesses in the sequence.

Out of the various possible ways to reorder the list after an operation, this problem focuses on transposing adjacent list elements—switching their positions in the list—with a unit cost for each transpose operation. You will show, by means of a potential function, that a particular heuristic for reordering the list, move-to-front, entails a total cost no worse than 4 times that of any other heuristic for maintaining the list order—even if the other heuristic knows the access sequence in advance! We call this type of analysis a *competitive analysis*.

For a heuristic  $H$  and a given initial ordering of the list, denote the access cost of sequence  $\sigma$  by  $C_H(\sigma)$ . Let  $m$  be the number of accesses in  $\sigma$ .

- a. Argue that if heuristic  $H$  does not know the access sequence in advance, then the worst-case cost for  $H$  on an access sequence  $\sigma$  is  $C_H(\sigma) = \Omega(mn)$ .

With the *move-to-front* heuristic, immediately after searching for an element  $x$ , we move  $x$  to the first position on the list (i.e., the front of the list).

Let  $\text{rank}_L(x)$  denote the rank of element  $x$  in list  $L$ , that is, the position of  $x$  in list  $L$ . For example, if  $x$  is the fourth element in  $L$ , then  $\text{rank}_L(x) = 4$ . Let  $c_i$  denote the cost of access  $\sigma_i$  using the move-to-front heuristic, which includes the cost of finding the element in the list and the cost of moving it to the front of the list by a series of transpositions of adjacent list elements.

- b. Show that if  $\sigma_i$  accesses element  $x$  in list  $L$  using the move-to-front heuristic, then  $c_i = 2 \cdot \text{rank}_L(x) - 1$ .

Now we compare move-to-front with any other heuristic  $H$  that processes an access sequence according to the two properties above. Heuristic  $H$  may transpose

elements in the list in any way it wants, and it might even know the entire access sequence in advance.

Let  $L_i$  be the list after access  $\sigma_i$  using move-to-front, and let  $L_i^*$  be the list after access  $\sigma_i$  using heuristic H. We denote the cost of access  $\sigma_i$  by  $c_i$  for move-to-front and by  $c_i^*$  for heuristic H. Suppose that heuristic H performs  $t_i^*$  transpositions during access  $\sigma_i$ .

c. In part (b), you showed that  $c_i = 2 \cdot \text{rank}_{L_{i-1}}(x) - 1$ . Now show that  $c_i^* = \text{rank}_{L_{i-1}^*}(x) + t_i^*$ .

We define an ***inversion*** in list  $L_i$  as a pair of elements  $y$  and  $z$  such that  $y$  precedes  $z$  in  $L_i$  and  $z$  precedes  $y$  in list  $L_i^*$ . Suppose that list  $L_i$  has  $q_i$  inversions after processing the access sequence  $\langle \sigma_1, \sigma_2, \dots, \sigma_i \rangle$ . Then, we define a potential function  $\Phi$  that maps  $L_i$  to a real number by  $\Phi(L_i) = 2q_i$ . For example, if  $L_i$  has the elements  $\langle e, c, a, d, b \rangle$  and  $L_i^*$  has the elements  $\langle c, a, b, d, e \rangle$ , then  $L_i$  has 5 inversions  $((e, c), (e, a), (e, d), (e, b), (d, b))$ , and so  $\Phi(L_i) = 10$ . Observe that  $\Phi(L_i) \geq 0$  for all  $i$  and that, if move-to-front and heuristic H start with the same list  $L_0$ , then  $\Phi(L_0) = 0$ .

d. Argue that a transposition either increases the potential by 2 or decreases the potential by 2.

Suppose that access  $\sigma_i$  finds the element  $x$ . To understand how the potential changes due to  $\sigma_i$ , let us partition the elements other than  $x$  into four sets, depending on where they are in the lists just before the  $i$ th access:

- Set  $A$  consists of elements that precede  $x$  in both  $L_{i-1}$  and  $L_{i-1}^*$ .
  - Set  $B$  consists of elements that precede  $x$  in  $L_{i-1}$  and follow  $x$  in  $L_{i-1}^*$ .
  - Set  $C$  consists of elements that follow  $x$  in  $L_{i-1}$  and precede  $x$  in  $L_{i-1}^*$ .
  - Set  $D$  consists of elements that follow  $x$  in both  $L_{i-1}$  and  $L_{i-1}^*$ .
- e. Argue that  $\text{rank}_{L_{i-1}}(x) = |A| + |B| + 1$  and  $\text{rank}_{L_{i-1}^*}(x) = |A| + |C| + 1$ .

f. Show that access  $\sigma_i$  causes a change in potential of

$$\Phi(L_i) - \Phi(L_{i-1}) \leq 2(|A| - |B| + t_i^*),$$

where, as before, heuristic H performs  $t_i^*$  transpositions during access  $\sigma_i$ .

Define the amortized cost  $\hat{c}_i$  of access  $\sigma_i$  by  $\hat{c}_i = c_i + \Phi(L_i) - \Phi(L_{i-1})$ .

g. Show that the amortized cost  $\hat{c}_i$  of access  $\sigma_i$  is bounded from above by  $4c_i^*$ .

h. Conclude that the cost  $C_{\text{MTF}}(\sigma)$  of access sequence  $\sigma$  with move-to-front is at most 4 times the cost  $C_H(\sigma)$  of  $\sigma$  with any other heuristic H, assuming that both heuristics start with the same list.

---

## Chapter notes

Aho, Hopcroft, and Ullman [5] used aggregate analysis to determine the running time of operations on a disjoint-set forest; we shall analyze this data structure using the potential method in Chapter 21. Tarjan [331] surveys the accounting and potential methods of amortized analysis and presents several applications. He attributes the accounting method to several authors, including M. R. Brown, R. E. Tarjan, S. Huddleston, and K. Mehlhorn. He attributes the potential method to D. D. Sleator. The term “amortized” is due to D. D. Sleator and R. E. Tarjan.

Potential functions are also useful for proving lower bounds for certain types of problems. For each configuration of the problem, we define a potential function that maps the configuration to a real number. Then we determine the potential  $\Phi_{\text{init}}$  of the initial configuration, the potential  $\Phi_{\text{final}}$  of the final configuration, and the maximum change in potential  $\Delta\Phi_{\text{max}}$  due to any step. The number of steps must therefore be at least  $|\Phi_{\text{final}} - \Phi_{\text{init}}| / |\Delta\Phi_{\text{max}}|$ . Examples of potential functions to prove lower bounds in I/O complexity appear in works by Cormen, Sundquist, and Wisniewski [79]; Floyd [107]; and Aggarwal and Vitter [3]. Krumme, Cybenko, and Venkataraman [221] applied potential functions to prove lower bounds on *gossiping*: communicating a unique item from each vertex in a graph to every other vertex.

The move-to-front heuristic from Problem 17-5 works quite well in practice. Moreover, if we recognize that when we find an element, we can splice it out of its position in the list and relocate it to the front of the list in constant time, we can show that the cost of move-to-front is at most twice the cost of any other heuristic including, again, one that knows the entire access sequence in advance.



---

## *V   Advanced Data Structures*

---

## Introduction

This part returns to studying data structures that support operations on dynamic sets, but at a more advanced level than Part III. Two of the chapters, for example, make extensive use of the amortized analysis techniques we saw in Chapter 17.

Chapter 18 presents B-trees, which are balanced search trees specifically designed to be stored on disks. Because disks operate much more slowly than random-access memory, we measure the performance of B-trees not only by how much computing time the dynamic-set operations consume but also by how many disk accesses they perform. For each B-tree operation, the number of disk accesses increases with the height of the B-tree, but B-tree operations keep the height low.

Chapter 19 gives an implementation of a mergeable heap, which supports the operations INSERT, MINIMUM, EXTRACT-MIN, and UNION.<sup>1</sup> The UNION operation unites, or merges, two heaps. Fibonacci heaps—the data structure in Chapter 19—also support the operations DELETE and DECREASE-KEY. We use amortized time bounds to measure the performance of Fibonacci heaps. The operations INSERT, MINIMUM, and UNION take only  $O(1)$  actual and amortized time on Fibonacci heaps, and the operations EXTRACT-MIN and DELETE take  $O(\lg n)$  amortized time. The most significant advantage of Fibonacci heaps, however, is that DECREASE-KEY takes only  $O(1)$  amortized time. Because the DECREASE-

---

<sup>1</sup>As in Problem 10-2, we have defined a mergeable heap to support MINIMUM and EXTRACT-MIN, and so we can also refer to it as a *mergeable min-heap*. Alternatively, if it supported MAXIMUM and EXTRACT-MAX, it would be a *mergeable max-heap*. Unless we specify otherwise, mergeable heaps will be by default mergeable min-heaps.

KEY operation takes constant amortized time, Fibonacci heaps are key components of some of the asymptotically fastest algorithms to date for graph problems.

Noting that we can beat the  $\Omega(n \lg n)$  lower bound for sorting when the keys are integers in a restricted range, Chapter 20 asks whether we can design a data structure that supports the dynamic-set operations SEARCH, INSERT, DELETE, MINIMUM, MAXIMUM, SUCCESSOR, and PREDECESSOR in  $o(\lg n)$  time when the keys are integers in a restricted range. The answer turns out to be that we can, by using a recursive data structure known as a van Emde Boas tree. If the keys are unique integers drawn from the set  $\{0, 1, 2, \dots, u - 1\}$ , where  $u$  is an exact power of 2, then van Emde Boas trees support each of the above operations in  $O(\lg \lg u)$  time.

Finally, Chapter 21 presents data structures for disjoint sets. We have a universe of  $n$  elements that are partitioned into dynamic sets. Initially, each element belongs to its own singleton set. The operation UNION unites two sets, and the query FIND-SET identifies the unique set that contains a given element at the moment. By representing each set as a simple rooted tree, we obtain surprisingly fast operations: a sequence of  $m$  operations runs in  $O(m \alpha(n))$  time, where  $\alpha(n)$  is an incredibly slowly growing function— $\alpha(n)$  is at most 4 in any conceivable application. The amortized analysis that proves this time bound is as complex as the data structure is simple.

The topics covered in this part are by no means the only examples of “advanced” data structures. Other advanced data structures include the following:

- **Dynamic trees**, introduced by Sleator and Tarjan [319] and discussed by Tarjan [330], maintain a forest of disjoint rooted trees. Each edge in each tree has a real-valued cost. Dynamic trees support queries to find parents, roots, edge costs, and the minimum edge cost on a simple path from a node up to a root. Trees may be manipulated by cutting edges, updating all edge costs on a simple path from a node up to a root, linking a root into another tree, and making a node the root of the tree it appears in. One implementation of dynamic trees gives an  $O(\lg n)$  amortized time bound for each operation; a more complicated implementation yields  $O(\lg n)$  worst-case time bounds. Dynamic trees are used in some of the asymptotically fastest network-flow algorithms.
- **Splay trees**, developed by Sleator and Tarjan [320] and, again, discussed by Tarjan [330], are a form of binary search tree on which the standard search-tree operations run in  $O(\lg n)$  amortized time. One application of splay trees simplifies dynamic trees.
- **Persistent** data structures allow queries, and sometimes updates as well, on past versions of a data structure. Driscoll, Sarnak, Sleator, and Tarjan [97] present techniques for making linked data structures persistent with only a small time

and space cost. Problem 13-1 gives a simple example of a persistent dynamic set.

- As in Chapter 20, several data structures allow a faster implementation of dictionary operations (INSERT, DELETE, and SEARCH) for a restricted universe of keys. By taking advantage of these restrictions, they are able to achieve better worst-case asymptotic running times than comparison-based data structures. Fredman and Willard introduced *fusion trees* [115], which were the first data structure to allow faster dictionary operations when the universe is restricted to integers. They showed how to implement these operations in  $O(\lg n / \lg \lg n)$  time. Several subsequent data structures, including *exponential search trees* [16], have also given improved bounds on some or all of the dictionary operations and are mentioned in the chapter notes throughout this book.
- *Dynamic graph data structures* support various queries while allowing the structure of a graph to change through operations that insert or delete vertices or edges. Examples of the queries that they support include vertex connectivity [166], edge connectivity, minimum spanning trees [165], biconnectivity, and transitive closure [164].

Chapter notes throughout this book mention additional data structures.



B-trees are balanced search trees designed to work well on disks or other direct-access secondary storage devices. B-trees are similar to red-black trees (Chapter 13), but they are better at minimizing disk I/O operations. Many database systems use B-trees, or variants of B-trees, to store information.

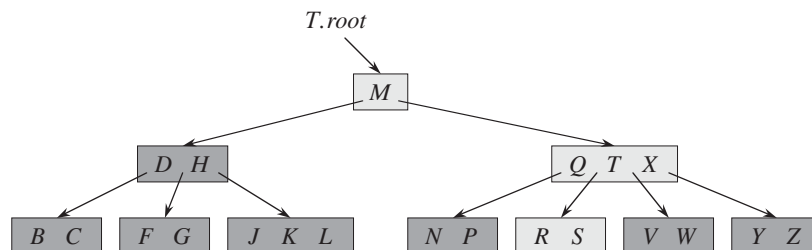
B-trees differ from red-black trees in that B-tree nodes may have many children, from a few to thousands. That is, the “branching factor” of a B-tree can be quite large, although it usually depends on characteristics of the disk unit used. B-trees are similar to red-black trees in that every  $n$ -node B-tree has height  $O(\lg n)$ . The exact height of a B-tree can be considerably less than that of a red-black tree, however, because its branching factor, and hence the base of the logarithm that expresses its height, can be much larger. Therefore, we can also use B-trees to implement many dynamic-set operations in time  $O(\lg n)$ .

B-trees generalize binary search trees in a natural manner. Figure 18.1 shows a simple B-tree. If an internal B-tree node  $x$  contains  $x.n$  keys, then  $x$  has  $x.n + 1$  children. The keys in node  $x$  serve as dividing points separating the range of keys handled by  $x$  into  $x.n + 1$  subranges, each handled by one child of  $x$ . When searching for a key in a B-tree, we make an  $(x.n + 1)$ -way decision based on comparisons with the  $x.n$  keys stored at node  $x$ . The structure of leaf nodes differs from that of internal nodes; we will examine these differences in Section 18.1.

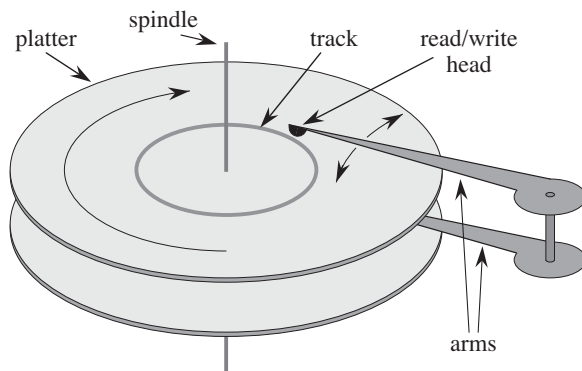
Section 18.1 gives a precise definition of B-trees and proves that the height of a B-tree grows only logarithmically with the number of nodes it contains. Section 18.2 describes how to search for a key and insert a key into a B-tree, and Section 18.3 discusses deletion. Before proceeding, however, we need to ask why we evaluate data structures designed to work on a disk differently from data structures designed to work in main random-access memory.

### Data structures on secondary storage

Computer systems take advantage of various technologies that provide memory capacity. The *primary memory* (or *main memory*) of a computer system normally



**Figure 18.1** A B-tree whose keys are the consonants of English. An internal node  $x$  containing  $x.n$  keys has  $x.n + 1$  children. All leaves are at the same depth in the tree. The lightly shaded nodes are examined in a search for the letter  $R$ .



**Figure 18.2** A typical disk drive. It comprises one or more platters (two platters are shown here) that rotate around a spindle. Each platter is read and written with a head at the end of an arm. Arms rotate around a common pivot axis. A track is the surface that passes beneath the read/write head when the head is stationary.

consists of silicon memory chips. This technology is typically more than an order of magnitude more expensive per bit stored than magnetic storage technology, such as tapes or disks. Most computer systems also have *secondary storage* based on magnetic disks; the amount of such secondary storage often exceeds the amount of primary memory by at least two orders of magnitude.

Figure 18.2 shows a typical disk drive. The drive consists of one or more *platters*, which rotate at a constant speed around a common *spindle*. A magnetizable material covers the surface of each platter. The drive reads and writes each platter by a *head* at the end of an *arm*. The arms can move their heads toward or away

from the spindle. When a given head is stationary, the surface that passes underneath it is called a *track*. Multiple platters increase only the disk drive's capacity and not its performance.

Although disks are cheaper and have higher capacity than main memory, they are much, much slower because they have moving mechanical parts.<sup>1</sup> The mechanical motion has two components: platter rotation and arm movement. As of this writing, commodity disks rotate at speeds of 5400–15,000 revolutions per minute (RPM). We typically see 15,000 RPM speeds in server-grade drives, 7200 RPM speeds in drives for desktops, and 5400 RPM speeds in drives for laptops. Although 7200 RPM may seem fast, one rotation takes 8.33 milliseconds, which is over 5 orders of magnitude longer than the 50 nanosecond access times (more or less) commonly found for silicon memory. In other words, if we have to wait a full rotation for a particular item to come under the read/write head, we could access main memory more than 100,000 times during that span. On average we have to wait for only half a rotation, but still, the difference in access times for silicon memory compared with disks is enormous. Moving the arms also takes some time. As of this writing, average access times for commodity disks are in the range of 8 to 11 milliseconds.

In order to amortize the time spent waiting for mechanical movements, disks access not just one item but several at a time. Information is divided into a number of equal-sized *pages* of bits that appear consecutively within tracks, and each disk read or write is of one or more entire pages. For a typical disk, a page might be  $2^{11}$  to  $2^{14}$  bytes in length. Once the read/write head is positioned correctly and the disk has rotated to the beginning of the desired page, reading or writing a magnetic disk is entirely electronic (aside from the rotation of the disk), and the disk can quickly read or write large amounts of data.

Often, accessing a page of information and reading it from a disk takes longer than examining all the information read. For this reason, in this chapter we shall look separately at the two principal components of the running time:

- the number of disk accesses, and
- the CPU (computing) time.

We measure the number of disk accesses in terms of the number of pages of information that need to be read from or written to the disk. We note that disk-access time is not constant—it depends on the distance between the current track and the desired track and also on the initial rotational position of the disk. We shall

---

<sup>1</sup>As of this writing, solid-state drives have recently come onto the consumer market. Although they are faster than mechanical disk drives, they cost more per gigabyte and have lower capacities than mechanical disk drives.

nonetheless use the number of pages read or written as a first-order approximation of the total time spent accessing the disk.

In a typical B-tree application, the amount of data handled is so large that all the data do not fit into main memory at once. The B-tree algorithms copy selected pages from disk into main memory as needed and write back onto disk the pages that have changed. B-tree algorithms keep only a constant number of pages in main memory at any time; thus, the size of main memory does not limit the size of B-trees that can be handled.

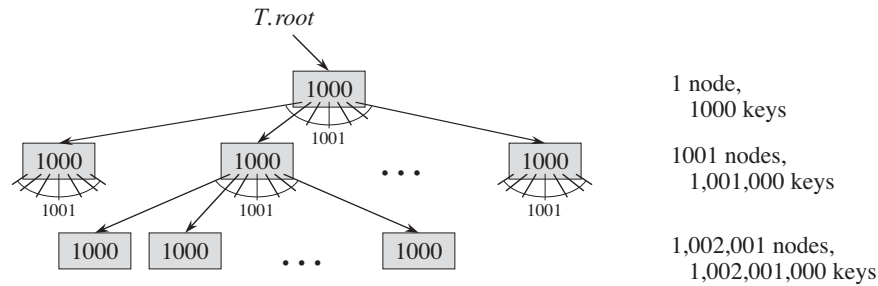
We model disk operations in our pseudocode as follows. Let  $x$  be a pointer to an object. If the object is currently in the computer's main memory, then we can refer to the attributes of the object as usual:  $x.key$ , for example. If the object referred to by  $x$  resides on disk, however, then we must perform the operation  $\text{DISK-READ}(x)$  to read object  $x$  into main memory before we can refer to its attributes. (We assume that if  $x$  is already in main memory, then  $\text{DISK-READ}(x)$  requires no disk accesses; it is a "no-op.") Similarly, the operation  $\text{DISK-WRITE}(x)$  is used to save any changes that have been made to the attributes of object  $x$ . That is, the typical pattern for working with an object is as follows:

```
 $x$  = a pointer to some object
DISK-READ( $x$ )
operations that access and/or modify the attributes of  $x$ 
DISK-WRITE( $x$ )           // omitted if no attributes of  $x$  were changed
other operations that access but do not modify attributes of  $x$ 
```

The system can keep only a limited number of pages in main memory at any one time. We shall assume that the system flushes from main memory pages no longer in use; our B-tree algorithms will ignore this issue.

Since in most systems the running time of a B-tree algorithm depends primarily on the number of  $\text{DISK-READ}$  and  $\text{DISK-WRITE}$  operations it performs, we typically want each of these operations to read or write as much information as possible. Thus, a B-tree node is usually as large as a whole disk page, and this size limits the number of children a B-tree node can have.

For a large B-tree stored on a disk, we often see branching factors between 50 and 2000, depending on the size of a key relative to the size of a page. A large branching factor dramatically reduces both the height of the tree and the number of disk accesses required to find any key. Figure 18.3 shows a B-tree with a branching factor of 1001 and height 2 that can store over one billion keys; nevertheless, since we can keep the root node permanently in main memory, we can find any key in this tree by making at most only two disk accesses.



**Figure 18.3** A B-tree of height 2 containing over one billion keys. Shown inside each node  $x$  is  $x.n$ , the number of keys in  $x$ . Each internal node and leaf contains 1000 keys. This B-tree has 1001 nodes at depth 1 and over one million leaves at depth 2.

## 18.1 Definition of B-trees

To keep things simple, we assume, as we have for binary search trees and red-black trees, that any “satellite information” associated with a key resides in the same node as the key. In practice, one might actually store with each key just a pointer to another disk page containing the satellite information for that key. The pseudocode in this chapter implicitly assumes that the satellite information associated with a key, or the pointer to such satellite information, travels with the key whenever the key is moved from node to node. A common variant on a B-tree, known as a  **$B^+$ -tree**, stores all the satellite information in the leaves and stores only keys and child pointers in the internal nodes, thus maximizing the branching factor of the internal nodes.

A **B-tree**  $T$  is a rooted tree (whose root is  $T.root$ ) having the following properties:

1. Every node  $x$  has the following attributes:
  - a.  $x.n$ , the number of keys currently stored in node  $x$ ,
  - b. the  $x.n$  keys themselves,  $x.key_1, x.key_2, \dots, x.key_{x.n}$ , stored in nondecreasing order, so that  $x.key_1 \leq x.key_2 \leq \dots \leq x.key_{x.n}$ ,
  - c.  $x.leaf$ , a boolean value that is TRUE if  $x$  is a leaf and FALSE if  $x$  is an internal node.
2. Each internal node  $x$  also contains  $x.n + 1$  pointers  $x.c_1, x.c_2, \dots, x.c_{x.n+1}$  to its children. Leaf nodes have no children, and so their  $c_i$  attributes are undefined.

3. The keys  $x.key_i$  separate the ranges of keys stored in each subtree: if  $k_i$  is any key stored in the subtree with root  $x.c_i$ , then

$$k_1 \leq x.key_1 \leq k_2 \leq x.key_2 \leq \cdots \leq x.key_{x.n} \leq k_{x.n+1} .$$

4. All leaves have the same depth, which is the tree's height  $h$ .
5. Nodes have lower and upper bounds on the number of keys they can contain. We express these bounds in terms of a fixed integer  $t \geq 2$  called the **minimum degree** of the B-tree:
- a. Every node other than the root must have at least  $t - 1$  keys. Every internal node other than the root thus has at least  $t$  children. If the tree is nonempty, the root must have at least one key.
  - b. Every node may contain at most  $2t - 1$  keys. Therefore, an internal node may have at most  $2t$  children. We say that a node is **full** if it contains exactly  $2t - 1$  keys.<sup>2</sup>

The simplest B-tree occurs when  $t = 2$ . Every internal node then has either 2, 3, or 4 children, and we have a **2-3-4 tree**. In practice, however, much larger values of  $t$  yield B-trees with smaller height.

### The height of a B-tree

The number of disk accesses required for most operations on a B-tree is proportional to the height of the B-tree. We now analyze the worst-case height of a B-tree.

#### **Theorem 18.1**

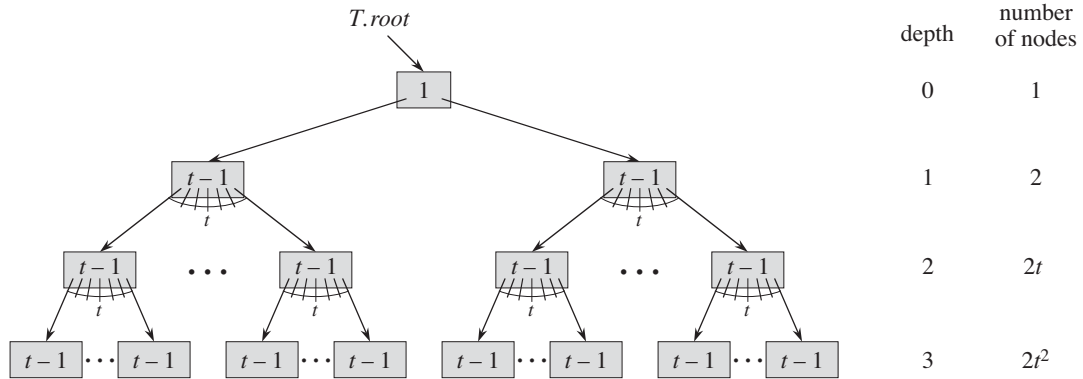
If  $n \geq 1$ , then for any  $n$ -key B-tree  $T$  of height  $h$  and minimum degree  $t \geq 2$ ,

$$h \leq \log_t \frac{n + 1}{2} .$$

**Proof** The root of a B-tree  $T$  contains at least one key, and all other nodes contain at least  $t - 1$  keys. Thus,  $T$ , whose height is  $h$ , has at least 2 nodes at depth 1, at least  $2t$  nodes at depth 2, at least  $2t^2$  nodes at depth 3, and so on, until at depth  $h$  it has at least  $2t^{h-1}$  nodes. Figure 18.4 illustrates such a tree for  $h = 3$ . Thus, the

---

<sup>2</sup>Another common variant on a B-tree, known as a **B\*-tree**, requires each internal node to be at least  $2/3$  full, rather than at least half full, as a B-tree requires.



**Figure 18.4** A B-tree of height 3 containing a minimum possible number of keys. Shown inside each node  $x$  is  $x.n$ .

number  $n$  of keys satisfies the inequality

$$\begin{aligned}
 n &\geq 1 + (t-1) \sum_{i=1}^h 2t^{i-1} \\
 &= 1 + 2(t-1) \left( \frac{t^h - 1}{t-1} \right) \\
 &= 2t^h - 1.
 \end{aligned}$$

By simple algebra, we get  $t^h \leq (n+1)/2$ . Taking base- $t$  logarithms of both sides proves the theorem. ■

Here we see the power of B-trees, as compared with red-black trees. Although the height of the tree grows as  $O(\lg n)$  in both cases (recall that  $t$  is a constant), for B-trees the base of the logarithm can be many times larger. Thus, B-trees save a factor of about  $\lg t$  over red-black trees in the number of nodes examined for most tree operations. Because we usually have to access the disk to examine an arbitrary node in a tree, B-trees avoid a substantial number of disk accesses.

## Exercises

### 18.1-1

Why don't we allow a minimum degree of  $t = 1$ ?

### 18.1-2

For what values of  $t$  is the tree of Figure 18.1 a legal B-tree?

**18.1-3**

Show all legal B-trees of minimum degree 2 that represent  $\{1, 2, 3, 4, 5\}$ .

**18.1-4**

As a function of the minimum degree  $t$ , what is the maximum number of keys that can be stored in a B-tree of height  $h$ ?

**18.1-5**

Describe the data structure that would result if each black node in a red-black tree were to absorb its red children, incorporating their children with its own.

---

## 18.2 Basic operations on B-trees

In this section, we present the details of the operations B-TREE-SEARCH, B-TREE-CREATE, and B-TREE-INSERT. In these procedures, we adopt two conventions:

- The root of the B-tree is always in main memory, so that we never need to perform a DISK-READ on the root; we do have to perform a DISK-WRITE of the root, however, whenever the root node is changed.
- Any nodes that are passed as parameters must already have had a DISK-READ operation performed on them.

The procedures we present are all “one-pass” algorithms that proceed downward from the root of the tree, without having to back up.

### Searching a B-tree

Searching a B-tree is much like searching a binary search tree, except that instead of making a binary, or “two-way,” branching decision at each node, we make a multiway branching decision according to the number of the node’s children. More precisely, at each internal node  $x$ , we make an  $(x.n + 1)$ -way branching decision.

B-TREE-SEARCH is a straightforward generalization of the TREE-SEARCH procedure defined for binary search trees. B-TREE-SEARCH takes as input a pointer to the root node  $x$  of a subtree and a key  $k$  to be searched for in that subtree. The top-level call is thus of the form B-TREE-SEARCH( $T.root, k$ ). If  $k$  is in the B-tree, B-TREE-SEARCH returns the ordered pair  $(y, i)$  consisting of a node  $y$  and an index  $i$  such that  $y.key_i = k$ . Otherwise, the procedure returns NIL.



```

B-TREE-SEARCH( $x, k$ )
1   $i = 1$ 
2  while  $i \leq x.n$  and  $k > x.key_i$ 
3       $i = i + 1$ 
4  if  $i \leq x.n$  and  $k == x.key_i$ 
5      return  $(x, i)$ 
6  elseif  $x.leaf$ 
7      return NIL
8  else DISK-READ( $x.c_i$ )
9      return B-TREE-SEARCH( $x.c_i, k$ )

```

Using a linear-search procedure, lines 1–3 find the smallest index  $i$  such that  $k \leq x.key_i$ , or else they set  $i$  to  $x.n + 1$ . Lines 4–5 check to see whether we have now discovered the key, returning if we have. Otherwise, lines 6–9 either terminate the search unsuccessfully (if  $x$  is a leaf) or recurse to search the appropriate subtree of  $x$ , after performing the necessary DISK-READ on that child.

Figure 18.1 illustrates the operation of B-TREE-SEARCH. The procedure examines the lightly shaded nodes during a search for the key  $R$ .

As in the TREE-SEARCH procedure for binary search trees, the nodes encountered during the recursion form a simple path downward from the root of the tree. The B-TREE-SEARCH procedure therefore accesses  $O(h) = O(\log_t n)$  disk pages, where  $h$  is the height of the B-tree and  $n$  is the number of keys in the B-tree. Since  $x.n < 2t$ , the **while** loop of lines 2–3 takes  $O(t)$  time within each node, and the total CPU time is  $O(th) = O(t \log_t n)$ .

### Creating an empty B-tree

To build a B-tree  $T$ , we first use B-TREE-CREATE to create an empty root node and then call B-TREE-INSERT to add new keys. Both of these procedures use an auxiliary procedure ALLOCATE-NODE, which allocates one disk page to be used as a new node in  $O(1)$  time. We can assume that a node created by ALLOCATE-NODE requires no DISK-READ, since there is as yet no useful information stored on the disk for that node.

```

B-TREE-CREATE( $T$ )
1   $x = \text{ALLOCATE-NODE}()$ 
2   $x.leaf = \text{TRUE}$ 
3   $x.n = 0$ 
4  DISK-WRITE( $x$ )
5   $T.root = x$ 

```

B-TREE-CREATE requires  $O(1)$  disk operations and  $O(1)$  CPU time.

### Inserting a key into a B-tree

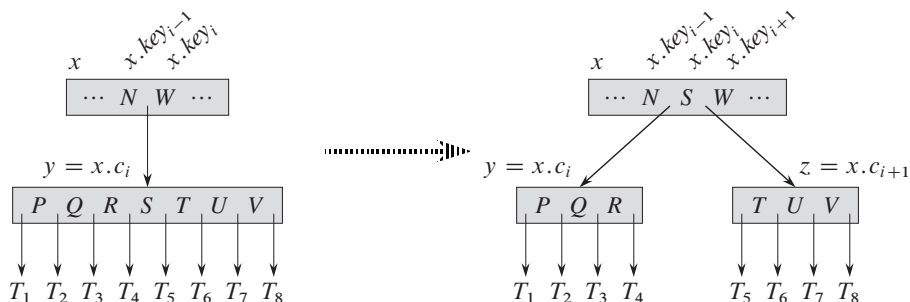
Inserting a key into a B-tree is significantly more complicated than inserting a key into a binary search tree. As with binary search trees, we search for the leaf position at which to insert the new key. With a B-tree, however, we cannot simply create a new leaf node and insert it, as the resulting tree would fail to be a valid B-tree. Instead, we insert the new key into an existing leaf node. Since we cannot insert a key into a leaf node that is full, we introduce an operation that *splits* a full node  $y$  (having  $2t - 1$  keys) around its *median key*  $y.key_t$  into two nodes having only  $t - 1$  keys each. The median key moves up into  $y$ 's parent to identify the dividing point between the two new trees. But if  $y$ 's parent is also full, we must split it before we can insert the new key, and thus we could end up splitting full nodes all the way up the tree.

As with a binary search tree, we can insert a key into a B-tree in a single pass down the tree from the root to a leaf. To do so, we do not wait to find out whether we will actually need to split a full node in order to do the insertion. Instead, as we travel down the tree searching for the position where the new key belongs, we split each full node we come to along the way (including the leaf itself). Thus whenever we want to split a full node  $y$ , we are assured that its parent is not full.

### Splitting a node in a B-tree

The procedure B-TREE-SPLIT-CHILD takes as input a *nonfull* internal node  $x$  (assumed to be in main memory) and an index  $i$  such that  $x.c_i$  (also assumed to be in main memory) is a *full* child of  $x$ . The procedure then splits this child in two and adjusts  $x$  so that it has an additional child. To split a full root, we will first make the root a child of a new empty root node, so that we can use B-TREE-SPLIT-CHILD. The tree thus grows in height by one; splitting is the only means by which the tree grows.

Figure 18.5 illustrates this process. We split the full node  $y = x.c_i$  about its median key  $S$ , which moves up into  $y$ 's parent node  $x$ . Those keys in  $y$  that are greater than the median key move into a new node  $z$ , which becomes a new child of  $x$ .



**Figure 18.5** Splitting a node with  $t = 4$ . Node  $y = x.c_i$  splits into two nodes,  $y$  and  $z$ , and the median key  $S$  of  $y$  moves up into  $y$ 's parent.

**B-TREE-SPLIT-CHILD( $x, i$ )**

```

1   $z = \text{ALLOCATE-NODE}()$ 
2   $y = x.c_i$ 
3   $z.\text{leaf} = y.\text{leaf}$ 
4   $z.n = t - 1$ 
5  for  $j = 1$  to  $t - 1$ 
6       $z.\text{key}_j = y.\text{key}_{j+t}$ 
7  if not  $y.\text{leaf}$ 
8      for  $j = 1$  to  $t$ 
9           $z.c_j = y.c_{j+t}$ 
10  $y.n = t - 1$ 
11 for  $j = x.n + 1$  downto  $i + 1$ 
12      $x.c_{j+1} = x.c_j$ 
13  $x.c_{i+1} = z$ 
14 for  $j = x.n$  downto  $i$ 
15      $x.\text{key}_{j+1} = x.\text{key}_j$ 
16  $x.\text{key}_i = y.\text{key}_t$ 
17  $x.n = x.n + 1$ 
18 DISK-WRITE( $y$ )
19 DISK-WRITE( $z$ )
20 DISK-WRITE( $x$ )

```

B-TREE-SPLIT-CHILD works by straightforward “cutting and pasting.” Here,  $x$  is the node being split, and  $y$  is  $x$ 's  $i$ th child (set in line 2). Node  $y$  originally has  $2t$  children ( $2t - 1$  keys) but is reduced to  $t$  children ( $t - 1$  keys) by this operation. Node  $z$  takes the  $t$  largest children ( $t - 1$  keys) from  $y$ , and  $z$  becomes a new child

of  $x$ , positioned just after  $y$  in  $x$ 's table of children. The median key of  $y$  moves up to become the key in  $x$  that separates  $y$  and  $z$ .

Lines 1–9 create node  $z$  and give it the largest  $t - 1$  keys and corresponding  $t$  children of  $y$ . Line 10 adjusts the key count for  $y$ . Finally, lines 11–17 insert  $z$  as a child of  $x$ , move the median key from  $y$  up to  $x$  in order to separate  $y$  from  $z$ , and adjust  $x$ 's key count. Lines 18–20 write out all modified disk pages. The CPU time used by B-TREE-SPLIT-CHILD is  $\Theta(t)$ , due to the loops on lines 5–6 and 8–9. (The other loops run for  $O(t)$  iterations.) The procedure performs  $O(1)$  disk operations.

### *Inserting a key into a B-tree in a single pass down the tree*

We insert a key  $k$  into a B-tree  $T$  of height  $h$  in a single pass down the tree, requiring  $O(h)$  disk accesses. The CPU time required is  $O(th) = O(t \log_t n)$ . The B-TREE-INSERT procedure uses B-TREE-SPLIT-CHILD to guarantee that the recursion never descends to a full node.

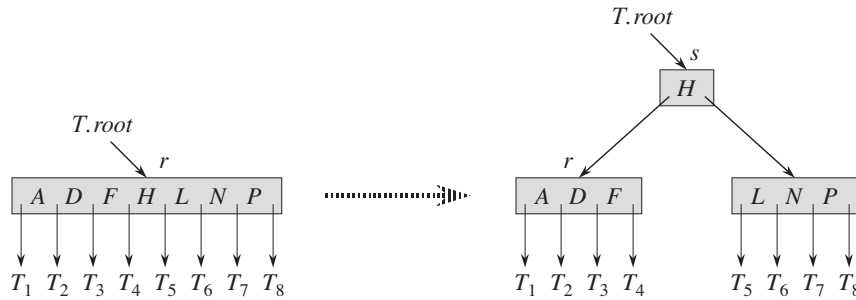
B-TREE-INSERT( $T, k$ )

```

1   $r = T.root$ 
2  if  $r.n == 2t - 1$ 
3       $s = \text{ALLOCATE-NODE}()$ 
4       $T.root = s$ 
5       $s.leaf = \text{FALSE}$ 
6       $s.n = 0$ 
7       $s.c_1 = r$ 
8      B-TREE-SPLIT-CHILD( $s, 1$ )
9      B-TREE-INSERT-NONFULL( $s, k$ )
10 else B-TREE-INSERT-NONFULL( $r, k$ )
```

Lines 3–9 handle the case in which the root node  $r$  is full: the root splits and a new node  $s$  (having two children) becomes the root. Splitting the root is the only way to increase the height of a B-tree. Figure 18.6 illustrates this case. Unlike a binary search tree, a B-tree increases in height at the top instead of at the bottom. The procedure finishes by calling B-TREE-INSERT-NONFULL to insert key  $k$  into the tree rooted at the nonfull root node. B-TREE-INSERT-NONFULL recurses as necessary down the tree, at all times guaranteeing that the node to which it recurses is not full by calling B-TREE-SPLIT-CHILD as necessary.

The auxiliary recursive procedure B-TREE-INSERT-NONFULL inserts key  $k$  into node  $x$ , which is assumed to be nonfull when the procedure is called. The operation of B-TREE-INSERT and the recursive operation of B-TREE-INSERT-NONFULL guarantee that this assumption is true.



**Figure 18.6** Splitting the root with  $t = 4$ . Root node  $r$  splits in two, and a new root node  $s$  is created. The new root contains the median key of  $r$  and has the two halves of  $r$  as children. The B-tree grows in height by one when the root is split.

B-TREE-INSERT-NONFULL( $x, k$ )

```

1   $i = x.n$ 
2  if  $x.leaf$ 
3      while  $i \geq 1$  and  $k < x.key_i$ 
4           $x.key_{i+1} = x.key_i$ 
5           $i = i - 1$ 
6       $x.key_{i+1} = k$ 
7       $x.n = x.n + 1$ 
8      DISK-WRITE( $x$ )
9  else while  $i \geq 1$  and  $k < x.key_i$ 
10      $i = i - 1$ 
11      $i = i + 1$ 
12     DISK-READ( $x.c_i$ )
13     if  $x.c_i.n == 2t - 1$ 
14         B-TREE-SPLIT-CHILD( $x, i$ )
15         if  $k > x.key_i$ 
16              $i = i + 1$ 
17     B-TREE-INSERT-NONFULL( $x.c_i, k$ )

```

The B-TREE-INSERT-NONFULL procedure works as follows. Lines 3–8 handle the case in which  $x$  is a leaf node by inserting key  $k$  into  $x$ . If  $x$  is not a leaf node, then we must insert  $k$  into the appropriate leaf node in the subtree rooted at internal node  $x$ . In this case, lines 9–11 determine the child of  $x$  to which the recursion descends. Line 13 detects whether the recursion would descend to a full child, in which case line 14 uses B-TREE-SPLIT-CHILD to split that child into two nonfull children, and lines 15–16 determine which of the two children is now the

correct one to descend to. (Note that there is no need for a  $\text{DISK-READ}(x.c_i)$  after line 16 increments  $i$ , since the recursion will descend in this case to a child that was just created by  $\text{B-TREE-SPLIT-CHILD}$ .) The net effect of lines 13–16 is thus to guarantee that the procedure never recurses to a full node. Line 17 then recurses to insert  $k$  into the appropriate subtree. Figure 18.7 illustrates the various cases of inserting into a B-tree.

For a B-tree of height  $h$ ,  $\text{B-TREE-INSERT}$  performs  $O(h)$  disk accesses, since only  $O(1)$   $\text{DISK-READ}$  and  $\text{DISK-WRITE}$  operations occur between calls to  $\text{B-TREE-INSERT-NONFULL}$ . The total CPU time used is  $O(th) = O(t \log_t n)$ . Since  $\text{B-TREE-INSERT-NONFULL}$  is tail-recursive, we can alternatively implement it as a **while** loop, thereby demonstrating that the number of pages that need to be in main memory at any time is  $O(1)$ .

## Exercises

### 18.2-1

Show the results of inserting the keys

$F, S, Q, K, C, L, H, T, V, W, M, R, N, P, A, B, X, Y, D, Z, E$

in order into an empty B-tree with minimum degree 2. Draw only the configurations of the tree just before some node must split, and also draw the final configuration.

### 18.2-2

Explain under what circumstances, if any, redundant  $\text{DISK-READ}$  or  $\text{DISK-WRITE}$  operations occur during the course of executing a call to  $\text{B-TREE-INSERT}$ . (A redundant  $\text{DISK-READ}$  is a  $\text{DISK-READ}$  for a page that is already in memory. A redundant  $\text{DISK-WRITE}$  writes to disk a page of information that is identical to what is already stored there.)

### 18.2-3

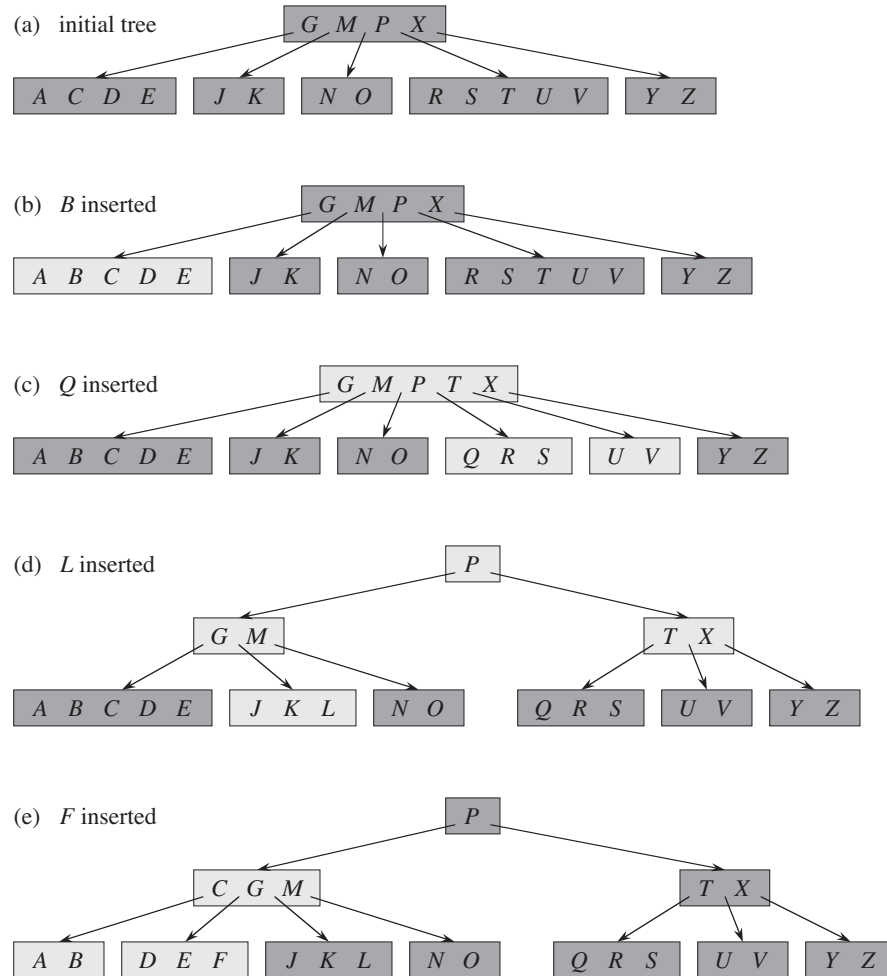
Explain how to find the minimum key stored in a B-tree and how to find the predecessor of a given key stored in a B-tree.

### 18.2-4 ★

Suppose that we insert the keys  $\{1, 2, \dots, n\}$  into an empty B-tree with minimum degree 2. How many nodes does the final B-tree have?

### 18.2-5

Since leaf nodes require no pointers to children, they could conceivably use a different (larger)  $t$  value than internal nodes for the same disk page size. Show how to modify the procedures for creating and inserting into a B-tree to handle this variation.



**Figure 18.7** Inserting keys into a B-tree. The minimum degree  $t$  for this B-tree is 3, so a node can hold at most 5 keys. Nodes that are modified by the insertion process are lightly shaded. (a) The initial tree for this example. (b) The result of inserting *B* into the initial tree; this is a simple insertion into a leaf node. (c) The result of inserting *Q* into the previous tree. The node *RSTUV* splits into two nodes containing *RS* and *UV*, the key *T* moves up to the root, and *Q* is inserted in the leftmost of the two halves (the *RS* node). (d) The result of inserting *L* into the previous tree. The root splits right away, since it is full, and the B-tree grows in height by one. Then *L* is inserted into the leaf containing *JK*. (e) The result of inserting *F* into the previous tree. The node *ABCDE* splits before *F* is inserted into the rightmost of the two halves (the *DE* node).

**18.2-6**

Suppose that we were to implement B-TREE-SEARCH to use binary search rather than linear search within each node. Show that this change makes the CPU time required  $O(\lg n)$ , independently of how  $t$  might be chosen as a function of  $n$ .

**18.2-7**

Suppose that disk hardware allows us to choose the size of a disk page arbitrarily, but that the time it takes to read the disk page is  $a + bt$ , where  $a$  and  $b$  are specified constants and  $t$  is the minimum degree for a B-tree using pages of the selected size. Describe how to choose  $t$  so as to minimize (approximately) the B-tree search time. Suggest an optimal value of  $t$  for the case in which  $a = 5$  milliseconds and  $b = 10$  microseconds.

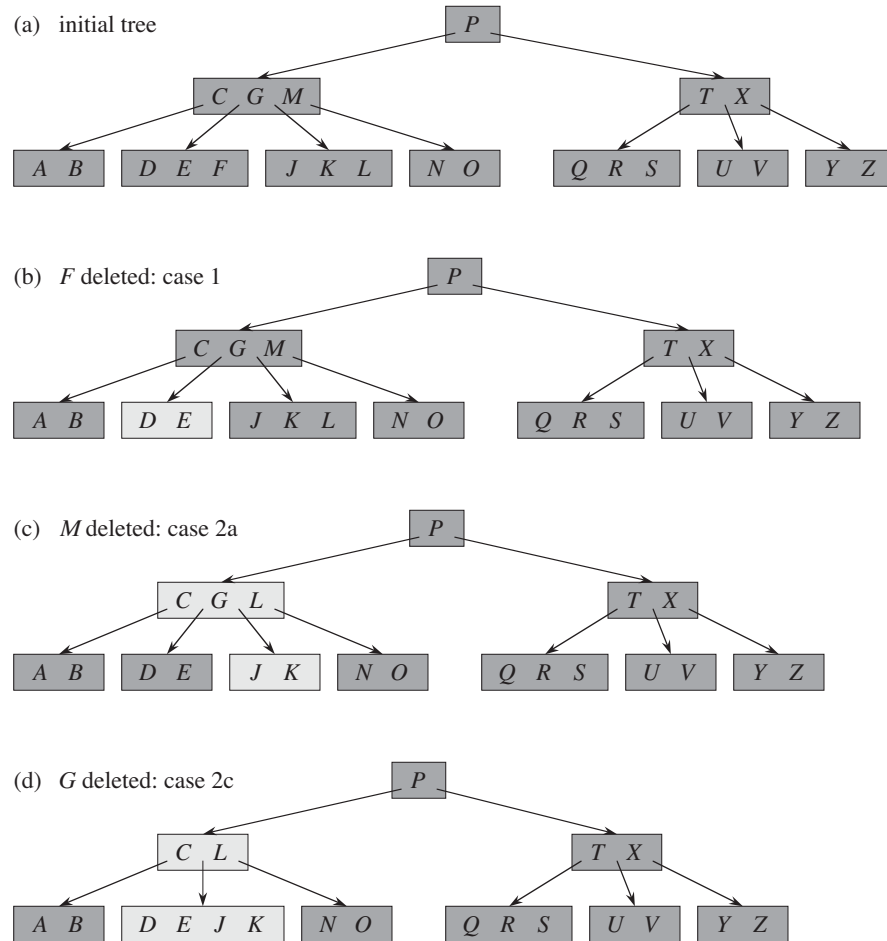
---

## 18.3 Deleting a key from a B-tree

Deletion from a B-tree is analogous to insertion but a little more complicated, because we can delete a key from any node—not just a leaf—and when we delete a key from an internal node, we will have to rearrange the node’s children. As in insertion, we must guard against deletion producing a tree whose structure violates the B-tree properties. Just as we had to ensure that a node didn’t get too big due to insertion, we must ensure that a node doesn’t get too small during deletion (except that the root is allowed to have fewer than the minimum number  $t - 1$  of keys). Just as a simple insertion algorithm might have to back up if a node on the path to where the key was to be inserted was full, a simple approach to deletion might have to back up if a node (other than the root) along the path to where the key is to be deleted has the minimum number of keys.

The procedure B-TREE-DELETE deletes the key  $k$  from the subtree rooted at  $x$ . We design this procedure to guarantee that whenever it calls itself recursively on a node  $x$ , the number of keys in  $x$  is at least the minimum degree  $t$ . Note that this condition requires one more key than the minimum required by the usual B-tree conditions, so that sometimes a key may have to be moved into a child node before recursion descends to that child. This strengthened condition allows us to delete a key from the tree in one downward pass without having to “back up” (with one exception, which we’ll explain). You should interpret the following specification for deletion from a B-tree with the understanding that if the root node  $x$  ever becomes an internal node having no keys (this situation can occur in cases 2c and 3b on pages 501–502), then we delete  $x$ , and  $x$ ’s only child  $x.c_1$  becomes the new root of the tree, decreasing the height of the tree by one and preserving the property that the root of the tree contains at least one key (unless the tree is empty).

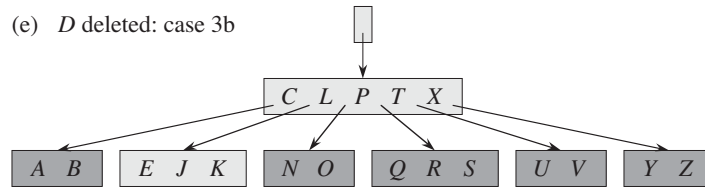




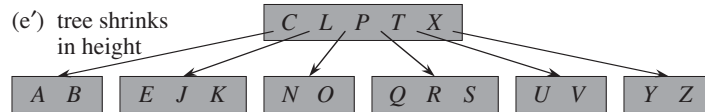
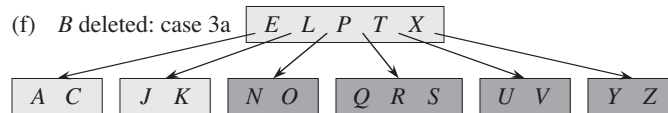
**Figure 18.8** Deleting keys from a B-tree. The minimum degree for this B-tree is  $t = 3$ , so a node (other than the root) cannot have fewer than 2 keys. Nodes that are modified are lightly shaded. (a) The B-tree of Figure 18.7(e). (b) Deletion of  $F$ . This is case 1: simple deletion from a leaf. (c) Deletion of  $M$ . This is case 2a: the predecessor  $L$  of  $M$  moves up to take  $M$ 's position. (d) Deletion of  $G$ . This is case 2c: we push  $G$  down to make node  $DEGJK$  and then delete  $G$  from this leaf (case 1).

We sketch how deletion works instead of presenting the pseudocode. Figure 18.8 illustrates the various cases of deleting keys from a B-tree.

1. If the key  $k$  is in node  $x$  and  $x$  is a leaf, delete the key  $k$  from  $x$ .
2. If the key  $k$  is in node  $x$  and  $x$  is an internal node, do the following:

(e)  $D$  deleted: case 3b

(e') tree shrinks in height

(f)  $B$  deleted: case 3a

**Figure 18.8, continued** (e) Deletion of  $D$ . This is case 3b: the recursion cannot descend to node  $CL$  because it has only 2 keys, so we push  $P$  down and merge it with  $CL$  and  $TX$  to form  $CLPTX$ ; then we delete  $D$  from a leaf (case 1). (e') After (e), we delete the root and the tree shrinks in height by one. (f) Deletion of  $B$ . This is case 3a:  $C$  moves to fill  $B$ 's position and  $E$  moves to fill  $C$ 's position.

- a. If the child  $y$  that precedes  $k$  in node  $x$  has at least  $t$  keys, then find the predecessor  $k'$  of  $k$  in the subtree rooted at  $y$ . Recursively delete  $k'$ , and replace  $k$  by  $k'$  in  $x$ . (We can find  $k'$  and delete it in a single downward pass.)
  - b. If  $y$  has fewer than  $t$  keys, then, symmetrically, examine the child  $z$  that follows  $k$  in node  $x$ . If  $z$  has at least  $t$  keys, then find the successor  $k'$  of  $k$  in the subtree rooted at  $z$ . Recursively delete  $k'$ , and replace  $k$  by  $k'$  in  $x$ . (We can find  $k'$  and delete it in a single downward pass.)
  - c. Otherwise, if both  $y$  and  $z$  have only  $t - 1$  keys, merge  $k$  and all of  $z$  into  $y$ , so that  $x$  loses both  $k$  and the pointer to  $z$ , and  $y$  now contains  $2t - 1$  keys. Then free  $z$  and recursively delete  $k$  from  $y$ .
3. If the key  $k$  is not present in internal node  $x$ , determine the root  $x.c_i$  of the appropriate subtree that must contain  $k$ , if  $k$  is in the tree at all. If  $x.c_i$  has only  $t - 1$  keys, execute step 3a or 3b as necessary to guarantee that we descend to a node containing at least  $t$  keys. Then finish by recursing on the appropriate child of  $x$ .

- a. If  $x.c_i$  has only  $t - 1$  keys but has an immediate sibling with at least  $t$  keys, give  $x.c_i$  an extra key by moving a key from  $x$  down into  $x.c_i$ , moving a key from  $x.c_i$ 's immediate left or right sibling up into  $x$ , and moving the appropriate child pointer from the sibling into  $x.c_i$ .
- b. If  $x.c_i$  and both of  $x.c_i$ 's immediate siblings have  $t - 1$  keys, merge  $x.c_i$  with one sibling, which involves moving a key from  $x$  down into the new merged node to become the median key for that node.

Since most of the keys in a B-tree are in the leaves, we may expect that in practice, deletion operations are most often used to delete keys from leaves. The B-TREE-DELETE procedure then acts in one downward pass through the tree, without having to back up. When deleting a key in an internal node, however, the procedure makes a downward pass through the tree but may have to return to the node from which the key was deleted to replace the key with its predecessor or successor (cases 2a and 2b).

Although this procedure seems complicated, it involves only  $O(h)$  disk operations for a B-tree of height  $h$ , since only  $O(1)$  calls to DISK-READ and DISK-WRITE are made between recursive invocations of the procedure. The CPU time required is  $O(th) = O(t \log_t n)$ .

## Exercises

### 18.3-1

Show the results of deleting  $C$ ,  $P$ , and  $V$ , in order, from the tree of Figure 18.8(f).

### 18.3-2

Write pseudocode for B-TREE-DELETE.

---

## Problems

### 18-1 Stacks on secondary storage

Consider implementing a stack in a computer that has a relatively small amount of fast primary memory and a relatively large amount of slower disk storage. The operations PUSH and POP work on single-word values. The stack we wish to support can grow to be much larger than can fit in memory, and thus most of it must be stored on disk.

A simple, but inefficient, stack implementation keeps the entire stack on disk. We maintain in memory a stack pointer, which is the disk address of the top element on the stack. If the pointer has value  $p$ , the top element is the  $(p \bmod m)$ th word on page  $\lfloor p/m \rfloor$  of the disk, where  $m$  is the number of words per page.

To implement the PUSH operation, we increment the stack pointer, read the appropriate page into memory from disk, copy the element to be pushed to the appropriate word on the page, and write the page back to disk. A POP operation is similar. We decrement the stack pointer, read in the appropriate page from disk, and return the top of the stack. We need not write back the page, since it was not modified.

Because disk operations are relatively expensive, we count two costs for any implementation: the total number of disk accesses and the total CPU time. Any disk access to a page of  $m$  words incurs charges of one disk access and  $\Theta(m)$  CPU time.

- a.* Asymptotically, what is the worst-case number of disk accesses for  $n$  stack operations using this simple implementation? What is the CPU time for  $n$  stack operations? (Express your answer in terms of  $m$  and  $n$  for this and subsequent parts.)

Now consider a stack implementation in which we keep one page of the stack in memory. (We also maintain a small amount of memory to keep track of which page is currently in memory.) We can perform a stack operation only if the relevant disk page resides in memory. If necessary, we can write the page currently in memory to the disk and read in the new page from the disk to memory. If the relevant disk page is already in memory, then no disk accesses are required.

- b.* What is the worst-case number of disk accesses required for  $n$  PUSH operations? What is the CPU time?
- c.* What is the worst-case number of disk accesses required for  $n$  stack operations? What is the CPU time?

Suppose that we now implement the stack by keeping two pages in memory (in addition to a small number of words for bookkeeping).

- d.* Describe how to manage the stack pages so that the amortized number of disk accesses for any stack operation is  $O(1/m)$  and the amortized CPU time for any stack operation is  $O(1)$ .

### 18-2 Joining and splitting 2-3-4 trees

The *join* operation takes two dynamic sets  $S'$  and  $S''$  and an element  $x$  such that for any  $x' \in S'$  and  $x'' \in S''$ , we have  $x'.key < x.key < x''.key$ . It returns a set  $S = S' \cup \{x\} \cup S''$ . The *split* operation is like an “inverse” join: given a dynamic set  $S$  and an element  $x \in S$ , it creates a set  $S'$  that consists of all elements in  $S - \{x\}$  whose keys are less than  $x.key$  and a set  $S''$  that consists of all elements in  $S - \{x\}$  whose keys are greater than  $x.key$ . In this problem, we investigate

how to implement these operations on 2-3-4 trees. We assume for convenience that elements consist only of keys and that all key values are distinct.

- a. Show how to maintain, for every node  $x$  of a 2-3-4 tree, the height of the subtree rooted at  $x$  as an attribute  $x.height$ . Make sure that your implementation does not affect the asymptotic running times of searching, insertion, and deletion.
- b. Show how to implement the join operation. Given two 2-3-4 trees  $T'$  and  $T''$  and a key  $k$ , the join operation should run in  $O(1 + |h' - h''|)$  time, where  $h'$  and  $h''$  are the heights of  $T'$  and  $T''$ , respectively.
- c. Consider the simple path  $p$  from the root of a 2-3-4 tree  $T$  to a given key  $k$ , the set  $S'$  of keys in  $T$  that are less than  $k$ , and the set  $S''$  of keys in  $T$  that are greater than  $k$ . Show that  $p$  breaks  $S'$  into a set of trees  $\{T'_0, T'_1, \dots, T'_m\}$  and a set of keys  $\{k'_1, k'_2, \dots, k'_m\}$ , where, for  $i = 1, 2, \dots, m$ , we have  $y < k'_i < z$  for any keys  $y \in T'_{i-1}$  and  $z \in T'_i$ . What is the relationship between the heights of  $T'_{i-1}$  and  $T'_i$ ? Describe how  $p$  breaks  $S''$  into sets of trees and keys.
- d. Show how to implement the split operation on  $T$ . Use the join operation to assemble the keys in  $S'$  into a single 2-3-4 tree  $T'$  and the keys in  $S''$  into a single 2-3-4 tree  $T''$ . The running time of the split operation should be  $O(\lg n)$ , where  $n$  is the number of keys in  $T$ . (*Hint:* The costs for joining should telescope.)

---

## Chapter notes

Knuth [211], Aho, Hopcroft, and Ullman [5], and Sedgewick [306] give further discussions of balanced-tree schemes and B-trees. Comer [74] provides a comprehensive survey of B-trees. Guibas and Sedgewick [155] discuss the relationships among various kinds of balanced-tree schemes, including red-black trees and 2-3-4 trees.

In 1970, J. E. Hopcroft invented 2-3 trees, a precursor to B-trees and 2-3-4 trees, in which every internal node has either two or three children. Bayer and McCreight [35] introduced B-trees in 1972; they did not explain their choice of name.

Bender, Demaine, and Farach-Colton [40] studied how to make B-trees perform well in the presence of memory-hierarchy effects. Their *cache-oblivious* algorithms work efficiently without explicitly knowing the data transfer sizes within the memory hierarchy.

---

## 19 Fibonacci Heaps

The Fibonacci heap data structure serves a dual purpose. First, it supports a set of operations that constitutes what is known as a “mergeable heap.” Second, several Fibonacci-heap operations run in constant amortized time, which makes this data structure well suited for applications that invoke these operations frequently.

### Mergeable heaps

A *mergeable heap* is any data structure that supports the following five operations, in which each element has a *key*:

MAKE-HEAP() creates and returns a new heap containing no elements.

INSERT( $H, x$ ) inserts element  $x$ , whose *key* has already been filled in, into heap  $H$ .

MINIMUM( $H$ ) returns a pointer to the element in heap  $H$  whose key is minimum.

EXTRACT-MIN( $H$ ) deletes the element from heap  $H$  whose key is minimum, returning a pointer to the element.

UNION( $H_1, H_2$ ) creates and returns a new heap that contains all the elements of heaps  $H_1$  and  $H_2$ . Heaps  $H_1$  and  $H_2$  are “destroyed” by this operation.

In addition to the mergeable-heap operations above, Fibonacci heaps also support the following two operations:

DECREASE-KEY( $H, x, k$ ) assigns to element  $x$  within heap  $H$  the new key value  $k$ , which we assume to be no greater than its current key value.<sup>1</sup>

DELETE( $H, x$ ) deletes element  $x$  from heap  $H$ .

---

<sup>1</sup>As mentioned in the introduction to Part V, our default mergeable heaps are mergeable min-heaps, and so the operations MINIMUM, EXTRACT-MIN, and DECREASE-KEY apply. Alternatively, we could define a *mergeable max-heap* with the operations MAXIMUM, EXTRACT-MAX, and INCREASE-KEY.

Procedure	Binary heap (worst-case)	Fibonacci heap (amortized)
MAKE-HEAP	$\Theta(1)$	$\Theta(1)$
INSERT	$\Theta(\lg n)$	$\Theta(1)$
MINIMUM	$\Theta(1)$	$\Theta(1)$
EXTRACT-MIN	$\Theta(\lg n)$	$O(\lg n)$
UNION	$\Theta(n)$	$\Theta(1)$
DECREASE-KEY	$\Theta(\lg n)$	$\Theta(1)$
DELETE	$\Theta(\lg n)$	$O(\lg n)$

**Figure 19.1** Running times for operations on two implementations of mergeable heaps. The number of items in the heap(s) at the time of an operation is denoted by  $n$ .

As the table in Figure 19.1 shows, if we don't need the UNION operation, ordinary binary heaps, as used in heapsort (Chapter 6), work fairly well. Operations other than UNION run in worst-case time  $O(\lg n)$  on a binary heap. If we need to support the UNION operation, however, binary heaps perform poorly. By concatenating the two arrays that hold the binary heaps to be merged and then running BUILD-MIN-HEAP (see Section 6.3), the UNION operation takes  $\Theta(n)$  time in the worst case.

Fibonacci heaps, on the other hand, have better asymptotic time bounds than binary heaps for the INSERT, UNION, and DECREASE-KEY operations, and they have the same asymptotic running times for the remaining operations. Note, however, that the running times for Fibonacci heaps in Figure 19.1 are amortized time bounds, not worst-case per-operation time bounds. The UNION operation takes only constant amortized time in a Fibonacci heap, which is significantly better than the linear worst-case time required in a binary heap (assuming, of course, that an amortized time bound suffices).

### Fibonacci heaps in theory and practice

From a theoretical standpoint, Fibonacci heaps are especially desirable when the number of EXTRACT-MIN and DELETE operations is small relative to the number of other operations performed. This situation arises in many applications. For example, some algorithms for graph problems may call DECREASE-KEY once per edge. For dense graphs, which have many edges, the  $\Theta(1)$  amortized time of each call of DECREASE-KEY adds up to a big improvement over the  $\Theta(\lg n)$  worst-case time of binary heaps. Fast algorithms for problems such as computing minimum spanning trees (Chapter 23) and finding single-source shortest paths (Chapter 24) make essential use of Fibonacci heaps.

From a practical point of view, however, the constant factors and programming complexity of Fibonacci heaps make them less desirable than ordinary binary (or  $k$ -ary) heaps for most applications, except for certain applications that manage large amounts of data. Thus, Fibonacci heaps are predominantly of theoretical interest. If a much simpler data structure with the same amortized time bounds as Fibonacci heaps were developed, it would be of practical use as well.

Both binary heaps and Fibonacci heaps are inefficient in how they support the operation SEARCH; it can take a while to find an element with a given key. For this reason, operations such as DECREASE-KEY and DELETE that refer to a given element require a pointer to that element as part of their input. As in our discussion of priority queues in Section 6.5, when we use a mergeable heap in an application, we often store a handle to the corresponding application object in each mergeable-heap element, as well as a handle to the corresponding mergeable-heap element in each application object. The exact nature of these handles depends on the application and its implementation.

Like several other data structures that we have seen, Fibonacci heaps are based on rooted trees. We represent each element by a node within a tree, and each node has a *key* attribute. For the remainder of this chapter, we shall use the term “node” instead of “element.” We shall also ignore issues of allocating nodes prior to insertion and freeing nodes following deletion, assuming instead that the code calling the heap procedures deals with these details.

Section 19.1 defines Fibonacci heaps, discusses how we represent them, and presents the potential function used for their amortized analysis. Section 19.2 shows how to implement the mergeable-heap operations and achieve the amortized time bounds shown in Figure 19.1. The remaining two operations, DECREASE-KEY and DELETE, form the focus of Section 19.3. Finally, Section 19.4 finishes a key part of the analysis and also explains the curious name of the data structure.

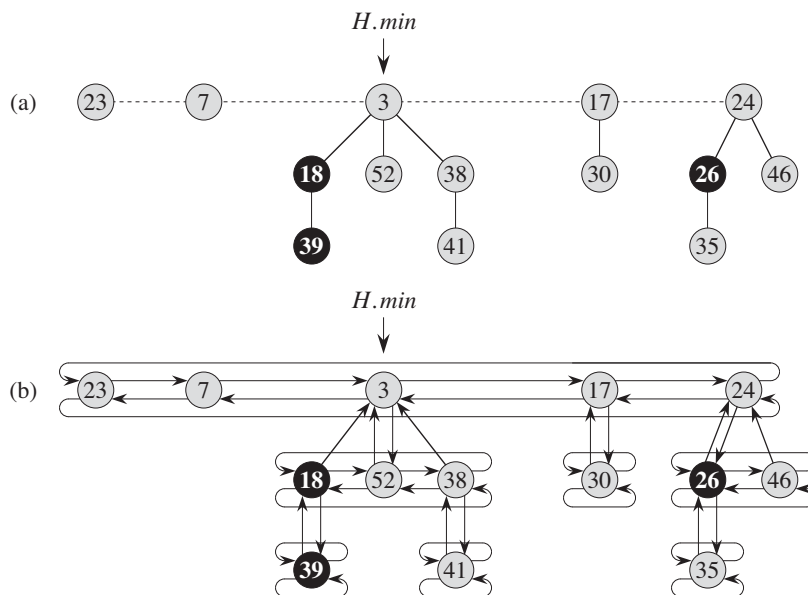
---

## 19.1 Structure of Fibonacci heaps

A **Fibonacci heap** is a collection of rooted trees that are *min-heap ordered*. That is, each tree obeys the *min-heap property*: the key of a node is greater than or equal to the key of its parent. Figure 19.2(a) shows an example of a Fibonacci heap.

As Figure 19.2(b) shows, each node  $x$  contains a pointer  $x.p$  to its parent and a pointer  $x.child$  to any one of its children. The children of  $x$  are linked together in a circular, doubly linked list, which we call the *child list* of  $x$ . Each child  $y$  in a child list has pointers  $y.left$  and  $y.right$  that point to  $y$ 's left and right siblings, respectively. If node  $y$  is an only child, then  $y.left = y.right = y$ . Siblings may appear in a child list in any order.





**Figure 19.2** (a) A Fibonacci heap consisting of five min-heap-ordered trees and 14 nodes. The dashed line indicates the root list. The minimum node of the heap is the node containing the key 3. Black nodes are marked. The potential of this particular Fibonacci heap is  $5 + 2 \cdot 3 = 11$ . (b) A more complete representation showing pointers  $p$  (up arrows),  $child$  (down arrows), and  $left$  and  $right$  (sideways arrows). The remaining figures in this chapter omit these details, since all the information shown here can be determined from what appears in part (a).

Circular, doubly linked lists (see Section 10.2) have two advantages for use in Fibonacci heaps. First, we can insert a node into any location or remove a node from anywhere in a circular, doubly linked list in  $O(1)$  time. Second, given two such lists, we can concatenate them (or “splice” them together) into one circular, doubly linked list in  $O(1)$  time. In the descriptions of Fibonacci heap operations, we shall refer to these operations informally, letting you fill in the details of their implementations if you wish.

Each node has two other attributes. We store the number of children in the child list of node  $x$  in  $x.degree$ . The boolean-valued attribute  $x.mark$  indicates whether node  $x$  has lost a child since the last time  $x$  was made the child of another node. Newly created nodes are unmarked, and a node  $x$  becomes unmarked whenever it is made the child of another node. Until we look at the DECREASE-KEY operation in Section 19.3, we will just set all  $mark$  attributes to FALSE.

We access a given Fibonacci heap  $H$  by a pointer  $H.min$  to the root of a tree containing the minimum key; we call this node the **minimum node** of the Fibonacci

heap. If more than one root has a key with the minimum value, then any such root may serve as the minimum node. When a Fibonacci heap  $H$  is empty,  $H.min$  is NIL.

The roots of all the trees in a Fibonacci heap are linked together using their *left* and *right* pointers into a circular, doubly linked list called the **root list** of the Fibonacci heap. The pointer  $H.min$  thus points to the node in the root list whose key is minimum. Trees may appear in any order within a root list.

We rely on one other attribute for a Fibonacci heap  $H$ :  $H.n$ , the number of nodes currently in  $H$ .

### Potential function

As mentioned, we shall use the potential method of Section 17.3 to analyze the performance of Fibonacci heap operations. For a given Fibonacci heap  $H$ , we indicate by  $t(H)$  the number of trees in the root list of  $H$  and by  $m(H)$  the number of marked nodes in  $H$ . We then define the potential  $\Phi(H)$  of Fibonacci heap  $H$  by

$$\Phi(H) = t(H) + 2m(H) . \quad (19.1)$$

(We will gain some intuition for this potential function in Section 19.3.) For example, the potential of the Fibonacci heap shown in Figure 19.2 is  $5 + 2 \cdot 3 = 11$ . The potential of a set of Fibonacci heaps is the sum of the potentials of its constituent Fibonacci heaps. We shall assume that a unit of potential can pay for a constant amount of work, where the constant is sufficiently large to cover the cost of any of the specific constant-time pieces of work that we might encounter.

We assume that a Fibonacci heap application begins with no heaps. The initial potential, therefore, is 0, and by equation (19.1), the potential is nonnegative at all subsequent times. From equation (17.3), an upper bound on the total amortized cost provides an upper bound on the total actual cost for the sequence of operations.

### Maximum degree

The amortized analyses we shall perform in the remaining sections of this chapter assume that we know an upper bound  $D(n)$  on the maximum degree of any node in an  $n$ -node Fibonacci heap. We won't prove it, but when only the mergeable-heap operations are supported,  $D(n) \leq \lfloor \lg n \rfloor$ . (Problem 19-2(d) asks you to prove this property.) In Sections 19.3 and 19.4, we shall show that when we support DECREASE-KEY and DELETE as well,  $D(n) = O(\lg n)$ .

## 19.2 Mergeable-heap operations

The mergeable-heap operations on Fibonacci heaps delay work as long as possible. The various operations have performance trade-offs. For example, we insert a node by adding it to the root list, which takes just constant time. If we were to start with an empty Fibonacci heap and then insert  $k$  nodes, the Fibonacci heap would consist of just a root list of  $k$  nodes. The trade-off is that if we then perform an EXTRACT-MIN operation on Fibonacci heap  $H$ , after removing the node that  $H.min$  points to, we would have to look through each of the remaining  $k - 1$  nodes in the root list to find the new minimum node. As long as we have to go through the entire root list during the EXTRACT-MIN operation, we also consolidate nodes into min-heap-ordered trees to reduce the size of the root list. We shall see that, no matter what the root list looks like before a EXTRACT-MIN operation, afterward each node in the root list has a degree that is unique within the root list, which leads to a root list of size at most  $D(n) + 1$ .

### Creating a new Fibonacci heap

To make an empty Fibonacci heap, the MAKE-FIB-HEAP procedure allocates and returns the Fibonacci heap object  $H$ , where  $H.n = 0$  and  $H.min = \text{NIL}$ ; there are no trees in  $H$ . Because  $t(H) = 0$  and  $m(H) = 0$ , the potential of the empty Fibonacci heap is  $\Phi(H) = 0$ . The amortized cost of MAKE-FIB-HEAP is thus equal to its  $O(1)$  actual cost.

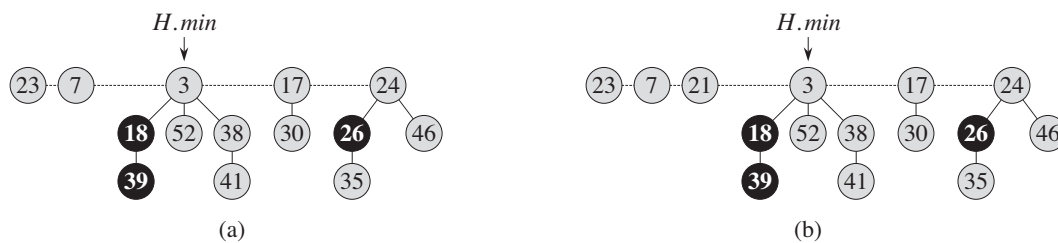
### Inserting a node

The following procedure inserts node  $x$  into Fibonacci heap  $H$ , assuming that the node has already been allocated and that  $x.key$  has already been filled in.

FIB-HEAP-INSERT( $H, x$ )

```

1   $x.degree = 0$ 
2   $x.p = \text{NIL}$ 
3   $x.child = \text{NIL}$ 
4   $x.mark = \text{FALSE}$ 
5  if  $H.min == \text{NIL}$ 
6      create a root list for  $H$  containing just  $x$ 
7       $H.min = x$ 
8  else insert  $x$  into  $H$ 's root list
9      if  $x.key < H.min.key$ 
10          $H.min = x$ 
11   $H.n = H.n + 1$ 
```



**Figure 19.3** Inserting a node into a Fibonacci heap. **(a)** A Fibonacci heap  $H$ . **(b)** Fibonacci heap  $H$  after inserting the node with key 21. The node becomes its own min-heap-ordered tree and is then added to the root list, becoming the left sibling of the root.

Lines 1–4 initialize some of the structural attributes of node  $x$ . Line 5 tests to see whether Fibonacci heap  $H$  is empty. If it is, then lines 6–7 make  $x$  be the only node in  $H$ 's root list and set  $H.min$  to point to  $x$ . Otherwise, lines 8–10 insert  $x$  into  $H$ 's root list and update  $H.min$  if necessary. Finally, line 11 increments  $H.n$  to reflect the addition of the new node. Figure 19.3 shows a node with key 21 inserted into the Fibonacci heap of Figure 19.2.

To determine the amortized cost of FIB-HEAP-INSERT, let  $H$  be the input Fibonacci heap and  $H'$  be the resulting Fibonacci heap. Then,  $t(H') = t(H) + 1$  and  $m(H') = m(H)$ , and the increase in potential is

$$((t(H) + 1) + 2m(H)) - (t(H) + 2m(H)) = 1.$$

Since the actual cost is  $O(1)$ , the amortized cost is  $O(1) + 1 = O(1)$ .

### Finding the minimum node

The minimum node of a Fibonacci heap  $H$  is given by the pointer  $H.min$ , so we can find the minimum node in  $O(1)$  actual time. Because the potential of  $H$  does not change, the amortized cost of this operation is equal to its  $O(1)$  actual cost.

### Uniting two Fibonacci heaps

The following procedure unites Fibonacci heaps  $H_1$  and  $H_2$ , destroying  $H_1$  and  $H_2$  in the process. It simply concatenates the root lists of  $H_1$  and  $H_2$  and then determines the new minimum node. Afterward, the objects representing  $H_1$  and  $H_2$  will never be used again.

FIB-HEAP-UNION( $H_1, H_2$ )

```

1   $H = \text{MAKE-FIB-HEAP}()$ 
2   $H.min = H_1.min$ 
3  concatenate the root list of  $H_2$  with the root list of  $H$ 
4  if ( $H_1.min == \text{NIL}$ ) or ( $H_2.min \neq \text{NIL}$  and  $H_2.min.key < H_1.min.key$ )
5       $H.min = H_2.min$ 
6   $H.n = H_1.n + H_2.n$ 
7  return  $H$ 

```

Lines 1–3 concatenate the root lists of  $H_1$  and  $H_2$  into a new root list  $H$ . Lines 2, 4, and 5 set the minimum node of  $H$ , and line 6 sets  $H.n$  to the total number of nodes. Line 7 returns the resulting Fibonacci heap  $H$ . As in the FIB-HEAP-INSERT procedure, all roots remain roots.

The change in potential is

$$\begin{aligned}
 \Phi(H) - (\Phi(H_1) + \Phi(H_2)) \\
 &= (t(H) + 2m(H)) - ((t(H_1) + 2m(H_1)) + (t(H_2) + 2m(H_2))) \\
 &= 0,
 \end{aligned}$$

because  $t(H) = t(H_1) + t(H_2)$  and  $m(H) = m(H_1) + m(H_2)$ . The amortized cost of FIB-HEAP-UNION is therefore equal to its  $O(1)$  actual cost.

### Extracting the minimum node

The process of extracting the minimum node is the most complicated of the operations presented in this section. It is also where the delayed work of consolidating trees in the root list finally occurs. The following pseudocode extracts the minimum node. The code assumes for convenience that when a node is removed from a linked list, pointers remaining in the list are updated, but pointers in the extracted node are left unchanged. It also calls the auxiliary procedure CONSOLIDATE, which we shall see shortly.

FIB-HEAP-EXTRACT-MIN( $H$ )

```

1   $z = H.min$ 
2  if  $z \neq \text{NIL}$ 
3      for each child  $x$  of  $z$ 
4          add  $x$  to the root list of  $H$ 
5           $x.p = \text{NIL}$ 
6      remove  $z$  from the root list of  $H$ 
7      if  $z == z.right$ 
8           $H.min = \text{NIL}$ 
9      else  $H.min = z.right$ 
10     CONSOLIDATE( $H$ )
11      $H.n = H.n - 1$ 
12 return  $z$ 

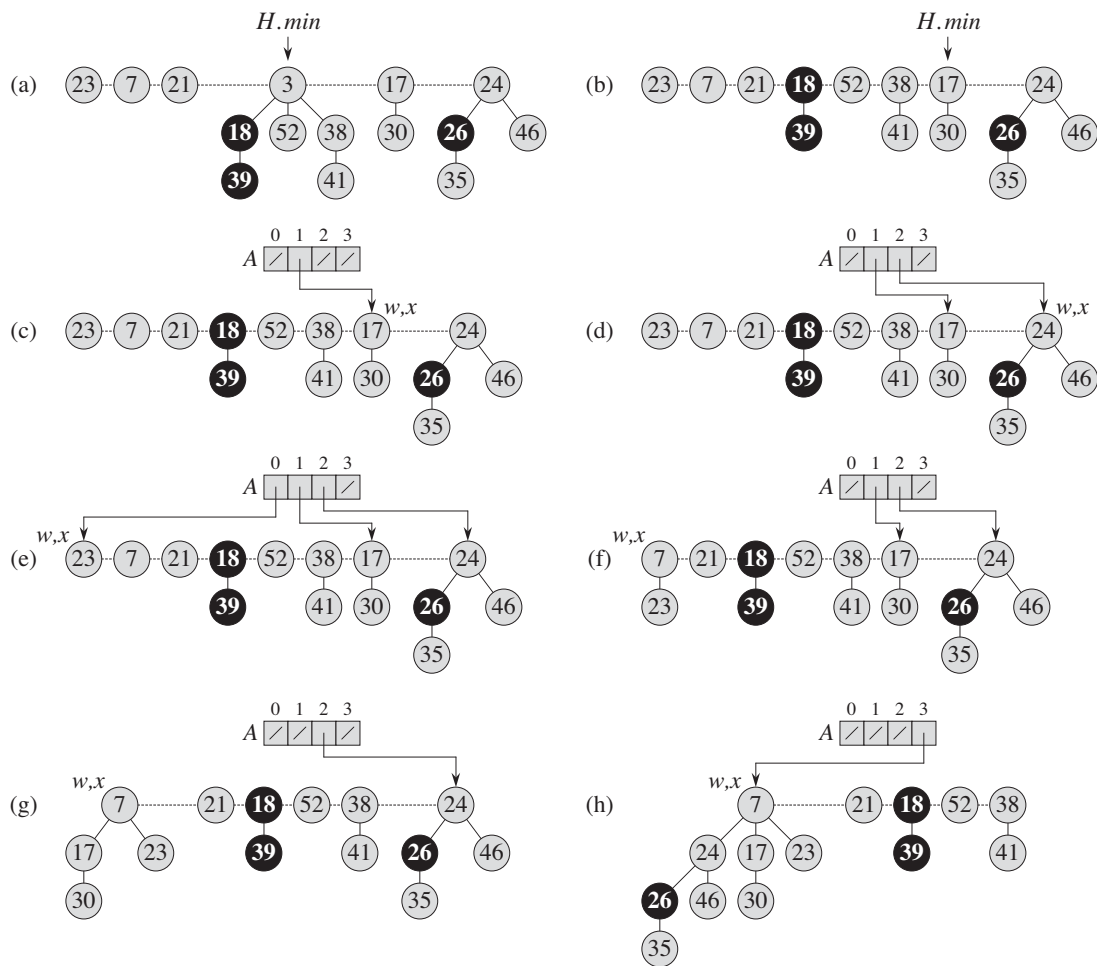
```

As Figure 19.4 illustrates, FIB-HEAP-EXTRACT-MIN works by first making a root out of each of the minimum node's children and removing the minimum node from the root list. It then consolidates the root list by linking roots of equal degree until at most one root remains of each degree.

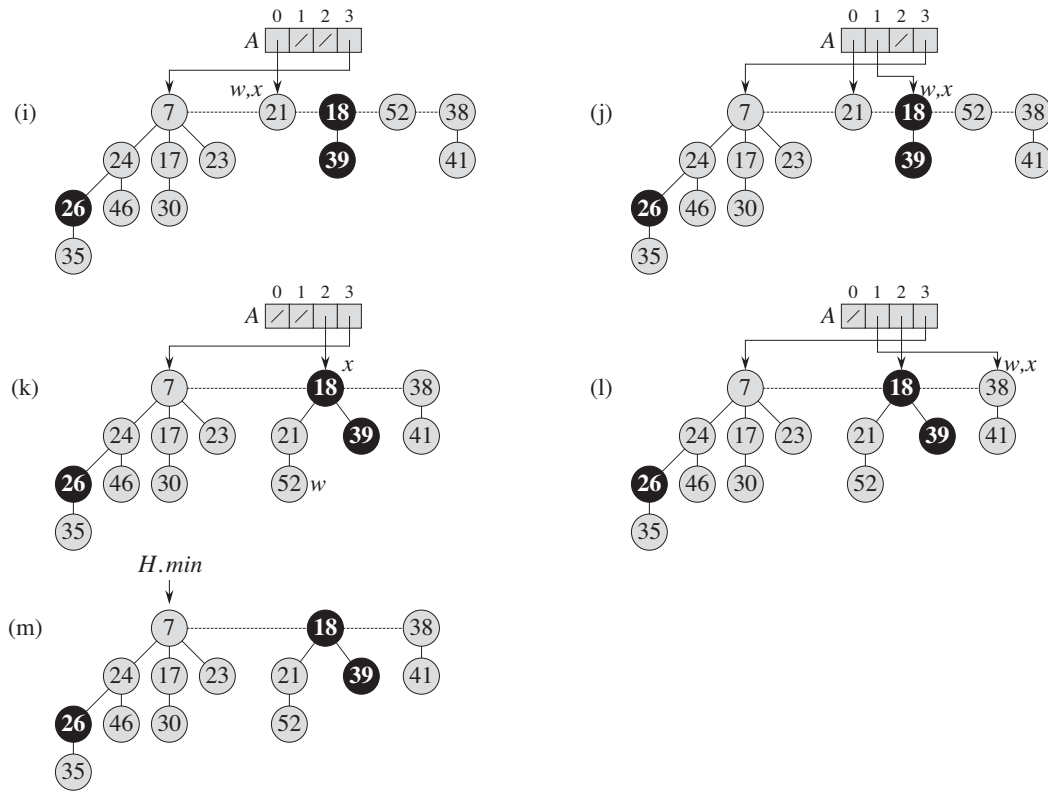
We start in line 1 by saving a pointer  $z$  to the minimum node; the procedure returns this pointer at the end. If  $z$  is NIL, then Fibonacci heap  $H$  is already empty and we are done. Otherwise, we delete node  $z$  from  $H$  by making all of  $z$ 's children roots of  $H$  in lines 3–5 (putting them into the root list) and removing  $z$  from the root list in line 6. If  $z$  is its own right sibling after line 6, then  $z$  was the only node on the root list and it had no children, so all that remains is to make the Fibonacci heap empty in line 8 before returning  $z$ . Otherwise, we set the pointer  $H.min$  into the root list to point to a root other than  $z$  (in this case,  $z$ 's right sibling), which is not necessarily going to be the new minimum node when FIB-HEAP-EXTRACT-MIN is done. Figure 19.4(b) shows the Fibonacci heap of Figure 19.4(a) after executing line 9.

The next step, in which we reduce the number of trees in the Fibonacci heap, is **consolidating** the root list of  $H$ , which the call CONSOLIDATE( $H$ ) accomplishes. Consolidating the root list consists of repeatedly executing the following steps until every root in the root list has a distinct *degree* value:

1. Find two roots  $x$  and  $y$  in the root list with the same degree. Without loss of generality, let  $x.key \leq y.key$ .
2. **Link**  $y$  to  $x$ : remove  $y$  from the root list, and make  $y$  a child of  $x$  by calling the FIB-HEAP-LINK procedure. This procedure increments the attribute  $x.degree$  and clears the mark on  $y$ .



**Figure 19.4** The action of FIB-HEAP-EXTRACT-MIN. (a) A Fibonacci heap  $H$ . (b) The situation after removing the minimum node  $z$  from the root list and adding its children to the root list. (c)–(e) The array  $A$  and the trees after each of the first three iterations of the **for** loop of lines 4–14 of the procedure CONSOLIDATE. The procedure processes the root list by starting at the node pointed to by  $H.min$  and following *right* pointers. Each part shows the values of  $w$  and  $x$  at the end of an iteration. (f)–(h) The next iteration of the **for** loop, with the values of  $w$  and  $x$  shown at the end of each iteration of the **while** loop of lines 7–13. Part (f) shows the situation after the first time through the **while** loop. The node with key 23 has been linked to the node with key 7, which  $x$  now points to. In part (g), the node with key 17 has been linked to the node with key 7, which  $x$  still points to. In part (h), the node with key 24 has been linked to the node with key 7. Since no node was previously pointed to by  $A[3]$ , at the end of the **for** loop iteration,  $A[3]$  is set to point to the root of the resulting tree.



**Figure 19.4, continued** (i)–(l) The situation after each of the next four iterations of the **for** loop. (m) Fibonacci heap  $H$  after reconstructing the root list from the array  $A$  and determining the new  $H.min$  pointer.

The procedure CONSOLIDATE uses an auxiliary array  $A[0..D(H.n)]$  to keep track of roots according to their degrees. If  $A[i] = y$ , then  $y$  is currently a root with  $y.degree = i$ . Of course, in order to allocate the array we have to know how to calculate the upper bound  $D(H.n)$  on the maximum degree, but we will see how to do so in Section 19.4.



CONSOLIDATE( $H$ )

```

1  let  $A[0 \dots D(H.n)]$  be a new array
2  for  $i = 0$  to  $D(H.n)$ 
3       $A[i] = \text{NIL}$ 
4  for each node  $w$  in the root list of  $H$ 
5       $x = w$ 
6       $d = x.\text{degree}$ 
7      while  $A[d] \neq \text{NIL}$ 
8           $y = A[d]$            // another node with the same degree as  $x$ 
9          if  $x.\text{key} > y.\text{key}$ 
10             exchange  $x$  with  $y$ 
11             FIB-HEAP-LINK( $H, y, x$ )
12              $A[d] = \text{NIL}$ 
13              $d = d + 1$ 
14       $A[d] = x$ 
15   $H.\text{min} = \text{NIL}$ 
16  for  $i = 0$  to  $D(H.n)$ 
17      if  $A[i] \neq \text{NIL}$ 
18          if  $H.\text{min} == \text{NIL}$ 
19              create a root list for  $H$  containing just  $A[i]$ 
20               $H.\text{min} = A[i]$ 
21          else insert  $A[i]$  into  $H$ 's root list
22              if  $A[i].\text{key} < H.\text{min}.\text{key}$ 
23                   $H.\text{min} = A[i]$ 

```

FIB-HEAP-LINK( $H, y, x$ )

```

1  remove  $y$  from the root list of  $H$ 
2  make  $y$  a child of  $x$ , incrementing  $x.\text{degree}$ 
3   $y.\text{mark} = \text{FALSE}$ 

```

In detail, the CONSOLIDATE procedure works as follows. Lines 1–3 allocate and initialize the array  $A$  by making each entry NIL. The **for** loop of lines 4–14 processes each root  $w$  in the root list. As we link roots together,  $w$  may be linked to some other node and no longer be a root. Nevertheless,  $w$  is always in a tree rooted at some node  $x$ , which may or may not be  $w$  itself. Because we want at most one root with each degree, we look in the array  $A$  to see whether it contains a root  $y$  with the same degree as  $x$ . If it does, then we link the roots  $x$  and  $y$  but guaranteeing that  $x$  remains a root after linking. That is, we link  $y$  to  $x$  after first exchanging the pointers to the two roots if  $y$ 's key is smaller than  $x$ 's key. After we link  $y$  to  $x$ , the degree of  $x$  has increased by 1, and so we continue this process, linking  $x$  and another root whose degree equals  $x$ 's new degree, until no other root

that we have processed has the same degree as  $x$ . We then set the appropriate entry of  $A$  to point to  $x$ , so that as we process roots later on, we have recorded that  $x$  is the unique root of its degree that we have already processed. When this **for** loop terminates, at most one root of each degree will remain, and the array  $A$  will point to each remaining root.

The **while** loop of lines 7–13 repeatedly links the root  $x$  of the tree containing node  $w$  to another tree whose root has the same degree as  $x$ , until no other root has the same degree. This **while** loop maintains the following invariant:

At the start of each iteration of the **while** loop,  $d = x.degree$ .

We use this loop invariant as follows:

**Initialization:** Line 6 ensures that the loop invariant holds the first time we enter the loop.

**Maintenance:** In each iteration of the **while** loop,  $A[d]$  points to some root  $y$ . Because  $d = x.degree = y.degree$ , we want to link  $x$  and  $y$ . Whichever of  $x$  and  $y$  has the smaller key becomes the parent of the other as a result of the link operation, and so lines 9–10 exchange the pointers to  $x$  and  $y$  if necessary. Next, we link  $y$  to  $x$  by the call `FIB-HEAP-LINK( $H, y, x$ )` in line 11. This call increments  $x.degree$  but leaves  $y.degree$  as  $d$ . Node  $y$  is no longer a root, and so line 12 removes the pointer to it in array  $A$ . Because the call of `FIB-HEAP-LINK` increments the value of  $x.degree$ , line 13 restores the invariant that  $d = x.degree$ .

**Termination:** We repeat the **while** loop until  $A[d] = \text{NIL}$ , in which case there is no other root with the same degree as  $x$ .

After the **while** loop terminates, we set  $A[d]$  to  $x$  in line 14 and perform the next iteration of the **for** loop.

Figures 19.4(c)–(e) show the array  $A$  and the resulting trees after the first three iterations of the **for** loop of lines 4–14. In the next iteration of the **for** loop, three links occur; their results are shown in Figures 19.4(f)–(h). Figures 19.4(i)–(l) show the result of the next four iterations of the **for** loop.

All that remains is to clean up. Once the **for** loop of lines 4–14 completes, line 15 empties the root list, and lines 16–23 reconstruct it from the array  $A$ . The resulting Fibonacci heap appears in Figure 19.4(m). After consolidating the root list, `FIB-HEAP-EXTRACT-MIN` finishes up by decrementing  $H.n$  in line 11 and returning a pointer to the deleted node  $z$  in line 12.

We are now ready to show that the amortized cost of extracting the minimum node of an  $n$ -node Fibonacci heap is  $O(D(n))$ . Let  $H$  denote the Fibonacci heap just prior to the `FIB-HEAP-EXTRACT-MIN` operation.

We start by accounting for the actual cost of extracting the minimum node. An  $O(D(n))$  contribution comes from `FIB-HEAP-EXTRACT-MIN` processing at

most  $D(n)$  children of the minimum node and from the work in lines 2–3 and 16–23 of CONSOLIDATE. It remains to analyze the contribution from the **for** loop of lines 4–14 in CONSOLIDATE, for which we use an aggregate analysis. The size of the root list upon calling CONSOLIDATE is at most  $D(n) + t(H) - 1$ , since it consists of the original  $t(H)$  root-list nodes, minus the extracted root node, plus the children of the extracted node, which number at most  $D(n)$ . Within a given iteration of the **for** loop of lines 4–14, the number of iterations of the **while** loop of lines 7–13 depends on the root list. But we know that every time through the **while** loop, one of the roots is linked to another, and thus the total number of iterations of the **while** loop over all iterations of the **for** loop is at most the number of roots in the root list. Hence, the total amount of work performed in the **for** loop is at most proportional to  $D(n) + t(H)$ . Thus, the total actual work in extracting the minimum node is  $O(D(n) + t(H))$ .

The potential before extracting the minimum node is  $t(H) + 2m(H)$ , and the potential afterward is at most  $(D(n) + 1) + 2m(H)$ , since at most  $D(n) + 1$  roots remain and no nodes become marked during the operation. The amortized cost is thus at most

$$\begin{aligned} & O(D(n) + t(H)) + ((D(n) + 1) + 2m(H)) - (t(H) + 2m(H)) \\ &= O(D(n)) + O(t(H)) - t(H) \\ &= O(D(n)) , \end{aligned}$$

since we can scale up the units of potential to dominate the constant hidden in  $O(t(H))$ . Intuitively, the cost of performing each link is paid for by the reduction in potential due to the link's reducing the number of roots by one. We shall see in Section 19.4 that  $D(n) = O(\lg n)$ , so that the amortized cost of extracting the minimum node is  $O(\lg n)$ .

## Exercises

### 19.2-1

Show the Fibonacci heap that results from calling FIB-HEAP-EXTRACT-MIN on the Fibonacci heap shown in Figure 19.4(m).

---

## 19.3 Decreasing a key and deleting a node

In this section, we show how to decrease the key of a node in a Fibonacci heap in  $O(1)$  amortized time and how to delete any node from an  $n$ -node Fibonacci heap in  $O(D(n))$  amortized time. In Section 19.4, we will show that the maxi-

imum degree  $D(n)$  is  $O(\lg n)$ , which will imply that FIB-HEAP-EXTRACT-MIN and FIB-HEAP-DELETE run in  $O(\lg n)$  amortized time.

### Decreasing a key

In the following pseudocode for the operation FIB-HEAP-DECREASE-KEY, we assume as before that removing a node from a linked list does not change any of the structural attributes in the removed node.

FIB-HEAP-DECREASE-KEY( $H, x, k$ )

```

1  if  $k > x.key$ 
2      error “new key is greater than current key”
3   $x.key = k$ 
4   $y = x.p$ 
5  if  $y \neq \text{NIL}$  and  $x.key < y.key$ 
6      CUT( $H, x, y$ )
7      CASCADING-CUT( $H, y$ )
8  if  $x.key < H.min.key$ 
9       $H.min = x$ 
```

CUT( $H, x, y$ )

```

1  remove  $x$  from the child list of  $y$ , decrementing  $y.degree$ 
2  add  $x$  to the root list of  $H$ 
3   $x.p = \text{NIL}$ 
4   $x.mark = \text{FALSE}$ 
```

CASCADING-CUT( $H, y$ )

```

1   $z = y.p$ 
2  if  $z \neq \text{NIL}$ 
3      if  $y.mark == \text{FALSE}$ 
4           $y.mark = \text{TRUE}$ 
5      else CUT( $H, y, z$ )
6      CASCADING-CUT( $H, z$ )
```

The FIB-HEAP-DECREASE-KEY procedure works as follows. Lines 1–3 ensure that the new key is no greater than the current key of  $x$  and then assign the new key to  $x$ . If  $x$  is a root or if  $x.key \geq y.key$ , where  $y$  is  $x$ ’s parent, then no structural changes need occur, since min-heap order has not been violated. Lines 4–5 test for this condition.

If min-heap order has been violated, many changes may occur. We start by **cutting**  $x$  in line 6. The CUT procedure “cuts” the link between  $x$  and its parent  $y$ , making  $x$  a root.

We use the *mark* attributes to obtain the desired time bounds. They record a little piece of the history of each node. Suppose that the following events have happened to node  $x$ :

1. at some time,  $x$  was a root,
2. then  $x$  was linked to (made the child of) another node,
3. then two children of  $x$  were removed by cuts.

As soon as the second child has been lost, we cut  $x$  from its parent, making it a new root. The attribute  $x.mark$  is TRUE if steps 1 and 2 have occurred and one child of  $x$  has been cut. The CUT procedure, therefore, clears  $x.mark$  in line 4, since it performs step 1. (We can now see why line 3 of FIB-HEAP-LINK clears  $y.mark$ : node  $y$  is being linked to another node, and so step 2 is being performed. The next time a child of  $y$  is cut,  $y.mark$  will be set to TRUE.)

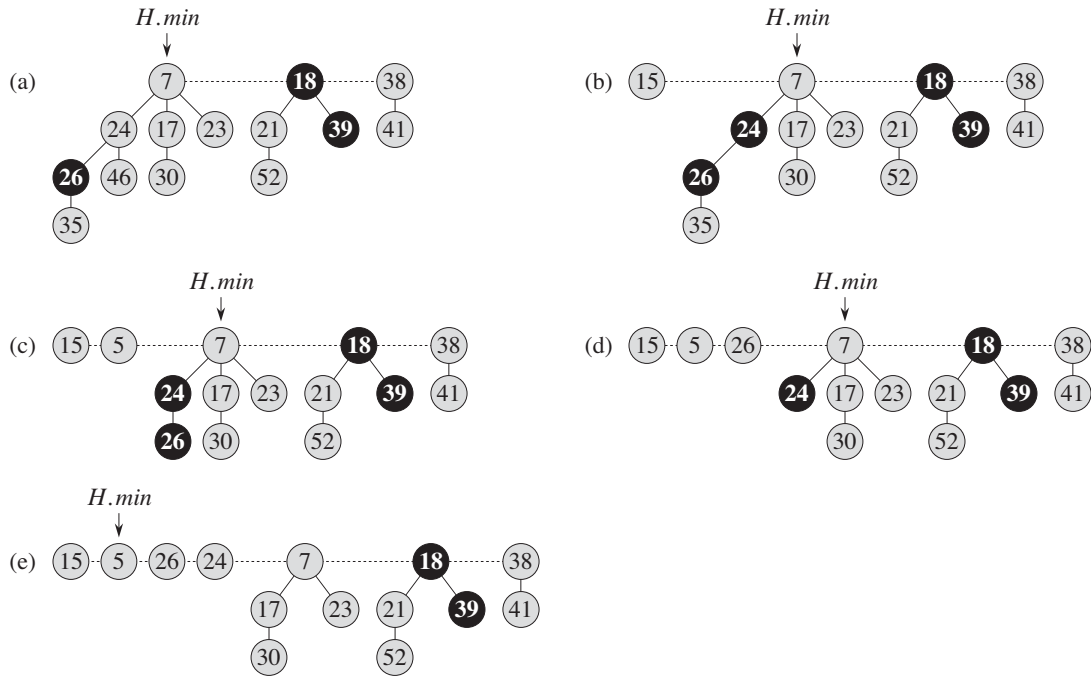
We are not yet done, because  $x$  might be the second child cut from its parent  $y$  since the time that  $y$  was linked to another node. Therefore, line 7 of FIB-HEAP-DECREASE-KEY attempts to perform a *cascading-cut* operation on  $y$ . If  $y$  is a root, then the test in line 2 of CASCADING-CUT causes the procedure to just return. If  $y$  is unmarked, the procedure marks it in line 4, since its first child has just been cut, and returns. If  $y$  is marked, however, it has just lost its second child;  $y$  is cut in line 5, and CASCADING-CUT calls itself recursively in line 6 on  $y$ 's parent  $z$ . The CASCADING-CUT procedure recurses its way up the tree until it finds either a root or an unmarked node.

Once all the cascading cuts have occurred, lines 8–9 of FIB-HEAP-DECREASE-KEY finish up by updating  $H.min$  if necessary. The only node whose key changed was the node  $x$  whose key decreased. Thus, the new minimum node is either the original minimum node or node  $x$ .

Figure 19.5 shows the execution of two calls of FIB-HEAP-DECREASE-KEY, starting with the Fibonacci heap shown in Figure 19.5(a). The first call, shown in Figure 19.5(b), involves no cascading cuts. The second call, shown in Figures 19.5(c)–(e), invokes two cascading cuts.

We shall now show that the amortized cost of FIB-HEAP-DECREASE-KEY is only  $O(1)$ . We start by determining its actual cost. The FIB-HEAP-DECREASE-KEY procedure takes  $O(1)$  time, plus the time to perform the cascading cuts. Suppose that a given invocation of FIB-HEAP-DECREASE-KEY results in  $c$  calls of CASCADING-CUT (the call made from line 7 of FIB-HEAP-DECREASE-KEY followed by  $c - 1$  recursive calls of CASCADING-CUT). Each call of CASCADING-CUT takes  $O(1)$  time exclusive of recursive calls. Thus, the actual cost of FIB-HEAP-DECREASE-KEY, including all recursive calls, is  $O(c)$ .

We next compute the change in potential. Let  $H$  denote the Fibonacci heap just prior to the FIB-HEAP-DECREASE-KEY operation. The call to CUT in line 6 of



**Figure 19.5** Two calls of FIB-HEAP-DECREASE-KEY. (a) The initial Fibonacci heap. (b) The node with key 46 has its key decreased to 15. The node becomes a root, and its parent (with key 24), which had previously been unmarked, becomes marked. (c)–(e) The node with key 35 has its key decreased to 5. In part (c), the node, now with key 5, becomes a root. Its parent, with key 26, is marked, so a cascading cut occurs. The node with key 26 is cut from its parent and made an unmarked root in (d). Another cascading cut occurs, since the node with key 24 is marked as well. This node is cut from its parent and made an unmarked root in part (e). The cascading cuts stop at this point, since the node with key 7 is a root. (Even if this node were not a root, the cascading cuts would stop, since it is unmarked.) Part (e) shows the result of the FIB-HEAP-DECREASE-KEY operation, with *H.min* pointing to the new minimum node.

FIB-HEAP-DECREASE-KEY creates a new tree rooted at node  $x$  and clears  $x$ 's mark bit (which may have already been FALSE). Each call of CASCADING-CUT, except for the last one, cuts a marked node and clears the mark bit. Afterward, the Fibonacci heap contains  $t(H) + c$  trees (the original  $t(H)$  trees,  $c - 1$  trees produced by cascading cuts, and the tree rooted at  $x$ ) and at most  $m(H) - c + 2$  marked nodes ( $c - 1$  were unmarked by cascading cuts and the last call of CASCADING-CUT may have marked a node). The change in potential is therefore at most

$$((t(H) + c) + 2(m(H) - c + 2)) - (t(H) + 2m(H)) = 4 - c.$$

Thus, the amortized cost of FIB-HEAP-DECREASE-KEY is at most

$$O(c) + 4 - c = O(1) ,$$

since we can scale up the units of potential to dominate the constant hidden in  $O(c)$ .

You can now see why we defined the potential function to include a term that is twice the number of marked nodes. When a marked node  $y$  is cut by a cascading cut, its mark bit is cleared, which reduces the potential by 2. One unit of potential pays for the cut and the clearing of the mark bit, and the other unit compensates for the unit increase in potential due to node  $y$  becoming a root.

### Deleting a node

The following pseudocode deletes a node from an  $n$ -node Fibonacci heap in  $O(D(n))$  amortized time. We assume that there is no key value of  $-\infty$  currently in the Fibonacci heap.

FIB-HEAP-DELETE( $H, x$ )

- 1 FIB-HEAP-DECREASE-KEY( $H, x, -\infty$ )
- 2 FIB-HEAP-EXTRACT-MIN( $H$ )

FIB-HEAP-DELETE makes  $x$  become the minimum node in the Fibonacci heap by giving it a uniquely small key of  $-\infty$ . The FIB-HEAP-EXTRACT-MIN procedure then removes node  $x$  from the Fibonacci heap. The amortized time of FIB-HEAP-DELETE is the sum of the  $O(1)$  amortized time of FIB-HEAP-DECREASE-KEY and the  $O(D(n))$  amortized time of FIB-HEAP-EXTRACT-MIN. Since we shall see in Section 19.4 that  $D(n) = O(\lg n)$ , the amortized time of FIB-HEAP-DELETE is  $O(\lg n)$ .

### Exercises

#### 19.3-1

Suppose that a root  $x$  in a Fibonacci heap is marked. Explain how  $x$  came to be a marked root. Argue that it doesn't matter to the analysis that  $x$  is marked, even though it is not a root that was first linked to another node and then lost one child.

#### 19.3-2

Justify the  $O(1)$  amortized time of FIB-HEAP-DECREASE-KEY as an average cost per operation by using aggregate analysis.

---

## 19.4 Bounding the maximum degree

To prove that the amortized time of FIB-HEAP-EXTRACT-MIN and FIB-HEAP-DELETE is  $O(\lg n)$ , we must show that the upper bound  $D(n)$  on the degree of any node of an  $n$ -node Fibonacci heap is  $O(\lg n)$ . In particular, we shall show that  $D(n) \leq \lfloor \log_\phi n \rfloor$ , where  $\phi$  is the golden ratio, defined in equation (3.24) as

$$\phi = (1 + \sqrt{5})/2 = 1.61803 \dots$$

The key to the analysis is as follows. For each node  $x$  within a Fibonacci heap, define  $\text{size}(x)$  to be the number of nodes, including  $x$  itself, in the subtree rooted at  $x$ . (Note that  $x$  need not be in the root list—it can be any node at all.) We shall show that  $\text{size}(x)$  is exponential in  $x.\text{degree}$ . Bear in mind that  $x.\text{degree}$  is always maintained as an accurate count of the degree of  $x$ .

### Lemma 19.1

Let  $x$  be any node in a Fibonacci heap, and suppose that  $x.\text{degree} = k$ . Let  $y_1, y_2, \dots, y_k$  denote the children of  $x$  in the order in which they were linked to  $x$ , from the earliest to the latest. Then,  $y_1.\text{degree} \geq 0$  and  $y_i.\text{degree} \geq i - 2$  for  $i = 2, 3, \dots, k$ .

**Proof** Obviously,  $y_1.\text{degree} \geq 0$ .

For  $i \geq 2$ , we note that when  $y_i$  was linked to  $x$ , all of  $y_1, y_2, \dots, y_{i-1}$  were children of  $x$ , and so we must have had  $x.\text{degree} \geq i - 1$ . Because node  $y_i$  is linked to  $x$  (by CONSOLIDATE) only if  $x.\text{degree} = y_i.\text{degree}$ , we must have also had  $y_i.\text{degree} \geq i - 1$  at that time. Since then, node  $y_i$  has lost at most one child, since it would have been cut from  $x$  (by CASCADING-CUT) if it had lost two children. We conclude that  $y_i.\text{degree} \geq i - 2$ . ■

We finally come to the part of the analysis that explains the name “Fibonacci heaps.” Recall from Section 3.2 that for  $k = 0, 1, 2, \dots$ , the  $k$ th Fibonacci number is defined by the recurrence

$$F_k = \begin{cases} 0 & \text{if } k = 0, \\ 1 & \text{if } k = 1, \\ F_{k-1} + F_{k-2} & \text{if } k \geq 2. \end{cases}$$

The following lemma gives another way to express  $F_k$ .



**Lemma 19.2**

For all integers  $k \geq 0$ ,

$$F_{k+2} = 1 + \sum_{i=0}^k F_i .$$

**Proof** The proof is by induction on  $k$ . When  $k = 0$ ,

$$\begin{aligned} 1 + \sum_{i=0}^0 F_i &= 1 + F_0 \\ &= 1 + 0 \\ &= F_2 . \end{aligned}$$

We now assume the inductive hypothesis that  $F_{k+1} = 1 + \sum_{i=0}^{k-1} F_i$ , and we have

$$\begin{aligned} F_{k+2} &= F_k + F_{k+1} \\ &= F_k + \left( 1 + \sum_{i=0}^{k-1} F_i \right) \\ &= 1 + \sum_{i=0}^k F_i . \end{aligned} \quad \blacksquare$$

**Lemma 19.3**

For all integers  $k \geq 0$ , the  $(k + 2)$ nd Fibonacci number satisfies  $F_{k+2} \geq \phi^k$ .

**Proof** The proof is by induction on  $k$ . The base cases are for  $k = 0$  and  $k = 1$ . When  $k = 0$  we have  $F_2 = 1 = \phi^0$ , and when  $k = 1$  we have  $F_3 = 2 > 1.619 > \phi^1$ . The inductive step is for  $k \geq 2$ , and we assume that  $F_{i+2} > \phi^i$  for  $i = 0, 1, \dots, k-1$ . Recall that  $\phi$  is the positive root of equation (3.23),  $x^2 = x + 1$ . Thus, we have

$$\begin{aligned} F_{k+2} &= F_{k+1} + F_k \\ &\geq \phi^{k-1} + \phi^{k-2} \quad (\text{by the inductive hypothesis}) \\ &= \phi^{k-2}(\phi + 1) \\ &= \phi^{k-2} \cdot \phi^2 \quad (\text{by equation (3.23)}) \\ &= \phi^k . \end{aligned} \quad \blacksquare$$

The following lemma and its corollary complete the analysis.

**Lemma 19.4**

Let  $x$  be any node in a Fibonacci heap, and let  $k = x.degree$ . Then  $size(x) \geq F_{k+2} \geq \phi^k$ , where  $\phi = (1 + \sqrt{5})/2$ .

**Proof** Let  $s_k$  denote the minimum possible size of any node of degree  $k$  in any Fibonacci heap. Trivially,  $s_0 = 1$  and  $s_1 = 2$ . The number  $s_k$  is at most  $size(x)$  and, because adding children to a node cannot decrease the node's size, the value of  $s_k$  increases monotonically with  $k$ . Consider some node  $z$ , in any Fibonacci heap, such that  $z.degree = k$  and  $size(z) = s_k$ . Because  $s_k \leq size(x)$ , we compute a lower bound on  $size(x)$  by computing a lower bound on  $s_k$ . As in Lemma 19.1, let  $y_1, y_2, \dots, y_k$  denote the children of  $z$  in the order in which they were linked to  $z$ . To bound  $s_k$ , we count one for  $z$  itself and one for the first child  $y_1$  (for which  $size(y_1) \geq 1$ ), giving

$$\begin{aligned} size(x) &\geq s_k \\ &\geq 2 + \sum_{i=2}^k s_{y_i.degree} \\ &\geq 2 + \sum_{i=2}^k s_{i-2}, \end{aligned}$$

where the last line follows from Lemma 19.1 (so that  $y_i.degree \geq i - 2$ ) and the monotonicity of  $s_k$  (so that  $s_{y_i.degree} \geq s_{i-2}$ ).

We now show by induction on  $k$  that  $s_k \geq F_{k+2}$  for all nonnegative integers  $k$ . The bases, for  $k = 0$  and  $k = 1$ , are trivial. For the inductive step, we assume that  $k \geq 2$  and that  $s_i \geq F_{i+2}$  for  $i = 0, 1, \dots, k - 1$ . We have

$$\begin{aligned} s_k &\geq 2 + \sum_{i=2}^k s_{i-2} \\ &\geq 2 + \sum_{i=2}^k F_i \\ &= 1 + \sum_{i=0}^k F_i \\ &= F_{k+2} && \text{(by Lemma 19.2)} \\ &\geq \phi^k && \text{(by Lemma 19.3) .} \end{aligned}$$

Thus, we have shown that  $size(x) \geq s_k \geq F_{k+2} \geq \phi^k$ . ■

**Corollary 19.5**

The maximum degree  $D(n)$  of any node in an  $n$ -node Fibonacci heap is  $O(\lg n)$ .

**Proof** Let  $x$  be any node in an  $n$ -node Fibonacci heap, and let  $k = x.\text{degree}$ . By Lemma 19.4, we have  $n \geq \text{size}(x) \geq \phi^k$ . Taking base- $\phi$  logarithms gives us  $k \leq \log_\phi n$ . (In fact, because  $k$  is an integer,  $k \leq \lfloor \log_\phi n \rfloor$ .) The maximum degree  $D(n)$  of any node is thus  $O(\lg n)$ . ■

**Exercises****19.4-1**

Professor Pinocchio claims that the height of an  $n$ -node Fibonacci heap is  $O(\lg n)$ . Show that the professor is mistaken by exhibiting, for any positive integer  $n$ , a sequence of Fibonacci-heap operations that creates a Fibonacci heap consisting of just one tree that is a linear chain of  $n$  nodes.

**19.4-2**

Suppose we generalize the cascading-cut rule to cut a node  $x$  from its parent as soon as it loses its  $k$ th child, for some integer constant  $k$ . (The rule in Section 19.3 uses  $k = 2$ .) For what values of  $k$  is  $D(n) = O(\lg n)$ ?

---

**Problems**
**19-1 Alternative implementation of deletion**

Professor Pisano has proposed the following variant of the FIB-HEAP-DELETE procedure, claiming that it runs faster when the node being deleted is not the node pointed to by  $H.\text{min}$ .

PISANO-DELETE( $H, x$ )

```

1  if  $x == H.\text{min}$ 
2      FIB-HEAP-EXTRACT-MIN( $H$ )
3  else  $y = x.p$ 
4      if  $y \neq \text{NIL}$ 
5          CUT( $H, x, y$ )
6          CASCADING-CUT( $H, y$ )
7      add  $x$ 's child list to the root list of  $H$ 
8      remove  $x$  from the root list of  $H$ 
```

- a. The professor's claim that this procedure runs faster is based partly on the assumption that line 7 can be performed in  $O(1)$  actual time. What is wrong with this assumption?
- b. Give a good upper bound on the actual time of PISANO-DELETE when  $x$  is not  $H.min$ . Your bound should be in terms of  $x.degree$  and the number  $c$  of calls to the CASCADING-CUT procedure.
- c. Suppose that we call PISANO-DELETE( $H, x$ ), and let  $H'$  be the Fibonacci heap that results. Assuming that node  $x$  is not a root, bound the potential of  $H'$  in terms of  $x.degree$ ,  $c$ ,  $t(H)$ , and  $m(H)$ .
- d. Conclude that the amortized time for PISANO-DELETE is asymptotically no better than for FIB-HEAP-DELETE, even when  $x \neq H.min$ .

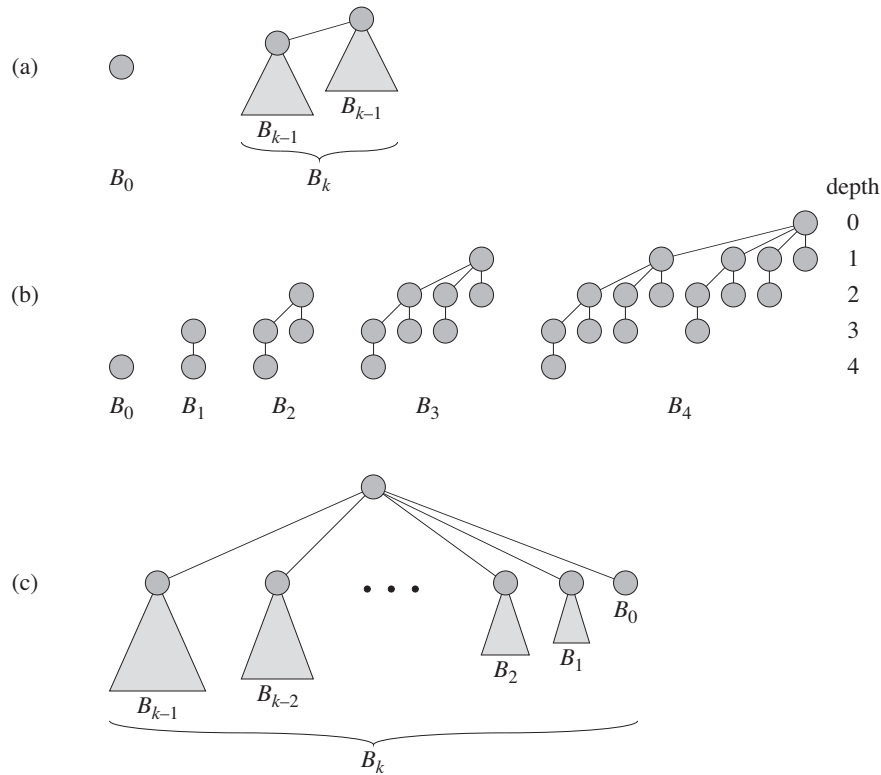
### 19-2 Binomial trees and binomial heaps

The **binomial tree**  $B_k$  is an ordered tree (see Section B.5.2) defined recursively. As shown in Figure 19.6(a), the binomial tree  $B_0$  consists of a single node. The binomial tree  $B_k$  consists of two binomial trees  $B_{k-1}$  that are linked together so that the root of one is the leftmost child of the root of the other. Figure 19.6(b) shows the binomial trees  $B_0$  through  $B_4$ .

- a. Show that for the binomial tree  $B_k$ ,
  1. there are  $2^k$  nodes,
  2. the height of the tree is  $k$ ,
  3. there are exactly  $\binom{k}{i}$  nodes at depth  $i$  for  $i = 0, 1, \dots, k$ , and
  4. the root has degree  $k$ , which is greater than that of any other node; moreover, as Figure 19.6(c) shows, if we number the children of the root from left to right by  $k-1, k-2, \dots, 0$ , then child  $i$  is the root of a subtree  $B_i$ .

A **binomial heap**  $H$  is a set of binomial trees that satisfies the following properties:

1. Each node has a *key* (like a Fibonacci heap).
2. Each binomial tree in  $H$  obeys the min-heap property.
3. For any nonnegative integer  $k$ , there is at most one binomial tree in  $H$  whose root has degree  $k$ .
- b. Suppose that a binomial heap  $H$  has a total of  $n$  nodes. Discuss the relationship between the binomial trees that  $H$  contains and the binary representation of  $n$ . Conclude that  $H$  consists of at most  $\lfloor \lg n \rfloor + 1$  binomial trees.



**Figure 19.6** (a) The recursive definition of the binomial tree  $B_k$ . Triangles represent rooted subtrees. (b) The binomial trees  $B_0$  through  $B_4$ . Node depths in  $B_4$  are shown. (c) Another way of looking at the binomial tree  $B_k$ .

Suppose that we represent a binomial heap as follows. The left-child, right-sibling scheme of Section 10.4 represents each binomial tree within a binomial heap. Each node contains its key; pointers to its parent, to its leftmost child, and to the sibling immediately to its right (these pointers are NIL when appropriate); and its degree (as in Fibonacci heaps, how many children it has). The roots form a singly linked root list, ordered by the degrees of the roots (from low to high), and we access the binomial heap by a pointer to the first node on the root list.

- c. Complete the description of how to represent a binomial heap (i.e., name the attributes, describe when attributes have the value NIL, and define how the root list is organized), and show how to implement the same seven operations on binomial heaps as this chapter implemented on Fibonacci heaps. Each operation should run in  $O(\lg n)$  worst-case time, where  $n$  is the number of nodes in

the binomial heap (or in the case of the UNION operation, in the two binomial heaps that are being united). The MAKE-HEAP operation should take constant time.

- d. Suppose that we were to implement only the mergeable-heap operations on a Fibonacci heap (i.e., we do not implement the DECREASE-KEY or DELETE operations). How would the trees in a Fibonacci heap resemble those in a binomial heap? How would they differ? Show that the maximum degree in an  $n$ -node Fibonacci heap would be at most  $\lfloor \lg n \rfloor$ .
- e. Professor McGee has devised a new data structure based on Fibonacci heaps. A McGee heap has the same structure as a Fibonacci heap and supports just the mergeable-heap operations. The implementations of the operations are the same as for Fibonacci heaps, except that insertion and union consolidate the root list as their last step. What are the worst-case running times of operations on McGee heaps?

### 19-3 More Fibonacci-heap operations

We wish to augment a Fibonacci heap  $H$  to support two new operations without changing the amortized running time of any other Fibonacci-heap operations.

- a. The operation FIB-HEAP-CHANGE-KEY( $H, x, k$ ) changes the key of node  $x$  to the value  $k$ . Give an efficient implementation of FIB-HEAP-CHANGE-KEY, and analyze the amortized running time of your implementation for the cases in which  $k$  is greater than, less than, or equal to  $x.key$ .
- b. Give an efficient implementation of FIB-HEAP-PRUNE( $H, r$ ), which deletes  $q = \min(r, H.n)$  nodes from  $H$ . You may choose any  $q$  nodes to delete. Analyze the amortized running time of your implementation. (*Hint:* You may need to modify the data structure and potential function.)

### 19-4 2-3-4 heaps

Chapter 18 introduced the 2-3-4 tree, in which every internal node (other than possibly the root) has two, three, or four children and all leaves have the same depth. In this problem, we shall implement **2-3-4 heaps**, which support the mergeable-heap operations.

The 2-3-4 heaps differ from 2-3-4 trees in the following ways. In 2-3-4 heaps, only leaves store keys, and each leaf  $x$  stores exactly one key in the attribute  $x.key$ . The keys in the leaves may appear in any order. Each internal node  $x$  contains a value  $x.small$  that is equal to the smallest key stored in any leaf in the subtree rooted at  $x$ . The root  $r$  contains an attribute  $r.height$  that gives the height of the

tree. Finally, 2-3-4 heaps are designed to be kept in main memory, so that disk reads and writes are not needed.

Implement the following 2-3-4 heap operations. In parts (a)–(e), each operation should run in  $O(\lg n)$  time on a 2-3-4 heap with  $n$  elements. The UNION operation in part (f) should run in  $O(\lg n)$  time, where  $n$  is the number of elements in the two input heaps.

- a.* MINIMUM, which returns a pointer to the leaf with the smallest key.
- b.* DECREASE-KEY, which decreases the key of a given leaf  $x$  to a given value  $k \leq x.key$ .
- c.* INSERT, which inserts leaf  $x$  with key  $k$ .
- d.* DELETE, which deletes a given leaf  $x$ .
- e.* EXTRACT-MIN, which extracts the leaf with the smallest key.
- f.* UNION, which unites two 2-3-4 heaps, returning a single 2-3-4 heap and destroying the input heaps.

---

## Chapter notes

Fredman and Tarjan [114] introduced Fibonacci heaps. Their paper also describes the application of Fibonacci heaps to the problems of single-source shortest paths, all-pairs shortest paths, weighted bipartite matching, and the minimum-spanning-tree problem.

Subsequently, Driscoll, Gabow, Shrairman, and Tarjan [96] developed “relaxed heaps” as an alternative to Fibonacci heaps. They devised two varieties of relaxed heaps. One gives the same amortized time bounds as Fibonacci heaps. The other allows DECREASE-KEY to run in  $O(1)$  worst-case (not amortized) time and EXTRACT-MIN and DELETE to run in  $O(\lg n)$  worst-case time. Relaxed heaps also have some advantages over Fibonacci heaps in parallel algorithms.

See also the chapter notes for Chapter 6 for other data structures that support fast DECREASE-KEY operations when the sequence of values returned by EXTRACT-MIN calls are monotonically increasing over time and the data are integers in a specific range.

In previous chapters, we saw data structures that support the operations of a priority queue—binary heaps in Chapter 6, red-black trees in Chapter 13,<sup>1</sup> and Fibonacci heaps in Chapter 19. In each of these data structures, at least one important operation took  $O(\lg n)$  time, either worst case or amortized. In fact, because each of these data structures bases its decisions on comparing keys, the  $\Omega(n \lg n)$  lower bound for sorting in Section 8.1 tells us that at least one operation will have to take  $\Omega(\lg n)$  time. Why? If we could perform both the INSERT and EXTRACT-MIN operations in  $o(\lg n)$  time, then we could sort  $n$  keys in  $o(n \lg n)$  time by first performing  $n$  INSERT operations, followed by  $n$  EXTRACT-MIN operations.

We saw in Chapter 8, however, that sometimes we can exploit additional information about the keys to sort in  $o(n \lg n)$  time. In particular, with counting sort we can sort  $n$  keys, each an integer in the range 0 to  $k$ , in time  $\Theta(n + k)$ , which is  $\Theta(n)$  when  $k = O(n)$ .

Since we can circumvent the  $\Omega(n \lg n)$  lower bound for sorting when the keys are integers in a bounded range, you might wonder whether we can perform each of the priority-queue operations in  $o(\lg n)$  time in a similar scenario. In this chapter, we shall see that we can: van Emde Boas trees support the priority-queue operations, and a few others, each in  $O(\lg \lg n)$  worst-case time. The hitch is that the keys must be integers in the range 0 to  $n - 1$ , with no duplicates allowed.

Specifically, van Emde Boas trees support each of the dynamic set operations listed on page 230—SEARCH, INSERT, DELETE, MINIMUM, MAXIMUM, SUCCESSOR, and PREDECESSOR—in  $O(\lg \lg n)$  time. In this chapter, we will omit discussion of satellite data and focus only on storing keys. Because we concentrate on keys and disallow duplicate keys to be stored, instead of describing the SEARCH

---

<sup>1</sup>Chapter 13 does not explicitly discuss how to implement EXTRACT-MIN and DECREASE-KEY, but we can easily build these operations for any data structure that supports MINIMUM, DELETE, and INSERT.



operation, we will implement the simpler operation  $\text{MEMBER}(S, x)$ , which returns a boolean indicating whether the value  $x$  is currently in dynamic set  $S$ .

So far, we have used the parameter  $n$  for two distinct purposes: the number of elements in the dynamic set, and the range of the possible values. To avoid any further confusion, from here on we will use  $n$  to denote the number of elements currently in the set and  $u$  as the range of possible values, so that each van Emde Boas tree operation runs in  $O(\lg \lg u)$  time. We call the set  $\{0, 1, 2, \dots, u-1\}$  the *universe* of values that can be stored and  $u$  the *universe size*. We assume throughout this chapter that  $u$  is an exact power of 2, i.e.,  $u = 2^k$  for some integer  $k \geq 1$ .

Section 20.1 starts us out by examining some simple approaches that will get us going in the right direction. We enhance these approaches in Section 20.2, introducing proto van Emde Boas structures, which are recursive but do not achieve our goal of  $O(\lg \lg u)$ -time operations. Section 20.3 modifies proto van Emde Boas structures to develop van Emde Boas trees, and it shows how to implement each operation in  $O(\lg \lg u)$  time.

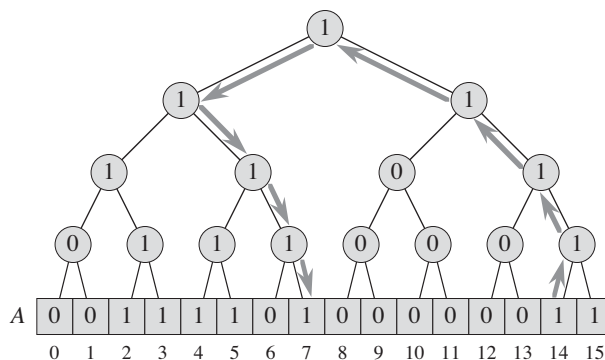
---

## 20.1 Preliminary approaches

In this section, we shall examine various approaches for storing a dynamic set. Although none will achieve the  $O(\lg \lg u)$  time bounds that we desire, we will gain insights that will help us understand van Emde Boas trees when we see them later in this chapter.

### Direct addressing

Direct addressing, as we saw in Section 11.1, provides the simplest approach to storing a dynamic set. Since in this chapter we are concerned only with storing keys, we can simplify the direct-addressing approach to store the dynamic set as a bit vector, as discussed in Exercise 11.1-2. To store a dynamic set of values from the universe  $\{0, 1, 2, \dots, u-1\}$ , we maintain an array  $A[0..u-1]$  of  $u$  bits. The entry  $A[x]$  holds a 1 if the value  $x$  is in the dynamic set, and it holds a 0 otherwise. Although we can perform each of the INSERT, DELETE, and MEMBER operations in  $O(1)$  time with a bit vector, the remaining operations—MINIMUM, MAXIMUM, SUCCESSOR, and PREDECESSOR—each take  $\Theta(u)$  time in the worst case because



**Figure 20.1** A binary tree of bits superimposed on top of a bit vector representing the set  $\{2, 3, 4, 5, 7, 14, 15\}$  when  $u = 16$ . Each internal node contains a 1 if and only if some leaf in its subtree contains a 1. The arrows show the path followed to determine the predecessor of 14 in the set.

we might have to scan through  $\Theta(u)$  elements.<sup>2</sup> For example, if a set contains only the values 0 and  $u - 1$ , then to find the successor of 0, we would have to scan entries 1 through  $u - 2$  before finding a 1 in  $A[u - 1]$ .

### Superimposing a binary tree structure

We can short-cut long scans in the bit vector by superimposing a binary tree of bits on top of it. Figure 20.1 shows an example. The entries of the bit vector form the leaves of the binary tree, and each internal node contains a 1 if and only if any leaf in its subtree contains a 1. In other words, the bit stored in an internal node is the logical-or of its two children.

The operations that took  $\Theta(u)$  worst-case time with an unadorned bit vector now use the tree structure:

- To find the minimum value in the set, start at the root and head down toward the leaves, always taking the leftmost node containing a 1.
- To find the maximum value in the set, start at the root and head down toward the leaves, always taking the rightmost node containing a 1.

<sup>2</sup>We assume throughout this chapter that MINIMUM and MAXIMUM return NIL if the dynamic set is empty and that SUCCESSOR and PREDECESSOR return NIL if the element they are given has no successor or predecessor, respectively.

- To find the successor of  $x$ , start at the leaf indexed by  $x$ , and head up toward the root until we enter a node from the left and this node has a 1 in its right child  $z$ . Then head down through node  $z$ , always taking the leftmost node containing a 1 (i.e., find the minimum value in the subtree rooted at the right child  $z$ ).
- To find the predecessor of  $x$ , start at the leaf indexed by  $x$ , and head up toward the root until we enter a node from the right and this node has a 1 in its left child  $z$ . Then head down through node  $z$ , always taking the rightmost node containing a 1 (i.e., find the maximum value in the subtree rooted at the left child  $z$ ).

Figure 20.1 shows the path taken to find the predecessor, 7, of the value 14.

We also augment the INSERT and DELETE operations appropriately. When inserting a value, we store a 1 in each node on the simple path from the appropriate leaf up to the root. When deleting a value, we go from the appropriate leaf up to the root, recomputing the bit in each internal node on the path as the logical-or of its two children.

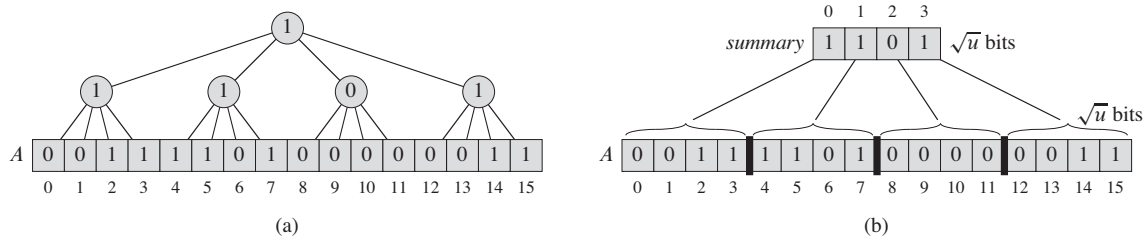
Since the height of the tree is  $\lg u$  and each of the above operations makes at most one pass up the tree and at most one pass down, each operation takes  $O(\lg u)$  time in the worst case.

This approach is only marginally better than just using a red-black tree. We can still perform the MEMBER operation in  $O(1)$  time, whereas searching a red-black tree takes  $O(\lg n)$  time. Then again, if the number  $n$  of elements stored is much smaller than the size  $u$  of the universe, a red-black tree would be faster for all the other operations.

### Superimposing a tree of constant height

What happens if we superimpose a tree with greater degree? Let us assume that the size of the universe is  $u = 2^{2k}$  for some integer  $k$ , so that  $\sqrt{u}$  is an integer. Instead of superimposing a binary tree on top of the bit vector, we superimpose a tree of degree  $\sqrt{u}$ . Figure 20.2(a) shows such a tree for the same bit vector as in Figure 20.1. The height of the resulting tree is always 2.

As before, each internal node stores the logical-or of the bits within its subtree, so that the  $\sqrt{u}$  internal nodes at depth 1 summarize each group of  $\sqrt{u}$  values. As Figure 20.2(b) demonstrates, we can think of these nodes as an array *summary*[0.. $\sqrt{u} - 1$ ], where *summary*[ $i$ ] contains a 1 if and only if the subarray  $A[i\sqrt{u}..(i+1)\sqrt{u} - 1]$  contains a 1. We call this  $\sqrt{u}$ -bit subarray of  $A$  the  $i$ th **cluster**. For a given value of  $x$ , the bit  $A[x]$  appears in cluster number  $\lfloor x/\sqrt{u} \rfloor$ . Now INSERT becomes an  $O(1)$ -time operation: to insert  $x$ , set both  $A[x]$  and *summary*[ $\lfloor x/\sqrt{u} \rfloor$ ] to 1. We can use the *summary* array to perform



**Figure 20.2** (a) A tree of degree  $\sqrt{u}$  superimposed on top of the same bit vector as in Figure 20.1. Each internal node stores the logical-or of the bits in its subtree. (b) A view of the same structure, but with the internal nodes at depth 1 treated as an array  $summary[0.. \sqrt{u} - 1]$ , where  $summary[i]$  is the logical-or of the subarray  $A[i\sqrt{u}.. (i+1)\sqrt{u} - 1]$ .

each of the operations MINIMUM, MAXIMUM, SUCCESSOR, PREDECESSOR, and DELETE in  $O(\sqrt{u})$  time:

- To find the minimum (maximum) value, find the leftmost (rightmost) entry in  $summary$  that contains a 1, say  $summary[i]$ , and then do a linear search within the  $i$ th cluster for the leftmost (rightmost) 1.
- To find the successor (predecessor) of  $x$ , first search to the right (left) within its cluster. If we find a 1, that position gives the result. Otherwise, let  $i = \lfloor x/\sqrt{u} \rfloor$  and search to the right (left) within the  $summary$  array from index  $i$ . The first position that holds a 1 gives the index of a cluster. Search within that cluster for the leftmost (rightmost) 1. That position holds the successor (predecessor).
- To delete the value  $x$ , let  $i = \lfloor x/\sqrt{u} \rfloor$ . Set  $A[x]$  to 0 and then set  $summary[i]$  to the logical-or of the bits in the  $i$ th cluster.

In each of the above operations, we search through at most two clusters of  $\sqrt{u}$  bits plus the  $summary$  array, and so each operation takes  $O(\sqrt{u})$  time.

At first glance, it seems as though we have made negative progress. Superimposing a binary tree gave us  $O(\lg u)$ -time operations, which are asymptotically faster than  $O(\sqrt{u})$  time. Using a tree of degree  $\sqrt{u}$  will turn out to be a key idea of van Emde Boas trees, however. We continue down this path in the next section.

## Exercises

### 20.1-1

Modify the data structures in this section to support duplicate keys.

**20.1-2**

Modify the data structures in this section to support keys that have associated satellite data.

**20.1-3**

Observe that, using the structures in this section, the way we find the successor and predecessor of a value  $x$  does not depend on whether  $x$  is in the set at the time. Show how to find the successor of  $x$  in a binary search tree when  $x$  is not stored in the tree.

**20.1-4**

Suppose that instead of superimposing a tree of degree  $\sqrt{u}$ , we were to superimpose a tree of degree  $u^{1/k}$ , where  $k > 1$  is a constant. What would be the height of such a tree, and how long would each of the operations take?

---

**20.2 A recursive structure**

In this section, we modify the idea of superimposing a tree of degree  $\sqrt{u}$  on top of a bit vector. In the previous section, we used a summary structure of size  $\sqrt{u}$ , with each entry pointing to another structure of size  $\sqrt{u}$ . Now, we make the structure recursive, shrinking the universe size by the square root at each level of recursion. Starting with a universe of size  $u$ , we make structures holding  $\sqrt{u} = u^{1/2}$  items, which themselves hold structures of  $u^{1/4}$  items, which hold structures of  $u^{1/8}$  items, and so on, down to a base size of 2.

For simplicity, in this section, we assume that  $u = 2^{2^k}$  for some integer  $k$ , so that  $u, u^{1/2}, u^{1/4}, \dots$  are integers. This restriction would be quite severe in practice, allowing only values of  $u$  in the sequence 2, 4, 16, 256, 65536,  $\dots$ . We shall see in the next section how to relax this assumption and assume only that  $u = 2^k$  for some integer  $k$ . Since the structure we examine in this section is only a precursor to the true van Emde Boas tree structure, we tolerate this restriction in favor of aiding our understanding.

Recalling that our goal is to achieve running times of  $O(\lg \lg u)$  for the operations, let's think about how we might obtain such running times. At the end of Section 4.3, we saw that by changing variables, we could show that the recurrence

$$T(n) = 2T(\lfloor \sqrt{n} \rfloor) + \lg n \quad (20.1)$$

has the solution  $T(n) = O(\lg n \lg \lg n)$ . Let's consider a similar, but simpler, recurrence:

$$T(u) = T(\sqrt{u}) + O(1). \quad (20.2)$$

If we use the same technique, changing variables, we can show that recurrence (20.2) has the solution  $T(u) = O(\lg \lg u)$ . Let  $m = \lg u$ , so that  $u = 2^m$  and we have

$$T(2^m) = T(2^{m/2}) + O(1) .$$

Now we rename  $S(m) = T(2^m)$ , giving the new recurrence

$$S(m) = S(m/2) + O(1) .$$

By case 2 of the master method, this recurrence has the solution  $S(m) = O(\lg m)$ . We change back from  $S(m)$  to  $T(u)$ , giving  $T(u) = T(2^m) = S(m) = O(\lg m) = O(\lg \lg u)$ .

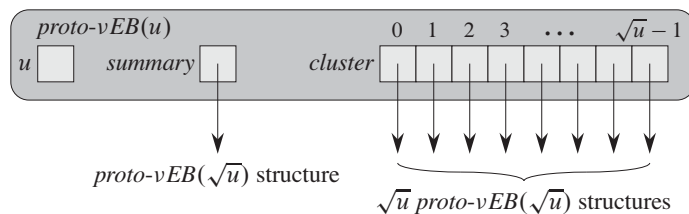
Recurrence (20.2) will guide our search for a data structure. We will design a recursive data structure that shrinks by a factor of  $\sqrt{u}$  in each level of its recursion. When an operation traverses this data structure, it will spend a constant amount of time at each level before recursing to the level below. Recurrence (20.2) will then characterize the running time of the operation.

Here is another way to think of how the term  $\lg \lg u$  ends up in the solution to recurrence (20.2). As we look at the universe size in each level of the recursive data structure, we see the sequence  $u, u^{1/2}, u^{1/4}, u^{1/8}, \dots$ . If we consider how many bits we need to store the universe size at each level, we need  $\lg u$  at the top level, and each level needs half the bits of the previous level. In general, if we start with  $b$  bits and halve the number of bits at each level, then after  $\lg b$  levels, we get down to just one bit. Since  $b = \lg u$ , we see that after  $\lg \lg u$  levels, we have a universe size of 2.

Looking back at the data structure in Figure 20.2, a given value  $x$  resides in cluster number  $\lfloor x/\sqrt{u} \rfloor$ . If we view  $x$  as a  $\lg u$ -bit binary integer, that cluster number,  $\lfloor x/\sqrt{u} \rfloor$ , is given by the most significant  $(\lg u)/2$  bits of  $x$ . Within its cluster,  $x$  appears in position  $x \bmod \sqrt{u}$ , which is given by the least significant  $(\lg u)/2$  bits of  $x$ . We will need to index in this way, and so let us define some functions that will help us do so:

$$\begin{aligned} \text{high}(x) &= \lfloor x/\sqrt{u} \rfloor , \\ \text{low}(x) &= x \bmod \sqrt{u} , \\ \text{index}(x, y) &= x\sqrt{u} + y . \end{aligned}$$

The function  $\text{high}(x)$  gives the most significant  $(\lg u)/2$  bits of  $x$ , producing the number of  $x$ 's cluster. The function  $\text{low}(x)$  gives the least significant  $(\lg u)/2$  bits of  $x$  and provides  $x$ 's position within its cluster. The function  $\text{index}(x, y)$  builds an element number from  $x$  and  $y$ , treating  $x$  as the most significant  $(\lg u)/2$  bits of the element number and  $y$  as the least significant  $(\lg u)/2$  bits. We have the identity  $x = \text{index}(\text{high}(x), \text{low}(x))$ . The value of  $u$  used by each of these functions will



**Figure 20.3** The information in a *proto-vEB(u)* structure when  $u \geq 4$ . The structure contains the universe size  $u$ , a pointer *summary* to a *proto-vEB(sqrt(u))* structure, and an array *cluster*[0 ..  $\sqrt{u}-1$ ] of  $\sqrt{u}$  pointers to *proto-vEB(sqrt(u))* structures.

always be the universe size of the data structure in which we call the function, which changes as we descend into the recursive structure.

### 20.2.1 Proto van Emde Boas structures

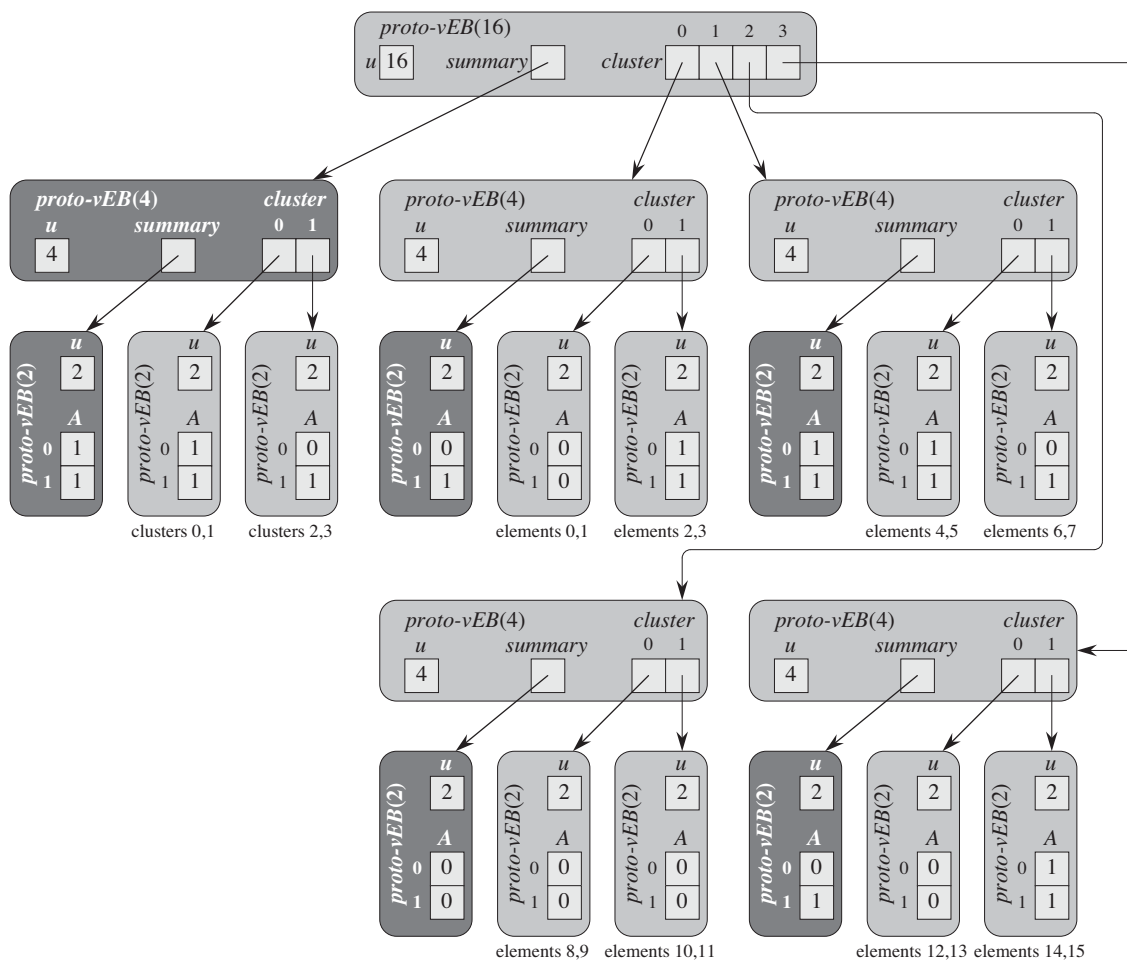
Taking our cue from recurrence (20.2), let us design a recursive data structure to support the operations. Although this data structure will fail to achieve our goal of  $O(\lg \lg u)$  time for some operations, it serves as a basis for the van Emde Boas tree structure that we will see in Section 20.3.

For the universe  $\{0, 1, 2, \dots, u-1\}$ , we define a **proto van Emde Boas structure**, or **proto-vEB structure**, which we denote as *proto-vEB(u)*, recursively as follows. Each *proto-vEB(u)* structure contains an attribute  $u$  giving its universe size. In addition, it contains the following:

- If  $u = 2$ , then it is the base size, and it contains an array  $A[0..1]$  of two bits.
- Otherwise,  $u = 2^{2^k}$  for some integer  $k \geq 1$ , so that  $u \geq 4$ . In addition to the universe size  $u$ , the data structure *proto-vEB(u)* contains the following attributes, illustrated in Figure 20.3:
  - a pointer named *summary* to a *proto-vEB(sqrt(u))* structure and
  - an array *cluster*[0 ..  $\sqrt{u}-1$ ] of  $\sqrt{u}$  pointers, each to a *proto-vEB(sqrt(u))* structure.

The element  $x$ , where  $0 \leq x < u$ , is recursively stored in the cluster numbered  $\text{high}(x)$  as element  $\text{low}(x)$  within that cluster.

In the two-level structure of the previous section, each node stores a summary array of size  $\sqrt{u}$ , in which each entry contains a bit. From the index of each entry, we can compute the starting index of the subarray of size  $\sqrt{u}$  that the bit summarizes. In the proto-vEB structure, we use explicit pointers rather than index



**Figure 20.4** A *proto-vEB(16)* structure representing the set {2, 3, 4, 5, 7, 14, 15}. It points to four *proto-vEB(4)* structures in *cluster*[0..3], and to a summary structure, which is also a *proto-vEB(4)*. Each *proto-vEB(4)* structure points to two *proto-vEB(2)* structures in *cluster*[0..1], and to a *proto-vEB(2)* summary. Each *proto-vEB(2)* structure contains just an array *A*[0..1] of two bits. The *proto-vEB(2)* structures above “elements *i*, *j*” store bits *i* and *j* of the actual dynamic set, and the *proto-vEB(2)* structures above “clusters *i*, *j*” store the summary bits for clusters *i* and *j* in the top-level *proto-vEB(16)* structure. For clarity, heavy shading indicates the top level of a *proto-vEB* structure that stores summary information for its parent structure; such a *proto-vEB* structure is otherwise identical to any other *proto-vEB* structure with the same universe size.



calculations. The array *summary* contains the summary bits stored recursively in a proto-vEB structure, and the array *cluster* contains  $\sqrt{u}$  pointers.

Figure 20.4 shows a fully expanded *proto-vEB*(16) structure representing the set  $\{2, 3, 4, 5, 7, 14, 15\}$ . If the value  $i$  is in the proto-vEB structure pointed to by *summary*, then the  $i$ th cluster contains some value in the set being represented. As in the tree of constant height, *cluster*[ $i$ ] represents the values  $i\sqrt{u}$  through  $(i + 1)\sqrt{u} - 1$ , which form the  $i$ th cluster.

At the base level, the elements of the actual dynamic sets are stored in some of the *proto-vEB*(2) structures, and the remaining *proto-vEB*(2) structures store summary bits. Beneath each of the non-summary base structures, the figure indicates which bits it stores. For example, the *proto-vEB*(2) structure labeled “elements 6,7” stores bit 6 (0, since element 6 is not in the set) in its  $A[0]$  and bit 7 (1, since element 7 is in the set) in its  $A[1]$ .

Like the clusters, each summary is just a dynamic set with universe size  $\sqrt{u}$ , and so we represent each summary as a *proto-vEB*( $\sqrt{u}$ ) structure. The four summary bits for the main *proto-vEB*(16) structure are in the leftmost *proto-vEB*(4) structure, and they ultimately appear in two *proto-vEB*(2) structures. For example, the *proto-vEB*(2) structure labeled “clusters 2,3” has  $A[0] = 0$ , indicating that cluster 2 of the *proto-vEB*(16) structure (containing elements 8, 9, 10, 11) is all 0, and  $A[1] = 1$ , telling us that cluster 3 (containing elements 12, 13, 14, 15) has at least one 1. Each *proto-vEB*(4) structure points to its own summary, which is itself stored as a *proto-vEB*(2) structure. For example, look at the *proto-vEB*(2) structure just to the left of the one labeled “elements 0,1.” Because its  $A[0]$  is 0, it tells us that the “elements 0,1” structure is all 0, and because its  $A[1]$  is 1, we know that the “elements 2,3” structure contains at least one 1.

## 20.2.2 Operations on a proto van Emde Boas structure

We shall now describe how to perform operations on a proto-vEB structure. We first examine the query operations—MEMBER, MINIMUM, MAXIMUM, and SUCCESSOR—which do not change the proto-vEB structure. We then discuss INSERT and DELETE. We leave MAXIMUM and PREDECESSOR, which are symmetric to MINIMUM and SUCCESSOR, respectively, as Exercise 20.2-1.

Each of the MEMBER, SUCCESSOR, PREDECESSOR, INSERT, and DELETE operations takes a parameter  $x$ , along with a proto-vEB structure  $V$ . Each of these operations assumes that  $0 \leq x < V.u$ .

### Determining whether a value is in the set

To perform MEMBER( $x$ ), we need to find the bit corresponding to  $x$  within the appropriate *proto-vEB*(2) structure. We can do so in  $O(\lg \lg u)$  time, bypassing

the *summary* structures altogether. The following procedure takes a *proto-vEB* structure  $V$  and a value  $x$ , and it returns a bit indicating whether  $x$  is in the dynamic set held by  $V$ .

```

PROTO-VEB-MEMBER( $V, x$ )
1  if  $V.u == 2$ 
2      return  $V.A[x]$ 
3  else return PROTO-VEB-MEMBER( $V.cluster[high(x)], low(x)$ )

```

The PROTO-VEB-MEMBER procedure works as follows. Line 1 tests whether we are in a base case, where  $V$  is a *proto-vEB*(2) structure. Line 2 handles the base case, simply returning the appropriate bit of array  $A$ . Line 3 deals with the recursive case, “drilling down” into the appropriate smaller *proto-vEB* structure. The value  $high(x)$  says which *proto-vEB*( $\sqrt{u}$ ) structure we visit, and  $low(x)$  determines which element within that *proto-vEB*( $\sqrt{u}$ ) structure we are querying.

Let’s see what happens when we call PROTO-VEB-MEMBER( $V, 6$ ) on the *proto-vEB*(16) structure in Figure 20.4. Since  $high(6) = 1$  when  $u = 16$ , we recurse into the *proto-vEB*(4) structure in the upper right, and we ask about element  $low(6) = 2$  of that structure. In this recursive call,  $u = 4$ , and so we recurse again. With  $u = 4$ , we have  $high(2) = 1$  and  $low(2) = 0$ , and so we ask about element 0 of the *proto-vEB*(2) structure in the upper right. This recursive call turns out to be a base case, and so it returns  $A[0] = 0$  back up through the chain of recursive calls. Thus, we get the result that PROTO-VEB-MEMBER( $V, 6$ ) returns 0, indicating that 6 is not in the set.

To determine the running time of PROTO-VEB-MEMBER, let  $T(u)$  denote its running time on a *proto-vEB*( $u$ ) structure. Each recursive call takes constant time, not including the time taken by the recursive calls that it makes. When PROTO-VEB-MEMBER makes a recursive call, it makes a call on a *proto-vEB*( $\sqrt{u}$ ) structure. Thus, we can characterize the running time by the recurrence  $T(u) = T(\sqrt{u}) + O(1)$ , which we have already seen as recurrence (20.2). Its solution is  $T(u) = O(\lg \lg u)$ , and so we conclude that PROTO-VEB-MEMBER runs in time  $O(\lg \lg u)$ .

### Finding the minimum element

Now we examine how to perform the MINIMUM operation. The procedure PROTO-VEB-MINIMUM( $V$ ) returns the minimum element in the *proto-vEB* structure  $V$ , or NIL if  $V$  represents an empty set.

PROTO-VEB-MINIMUM( $V$ )

```

1  if  $V.u == 2$ 
2      if  $V.A[0] == 1$ 
3          return 0
4      elseif  $V.A[1] == 1$ 
5          return 1
6      else return NIL
7  else  $min\_cluster = \text{PROTO-VEB-MINIMUM}(V.summary)$ 
8      if  $min\_cluster == \text{NIL}$ 
9          return NIL
10     else  $offset = \text{PROTO-VEB-MINIMUM}(V.cluster[min\_cluster])$ 
11     return  $\text{index}(min\_cluster, offset)$ 

```

This procedure works as follows. Line 1 tests for the base case, which lines 2–6 handle by brute force. Lines 7–11 handle the recursive case. First, line 7 finds the number of the first cluster that contains an element of the set. It does so by recursively calling PROTO-VEB-MINIMUM on  $V.summary$ , which is a  $proto\text{-}vEB(\sqrt{u})$  structure. Line 7 assigns this cluster number to the variable  $min\_cluster$ . If the set is empty, then the recursive call returned NIL, and line 9 returns NIL. Otherwise, the minimum element of the set is somewhere in cluster number  $min\_cluster$ . The recursive call in line 10 finds the offset within the cluster of the minimum element in this cluster. Finally, line 11 constructs the value of the minimum element from the cluster number and offset, and it returns this value.

Although querying the summary information allows us to quickly find the cluster containing the minimum element, because this procedure makes two recursive calls on  $proto\text{-}vEB(\sqrt{u})$  structures, it does not run in  $O(\lg \lg u)$  time in the worst case. Letting  $T(u)$  denote the worst-case time for PROTO-VEB-MINIMUM on a  $proto\text{-}vEB(u)$  structure, we have the recurrence

$$T(u) = 2T(\sqrt{u}) + O(1). \quad (20.3)$$

Again, we use a change of variables to solve this recurrence, letting  $m = \lg u$ , which gives

$$T(2^m) = 2T(2^{m/2}) + O(1).$$

Renaming  $S(m) = T(2^m)$  gives

$$S(m) = 2S(m/2) + O(1),$$

which, by case 1 of the master method, has the solution  $S(m) = \Theta(m)$ . By changing back from  $S(m)$  to  $T(u)$ , we have that  $T(u) = T(2^m) = S(m) = \Theta(m) = \Theta(\lg u)$ . Thus, we see that because of the second recursive call, PROTO-VEB-MINIMUM runs in  $\Theta(\lg u)$  time rather than the desired  $O(\lg \lg u)$  time.

### Finding the successor

The SUCCESSOR operation is even worse. In the worst case, it makes two recursive calls, along with a call to PROTO-VEB-MINIMUM. The procedure PROTO-VEB-SUCCESSOR( $V, x$ ) returns the smallest element in the proto-vEB structure  $V$  that is greater than  $x$ , or NIL if no element in  $V$  is greater than  $x$ . It does not require  $x$  to be a member of the set, but it does assume that  $0 \leq x < V.u$ .

PROTO-VEB-SUCCESSOR( $V, x$ )

```

1  if  $V.u == 2$ 
2      if  $x == 0$  and  $V.A[1] == 1$ 
3          return 1
4      else return NIL
5  else  $offset = \text{PROTO-VEB-SUCCESSOR}(V.cluster[high(x)], low(x))$ 
6      if  $offset \neq \text{NIL}$ 
7          return  $\text{index}(high(x), offset)$ 
8      else  $succ-cluster = \text{PROTO-VEB-SUCCESSOR}(V.summary, high(x))$ 
9          if  $succ-cluster == \text{NIL}$ 
10             return NIL
11         else  $offset = \text{PROTO-VEB-MINIMUM}(V.cluster[succ-cluster])$ 
12         return  $\text{index}(succ-cluster, offset)$ 
```

The PROTO-VEB-SUCCESSOR procedure works as follows. As usual, line 1 tests for the base case, which lines 2–4 handle by brute force: the only way that  $x$  can have a successor within a *proto-vEB*(2) structure is when  $x = 0$  and  $A[1]$  is 1. Lines 5–12 handle the recursive case. Line 5 searches for a successor to  $x$  within  $x$ 's cluster, assigning the result to *offset*. Line 6 determines whether  $x$  has a successor within its cluster; if it does, then line 7 computes and returns the value of this successor. Otherwise, we have to search in other clusters. Line 8 assigns to *succ-cluster* the number of the next nonempty cluster, using the summary information to find it. Line 9 tests whether *succ-cluster* is NIL, with line 10 returning NIL if all succeeding clusters are empty. If *succ-cluster* is non-NIL, line 11 assigns the first element within that cluster to *offset*, and line 12 computes and returns the minimum element in that cluster.

In the worst case, PROTO-VEB-SUCCESSOR calls itself recursively twice on *proto-vEB*( $\sqrt{u}$ ) structures, and it makes one call to PROTO-VEB-MINIMUM on a *proto-vEB*( $\sqrt{u}$ ) structure. Thus, the recurrence for the worst-case running time  $T(u)$  of PROTO-VEB-SUCCESSOR is

$$\begin{aligned}
 T(u) &= 2T(\sqrt{u}) + \Theta(\lg \sqrt{u}) \\
 &= 2T(\sqrt{u}) + \Theta(\lg u) .
 \end{aligned}$$

We can employ the same technique that we used for recurrence (20.1) to show that this recurrence has the solution  $T(u) = \Theta(\lg u \lg \lg u)$ . Thus, **PROTO-VEB-SUCCESSOR** is asymptotically slower than **PROTO-VEB-MINIMUM**.

### Inserting an element

To insert an element, we need to insert it into the appropriate cluster and also set the summary bit for that cluster to 1. The procedure **PROTO-VEB-INSERT**( $V, x$ ) inserts the value  $x$  into the proto-veb structure  $V$ .

**PROTO-VEB-INSERT**( $V, x$ )

```

1  if  $V.u == 2$ 
2       $V.A[x] = 1$ 
3  else PROTO-VEB-INSERT( $V.cluster[high(x)], low(x)$ )
4      PROTO-VEB-INSERT( $V.summary, high(x)$ )
```

In the base case, line 2 sets the appropriate bit in the array  $A$  to 1. In the recursive case, the recursive call in line 3 inserts  $x$  into the appropriate cluster, and line 4 sets the summary bit for that cluster to 1.

Because **PROTO-VEB-INSERT** makes two recursive calls in the worst case, recurrence (20.3) characterizes its running time. Hence, **PROTO-VEB-INSERT** runs in  $\Theta(\lg u)$  time.

### Deleting an element

The **DELETE** operation is more complicated than insertion. Whereas we can always set a summary bit to 1 when inserting, we cannot always reset the same summary bit to 0 when deleting. We need to determine whether any bit in the appropriate cluster is 1. As we have defined proto-veb structures, we would have to examine all  $\sqrt{u}$  bits within a cluster to determine whether any of them are 1. Alternatively, we could add an attribute  $n$  to the proto-veb structure, counting how many elements it has. We leave implementation of **PROTO-VEB-DELETE** as Exercises 20.2-2 and 20.2-3.

Clearly, we need to modify the proto-veb structure to get each operation down to making at most one recursive call. We will see in the next section how to do so.

### Exercises

#### 20.2-1

Write pseudocode for the procedures **PROTO-VEB-MAXIMUM** and **PROTO-VEB-PREDECESSOR**.

**20.2-2**

Write pseudocode for `PROTO-VEB-DELETE`. It should update the appropriate summary bit by scanning the related bits within the cluster. What is the worst-case running time of your procedure?

**20.2-3**

Add the attribute  $n$  to each proto-veb structure, giving the number of elements currently in the set it represents, and write pseudocode for `PROTO-VEB-DELETE` that uses the attribute  $n$  to decide when to reset summary bits to 0. What is the worst-case running time of your procedure? What other procedures need to change because of the new attribute? Do these changes affect their running times?

**20.2-4**

Modify the proto-veb structure to support duplicate keys.

**20.2-5**

Modify the proto-veb structure to support keys that have associated satellite data.

**20.2-6**

Write pseudocode for a procedure that creates a *proto-veb*( $u$ ) structure.

**20.2-7**

Argue that if line 9 of `PROTO-VEB-MINIMUM` is executed, then the proto-veb structure is empty.

**20.2-8**

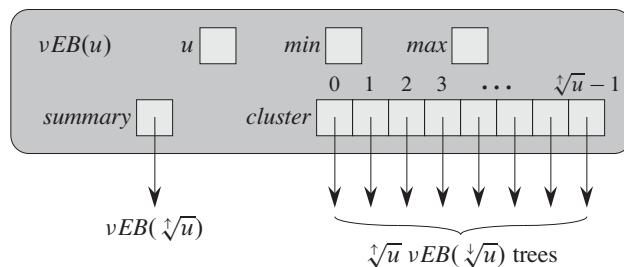
Suppose that we designed a proto-veb structure in which each *cluster* array had only  $u^{1/4}$  elements. What would the running times of each operation be?

---

## 20.3 The van Emde Boas tree

The proto-veb structure of the previous section is close to what we need to achieve  $O(\lg \lg u)$  running times. It falls short because we have to recurse too many times in most of the operations. In this section, we shall design a data structure that is similar to the proto-veb structure but stores a little more information, thereby removing the need for some of the recursion.

In Section 20.2, we observed that the assumption that we made about the universe size—that  $u = 2^{2^k}$  for some integer  $k$ —is unduly restrictive, confining the possible values of  $u$  to an overly sparse set. From this point on, therefore, we will allow the universe size  $u$  to be any exact power of 2, and when  $\sqrt{u}$  is not an inte-



**Figure 20.5** The information in a  $vEB(u)$  tree when  $u > 2$ . The structure contains the universe size  $u$ , elements  $min$  and  $max$ , a pointer  $summary$  to a  $vEB(\sqrt[4]{u})$  tree, and an array  $cluster[0.. \sqrt[4]{u}-1]$  of  $\sqrt[4]{u}$  pointers to  $vEB(\sqrt[4]{u})$  trees.

ger—that is, if  $u$  is an odd power of 2 ( $u = 2^{2k+1}$  for some integer  $k \geq 0$ )—then we will divide the  $\lg u$  bits of a number into the most significant  $\lceil (\lg u)/2 \rceil$  bits and the least significant  $\lfloor (\lg u)/2 \rfloor$  bits. For convenience, we denote  $2^{\lceil (\lg u)/2 \rceil}$  (the “upper square root” of  $u$ ) by  $\sqrt[4]{u}$  and  $2^{\lfloor (\lg u)/2 \rfloor}$  (the “lower square root” of  $u$ ) by  $\sqrt[4]{u}$ , so that  $u = \sqrt[4]{u} \cdot \sqrt[4]{u}$  and, when  $u$  is an even power of 2 ( $u = 2^{2k}$  for some integer  $k$ ),  $\sqrt[4]{u} = \sqrt[4]{u} = \sqrt{u}$ . Because we now allow  $u$  to be an odd power of 2, we must redefine our helpful functions from Section 20.2:

$$\begin{aligned} \text{high}(x) &= \lfloor x / \sqrt[4]{u} \rfloor, \\ \text{low}(x) &= x \bmod \sqrt[4]{u}, \\ \text{index}(x, y) &= x \sqrt[4]{u} + y. \end{aligned}$$

### 20.3.1 van Emde Boas trees

The *van Emde Boas tree*, or *vEB tree*, modifies the proto-vEB structure. We denote a vEB tree with a universe size of  $u$  as  $vEB(u)$  and, unless  $u$  equals the base size of 2, the attribute  $summary$  points to a  $vEB(\sqrt[4]{u})$  tree and the array  $cluster[0.. \sqrt[4]{u}-1]$  points to  $\sqrt[4]{u}$   $vEB(\sqrt[4]{u})$  trees. As Figure 20.5 illustrates, a vEB tree contains two attributes not found in a proto-vEB structure:

- $min$  stores the minimum element in the vEB tree, and
- $max$  stores the maximum element in the vEB tree.

Furthermore, the element stored in  $min$  does not appear in any of the recursive  $vEB(\sqrt[4]{u})$  trees that the  $cluster$  array points to. The elements stored in a  $vEB(u)$  tree  $V$ , therefore, are  $V.min$  plus all the elements recursively stored in the  $vEB(\sqrt[4]{u})$  trees pointed to by  $V.cluster[0.. \sqrt[4]{u}-1]$ . Note that when a vEB tree contains two or more elements, we treat  $min$  and  $max$  differently: the element

stored in *min* does not appear in any of the clusters, but the element stored in *max* does.

Since the base size is 2, a  $vEB(2)$  tree does not need the array  $A$  that the corresponding *proto- $vEB(2)$*  structure has. Instead, we can determine its elements from its *min* and *max* attributes. In a  $vEB$  tree with no elements, regardless of its universe size  $u$ , both *min* and *max* are NIL.

Figure 20.6 shows a  $vEB(16)$  tree  $V$  holding the set  $\{2, 3, 4, 5, 7, 14, 15\}$ . Because the smallest element is 2,  $V.min$  equals 2, and even though  $high(2) = 0$ , the element 2 does not appear in the  $vEB(4)$  tree pointed to by  $V.cluster[0]$ : notice that  $V.cluster[0].min$  equals 3, and so 2 is not in this  $vEB$  tree. Similarly, since  $V.cluster[0].min$  equals 3, and 2 and 3 are the only elements in  $V.cluster[0]$ , the  $vEB(2)$  clusters within  $V.cluster[0]$  are empty.

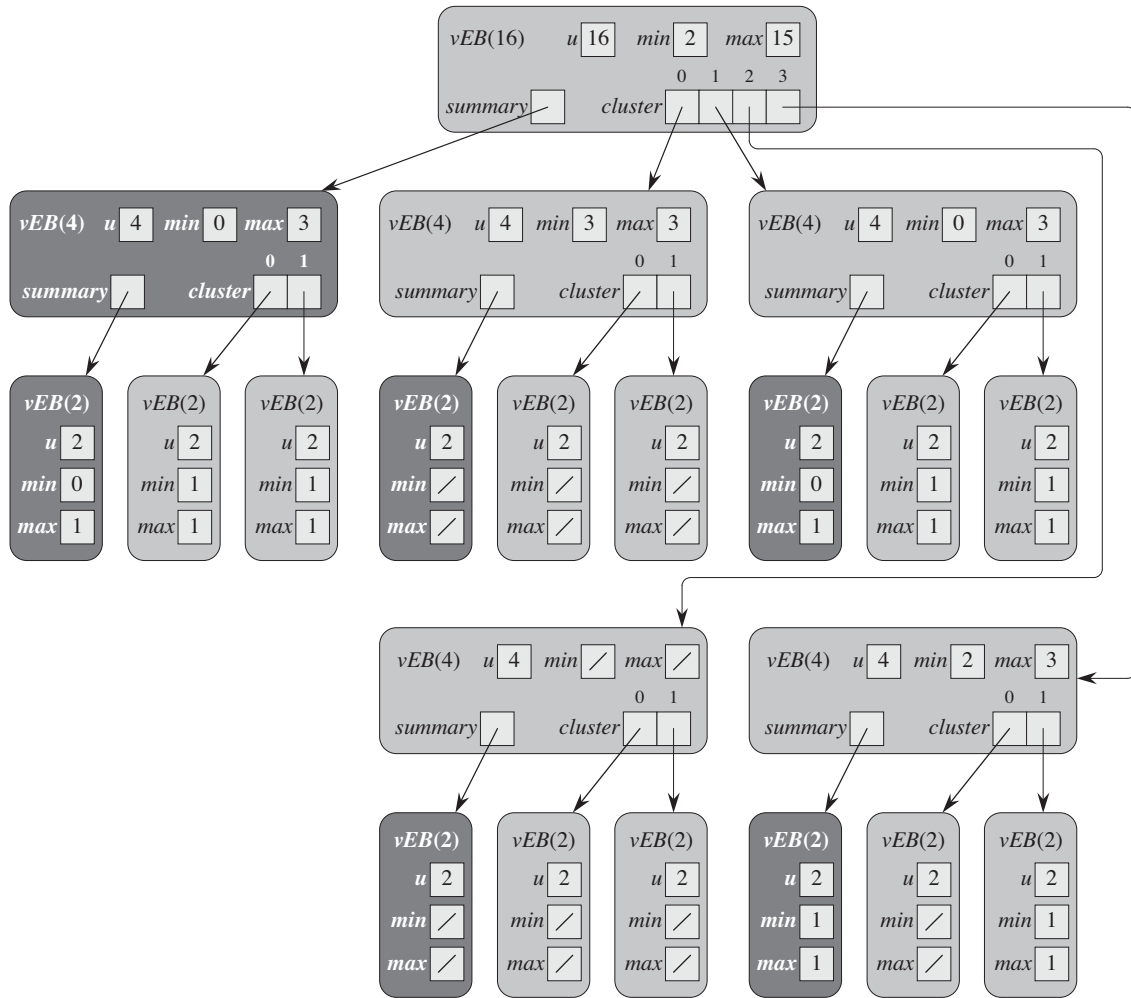
The *min* and *max* attributes will turn out to be key to reducing the number of recursive calls within the operations on  $vEB$  trees. These attributes will help us in four ways:

1. The MINIMUM and MAXIMUM operations do not even need to recurse, for they can just return the values of *min* or *max*.
2. The SUCCESSOR operation can avoid making a recursive call to determine whether the successor of a value  $x$  lies within  $high(x)$ . That is because  $x$ 's successor lies within its cluster if and only if  $x$  is strictly less than the *max* attribute of its cluster. A symmetric argument holds for PREDECESSOR and *min*.
3. We can tell whether a  $vEB$  tree has no elements, exactly one element, or at least two elements in constant time from its *min* and *max* values. This ability will help in the INSERT and DELETE operations. If *min* and *max* are both NIL, then the  $vEB$  tree has no elements. If *min* and *max* are non-NIL but are equal to each other, then the  $vEB$  tree has exactly one element. Otherwise, both *min* and *max* are non-NIL but are unequal, and the  $vEB$  tree has two or more elements.
4. If we know that a  $vEB$  tree is empty, we can insert an element into it by updating only its *min* and *max* attributes. Hence, we can insert into an empty  $vEB$  tree in constant time. Similarly, if we know that a  $vEB$  tree has only one element, we can delete that element in constant time by updating only *min* and *max*. These properties will allow us to cut short the chain of recursive calls.

Even if the universe size  $u$  is an odd power of 2, the difference in the sizes of the summary  $vEB$  tree and the clusters will not turn out to affect the asymptotic running times of the  $vEB$ -tree operations. The recursive procedures that implement the  $vEB$ -tree operations will all have running times characterized by the recurrence

$$T(u) \leq T(\lceil \sqrt{u} \rceil) + O(1). \quad (20.4)$$





**Figure 20.6** A  $vEB(16)$  tree corresponding to the proto- $vEB$  tree in Figure 20.4. It stores the set  $\{2, 3, 4, 5, 7, 14, 15\}$ . Slashes indicate NIL values. The value stored in the  $min$  attribute of a  $vEB$  tree does not appear in any of its clusters. Heavy shading serves the same purpose here as in Figure 20.4.

This recurrence looks similar to recurrence (20.2), and we will solve it in a similar fashion. Letting  $m = \lg u$ , we rewrite it as

$$T(2^m) \leq T(2^{\lceil m/2 \rceil}) + O(1) .$$

Noting that  $\lceil m/2 \rceil \leq 2m/3$  for all  $m \geq 2$ , we have

$$T(2^m) \leq T(2^{2m/3}) + O(1) .$$

Letting  $S(m) = T(2^m)$ , we rewrite this last recurrence as

$$S(m) \leq S(2m/3) + O(1) ,$$

which, by case 2 of the master method, has the solution  $S(m) = O(\lg m)$ . (In terms of the asymptotic solution, the fraction  $2/3$  does not make any difference compared with the fraction  $1/2$ , because when we apply the master method, we find that  $\log_{3/2} 1 = \log_2 1 = 0$ .) Thus, we have  $T(u) = T(2^m) = S(m) = O(\lg m) = O(\lg \lg u)$ .

Before using a van Emde Boas tree, we must know the universe size  $u$ , so that we can create a van Emde Boas tree of the appropriate size that initially represents an empty set. As Problem 20-1 asks you to show, the total space requirement of a van Emde Boas tree is  $O(u)$ , and it is straightforward to create an empty tree in  $O(u)$  time. In contrast, we can create an empty red-black tree in constant time. Therefore, we might not want to use a van Emde Boas tree when we perform only a small number of operations, since the time to create the data structure would exceed the time saved in the individual operations. This drawback is usually not significant, since we typically use a simple data structure, such as an array or linked list, to represent a set with only a few elements.

### 20.3.2 Operations on a van Emde Boas tree

We are now ready to see how to perform operations on a van Emde Boas tree. As we did for the proto van Emde Boas structure, we will consider the querying operations first, and then INSERT and DELETE. Due to the slight asymmetry between the minimum and maximum elements in a vEB tree—when a vEB tree contains at least two elements, the minimum element does not appear within a cluster but the maximum element does—we will provide pseudocode for all five querying operations. As in the operations on proto van Emde Boas structures, the operations here take parameters  $V$  and  $x$ , where  $V$  is a van Emde Boas tree and  $x$  is an element, assume that  $0 \leq x < V.u$ .

#### Finding the minimum and maximum elements

Because we store the minimum and maximum in the attributes *min* and *max*, two of the operations are one-liners, taking constant time:

VEB-TREE-MINIMUM( $V$ )

1   **return**  $V.min$

VEB-TREE-MAXIMUM( $V$ )

1   **return**  $V.max$

### Determining whether a value is in the set

The procedure VEB-TREE-MEMBER( $V, x$ ) has a recursive case like that of PROTO-VEB-MEMBER, but the base case is a little different. We also check directly whether  $x$  equals the minimum or maximum element. Since a vEB tree doesn't store bits as a proto-vEB structure does, we design VEB-TREE-MEMBER to return TRUE or FALSE rather than 1 or 0.

VEB-TREE-MEMBER( $V, x$ )

1   **if**  $x == V.min$  or  $x == V.max$

2       **return** TRUE

3   **elseif**  $V.u == 2$

4       **return** FALSE

5   **else return** VEB-TREE-MEMBER( $V.cluster[high(x)], low(x)$ )

Line 1 checks to see whether  $x$  equals either the minimum or maximum element. If it does, line 2 returns TRUE. Otherwise, line 3 tests for the base case. Since a  $vEB(2)$  tree has no elements other than those in  $min$  and  $max$ , if it is the base case, line 4 returns FALSE. The other possibility—it is not a base case and  $x$  equals neither  $min$  nor  $max$ —is handled by the recursive call in line 5.

Recurrence (20.4) characterizes the running time of the VEB-TREE-MEMBER procedure, and so this procedure takes  $O(\lg \lg u)$  time.

### Finding the successor and predecessor

Next we see how to implement the SUCCESSOR operation. Recall that the procedure PROTO-VEB-SUCCESSOR( $V, x$ ) could make two recursive calls: one to determine whether  $x$ 's successor resides in the same cluster as  $x$  and, if it does not, one to find the cluster containing  $x$ 's successor. Because we can access the maximum value in a vEB tree quickly, we can avoid making two recursive calls, and instead make one recursive call on either a cluster or on the summary, but not on both.

```

VEB-TREE-SUCCESSOR( $V, x$ )
1  if  $V.u == 2$ 
2      if  $x == 0$  and  $V.max == 1$ 
3          return 1
4      else return NIL
5  elseif  $V.min \neq \text{NIL}$  and  $x < V.min$ 
6      return  $V.min$ 
7  else  $max\text{-}low = \text{VEB-TREE-MAXIMUM}(V.cluster[high(x)])$ 
8      if  $max\text{-}low \neq \text{NIL}$  and  $low(x) < max\text{-}low$ 
9           $offset = \text{VEB-TREE-SUCCESSOR}(V.cluster[high(x)], low(x))$ 
10         return  $index(high(x), offset)$ 
11     else  $succ\text{-}cluster = \text{VEB-TREE-SUCCESSOR}(V.summary, high(x))$ 
12         if  $succ\text{-}cluster == \text{NIL}$ 
13             return NIL
14         else  $offset = \text{VEB-TREE-MINIMUM}(V.cluster[succ\text{-}cluster])$ 
15         return  $index(succ\text{-}cluster, offset)$ 

```

This procedure has six **return** statements and several cases. We start with the base case in lines 2–4, which returns 1 in line 3 if we are trying to find the successor of 0 and 1 is in the 2-element set; otherwise, the base case returns NIL in line 4.

If we are not in the base case, we next check in line 5 whether  $x$  is strictly less than the minimum element. If so, then we simply return the minimum element in line 6.

If we get to line 7, then we know that we are not in a base case and that  $x$  is greater than or equal to the minimum value in the vEB tree  $V$ . Line 7 assigns to  $max\text{-}low$  the maximum element in  $x$ 's cluster. If  $x$ 's cluster contains some element that is greater than  $x$ , then we know that  $x$ 's successor lies somewhere within  $x$ 's cluster. Line 8 tests for this condition. If  $x$ 's successor is within  $x$ 's cluster, then line 9 determines where in the cluster it is, and line 10 returns the successor in the same way as line 7 of **PROTO-VEB-SUCCESSOR**.

We get to line 11 if  $x$  is greater than or equal to the greatest element in its cluster. In this case, lines 11–15 find  $x$ 's successor in the same way as lines 8–12 of **PROTO-VEB-SUCCESSOR**.

It is easy to see how recurrence (20.4) characterizes the running time of **VEB-TREE-SUCCESSOR**. Depending on the result of the test in line 7, the procedure calls itself recursively in either line 9 (on a vEB tree with universe size  $\sqrt[4]{u}$ ) or line 11 (on a vEB tree with universe size  $\sqrt[4]{u}$ ). In either case, the one recursive call is on a vEB tree with universe size at most  $\sqrt[4]{u}$ . The remainder of the procedure, including the calls to **VEB-TREE-MINIMUM** and **VEB-TREE-MAXIMUM**, takes  $O(1)$  time. Hence, **VEB-TREE-SUCCESSOR** runs in  $O(\lg \lg u)$  worst-case time.

The `VEB-TREE-PREDECESSOR` procedure is symmetric to the `VEB-TREE-SUCCESSOR` procedure, but with one additional case:

```

VEB-TREE-PREDECESSOR( $V, x$ )
1  if  $V.u == 2$ 
2      if  $x == 1$  and  $V.min == 0$ 
3          return 0
4      else return NIL
5  elseif  $V.max \neq \text{NIL}$  and  $x > V.max$ 
6      return  $V.max$ 
7  else  $min-low = \text{VEB-TREE-MINIMUM}(V.cluster[high(x)])$ 
8      if  $min-low \neq \text{NIL}$  and  $low(x) > min-low$ 
9           $offset = \text{VEB-TREE-PREDECESSOR}(V.cluster[high(x)], low(x))$ 
10         return  $\text{index}(high(x), offset)$ 
11     else  $pred-cluster = \text{VEB-TREE-PREDECESSOR}(V.summary, high(x))$ 
12         if  $pred-cluster == \text{NIL}$ 
13             if  $V.min \neq \text{NIL}$  and  $x > V.min$ 
14                 return  $V.min$ 
15             else return NIL
16         else  $offset = \text{VEB-TREE-MAXIMUM}(V.cluster[pred-cluster])$ 
17         return  $\text{index}(pred-cluster, offset)$ 

```

Lines 13–14 form the additional case. This case occurs when  $x$ 's predecessor, if it exists, does not reside in  $x$ 's cluster. In `VEB-TREE-SUCCESSOR`, we were assured that if  $x$ 's successor resides outside of  $x$ 's cluster, then it must reside in a higher-numbered cluster. But if  $x$ 's predecessor is the minimum value in vEB tree  $V$ , then the successor resides in no cluster at all. Line 13 checks for this condition, and line 14 returns the minimum value as appropriate.

This extra case does not affect the asymptotic running time of `VEB-TREE-PREDECESSOR` when compared with `VEB-TREE-SUCCESSOR`, and so `VEB-TREE-PREDECESSOR` runs in  $O(\lg \lg u)$  worst-case time.

### Inserting an element

Now we examine how to insert an element into a vEB tree. Recall that `PROTO-VEB-INSERT` made two recursive calls: one to insert the element and one to insert the element's cluster number into the summary. The `VEB-TREE-INSERT` procedure will make only one recursive call. How can we get away with just one? When we insert an element, either the cluster that it goes into already has another element or it does not. If the cluster already has another element, then the cluster number is already in the summary, and so we do not need to make that recursive call. If

the cluster does not already have another element, then the element being inserted becomes the only element in the cluster, and we do not need to recurse to insert an element into an empty vEB tree:

VEB-EMPTY-TREE-INSERT( $V, x$ )

```
1   $V.min = x$ 
2   $V.max = x$ 
```

With this procedure in hand, here is the pseudocode for VEB-TREE-INSERT( $V, x$ ), which assumes that  $x$  is not already an element in the set represented by vEB tree  $V$ :

VEB-TREE-INSERT( $V, x$ )

```
1  if  $V.min == \text{NIL}$ 
2      VEB-EMPTY-TREE-INSERT( $V, x$ )
3  else if  $x < V.min$ 
4      exchange  $x$  with  $V.min$ 
5      if  $V.u > 2$ 
6          if VEB-TREE-MINIMUM( $V.cluster[\text{high}(x)]$ ) == NIL
7              VEB-TREE-INSERT( $V.summary, \text{high}(x)$ )
8              VEB-EMPTY-TREE-INSERT( $V.cluster[\text{high}(x)], \text{low}(x)$ )
9          else VEB-TREE-INSERT( $V.cluster[\text{high}(x)], \text{low}(x)$ )
10     if  $x > V.max$ 
11          $V.max = x$ 
```

This procedure works as follows. Line 1 tests whether  $V$  is an empty vEB tree and, if it is, then line 2 handles this easy case. Lines 3–11 assume that  $V$  is not empty, and therefore some element will be inserted into one of  $V$ 's clusters. But that element might not necessarily be the element  $x$  passed to VEB-TREE-INSERT. If  $x < min$ , as tested in line 3, then  $x$  needs to become the new  $min$ . We don't want to lose the original  $min$ , however, and so we need to insert it into one of  $V$ 's clusters. In this case, line 4 exchanges  $x$  with  $min$ , so that we insert the original  $min$  into one of  $V$ 's clusters.

We execute lines 6–9 only if  $V$  is not a base-case vEB tree. Line 6 determines whether the cluster that  $x$  will go into is currently empty. If so, then line 7 inserts  $x$ 's cluster number into the summary and line 8 handles the easy case of inserting  $x$  into an empty cluster. If  $x$ 's cluster is not currently empty, then line 9 inserts  $x$  into its cluster. In this case, we do not need to update the summary, since  $x$ 's cluster number is already a member of the summary.

Finally, lines 10–11 take care of updating  $max$  if  $x > max$ . Note that if  $V$  is a base-case vEB tree that is not empty, then lines 3–4 and 10–11 update  $min$  and  $max$  properly.

Once again, we can easily see how recurrence (20.4) characterizes the running time. Depending on the result of the test in line 6, either the recursive call in line 7 (run on a vEB tree with universe size  $\sqrt[4]{u}$ ) or the recursive call in line 9 (run on a vEB with universe size  $\sqrt[4]{u}$ ) executes. In either case, the one recursive call is on a vEB tree with universe size at most  $\sqrt[4]{u}$ . Because the remainder of VEB-TREE-INSERT takes  $O(1)$  time, recurrence (20.4) applies, and so the running time is  $O(\lg \lg u)$ .

### Deleting an element

Finally, we look at how to delete an element from a vEB tree. The procedure VEB-TREE-DELETE( $V, x$ ) assumes that  $x$  is currently an element in the set represented by the vEB tree  $V$ .

```

VEB-TREE-DELETE( $V, x$ )
1  if  $V.min == V.max$ 
2       $V.min = \text{NIL}$ 
3       $V.max = \text{NIL}$ 
4  elseif  $V.u == 2$ 
5      if  $x == 0$ 
6           $V.min = 1$ 
7      else  $V.min = 0$ 
8           $V.max = V.min$ 
9  else if  $x == V.min$ 
10      $first\_cluster = \text{VEB-TREE-MINIMUM}(V.summary)$ 
11      $x = \text{index}(first\_cluster,$ 
12          $\text{VEB-TREE-MINIMUM}(V.cluster[first\_cluster]))$ 
13      $V.min = x$ 
14      $\text{VEB-TREE-DELETE}(V.cluster[\text{high}(x)], \text{low}(x))$ 
15     if  $\text{VEB-TREE-MINIMUM}(V.cluster[\text{high}(x)]) == \text{NIL}$ 
16          $\text{VEB-TREE-DELETE}(V.summary, \text{high}(x))$ 
17     if  $x == V.max$ 
18          $summary\_max = \text{VEB-TREE-MAXIMUM}(V.summary)$ 
19         if  $summary\_max == \text{NIL}$ 
20              $V.max = V.min$ 
21         else  $V.max = \text{index}(summary\_max,$ 
22              $\text{VEB-TREE-MAXIMUM}(V.cluster[summary\_max]))$ 
23 elseif  $x == V.max$ 
24      $V.max = \text{index}(\text{high}(x),$ 
25          $\text{VEB-TREE-MAXIMUM}(V.cluster[\text{high}(x)]))$ 

```

The `VEB-TREE-DELETE` procedure works as follows. If the `vEB` tree  $V$  contains only one element, then it's just as easy to delete it as it was to insert an element into an empty `vEB` tree: just set  $\text{min}$  and  $\text{max}$  to `NIL`. Lines 1–3 handle this case. Otherwise,  $V$  has at least two elements. Line 4 tests whether  $V$  is a base-case `vEB` tree and, if so, lines 5–8 set  $\text{min}$  and  $\text{max}$  to the one remaining element.

Lines 9–22 assume that  $V$  has two or more elements and that  $u \geq 4$ . In this case, we will have to delete an element from a cluster. The element we delete from a cluster might not be  $x$ , however, because if  $x$  equals  $\text{min}$ , then once we have deleted  $x$ , some other element within one of  $V$ 's clusters becomes the new  $\text{min}$ , and we have to delete that other element from its cluster. If the test in line 9 reveals that we are in this case, then line 10 sets *first-cluster* to the number of the cluster that contains the lowest element other than  $\text{min}$ , and line 11 sets  $x$  to the value of the lowest element in that cluster. This element becomes the new  $\text{min}$  in line 12 and, because we set  $x$  to its value, it is the element that will be deleted from its cluster.

When we reach line 13, we know that we need to delete element  $x$  from its cluster, whether  $x$  was the value originally passed to `VEB-TREE-DELETE` or  $x$  is the element becoming the new minimum. Line 13 deletes  $x$  from its cluster. That cluster might now become empty, which line 14 tests, and if it does, then we need to remove  $x$ 's cluster number from the summary, which line 15 handles. After updating the summary, we might need to update  $\text{max}$ . Line 16 checks to see whether we are deleting the maximum element in  $V$  and, if we are, then line 17 sets *summary-max* to the number of the highest-numbered nonempty cluster. (The call `VEB-TREE-MAXIMUM( $V.\text{summary}$ )` works because we have already recursively called `VEB-TREE-DELETE` on  $V.\text{summary}$ , and therefore  $V.\text{summary.max}$  has already been updated as necessary.) If all of  $V$ 's clusters are empty, then the only remaining element in  $V$  is  $\text{min}$ ; line 18 checks for this case, and line 19 updates  $\text{max}$  appropriately. Otherwise, line 20 sets  $\text{max}$  to the maximum element in the highest-numbered cluster. (If this cluster is where the element has been deleted, we again rely on the recursive call in line 13 having already corrected that cluster's  $\text{max}$  attribute.)

Finally, we have to handle the case in which  $x$ 's cluster did not become empty due to  $x$  being deleted. Although we do not have to update the summary in this case, we might have to update  $\text{max}$ . Line 21 tests for this case, and if we have to update  $\text{max}$ , line 22 does so (again relying on the recursive call to have corrected  $\text{max}$  in the cluster).

Now we show that `VEB-TREE-DELETE` runs in  $O(\lg \lg u)$  time in the worst case. At first glance, you might think that recurrence (20.4) does not always apply, because a single call of `VEB-TREE-DELETE` can make two recursive calls: one on line 13 and one on line 15. Although the procedure can make both recursive calls, let's think about what happens when it does. In order for the recursive call on



line 15 to occur, the test on line 14 must show that  $x$ 's cluster is empty. The only way that  $x$ 's cluster can be empty is if  $x$  was the only element in its cluster when we made the recursive call on line 13. But if  $x$  was the only element in its cluster, then that recursive call took  $O(1)$  time, because it executed only lines 1–3. Thus, we have two mutually exclusive possibilities:

- The recursive call on line 13 took constant time.
- The recursive call on line 15 did not occur.

In either case, recurrence (20.4) characterizes the running time of VEB-TREE-DELETE, and hence its worst-case running time is  $O(\lg \lg u)$ .

## Exercises

### 20.3-1

Modify vEB trees to support duplicate keys.

### 20.3-2

Modify vEB trees to support keys that have associated satellite data.

### 20.3-3

Write pseudocode for a procedure that creates an empty van Emde Boas tree.

### 20.3-4

What happens if you call VEB-TREE-INSERT with an element that is already in the vEB tree? What happens if you call VEB-TREE-DELETE with an element that is not in the vEB tree? Explain why the procedures exhibit the behavior that they do. Show how to modify vEB trees and their operations so that we can check in constant time whether an element is present.

### 20.3-5

Suppose that instead of  $\sqrt[k]{u}$  clusters, each with universe size  $\sqrt[k]{u}$ , we constructed vEB trees to have  $u^{1/k}$  clusters, each with universe size  $u^{1-1/k}$ , where  $k > 1$  is a constant. If we were to modify the operations appropriately, what would be their running times? For the purpose of analysis, assume that  $u^{1/k}$  and  $u^{1-1/k}$  are always integers.

### 20.3-6

Creating a vEB tree with universe size  $u$  requires  $O(u)$  time. Suppose we wish to explicitly account for that time. What is the smallest number of operations  $n$  for which the amortized time of each operation in a vEB tree is  $O(\lg \lg u)$ ?

---

## Problems

### 20-1 Space requirements for van Emde Boas trees

This problem explores the space requirements for van Emde Boas trees and suggests a way to modify the data structure to make its space requirement depend on the number  $n$  of elements actually stored in the tree, rather than on the universe size  $u$ . For simplicity, assume that  $\sqrt{u}$  is always an integer.

- a. Explain why the following recurrence characterizes the space requirement  $P(u)$  of a van Emde Boas tree with universe size  $u$ :

$$P(u) = (\sqrt{u} + 1)P(\sqrt{u}) + \Theta(\sqrt{u}) . \quad (20.5)$$

- b. Prove that recurrence (20.5) has the solution  $P(u) = O(u)$ .

In order to reduce the space requirements, let us define a **reduced-space van Emde Boas tree**, or **RS-vEB tree**, as a vEB tree  $V$  but with the following changes:

- The attribute  $V.cluster$ , rather than being stored as a simple array of pointers to vEB trees with universe size  $\sqrt{u}$ , is a hash table (see Chapter 11) stored as a dynamic table (see Section 17.4). Corresponding to the array version of  $V.cluster$ , the hash table stores pointers to RS-vEB trees with universe size  $\sqrt{u}$ . To find the  $i$ th cluster, we look up the key  $i$  in the hash table, so that we can find the  $i$ th cluster by a single search in the hash table.
- The hash table stores only pointers to nonempty clusters. A search in the hash table for an empty cluster returns NIL, indicating that the cluster is empty.
- The attribute  $V.summary$  is NIL if all clusters are empty. Otherwise,  $V.summary$  points to an RS-vEB tree with universe size  $\sqrt{u}$ .

Because the hash table is implemented with a dynamic table, the space it requires is proportional to the number of nonempty clusters.

When we need to insert an element into an empty RS-vEB tree, we create the RS-vEB tree by calling the following procedure, where the parameter  $u$  is the universe size of the RS-vEB tree:

CREATE-NEW-RS-vEB-TREE( $u$ )

```

1  allocate a new vEB tree  $V$ 
2   $V.u = u$ 
3   $V.min = \text{NIL}$ 
4   $V.max = \text{NIL}$ 
5   $V.summary = \text{NIL}$ 
6  create  $V.cluster$  as an empty dynamic hash table
7  return  $V$ 
```

- c. Modify the `VEB-TREE-INSERT` procedure to produce pseudocode for the procedure `RS-VEB-TREE-INSERT( $V, x$ )`, which inserts  $x$  into the RS-`VEB` tree  $V$ , calling `CREATE-NEW-RS-VEB-TREE` as appropriate.
- d. Modify the `VEB-TREE-SUCCESSOR` procedure to produce pseudocode for the procedure `RS-VEB-TREE-SUCCESSOR( $V, x$ )`, which returns the successor of  $x$  in RS-`VEB` tree  $V$ , or `NIL` if  $x$  has no successor in  $V$ .
- e. Prove that, under the assumption of simple uniform hashing, your RS-`VEB-TREE-INSERT` and RS-`VEB-TREE-SUCCESSOR` procedures run in  $O(\lg \lg u)$  expected time.
- f. Assuming that elements are never deleted from a `VEB` tree, prove that the space requirement for the RS-`VEB` tree structure is  $O(n)$ , where  $n$  is the number of elements actually stored in the RS-`VEB` tree.
- g. RS-`VEB` trees have another advantage over `VEB` trees: they require less time to create. How long does it take to create an empty RS-`VEB` tree?

### 20-2 *y-fast tries*

This problem investigates D. Willard's "y-fast tries" which, like van Emde Boas trees, perform each of the operations `MEMBER`, `MINIMUM`, `MAXIMUM`, `PREDECESSOR`, and `SUCCESSOR` on elements drawn from a universe with size  $u$  in  $O(\lg \lg u)$  worst-case time. The `INSERT` and `DELETE` operations take  $O(\lg \lg u)$  amortized time. Like reduced-space van Emde Boas trees (see Problem 20-1), y-fast tries use only  $O(n)$  space to store  $n$  elements. The design of y-fast tries relies on perfect hashing (see Section 11.5).

As a preliminary structure, suppose that we create a perfect hash table containing not only every element in the dynamic set, but every prefix of the binary representation of every element in the set. For example, if  $u = 16$ , so that  $\lg u = 4$ , and  $x = 13$  is in the set, then because the binary representation of 13 is 1101, the perfect hash table would contain the strings 1, 11, 110, and 1101. In addition to the hash table, we create a doubly linked list of the elements currently in the set, in increasing order.

- a. How much space does this structure require?
- b. Show how to perform the `MINIMUM` and `MAXIMUM` operations in  $O(1)$  time; the `MEMBER`, `PREDECESSOR`, and `SUCCESSOR` operations in  $O(\lg \lg u)$  time; and the `INSERT` and `DELETE` operations in  $O(\lg u)$  time.

To reduce the space requirement to  $O(n)$ , we make the following changes to the data structure:

- We cluster the  $n$  elements into  $n / \lg u$  groups of size  $\lg u$ . (Assume for now that  $\lg u$  divides  $n$ .) The first group consists of the  $\lg u$  smallest elements in the set, the second group consists of the next  $\lg u$  smallest elements, and so on.
- We designate a “representative” value for each group. The representative of the  $i$ th group is at least as large as the largest element in the  $i$ th group, and it is smaller than every element of the  $(i + 1)$ st group. (The representative of the last group can be the maximum possible element  $u - 1$ .) Note that a representative might be a value not currently in the set.
- We store the  $\lg u$  elements of each group in a balanced binary search tree, such as a red-black tree. Each representative points to the balanced binary search tree for its group, and each balanced binary search tree points to its group’s representative.
- The perfect hash table stores only the representatives, which are also stored in a doubly linked list in increasing order.

We call this structure a *y-fast trie*.

- c.* Show that a *y-fast trie* requires only  $O(n)$  space to store  $n$  elements.
- d.* Show how to perform the MINIMUM and MAXIMUM operations in  $O(\lg \lg u)$  time with a *y-fast trie*.
- e.* Show how to perform the MEMBER operation in  $O(\lg \lg u)$  time.
- f.* Show how to perform the PREDECESSOR and SUCCESSOR operations in  $O(\lg \lg u)$  time.
- g.* Explain why the INSERT and DELETE operations take  $\Omega(\lg \lg u)$  time.
- h.* Show how to relax the requirement that each group in a *y-fast trie* has exactly  $\lg u$  elements to allow INSERT and DELETE to run in  $O(\lg \lg u)$  amortized time without affecting the asymptotic running times of the other operations.

---

## Chapter notes

The data structure in this chapter is named after P. van Emde Boas, who described an early form of the idea in 1975 [339]. Later papers by van Emde Boas [340] and van Emde Boas, Kaas, and Zijlstra [341] refined the idea and the exposition. Mehlhorn and Näher [252] subsequently extended the ideas to apply to universe

sizes that are prime. Mehlhorn's book [249] contains a slightly different treatment of van Emde Boas trees than the one in this chapter.

Using the ideas behind van Emde Boas trees, Dementiev et al. [83] developed a nonrecursive, three-level search tree that ran faster than van Emde Boas trees in their own experiments.

Wang and Lin [347] designed a hardware-pipelined version of van Emde Boas trees, which achieves constant amortized time per operation and uses  $O(\lg \lg u)$  stages in the pipeline.

A lower bound by Pătraşcu and Thorup [273, 274] for finding the predecessor shows that van Emde Boas trees are optimal for this operation, even if randomization is allowed.

---

## 21 Data Structures for Disjoint Sets

Some applications involve grouping  $n$  distinct elements into a collection of disjoint sets. These applications often need to perform two operations in particular: finding the unique set that contains a given element and uniting two sets. This chapter explores methods for maintaining a data structure that supports these operations.

Section 21.1 describes the operations supported by a disjoint-set data structure and presents a simple application. In Section 21.2, we look at a simple linked-list implementation for disjoint sets. Section 21.3 presents a more efficient representation using rooted trees. The running time using the tree representation is theoretically superlinear, but for all practical purposes it is linear. Section 21.4 defines and discusses a very quickly growing function and its very slowly growing inverse, which appears in the running time of operations on the tree-based implementation, and then, by a complex amortized analysis, proves an upper bound on the running time that is just barely superlinear.

---

### 21.1 Disjoint-set operations

A *disjoint-set data structure* maintains a collection  $\mathcal{S} = \{S_1, S_2, \dots, S_k\}$  of disjoint dynamic sets. We identify each set by a *representative*, which is some member of the set. In some applications, it doesn't matter which member is used as the representative; we care only that if we ask for the representative of a dynamic set twice without modifying the set between the requests, we get the same answer both times. Other applications may require a prespecified rule for choosing the representative, such as choosing the smallest member in the set (assuming, of course, that the elements can be ordered).

As in the other dynamic-set implementations we have studied, we represent each element of a set by an object. Letting  $x$  denote an object, we wish to support the following operations:

MAKE-SET( $x$ ) creates a new set whose only member (and thus representative) is  $x$ . Since the sets are disjoint, we require that  $x$  not already be in some other set.

UNION( $x, y$ ) unites the dynamic sets that contain  $x$  and  $y$ , say  $S_x$  and  $S_y$ , into a new set that is the union of these two sets. We assume that the two sets are disjoint prior to the operation. The representative of the resulting set is any member of  $S_x \cup S_y$ , although many implementations of UNION specifically choose the representative of either  $S_x$  or  $S_y$  as the new representative. Since we require the sets in the collection to be disjoint, conceptually we destroy sets  $S_x$  and  $S_y$ , removing them from the collection  $\mathcal{S}$ . In practice, we often absorb the elements of one of the sets into the other set.

FIND-SET( $x$ ) returns a pointer to the representative of the (unique) set containing  $x$ .

Throughout this chapter, we shall analyze the running times of disjoint-set data structures in terms of two parameters:  $n$ , the number of MAKE-SET operations, and  $m$ , the total number of MAKE-SET, UNION, and FIND-SET operations. Since the sets are disjoint, each UNION operation reduces the number of sets by one. After  $n - 1$  UNION operations, therefore, only one set remains. The number of UNION operations is thus at most  $n - 1$ . Note also that since the MAKE-SET operations are included in the total number of operations  $m$ , we have  $m \geq n$ . We assume that the  $n$  MAKE-SET operations are the first  $n$  operations performed.

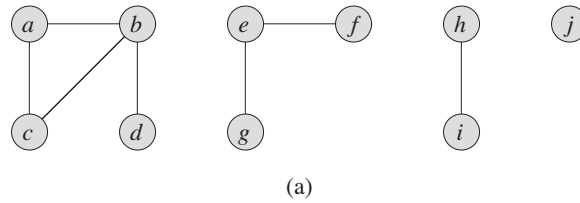
### An application of disjoint-set data structures

One of the many applications of disjoint-set data structures arises in determining the connected components of an undirected graph (see Section B.4). Figure 21.1(a), for example, shows a graph with four connected components.

The procedure CONNECTED-COMPONENTS that follows uses the disjoint-set operations to compute the connected components of a graph. Once CONNECTED-COMPONENTS has preprocessed the graph, the procedure SAME-COMPONENT answers queries about whether two vertices are in the same connected component.<sup>1</sup> (In pseudocode, we denote the set of vertices of a graph  $G$  by  $G.V$  and the set of edges by  $G.E$ .)

---

<sup>1</sup>When the edges of the graph are static—not changing over time—we can compute the connected components faster by using depth-first search (Exercise 22.3-12). Sometimes, however, the edges are added dynamically and we need to maintain the connected components as each edge is added. In this case, the implementation given here can be more efficient than running a new depth-first search for each new edge.



Edge processed	Collection of disjoint sets									
initial sets	{a}	{b}	{c}	{d}	{e}	{f}	{g}	{h}	{i}	{j}
(b,d)	{a}	{b,d}	{c}		{e}	{f}	{g}	{h}	{i}	{j}
(e,g)	{a}	{b,d}	{c}		{e,g}	{f}		{h}	{i}	{j}
(a,c)	{a,c}	{b,d}			{e,g}	{f}		{h}	{i}	{j}
(h,i)	{a,c}	{b,d}			{e,g}	{f}		{h,i}		{j}
(a,b)	{a,b,c,d}				{e,g}	{f}		{h,i}		{j}
(e,f)	{a,b,c,d}				{e,f,g}			{h,i}		{j}
(b,c)	{a,b,c,d}				{e,f,g}			{h,i}		{j}

(b)

**Figure 21.1** (a) A graph with four connected components:  $\{a, b, c, d\}$ ,  $\{e, f, g\}$ ,  $\{h, i\}$ , and  $\{j\}$ . (b) The collection of disjoint sets after processing each edge.

#### CONNECTED-COMPONENTS( $G$ )

```

1  for each vertex  $v \in G.V$ 
2      MAKE-SET( $v$ )
3  for each edge  $(u, v) \in G.E$ 
4      if FIND-SET( $u$ )  $\neq$  FIND-SET( $v$ )
5          UNION( $u, v$ )

```

#### SAME-COMPONENT( $u, v$ )

```

1  if FIND-SET( $u$ ) == FIND-SET( $v$ )
2      return TRUE
3  else return FALSE

```

The procedure **CONNECTED-COMPONENTS** initially places each vertex  $v$  in its own set. Then, for each edge  $(u, v)$ , it unites the sets containing  $u$  and  $v$ . By Exercise 21.1-2, after processing all the edges, two vertices are in the same connected component if and only if the corresponding objects are in the same set. Thus, **CONNECTED-COMPONENTS** computes sets in such a way that the procedure **SAME-COMPONENT** can determine whether two vertices are in the same con-



nected component. Figure 21.1(b) illustrates how CONNECTED-COMPONENTS computes the disjoint sets.

In an actual implementation of this connected-components algorithm, the representations of the graph and the disjoint-set data structure would need to reference each other. That is, an object representing a vertex would contain a pointer to the corresponding disjoint-set object, and vice versa. These programming details depend on the implementation language, and we do not address them further here.

## Exercises

### 21.1-1

Suppose that CONNECTED-COMPONENTS is run on the undirected graph  $G = (V, E)$ , where  $V = \{a, b, c, d, e, f, g, h, i, j, k\}$  and the edges of  $E$  are processed in the order  $(d, i), (f, k), (g, i), (b, g), (a, h), (i, j), (d, k), (b, j), (d, f), (g, j), (a, e)$ . List the vertices in each connected component after each iteration of lines 3–5.

### 21.1-2

Show that after all edges are processed by CONNECTED-COMPONENTS, two vertices are in the same connected component if and only if they are in the same set.

### 21.1-3

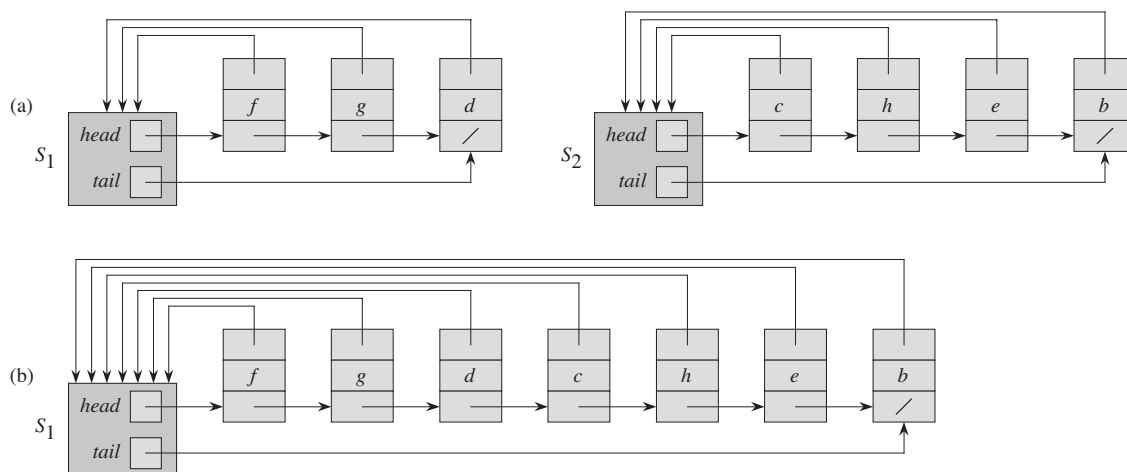
During the execution of CONNECTED-COMPONENTS on an undirected graph  $G = (V, E)$  with  $k$  connected components, how many times is FIND-SET called? How many times is UNION called? Express your answers in terms of  $|V|$ ,  $|E|$ , and  $k$ .

---

## 21.2 Linked-list representation of disjoint sets

Figure 21.2(a) shows a simple way to implement a disjoint-set data structure: each set is represented by its own linked list. The object for each set has attributes *head*, pointing to the first object in the list, and *tail*, pointing to the last object. Each object in the list contains a set member, a pointer to the next object in the list, and a pointer back to the set object. Within each linked list, the objects may appear in any order. The representative is the set member in the first object in the list.

With this linked-list representation, both MAKE-SET and FIND-SET are easy, requiring  $O(1)$  time. To carry out MAKE-SET( $x$ ), we create a new linked list whose only object is  $x$ . For FIND-SET( $x$ ), we just follow the pointer from  $x$  back to its set object and then return the member in the object that *head* points to. For example, in Figure 21.2(a), the call FIND-SET( $g$ ) would return  $f$ .



**Figure 21.2** (a) Linked-list representations of two sets. Set  $S_1$  contains members  $d$ ,  $f$ , and  $g$ , with representative  $f$ , and set  $S_2$  contains members  $b$ ,  $c$ ,  $e$ , and  $h$ , with representative  $c$ . Each object in the list contains a set member, a pointer to the next object in the list, and a pointer back to the set object. Each set object has pointers  $head$  and  $tail$  to the first and last objects, respectively. (b) The result of  $\text{UNION}(g, e)$ , which appends the linked list containing  $e$  to the linked list containing  $g$ . The representative of the resulting set is  $f$ . The set object for  $e$ 's list,  $S_2$ , is destroyed.

### A simple implementation of union

The simplest implementation of the  $\text{UNION}$  operation using the linked-list set representation takes significantly more time than  $\text{MAKE-SET}$  or  $\text{FIND-SET}$ . As Figure 21.2(b) shows, we perform  $\text{UNION}(x, y)$  by appending  $y$ 's list onto the end of  $x$ 's list. The representative of  $x$ 's list becomes the representative of the resulting set. We use the  $tail$  pointer for  $x$ 's list to quickly find where to append  $y$ 's list. Because all members of  $y$ 's list join  $x$ 's list, we can destroy the set object for  $y$ 's list. Unfortunately, we must update the pointer to the set object for each object originally on  $y$ 's list, which takes time linear in the length of  $y$ 's list. In Figure 21.2, for example, the operation  $\text{UNION}(g, e)$  causes pointers to be updated in the objects for  $b$ ,  $c$ ,  $e$ , and  $h$ .

In fact, we can easily construct a sequence of  $m$  operations on  $n$  objects that requires  $\Theta(n^2)$  time. Suppose that we have objects  $x_1, x_2, \dots, x_n$ . We execute the sequence of  $n$   $\text{MAKE-SET}$  operations followed by  $n - 1$   $\text{UNION}$  operations shown in Figure 21.3, so that  $m = 2n - 1$ . We spend  $\Theta(n)$  time performing the  $n$   $\text{MAKE-SET}$  operations. Because the  $i$ th  $\text{UNION}$  operation updates  $i$  objects, the total number of objects updated by all  $n - 1$   $\text{UNION}$  operations is

Operation	Number of objects updated
MAKE-SET( $x_1$ )	1
MAKE-SET( $x_2$ )	1
$\vdots$	$\vdots$
MAKE-SET( $x_n$ )	1
UNION( $x_2, x_1$ )	1
UNION( $x_3, x_2$ )	2
UNION( $x_4, x_3$ )	3
$\vdots$	$\vdots$
UNION( $x_n, x_{n-1}$ )	$n - 1$

**Figure 21.3** A sequence of  $2n - 1$  operations on  $n$  objects that takes  $\Theta(n^2)$  time, or  $\Theta(n)$  time per operation on average, using the linked-list set representation and the simple implementation of UNION.

$$\sum_{i=1}^{n-1} i = \Theta(n^2) .$$

The total number of operations is  $2n - 1$ , and so each operation on average requires  $\Theta(n)$  time. That is, the amortized time of an operation is  $\Theta(n)$ .

### A weighted-union heuristic

In the worst case, the above implementation of the UNION procedure requires an average of  $\Theta(n)$  time per call because we may be appending a longer list onto a shorter list; we must update the pointer to the set object for each member of the longer list. Suppose instead that each list also includes the length of the list (which we can easily maintain) and that we always append the shorter list onto the longer, breaking ties arbitrarily. With this simple **weighted-union heuristic**, a single UNION operation can still take  $\Omega(n)$  time if both sets have  $\Omega(n)$  members. As the following theorem shows, however, a sequence of  $m$  MAKE-SET, UNION, and FIND-SET operations,  $n$  of which are MAKE-SET operations, takes  $O(m + n \lg n)$  time.

#### **Theorem 21.1**

Using the linked-list representation of disjoint sets and the weighted-union heuristic, a sequence of  $m$  MAKE-SET, UNION, and FIND-SET operations,  $n$  of which are MAKE-SET operations, takes  $O(m + n \lg n)$  time.

**Proof** Because each UNION operation unites two disjoint sets, we perform at most  $n - 1$  UNION operations over all. We now bound the total time taken by these UNION operations. We start by determining, for each object, an upper bound on the number of times the object's pointer back to its set object is updated. Consider a particular object  $x$ . We know that each time  $x$ 's pointer was updated,  $x$  must have started in the smaller set. The first time  $x$ 's pointer was updated, therefore, the resulting set must have had at least 2 members. Similarly, the next time  $x$ 's pointer was updated, the resulting set must have had at least 4 members. Continuing on, we observe that for any  $k \leq n$ , after  $x$ 's pointer has been updated  $\lceil \lg k \rceil$  times, the resulting set must have at least  $k$  members. Since the largest set has at most  $n$  members, each object's pointer is updated at most  $\lceil \lg n \rceil$  times over all the UNION operations. Thus the total time spent updating object pointers over all UNION operations is  $O(n \lg n)$ . We must also account for updating the *tail* pointers and the list lengths, which take only  $\Theta(1)$  time per UNION operation. The total time spent in all UNION operations is thus  $O(n \lg n)$ .

The time for the entire sequence of  $m$  operations follows easily. Each MAKE-SET and FIND-SET operation takes  $O(1)$  time, and there are  $O(m)$  of them. The total time for the entire sequence is thus  $O(m + n \lg n)$ . ■

## Exercises

### 21.2-1

Write pseudocode for MAKE-SET, FIND-SET, and UNION using the linked-list representation and the weighted-union heuristic. Make sure to specify the attributes that you assume for set objects and list objects.

### 21.2-2

Show the data structure that results and the answers returned by the FIND-SET operations in the following program. Use the linked-list representation with the weighted-union heuristic.

```

1  for  $i = 1$  to 16
2      MAKE-SET( $x_i$ )
3  for  $i = 1$  to 15 by 2
4      UNION( $x_i, x_{i+1}$ )
5  for  $i = 1$  to 13 by 4
6      UNION( $x_i, x_{i+2}$ )
7  UNION( $x_1, x_5$ )
8  UNION( $x_{11}, x_{13}$ )
9  UNION( $x_1, x_{10}$ )
10 FIND-SET( $x_2$ )
11 FIND-SET( $x_9$ )

```

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.)

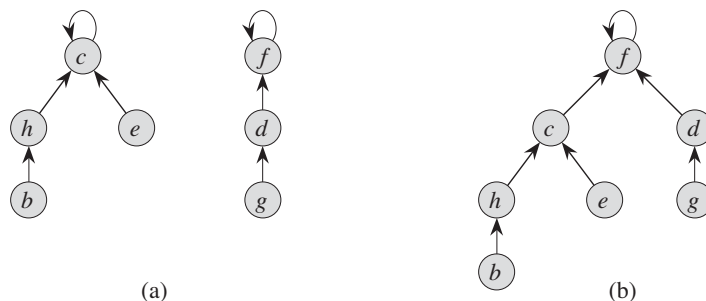
### 21.2-6

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  $\text{UNION}(e, g)$ .

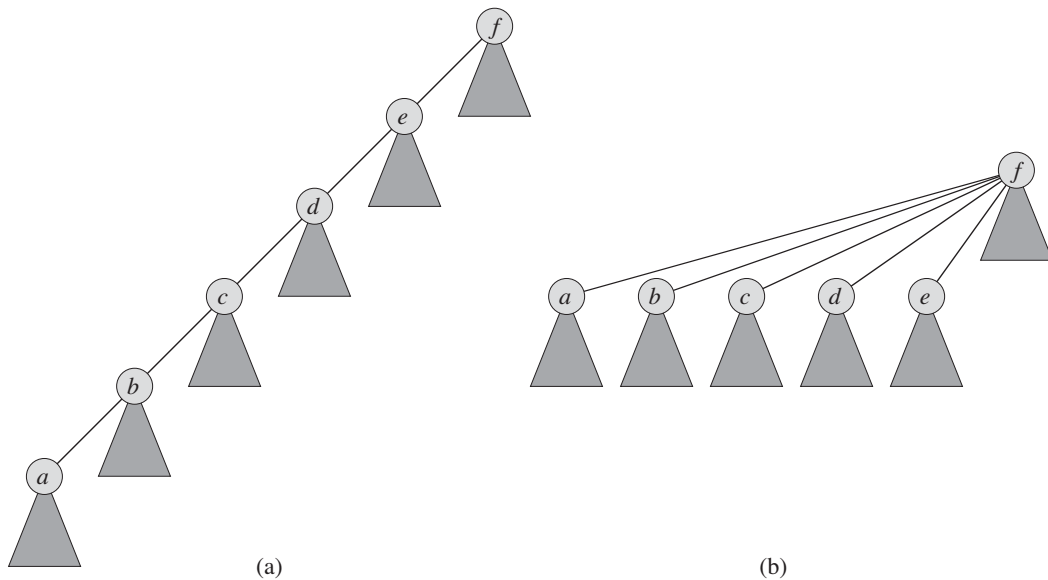
We perform the three disjoint-set operations as follows. A **MAKE-SET** operation simply creates a tree with just one node. We perform a **FIND-SET** operation by following parent pointers until we find the root of the tree. The nodes visited on this simple path toward the root constitute the *find path*. 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,

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

■

### Corollary 21.9

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

■

## 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  $iter(x)$  to increase (to at least  $i + 1$ ) or  $level(x)$  to increase (which occurs if  $iter(x)$  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   $u.\text{color} = \text{BLACK}$ 
8  for each node  $v$  such that  $\{u, v\} \in P$ 
9      if  $v.\text{color} == \text{BLACK}$ 
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.