

---

## *II   Sorting and Order Statistics*

---

## Introduction

This part presents several algorithms that solve the following *sorting problem*:

**Input:** A sequence of  $n$  numbers  $\langle a_1, a_2, \dots, a_n \rangle$ .

**Output:** A permutation (reordering)  $\langle a'_1, a'_2, \dots, a'_n \rangle$  of the input sequence such that  $a'_1 \leq a'_2 \leq \dots \leq a'_n$ .

The input sequence is usually an  $n$ -element array, although it may be represented in some other fashion, such as a linked list.

### The structure of the data

In practice, the numbers to be sorted are rarely isolated values. Each is usually part of a collection of data called a *record*. Each record contains a *key*, which is the value to be sorted. The remainder of the record consists of *satellite data*, which are usually carried around with the key. In practice, when a sorting algorithm permutes the keys, it must permute the satellite data as well. If each record includes a large amount of satellite data, we often permute an array of pointers to the records rather than the records themselves in order to minimize data movement.

In a sense, it is these implementation details that distinguish an algorithm from a full-blown program. A sorting algorithm describes the *method* by which we determine the sorted order, regardless of whether we are sorting individual numbers or large records containing many bytes of satellite data. Thus, when focusing on the problem of sorting, we typically assume that the input consists only of numbers. Translating an algorithm for sorting numbers into a program for sorting records

is conceptually straightforward, although in a given engineering situation other subtleties may make the actual programming task a challenge.

### **Why sorting?**

Many computer scientists consider sorting to be the most fundamental problem in the study of algorithms. There are several reasons:

- Sometimes an application inherently needs to sort information. For example, in order to prepare customer statements, banks need to sort checks by check number.
- Algorithms often use sorting as a key subroutine. For example, a program that renders graphical objects which are layered on top of each other might have to sort the objects according to an “above” relation so that it can draw these objects from bottom to top. We shall see numerous algorithms in this text that use sorting as a subroutine.
- We can draw from among a wide variety of sorting algorithms, and they employ a rich set of techniques. In fact, many important techniques used throughout algorithm design appear in the body of sorting algorithms that have been developed over the years. In this way, sorting is also a problem of historical interest.
- We can prove a nontrivial lower bound for sorting (as we shall do in Chapter 8). Our best upper bounds match the lower bound asymptotically, and so we know that our sorting algorithms are asymptotically optimal. Moreover, we can use the lower bound for sorting to prove lower bounds for certain other problems.
- Many engineering issues come to the fore when implementing sorting algorithms. The fastest sorting program for a particular situation may depend on many factors, such as prior knowledge about the keys and satellite data, the memory hierarchy (caches and virtual memory) of the host computer, and the software environment. Many of these issues are best dealt with at the algorithmic level, rather than by “tweaking” the code.

### **Sorting algorithms**

We introduced two algorithms that sort  $n$  real numbers in Chapter 2. Insertion sort takes  $\Theta(n^2)$  time in the worst case. Because its inner loops are tight, however, it is a fast in-place sorting algorithm for small input sizes. (Recall that a sorting algorithm sorts *in place* if only a constant number of elements of the input array are ever stored outside the array.) Merge sort has a better asymptotic running time,  $\Theta(n \lg n)$ , but the MERGE procedure it uses does not operate in place.

In this part, we shall introduce two more algorithms that sort arbitrary real numbers. Heapsort, presented in Chapter 6, sorts  $n$  numbers in place in  $O(n \lg n)$  time. It uses an important data structure, called a heap, with which we can also implement a priority queue.

Quicksort, in Chapter 7, also sorts  $n$  numbers in place, but its worst-case running time is  $\Theta(n^2)$ . Its expected running time is  $\Theta(n \lg n)$ , however, and it generally outperforms heapsort in practice. Like insertion sort, quicksort has tight code, and so the hidden constant factor in its running time is small. It is a popular algorithm for sorting large input arrays.

Insertion sort, merge sort, heapsort, and quicksort are all comparison sorts: they determine the sorted order of an input array by comparing elements. Chapter 8 begins by introducing the decision-tree model in order to study the performance limitations of comparison sorts. Using this model, we prove a lower bound of  $\Omega(n \lg n)$  on the worst-case running time of any comparison sort on  $n$  inputs, thus showing that heapsort and merge sort are asymptotically optimal comparison sorts.

Chapter 8 then goes on to show that we can beat this lower bound of  $\Omega(n \lg n)$  if we can gather information about the sorted order of the input by means other than comparing elements. The counting sort algorithm, for example, assumes that the input numbers are in the set  $\{0, 1, \dots, k\}$ . By using array indexing as a tool for determining relative order, counting sort can sort  $n$  numbers in  $\Theta(k + n)$  time. Thus, when  $k = O(n)$ , counting sort runs in time that is linear in the size of the input array. A related algorithm, radix sort, can be used to extend the range of counting sort. If there are  $n$  integers to sort, each integer has  $d$  digits, and each digit can take on up to  $k$  possible values, then radix sort can sort the numbers in  $\Theta(d(n + k))$  time. When  $d$  is a constant and  $k$  is  $O(n)$ , radix sort runs in linear time. A third algorithm, bucket sort, requires knowledge of the probabilistic distribution of numbers in the input array. It can sort  $n$  real numbers uniformly distributed in the half-open interval  $[0, 1)$  in average-case  $O(n)$  time.

The following table summarizes the running times of the sorting algorithms from Chapters 2 and 6–8. As usual,  $n$  denotes the number of items to sort. For counting sort, the items to sort are integers in the set  $\{0, 1, \dots, k\}$ . For radix sort, each item is a  $d$ -digit number, where each digit takes on  $k$  possible values. For bucket sort, we assume that the keys are real numbers uniformly distributed in the half-open interval  $[0, 1)$ . The rightmost column gives the average-case or expected running time, indicating which it gives when it differs from the worst-case running time. We omit the average-case running time of heapsort because we do not analyze it in this book.

Algorithm	Worst-case running time	Average-case/expected running time
Insertion sort	$\Theta(n^2)$	$\Theta(n^2)$
Merge sort	$\Theta(n \lg n)$	$\Theta(n \lg n)$
Heapsort	$O(n \lg n)$	—
Quicksort	$\Theta(n^2)$	$\Theta(n \lg n)$ (expected)
Counting sort	$\Theta(k + n)$	$\Theta(k + n)$
Radix sort	$\Theta(d(n + k))$	$\Theta(d(n + k))$
Bucket sort	$\Theta(n^2)$	$\Theta(n)$ (average-case)

### Order statistics

The  $i$ th order statistic of a set of  $n$  numbers is the  $i$ th smallest number in the set. We can, of course, select the  $i$ th order statistic by sorting the input and indexing the  $i$ th element of the output. With no assumptions about the input distribution, this method runs in  $\Omega(n \lg n)$  time, as the lower bound proved in Chapter 8 shows.

In Chapter 9, we show that we can find the  $i$ th smallest element in  $O(n)$  time, even when the elements are arbitrary real numbers. We present a randomized algorithm with tight pseudocode that runs in  $\Theta(n^2)$  time in the worst case, but whose expected running time is  $O(n)$ . We also give a more complicated algorithm that runs in  $O(n)$  worst-case time.

### Background

Although most of this part does not rely on difficult mathematics, some sections do require mathematical sophistication. In particular, analyses of quicksort, bucket sort, and the order-statistic algorithm use probability, which is reviewed in Appendix C, and the material on probabilistic analysis and randomized algorithms in Chapter 5. The analysis of the worst-case linear-time algorithm for order statistics involves somewhat more sophisticated mathematics than the other worst-case analyses in this part.

---

## 6 Heapsort

In this chapter, we introduce another sorting algorithm: heapsort. Like merge sort, but unlike insertion sort, heapsort's running time is  $O(n \lg n)$ . Like insertion sort, but unlike merge sort, heapsort sorts in place: only a constant number of array elements are stored outside the input array at any time. Thus, heapsort combines the better attributes of the two sorting algorithms we have already discussed.

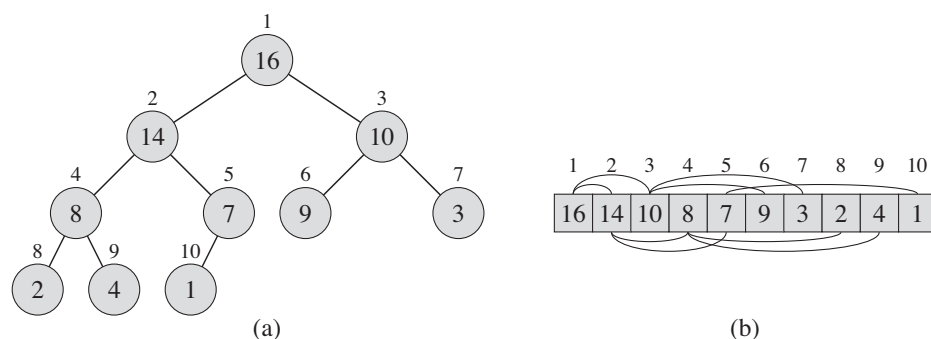
Heapsort also introduces another algorithm design technique: using a data structure, in this case one we call a “heap,” to manage information. Not only is the heap data structure useful for heapsort, but it also makes an efficient priority queue. The heap data structure will reappear in algorithms in later chapters.

The term “heap” was originally coined in the context of heapsort, but it has since come to refer to “garbage-collected storage,” such as the programming languages Java and Lisp provide. Our heap data structure is *not* garbage-collected storage, and whenever we refer to heaps in this book, we shall mean a data structure rather than an aspect of garbage collection.

---

### 6.1 Heaps

The (*binary*) *heap* data structure is an array object that we can view as a nearly complete binary tree (see Section B.5.3), as shown in Figure 6.1. Each node of the tree corresponds to an element of the array. The tree is completely filled on all levels except possibly the lowest, which is filled from the left up to a point. An array  $A$  that represents a heap is an object with two attributes:  $A.length$ , which (as usual) gives the number of elements in the array, and  $A.heap-size$ , which represents how many elements in the heap are stored within array  $A$ . That is, although  $A[1 \dots A.length]$  may contain numbers, only the elements in  $A[1 \dots A.heap-size]$ , where  $0 \leq A.heap-size \leq A.length$ , are valid elements of the heap. The root of the tree is  $A[1]$ , and given the index  $i$  of a node, we can easily compute the indices of its parent, left child, and right child:



**Figure 6.1** A max-heap viewed as (a) a binary tree and (b) an array. The number within the circle at each node in the tree is the value stored at that node. The number above a node is the corresponding index in the array. Above and below the array are lines showing parent-child relationships; parents are always to the left of their children. The tree has height three; the node at index 4 (with value 8) has height one.

PARENT( $i$ )

1 **return**  $\lfloor i/2 \rfloor$

LEFT( $i$ )

1 **return**  $2i$

RIGHT( $i$ )

1 **return**  $2i + 1$

On most computers, the LEFT procedure can compute  $2i$  in one instruction by simply shifting the binary representation of  $i$  left by one bit position. Similarly, the RIGHT procedure can quickly compute  $2i + 1$  by shifting the binary representation of  $i$  left by one bit position and then adding in a 1 as the low-order bit. The PARENT procedure can compute  $\lfloor i/2 \rfloor$  by shifting  $i$  right one bit position. Good implementations of heapsort often implement these procedures as “macros” or “in-line” procedures.

There are two kinds of binary heaps: max-heaps and min-heaps. In both kinds, the values in the nodes satisfy a *heap property*, the specifics of which depend on the kind of heap. In a *max-heap*, the *max-heap property* is that for every node  $i$  other than the root,

$$A[\text{PARENT}(i)] \geq A[i],$$

that is, the value of a node is at most the value of its parent. Thus, the largest element in a max-heap is stored at the root, and the subtree rooted at a node contains

values no larger than that contained at the node itself. A *min-heap* is organized in the opposite way; the *min-heap property* is that for every node  $i$  other than the root,

$$A[\text{PARENT}(i)] \leq A[i] .$$

The smallest element in a min-heap is at the root.

For the heapsort algorithm, we use max-heaps. Min-heaps commonly implement priority queues, which we discuss in Section 6.5. We shall be precise in specifying whether we need a max-heap or a min-heap for any particular application, and when properties apply to either max-heaps or min-heaps, we just use the term “heap.”

Viewing a heap as a tree, we define the *height* of a node in a heap to be the number of edges on the longest simple downward path from the node to a leaf, and we define the height of the heap to be the height of its root. Since a heap of  $n$  elements is based on a complete binary tree, its height is  $\Theta(\lg n)$  (see Exercise 6.1-2). We shall see that the basic operations on heaps run in time at most proportional to the height of the tree and thus take  $O(\lg n)$  time. The remainder of this chapter presents some basic procedures and shows how they are used in a sorting algorithm and a priority-queue data structure.

- The MAX-HEAPIFY procedure, which runs in  $O(\lg n)$  time, is the key to maintaining the max-heap property.
- The BUILD-MAX-HEAP procedure, which runs in linear time, produces a max-heap from an unordered input array.
- The HEAPSORT procedure, which runs in  $O(n \lg n)$  time, sorts an array in place.
- The MAX-HEAP-INSERT, HEAP-EXTRACT-MAX, HEAP-INCREASE-KEY, and HEAP-MAXIMUM procedures, which run in  $O(\lg n)$  time, allow the heap data structure to implement a priority queue.

## Exercises

### 6.1-1

What are the minimum and maximum numbers of elements in a heap of height  $h$ ?

### 6.1-2

Show that an  $n$ -element heap has height  $\lfloor \lg n \rfloor$ .

### 6.1-3

Show that in any subtree of a max-heap, the root of the subtree contains the largest value occurring anywhere in that subtree.



**6.1-4**

Where in a max-heap might the smallest element reside, assuming that all elements are distinct?

**6.1-5**

Is an array that is in sorted order a min-heap?

**6.1-6**

Is the array with values  $\langle 23, 17, 14, 6, 13, 10, 1, 5, 7, 12 \rangle$  a max-heap?

**6.1-7**

Show that, with the array representation for storing an  $n$ -element heap, the leaves are the nodes indexed by  $\lfloor n/2 \rfloor + 1, \lfloor n/2 \rfloor + 2, \dots, n$ .

---

## 6.2 Maintaining the heap property

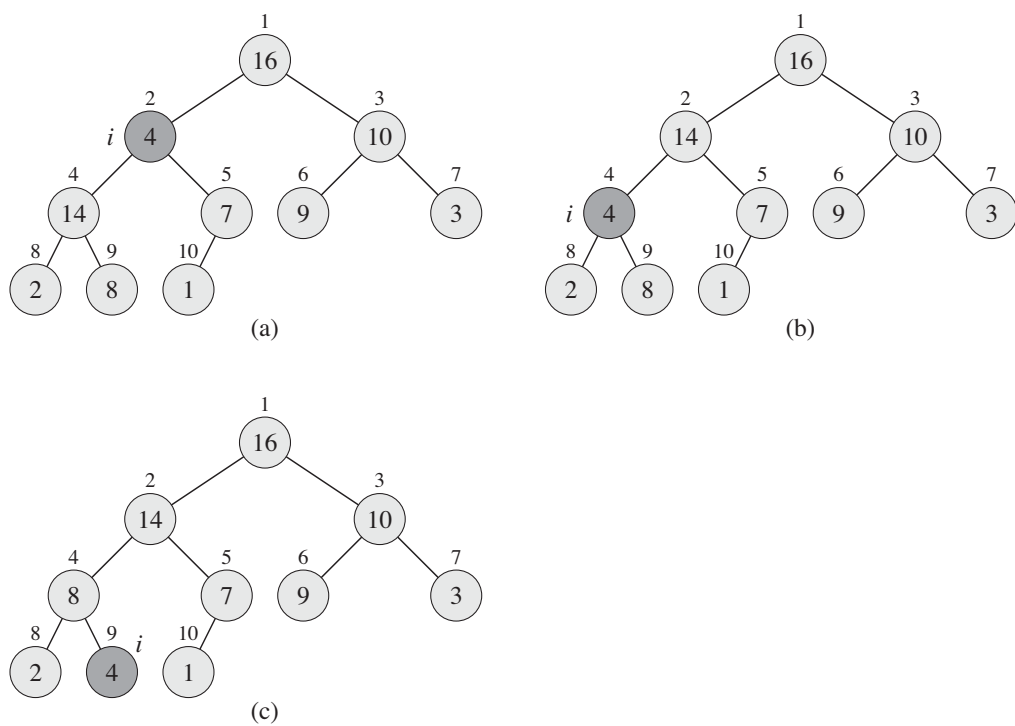
In order to maintain the max-heap property, we call the procedure MAX-HEAPIFY. Its inputs are an array  $A$  and an index  $i$  into the array. When it is called, MAX-HEAPIFY assumes that the binary trees rooted at  $\text{LEFT}(i)$  and  $\text{RIGHT}(i)$  are max-heaps, but that  $A[i]$  might be smaller than its children, thus violating the max-heap property. MAX-HEAPIFY lets the value at  $A[i]$  “float down” in the max-heap so that the subtree rooted at index  $i$  obeys the max-heap property.

MAX-HEAPIFY( $A, i$ )

```

1   $l = \text{LEFT}(i)$ 
2   $r = \text{RIGHT}(i)$ 
3  if  $l \leq A.\text{heap-size}$  and  $A[l] > A[i]$ 
4       $\text{largest} = l$ 
5  else  $\text{largest} = i$ 
6  if  $r \leq A.\text{heap-size}$  and  $A[r] > A[\text{largest}]$ 
7       $\text{largest} = r$ 
8  if  $\text{largest} \neq i$ 
9      exchange  $A[i]$  with  $A[\text{largest}]$ 
10     MAX-HEAPIFY( $A, \text{largest}$ )
```

Figure 6.2 illustrates the action of MAX-HEAPIFY. At each step, the largest of the elements  $A[i]$ ,  $A[\text{LEFT}(i)]$ , and  $A[\text{RIGHT}(i)]$  is determined, and its index is stored in  $\text{largest}$ . If  $A[i]$  is largest, then the subtree rooted at node  $i$  is already a max-heap and the procedure terminates. Otherwise, one of the two children has the largest element, and  $A[i]$  is swapped with  $A[\text{largest}]$ , which causes node  $i$  and its



**Figure 6.2** The action of  $\text{MAX-HEAPIFY}(A, 2)$ , where  $A.\text{heap-size} = 10$ . (a) The initial configuration, with  $A[2]$  at node  $i = 2$  violating the max-heap property since it is not larger than both children. The max-heap property is restored for node 2 in (b) by exchanging  $A[2]$  with  $A[4]$ , which destroys the max-heap property for node 4. The recursive call  $\text{MAX-HEAPIFY}(A, 4)$  now has  $i = 4$ . After swapping  $A[4]$  with  $A[9]$ , as shown in (c), node 4 is fixed up, and the recursive call  $\text{MAX-HEAPIFY}(A, 9)$  yields no further change to the data structure.

children to satisfy the max-heap property. The node indexed by *largest*, however, now has the original value  $A[i]$ , and thus the subtree rooted at *largest* might violate the max-heap property. Consequently, we call  $\text{MAX-HEAPIFY}$  recursively on that subtree.

The running time of  $\text{MAX-HEAPIFY}$  on a subtree of size  $n$  rooted at a given node  $i$  is the  $\Theta(1)$  time to fix up the relationships among the elements  $A[i]$ ,  $A[\text{LEFT}(i)]$ , and  $A[\text{RIGHT}(i)]$ , plus the time to run  $\text{MAX-HEAPIFY}$  on a subtree rooted at one of the children of node  $i$  (assuming that the recursive call occurs). The children's subtrees each have size at most  $2n/3$ —the worst case occurs when the bottom level of the tree is exactly half full—and therefore we can describe the running time of  $\text{MAX-HEAPIFY}$  by the recurrence

$$T(n) \leq T(2n/3) + \Theta(1).$$

The solution to this recurrence, by case 2 of the master theorem (Theorem 4.1), is  $T(n) = O(\lg n)$ . Alternatively, we can characterize the running time of MAX-HEAPIFY on a node of height  $h$  as  $O(h)$ .

## Exercises

### 6.2-1

Using Figure 6.2 as a model, illustrate the operation of MAX-HEAPIFY( $A, 3$ ) on the array  $A = \langle 27, 17, 3, 16, 13, 10, 1, 5, 7, 12, 4, 8, 9, 0 \rangle$ .

### 6.2-2

Starting with the procedure MAX-HEAPIFY, write pseudocode for the procedure MIN-HEAPIFY( $A, i$ ), which performs the corresponding manipulation on a min-heap. How does the running time of MIN-HEAPIFY compare to that of MAX-HEAPIFY?

### 6.2-3

What is the effect of calling MAX-HEAPIFY( $A, i$ ) when the element  $A[i]$  is larger than its children?

### 6.2-4

What is the effect of calling MAX-HEAPIFY( $A, i$ ) for  $i > A.heap-size/2$ ?

### 6.2-5

The code for MAX-HEAPIFY is quite efficient in terms of constant factors, except possibly for the recursive call in line 10, which might cause some compilers to produce inefficient code. Write an efficient MAX-HEAPIFY that uses an iterative control construct (a loop) instead of recursion.

### 6.2-6

Show that the worst-case running time of MAX-HEAPIFY on a heap of size  $n$  is  $\Omega(\lg n)$ . (*Hint:* For a heap with  $n$  nodes, give node values that cause MAX-HEAPIFY to be called recursively at every node on a simple path from the root down to a leaf.)

---

## 6.3 Building a heap

We can use the procedure MAX-HEAPIFY in a bottom-up manner to convert an array  $A[1..n]$ , where  $n = A.length$ , into a max-heap. By Exercise 6.1-7, the elements in the subarray  $A[(\lfloor n/2 \rfloor + 1) .. n]$  are all leaves of the tree, and so each is

a 1-element heap to begin with. The procedure BUILD-MAX-HEAP goes through the remaining nodes of the tree and runs MAX-HEAPIFY on each one.

```

BUILD-MAX-HEAP( $A$ )
1   $A.heap-size = A.length$ 
2  for  $i = \lfloor A.length/2 \rfloor$  downto 1
3      MAX-HEAPIFY( $A, i$ )

```

Figure 6.3 shows an example of the action of BUILD-MAX-HEAP.

To show why BUILD-MAX-HEAP works correctly, we use the following loop invariant:

At the start of each iteration of the **for** loop of lines 2–3, each node  $i + 1$ ,  $i + 2, \dots, n$  is the root of a max-heap.

We need to show that this invariant is true prior to the first loop iteration, that each iteration of the loop maintains the invariant, and that the invariant provides a useful property to show correctness when the loop terminates.

**Initialization:** Prior to the first iteration of the loop,  $i = \lfloor n/2 \rfloor$ . Each node  $\lfloor n/2 \rfloor + 1, \lfloor n/2 \rfloor + 2, \dots, n$  is a leaf and is thus the root of a trivial max-heap.

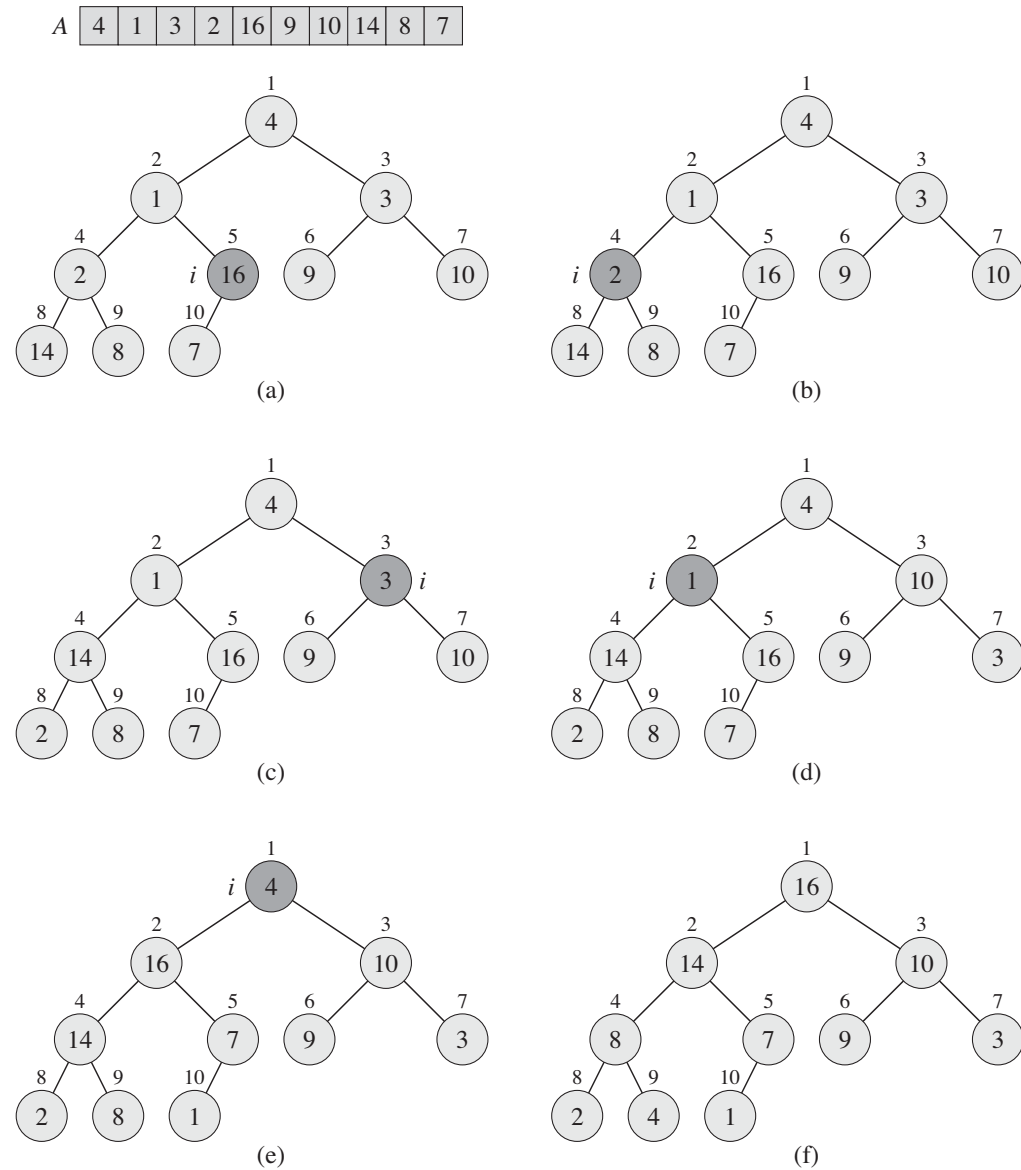
**Maintenance:** To see that each iteration maintains the loop invariant, observe that the children of node  $i$  are numbered higher than  $i$ . By the loop invariant, therefore, they are both roots of max-heaps. This is precisely the condition required for the call MAX-HEAPIFY( $A, i$ ) to make node  $i$  a max-heap root. Moreover, the MAX-HEAPIFY call preserves the property that nodes  $i + 1, i + 2, \dots, n$  are all roots of max-heaps. Decrementing  $i$  in the **for** loop update reestablishes the loop invariant for the next iteration.

**Termination:** At termination,  $i = 0$ . By the loop invariant, each node  $1, 2, \dots, n$  is the root of a max-heap. In particular, node 1 is.

We can compute a simple upper bound on the running time of BUILD-MAX-HEAP as follows. Each call to MAX-HEAPIFY costs  $O(\lg n)$  time, and BUILD-MAX-HEAP makes  $O(n)$  such calls. Thus, the running time is  $O(n \lg n)$ . This upper bound, though correct, is not asymptotically tight.

We can derive a tighter bound by observing that the time for MAX-HEAPIFY to run at a node varies with the height of the node in the tree, and the heights of most nodes are small. Our tighter analysis relies on the properties that an  $n$ -element heap has height  $\lceil \lg n \rceil$  (see Exercise 6.1-2) and at most  $\lceil n/2^{h+1} \rceil$  nodes of any height  $h$  (see Exercise 6.3-3).

The time required by MAX-HEAPIFY when called on a node of height  $h$  is  $O(h)$ , and so we can express the total cost of BUILD-MAX-HEAP as being bounded from above by



**Figure 6.3** The operation of BUILD-MAX-HEAP, showing the data structure before the call to MAX-HEAPIFY in line 3 of BUILD-MAX-HEAP. **(a)** A 10-element input array  $A$  and the binary tree it represents. The figure shows that the loop index  $i$  refers to node 5 before the call MAX-HEAPIFY( $A, i$ ). **(b)** The data structure that results. The loop index  $i$  for the next iteration refers to node 4. **(c)–(e)** Subsequent iterations of the **for** loop in BUILD-MAX-HEAP. Observe that whenever MAX-HEAPIFY is called on a node, the two subtrees of that node are both max-heaps. **(f)** The max-heap after BUILD-MAX-HEAP finishes.

$$\sum_{h=0}^{\lfloor \lg n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O\left(n \sum_{h=0}^{\lfloor \lg n \rfloor} \frac{h}{2^h}\right).$$

We evaluate the last summation by substituting  $x = 1/2$  in the formula (A.8), yielding

$$\begin{aligned} \sum_{h=0}^{\infty} \frac{h}{2^h} &= \frac{1/2}{(1 - 1/2)^2} \\ &= 2. \end{aligned}$$

Thus, we can bound the running time of BUILD-MAX-HEAP as

$$\begin{aligned} O\left(n \sum_{h=0}^{\lfloor \lg n \rfloor} \frac{h}{2^h}\right) &= O\left(n \sum_{h=0}^{\infty} \frac{h}{2^h}\right) \\ &= O(n). \end{aligned}$$

Hence, we can build a max-heap from an unordered array in linear time.

We can build a min-heap by the procedure BUILD-MIN-HEAP, which is the same as BUILD-MAX-HEAP but with the call to MAX-HEAPIFY in line 3 replaced by a call to MIN-HEAPIFY (see Exercise 6.2-2). BUILD-MIN-HEAP produces a min-heap from an unordered linear array in linear time.

## Exercises

### 6.3-1

Using Figure 6.3 as a model, illustrate the operation of BUILD-MAX-HEAP on the array  $A = \langle 5, 3, 17, 10, 84, 19, 6, 22, 9 \rangle$ .

### 6.3-2

Why do we want the loop index  $i$  in line 2 of BUILD-MAX-HEAP to decrease from  $\lfloor A.length/2 \rfloor$  to 1 rather than increase from 1 to  $\lfloor A.length/2 \rfloor$ ?

### 6.3-3

Show that there are at most  $\lceil n/2^{h+1} \rceil$  nodes of height  $h$  in any  $n$ -element heap.

---

## 6.4 The heapsort algorithm

The heapsort algorithm starts by using BUILD-MAX-HEAP to build a max-heap on the input array  $A[1..n]$ , where  $n = A.length$ . Since the maximum element of the array is stored at the root  $A[1]$ , we can put it into its correct final position

by exchanging it with  $A[n]$ . If we now discard node  $n$  from the heap—and we can do so by simply decrementing  $A.\text{heap-size}$ —we observe that the children of the root remain max-heaps, but the new root element might violate the max-heap property. All we need to do to restore the max-heap property, however, is call  $\text{MAX-HEAPIFY}(A, 1)$ , which leaves a max-heap in  $A[1..n-1]$ . The heapsort algorithm then repeats this process for the max-heap of size  $n-1$  down to a heap of size 2. (See Exercise 6.4-2 for a precise loop invariant.)

**HEAPSORT**( $A$ )

```

1  BUILD-MAX-HEAP( $A$ )
2  for  $i = A.\text{length}$  downto 2
3      exchange  $A[1]$  with  $A[i]$ 
4       $A.\text{heap-size} = A.\text{heap-size} - 1$ 
5      MAX-HEAPIFY( $A, 1$ )
```

Figure 6.4 shows an example of the operation of HEAPSORT after line 1 has built the initial max-heap. The figure shows the max-heap before the first iteration of the **for** loop of lines 2–5 and after each iteration.

The HEAPSORT procedure takes time  $O(n \lg n)$ , since the call to BUILD-MAX-HEAP takes time  $O(n)$  and each of the  $n-1$  calls to MAX-HEAPIFY takes time  $O(\lg n)$ .

## Exercises

### 6.4-1

Using Figure 6.4 as a model, illustrate the operation of HEAPSORT on the array  $A = \langle 5, 13, 2, 25, 7, 17, 20, 8, 4 \rangle$ .

### 6.4-2

Argue the correctness of HEAPSORT using the following loop invariant:

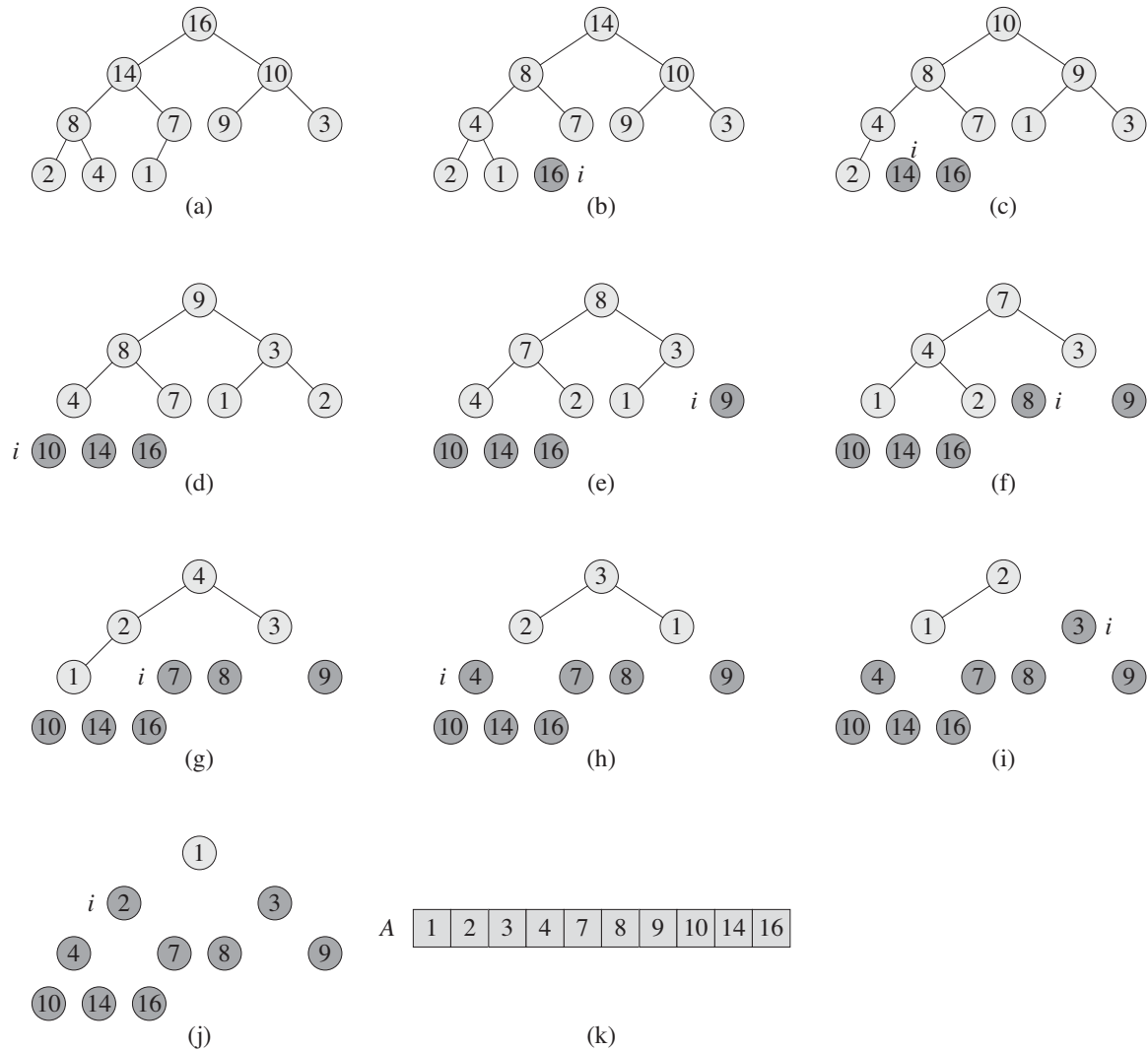
At the start of each iteration of the **for** loop of lines 2–5, the subarray  $A[1..i]$  is a max-heap containing the  $i$  smallest elements of  $A[1..n]$ , and the subarray  $A[i+1..n]$  contains the  $n-i$  largest elements of  $A[1..n]$ , sorted.

### 6.4-3

What is the running time of HEAPSORT on an array  $A$  of length  $n$  that is already sorted in increasing order? What about decreasing order?

### 6.4-4

Show that the worst-case running time of HEAPSORT is  $\Omega(n \lg n)$ .



**Figure 6.4** The operation of HEAPSORT. (a) The max-heap data structure just after BUILD-MAX-HEAP has built it in line 1. (b)–(j) The max-heap just after each call of MAX-HEAPIFY in line 5, showing the value of  $i$  at that time. Only lightly shaded nodes remain in the heap. (k) The resulting sorted array  $A$ .



**6.4-5 ★**

Show that when all elements are distinct, the best-case running time of HEAPSORT is  $\Omega(n \lg n)$ .

---

**6.5 Priority queues**

Heapsort is an excellent algorithm, but a good implementation of quicksort, presented in Chapter 7, usually beats it in practice. Nevertheless, the heap data structure itself has many uses. In this section, we present one of the most popular applications of a heap: as an efficient priority queue. As with heaps, priority queues come in two forms: max-priority queues and min-priority queues. We will focus here on how to implement max-priority queues, which are in turn based on max-heaps; Exercise 6.5-3 asks you to write the procedures for min-priority queues.

A **priority queue** is a data structure for maintaining a set  $S$  of elements, each with an associated value called a **key**. A **max-priority queue** supports the following operations:

INSERT( $S, x$ ) inserts the element  $x$  into the set  $S$ , which is equivalent to the operation  $S = S \cup \{x\}$ .

MAXIMUM( $S$ ) returns the element of  $S$  with the largest key.

EXTRACT-MAX( $S$ ) removes and returns the element of  $S$  with the largest key.

INCREASE-KEY( $S, x, k$ ) increases the value of element  $x$ 's key to the new value  $k$ , which is assumed to be at least as large as  $x$ 's current key value.

Among their other applications, we can use max-priority queues to schedule jobs on a shared computer. The max-priority queue keeps track of the jobs to be performed and their relative priorities. When a job is finished or interrupted, the scheduler selects the highest-priority job from among those pending by calling EXTRACT-MAX. The scheduler can add a new job to the queue at any time by calling INSERT.

Alternatively, a **min-priority queue** supports the operations INSERT, MINIMUM, EXTRACT-MIN, and DECREASE-KEY. A min-priority queue can be used in an event-driven simulator. The items in the queue are events to be simulated, each with an associated time of occurrence that serves as its key. The events must be simulated in order of their time of occurrence, because the simulation of an event can cause other events to be simulated in the future. The simulation program calls EXTRACT-MIN at each step to choose the next event to simulate. As new events are produced, the simulator inserts them into the min-priority queue by calling INSERT.

We shall see other uses for min-priority queues, highlighting the DECREASE-KEY operation, in Chapters 23 and 24.

Not surprisingly, we can use a heap to implement a priority queue. In a given application, such as job scheduling or event-driven simulation, elements of a priority queue correspond to objects in the application. We often need to determine which application object corresponds to a given priority-queue element, and vice versa. When we use a heap to implement a priority queue, therefore, we often need to store a *handle* to the corresponding application object in each heap element. The exact makeup of the handle (such as a pointer or an integer) depends on the application. Similarly, we need to store a handle to the corresponding heap element in each application object. Here, the handle would typically be an array index. Because heap elements change locations within the array during heap operations, an actual implementation, upon relocating a heap element, would also have to update the array index in the corresponding application object. Because the details of accessing application objects depend heavily on the application and its implementation, we shall not pursue them here, other than noting that in practice, these handles do need to be correctly maintained.

Now we discuss how to implement the operations of a max-priority queue. The procedure HEAP-MAXIMUM implements the MAXIMUM operation in  $\Theta(1)$  time.

HEAP-MAXIMUM( $A$ )

```
1 return  $A[1]$ 
```

The procedure HEAP-EXTRACT-MAX implements the EXTRACT-MAX operation. It is similar to the **for** loop body (lines 3–5) of the HEAPSORT procedure.

HEAP-EXTRACT-MAX( $A$ )

```
1 if  $A.heap-size < 1$ 
2   error “heap underflow”
3  $max = A[1]$ 
4  $A[1] = A[A.heap-size]$ 
5  $A.heap-size = A.heap-size - 1$ 
6 MAX-HEAPIFY( $A, 1$ )
7 return  $max$ 
```

The running time of HEAP-EXTRACT-MAX is  $O(\lg n)$ , since it performs only a constant amount of work on top of the  $O(\lg n)$  time for MAX-HEAPIFY.

The procedure HEAP-INCREASE-KEY implements the INCREASE-KEY operation. An index  $i$  into the array identifies the priority-queue element whose key we wish to increase. The procedure first updates the key of element  $A[i]$  to its new value. Because increasing the key of  $A[i]$  might violate the max-heap property,

the procedure then, in a manner reminiscent of the insertion loop (lines 5–7) of INSERTION-SORT from Section 2.1, traverses a simple path from this node toward the root to find a proper place for the newly increased key. As HEAP-INCREASE-KEY traverses this path, it repeatedly compares an element to its parent, exchanging their keys and continuing if the element's key is larger, and terminating if the element's key is smaller, since the max-heap property now holds. (See Exercise 6.5-5 for a precise loop invariant.)

HEAP-INCREASE-KEY( $A, i, key$ )

```

1  if  $key < A[i]$ 
2      error “new key is smaller than current key”
3   $A[i] = key$ 
4  while  $i > 1$  and  $A[PARENT(i)] < A[i]$ 
5      exchange  $A[i]$  with  $A[PARENT(i)]$ 
6       $i = PARENT(i)$ 
```

Figure 6.5 shows an example of a HEAP-INCREASE-KEY operation. The running time of HEAP-INCREASE-KEY on an  $n$ -element heap is  $O(\lg n)$ , since the path traced from the node updated in line 3 to the root has length  $O(\lg n)$ .

The procedure MAX-HEAP-INSERT implements the INSERT operation. It takes as an input the key of the new element to be inserted into max-heap  $A$ . The procedure first expands the max-heap by adding to the tree a new leaf whose key is  $-\infty$ . Then it calls HEAP-INCREASE-KEY to set the key of this new node to its correct value and maintain the max-heap property.

MAX-HEAP-INSERT( $A, key$ )

```

1   $A.heap-size = A.heap-size + 1$ 
2   $A[A.heap-size] = -\infty$ 
3  HEAP-INCREASE-KEY( $A, A.heap-size, key$ )
```

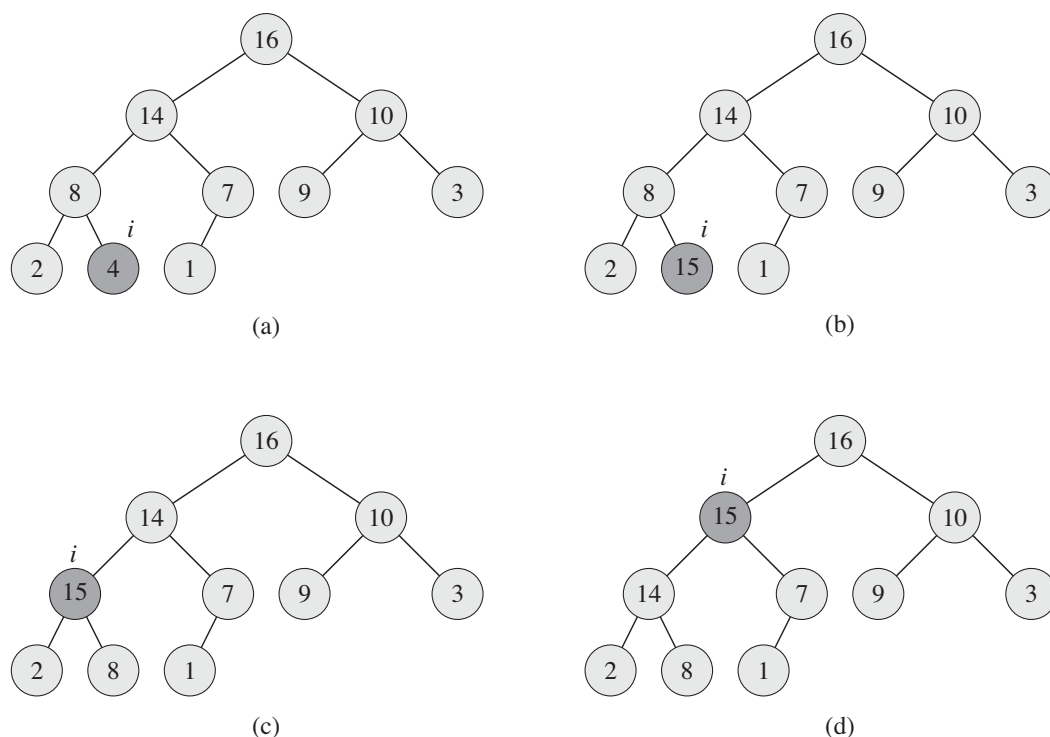
The running time of MAX-HEAP-INSERT on an  $n$ -element heap is  $O(\lg n)$ .

In summary, a heap can support any priority-queue operation on a set of size  $n$  in  $O(\lg n)$  time.

## Exercises

### 6.5-1

Illustrate the operation of HEAP-EXTRACT-MAX on the heap  $A = \langle 15, 13, 9, 5, 12, 8, 7, 4, 0, 6, 2, 1 \rangle$ .



**Figure 6.5** The operation of HEAP-INCREASE-KEY. **(a)** The max-heap of Figure 6.4(a) with a node whose index is  $i$  heavily shaded. **(b)** This node has its key increased to 15. **(c)** After one iteration of the **while** loop of lines 4–6, the node and its parent have exchanged keys, and the index  $i$  moves up to the parent. **(d)** The max-heap after one more iteration of the **while** loop. At this point,  $A[\text{PARENT}(i)] \geq A[i]$ . The max-heap property now holds and the procedure terminates.

### 6.5-2

Illustrate the operation of MAX-HEAP-INSERT( $A$ , 10) on the heap  $A = \langle 15, 13, 9, 5, 12, 8, 7, 4, 0, 6, 2, 1 \rangle$ .

### 6.5-3

Write pseudocode for the procedures HEAP-MINIMUM, HEAP-EXTRACT-MIN, HEAP-DECREASE-KEY, and MIN-HEAP-INSERT that implement a min-priority queue with a min-heap.

### 6.5-4

Why do we bother setting the key of the inserted node to  $-\infty$  in line 2 of MAX-HEAP-INSERT when the next thing we do is increase its key to the desired value?

**6.5-5**

Argue the correctness of HEAP-INCREASE-KEY using the following loop invariant:

At the start of each iteration of the **while** loop of lines 4–6, the subarray  $A[1 \dots A.heap-size]$  satisfies the max-heap property, except that there may be one violation:  $A[i]$  may be larger than  $A[PARENT(i)]$ .

You may assume that the subarray  $A[1 \dots A.heap-size]$  satisfies the max-heap property at the time HEAP-INCREASE-KEY is called.

**6.5-6**

Each exchange operation on line 5 of HEAP-INCREASE-KEY typically requires three assignments. Show how to use the idea of the inner loop of INSERTION-SORT to reduce the three assignments down to just one assignment.

**6.5-7**

Show how to implement a first-in, first-out queue with a priority queue. Show how to implement a stack with a priority queue. (Queues and stacks are defined in Section 10.1.)

**6.5-8**

The operation HEAP-DELETE( $A, i$ ) deletes the item in node  $i$  from heap  $A$ . Give an implementation of HEAP-DELETE that runs in  $O(\lg n)$  time for an  $n$ -element max-heap.

**6.5-9**

Give an  $O(n \lg k)$ -time algorithm to merge  $k$  sorted lists into one sorted list, where  $n$  is the total number of elements in all the input lists. (*Hint:* Use a min-heap for  $k$ -way merging.)

---

**Problems****6-1 Building a heap using insertion**

We can build a heap by repeatedly calling MAX-HEAP-INSERT to insert the elements into the heap. Consider the following variation on the BUILD-MAX-HEAP procedure:

BUILD-MAX-HEAP'(A)

```

1  A.heap-size = 1
2  for i = 2 to A.length
3      MAX-HEAP-INSERT(A, A[i])

```

- a. Do the procedures BUILD-MAX-HEAP and BUILD-MAX-HEAP' always create the same heap when run on the same input array? Prove that they do, or provide a counterexample.
- b. Show that in the worst case, BUILD-MAX-HEAP' requires  $\Theta(n \lg n)$  time to build an  $n$ -element heap.

### 6-2 Analysis of $d$ -ary heaps

A  **$d$ -ary heap** is like a binary heap, but (with one possible exception) non-leaf nodes have  $d$  children instead of 2 children.

- a. How would you represent a  $d$ -ary heap in an array?
- b. What is the height of a  $d$ -ary heap of  $n$  elements in terms of  $n$  and  $d$ ?
- c. Give an efficient implementation of EXTRACT-MAX in a  $d$ -ary max-heap. Analyze its running time in terms of  $d$  and  $n$ .
- d. Give an efficient implementation of INSERT in a  $d$ -ary max-heap. Analyze its running time in terms of  $d$  and  $n$ .
- e. Give an efficient implementation of INCREASE-KEY( $A, i, k$ ), which flags an error if  $k < A[i]$ , but otherwise sets  $A[i] = k$  and then updates the  $d$ -ary max-heap structure appropriately. Analyze its running time in terms of  $d$  and  $n$ .

### 6-3 Young tableaux

An  $m \times n$  **Young tableau** is an  $m \times n$  matrix such that the entries of each row are in sorted order from left to right and the entries of each column are in sorted order from top to bottom. Some of the entries of a Young tableau may be  $\infty$ , which we treat as nonexistent elements. Thus, a Young tableau can be used to hold  $r \leq mn$  finite numbers.

- a. Draw a  $4 \times 4$  Young tableau containing the elements  $\{9, 16, 3, 2, 4, 8, 5, 14, 12\}$ .
- b. Argue that an  $m \times n$  Young tableau  $Y$  is empty if  $Y[1, 1] = \infty$ . Argue that  $Y$  is full (contains  $mn$  elements) if  $Y[m, n] < \infty$ .

- c. Give an algorithm to implement EXTRACT-MIN on a nonempty  $m \times n$  Young tableau that runs in  $O(m + n)$  time. Your algorithm should use a recursive subroutine that solves an  $m \times n$  problem by recursively solving either an  $(m - 1) \times n$  or an  $m \times (n - 1)$  subproblem. (*Hint:* Think about MAX-HEAPIFY.) Define  $T(p)$ , where  $p = m + n$ , to be the maximum running time of EXTRACT-MIN on any  $m \times n$  Young tableau. Give and solve a recurrence for  $T(p)$  that yields the  $O(m + n)$  time bound.
- d. Show how to insert a new element into a nonfull  $m \times n$  Young tableau in  $O(m + n)$  time.
- e. Using no other sorting method as a subroutine, show how to use an  $n \times n$  Young tableau to sort  $n^2$  numbers in  $O(n^3)$  time.
- f. Give an  $O(m + n)$ -time algorithm to determine whether a given number is stored in a given  $m \times n$  Young tableau.

---

## Chapter notes

The heapsort algorithm was invented by Williams [357], who also described how to implement a priority queue with a heap. The BUILD-MAX-HEAP procedure was suggested by Floyd [106].

We use min-heaps to implement min-priority queues in Chapters 16, 23, and 24. We also give an implementation with improved time bounds for certain operations in Chapter 19 and, assuming that the keys are drawn from a bounded set of non-negative integers, Chapter 20.

If the data are  $b$ -bit integers, and the computer memory consists of addressable  $b$ -bit words, Fredman and Willard [115] showed how to implement MINIMUM in  $O(1)$  time and INSERT and EXTRACT-MIN in  $O(\sqrt{\lg n})$  time. Thorup [337] has improved the  $O(\sqrt{\lg n})$  bound to  $O(\lg \lg n)$  time. This bound uses an amount of space unbounded in  $n$ , but it can be implemented in linear space by using randomized hashing.

An important special case of priority queues occurs when the sequence of EXTRACT-MIN operations is **monotone**, that is, the values returned by successive EXTRACT-MIN operations are monotonically increasing over time. This case arises in several important applications, such as Dijkstra's single-source shortest-paths algorithm, which we discuss in Chapter 24, and in discrete-event simulation. For Dijkstra's algorithm it is particularly important that the DECREASE-KEY operation be implemented efficiently. For the monotone case, if the data are integers in the range  $1, 2, \dots, C$ , Ahuja, Mehlhorn, Orlin, and Tarjan [8] describe

how to implement EXTRACT-MIN and INSERT in  $O(\lg C)$  amortized time (see Chapter 17 for more on amortized analysis) and DECREASE-KEY in  $O(1)$  time, using a data structure called a radix heap. The  $O(\lg C)$  bound can be improved to  $O(\sqrt{\lg C})$  using Fibonacci heaps (see Chapter 19) in conjunction with radix heaps. Cherkassky, Goldberg, and Silverstein [65] further improved the bound to  $O(\lg^{1/3+\epsilon} C)$  expected time by combining the multilevel bucketing structure of Denardo and Fox [85] with the heap of Thorup mentioned earlier. Raman [291] further improved these results to obtain a bound of  $O(\min(\lg^{1/4+\epsilon} C, \lg^{1/3+\epsilon} n))$ , for any fixed  $\epsilon > 0$ .



The quicksort algorithm has a worst-case running time of  $\Theta(n^2)$  on an input array of  $n$  numbers. Despite this slow worst-case running time, quicksort is often the best practical choice for sorting because it is remarkably efficient on the average: its expected running time is  $\Theta(n \lg n)$ , and the constant factors hidden in the  $\Theta(n \lg n)$  notation are quite small. It also has the advantage of sorting in place (see page 17), and it works well even in virtual-memory environments.

Section 7.1 describes the algorithm and an important subroutine used by quicksort for partitioning. Because the behavior of quicksort is complex, we start with an intuitive discussion of its performance in Section 7.2 and postpone its precise analysis to the end of the chapter. Section 7.3 presents a version of quicksort that uses random sampling. This algorithm has a good expected running time, and no particular input elicits its worst-case behavior. Section 7.4 analyzes the randomized algorithm, showing that it runs in  $\Theta(n^2)$  time in the worst case and, assuming distinct elements, in expected  $O(n \lg n)$  time.

---

## 7.1 Description of quicksort

Quicksort, like merge sort, applies the divide-and-conquer paradigm introduced in Section 2.3.1. Here is the three-step divide-and-conquer process for sorting a typical subarray  $A[p \dots r]$ :

**Divide:** Partition (rearrange) the array  $A[p \dots r]$  into two (possibly empty) subarrays  $A[p \dots q - 1]$  and  $A[q + 1 \dots r]$  such that each element of  $A[p \dots q - 1]$  is less than or equal to  $A[q]$ , which is, in turn, less than or equal to each element of  $A[q + 1 \dots r]$ . Compute the index  $q$  as part of this partitioning procedure.

**Conquer:** Sort the two subarrays  $A[p \dots q - 1]$  and  $A[q + 1 \dots r]$  by recursive calls to quicksort.

**Combine:** Because the subarrays are already sorted, no work is needed to combine them: the entire array  $A[p \dots r]$  is now sorted.

The following procedure implements quicksort:

```

QUICKSORT( $A, p, r$ )
1  if  $p < r$ 
2       $q = \text{PARTITION}(A, p, r)$ 
3      QUICKSORT( $A, p, q - 1$ )
4      QUICKSORT( $A, q + 1, r$ )

```

To sort an entire array  $A$ , the initial call is  $\text{QUICKSORT}(A, 1, A.\text{length})$ .

### Partitioning the array

The key to the algorithm is the **PARTITION** procedure, which rearranges the subarray  $A[p \dots r]$  in place.

```

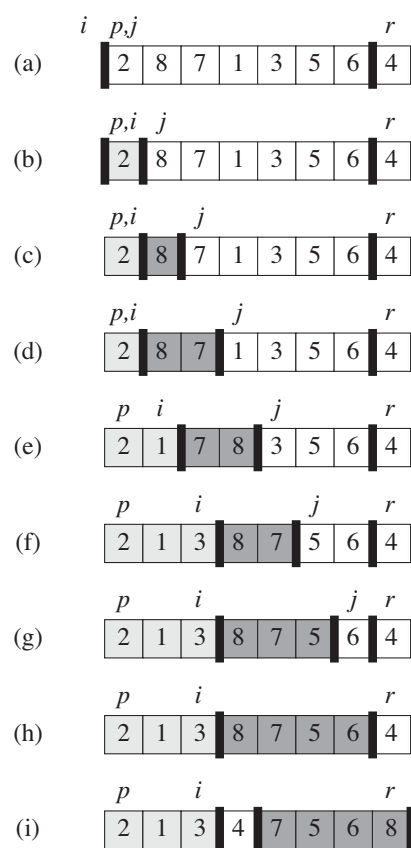
PARTITION( $A, p, r$ )
1   $x = A[r]$ 
2   $i = p - 1$ 
3  for  $j = p$  to  $r - 1$ 
4      if  $A[j] \leq x$ 
5           $i = i + 1$ 
6          exchange  $A[i]$  with  $A[j]$ 
7  exchange  $A[i + 1]$  with  $A[r]$ 
8  return  $i + 1$ 

```

Figure 7.1 shows how **PARTITION** works on an 8-element array. **PARTITION** always selects an element  $x = A[r]$  as a *pivot* element around which to partition the subarray  $A[p \dots r]$ . As the procedure runs, it partitions the array into four (possibly empty) regions. At the start of each iteration of the **for** loop in lines 3–6, the regions satisfy certain properties, shown in Figure 7.2. We state these properties as a loop invariant:

At the beginning of each iteration of the loop of lines 3–6, for any array index  $k$ ,

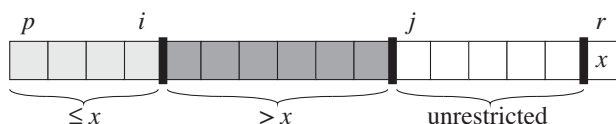
1. If  $p \leq k \leq i$ , then  $A[k] \leq x$ .
2. If  $i + 1 \leq k \leq j - 1$ , then  $A[k] > x$ .
3. If  $k = r$ , then  $A[k] = x$ .



**Figure 7.1** The operation of PARTITION on a sample array. Array entry  $A[r]$  becomes the pivot element  $x$ . Lightly shaded array elements are all in the first partition with values no greater than  $x$ . Heavily shaded elements are in the second partition with values greater than  $x$ . The unshaded elements have not yet been put in one of the first two partitions, and the final white element is the pivot  $x$ . (a) The initial array and variable settings. None of the elements have been placed in either of the first two partitions. (b) The value 2 is “swapped with itself” and put in the partition of smaller values. (c)–(d) The values 8 and 7 are added to the partition of larger values. (e) The values 1 and 8 are swapped, and the smaller partition grows. (f) The values 3 and 7 are swapped, and the smaller partition grows. (g)–(h) The larger partition grows to include 5 and 6, and the loop terminates. (i) In lines 7–8, the pivot element is swapped so that it lies between the two partitions.

The indices between  $j$  and  $r - 1$  are not covered by any of the three cases, and the values in these entries have no particular relationship to the pivot  $x$ .

We need to show that this loop invariant is true prior to the first iteration, that each iteration of the loop maintains the invariant, and that the invariant provides a useful property to show correctness when the loop terminates.



**Figure 7.2** The four regions maintained by the procedure PARTITION on a subarray  $A[p..r]$ . The values in  $A[p..i]$  are all less than or equal to  $x$ , the values in  $A[i+1..j-1]$  are all greater than  $x$ , and  $A[r] = x$ . The subarray  $A[j..r-1]$  can take on any values.

**Initialization:** Prior to the first iteration of the loop,  $i = p - 1$  and  $j = p$ . Because no values lie between  $p$  and  $i$  and no values lie between  $i + 1$  and  $j - 1$ , the first two conditions of the loop invariant are trivially satisfied. The assignment in line 1 satisfies the third condition.

**Maintenance:** As Figure 7.3 shows, we consider two cases, depending on the outcome of the test in line 4. Figure 7.3(a) shows what happens when  $A[j] > x$ ; the only action in the loop is to increment  $j$ . After  $j$  is incremented, condition 2 holds for  $A[j - 1]$  and all other entries remain unchanged. Figure 7.3(b) shows what happens when  $A[j] \leq x$ ; the loop increments  $i$ , swaps  $A[i]$  and  $A[j]$ , and then increments  $j$ . Because of the swap, we now have that  $A[i] \leq x$ , and condition 1 is satisfied. Similarly, we also have that  $A[j - 1] > x$ , since the item that was swapped into  $A[j - 1]$  is, by the loop invariant, greater than  $x$ .

**Termination:** At termination,  $j = r$ . Therefore, every entry in the array is in one of the three sets described by the invariant, and we have partitioned the values in the array into three sets: those less than or equal to  $x$ , those greater than  $x$ , and a singleton set containing  $x$ .

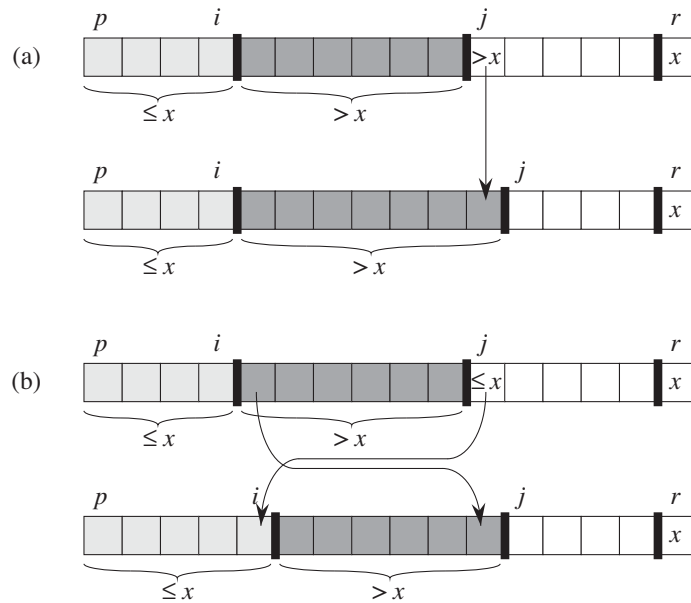
The final two lines of PARTITION finish up by swapping the pivot element with the leftmost element greater than  $x$ , thereby moving the pivot into its correct place in the partitioned array, and then returning the pivot's new index. The output of PARTITION now satisfies the specifications given for the divide step. In fact, it satisfies a slightly stronger condition: after line 2 of QUICKSORT,  $A[q]$  is strictly less than every element of  $A[q + 1..r]$ .

The running time of PARTITION on the subarray  $A[p..r]$  is  $\Theta(n)$ , where  $n = r - p + 1$  (see Exercise 7.1-3).

## Exercises

### 7.1-1

Using Figure 7.1 as a model, illustrate the operation of PARTITION on the array  $A = \langle 13, 19, 9, 5, 12, 8, 7, 4, 21, 2, 6, 11 \rangle$ .



**Figure 7.3** The two cases for one iteration of procedure PARTITION. (a) If  $A[j] > x$ , the only action is to increment  $j$ , which maintains the loop invariant. (b) If  $A[j] \leq x$ , index  $i$  is incremented,  $A[i]$  and  $A[j]$  are swapped, and then  $j$  is incremented. Again, the loop invariant is maintained.

### 7.1-2

What value of  $q$  does PARTITION return when all elements in the array  $A[p..r]$  have the same value? Modify PARTITION so that  $q = \lfloor (p + r)/2 \rfloor$  when all elements in the array  $A[p..r]$  have the same value.

### 7.1-3

Give a brief argument that the running time of PARTITION on a subarray of size  $n$  is  $\Theta(n)$ .

### 7.1-4

How would you modify QUICKSORT to sort into nonincreasing order?

## 7.2 Performance of quicksort

The running time of quicksort depends on whether the partitioning is balanced or unbalanced, which in turn depends on which elements are used for partitioning. If the partitioning is balanced, the algorithm runs asymptotically as fast as merge

sort. If the partitioning is unbalanced, however, it can run asymptotically as slowly as insertion sort. In this section, we shall informally investigate how quicksort performs under the assumptions of balanced versus unbalanced partitioning.

### Worst-case partitioning

The worst-case behavior for quicksort occurs when the partitioning routine produces one subproblem with  $n - 1$  elements and one with 0 elements. (We prove this claim in Section 7.4.1.) Let us assume that this unbalanced partitioning arises in each recursive call. The partitioning costs  $\Theta(n)$  time. Since the recursive call on an array of size 0 just returns,  $T(0) = \Theta(1)$ , and the recurrence for the running time is

$$\begin{aligned} T(n) &= T(n-1) + T(0) + \Theta(n) \\ &= T(n-1) + \Theta(n) . \end{aligned}$$

Intuitively, if we sum the costs incurred at each level of the recursion, we get an arithmetic series (equation (A.2)), which evaluates to  $\Theta(n^2)$ . Indeed, it is straightforward to use the substitution method to prove that the recurrence  $T(n) = T(n-1) + \Theta(n)$  has the solution  $T(n) = \Theta(n^2)$ . (See Exercise 7.2-1.)

Thus, if the partitioning is maximally unbalanced at every recursive level of the algorithm, the running time is  $\Theta(n^2)$ . Therefore the worst-case running time of quicksort is no better than that of insertion sort. Moreover, the  $\Theta(n^2)$  running time occurs when the input array is already completely sorted—a common situation in which insertion sort runs in  $O(n)$  time.

### Best-case partitioning

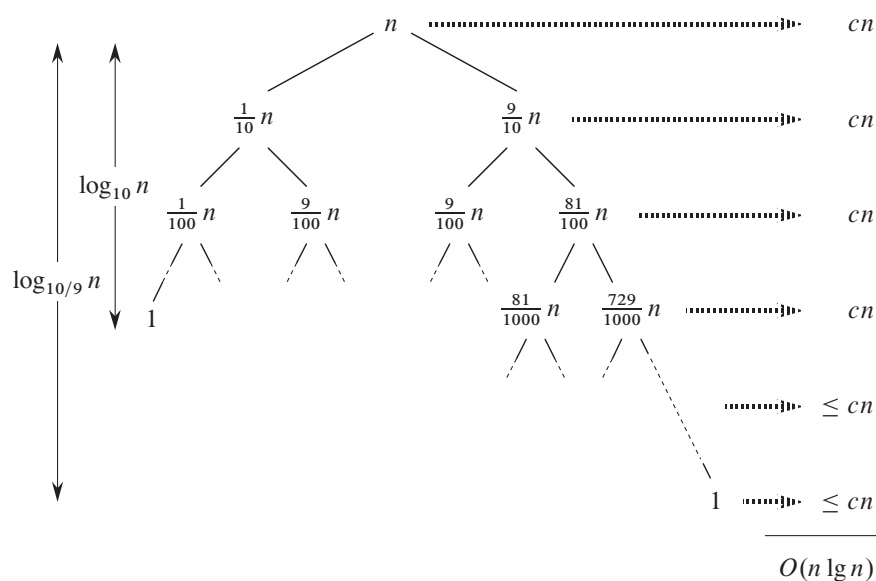
In the most even possible split, PARTITION produces two subproblems, each of size no more than  $n/2$ , since one is of size  $\lfloor n/2 \rfloor$  and one of size  $\lceil n/2 \rceil - 1$ . In this case, quicksort runs much faster. The recurrence for the running time is then

$$T(n) = 2T(n/2) + \Theta(n) ,$$

where we tolerate the sloppiness from ignoring the floor and ceiling and from subtracting 1. By case 2 of the master theorem (Theorem 4.1), this recurrence has the solution  $T(n) = \Theta(n \lg n)$ . By equally balancing the two sides of the partition at every level of the recursion, we get an asymptotically faster algorithm.

### Balanced partitioning

The average-case running time of quicksort is much closer to the best case than to the worst case, as the analyses in Section 7.4 will show. The key to understand-



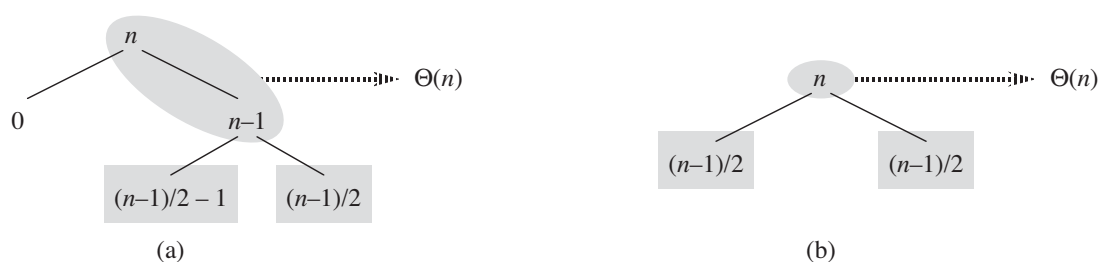
**Figure 7.4** A recursion tree for QUICKSORT in which PARTITION always produces a 9-to-1 split, yielding a running time of  $O(n \lg n)$ . Nodes show subproblem sizes, with per-level costs on the right. The per-level costs include the constant  $c$  implicit in the  $\Theta(n)$  term.

ing why is to understand how the balance of the partitioning is reflected in the recurrence that describes the running time.

Suppose, for example, that the partitioning algorithm always produces a 9-to-1 proportional split, which at first blush seems quite unbalanced. We then obtain the recurrence

$$T(n) = T(9n/10) + T(n/10) + cn ,$$

on the running time of quicksort, where we have explicitly included the constant  $c$  hidden in the  $\Theta(n)$  term. Figure 7.4 shows the recursion tree for this recurrence. Notice that every level of the tree has cost  $cn$ , until the recursion reaches a boundary condition at depth  $\log_{10} n = \Theta(\lg n)$ , and then the levels have cost at most  $cn$ . The recursion terminates at depth  $\log_{10/9} n = \Theta(\lg n)$ . The total cost of quicksort is therefore  $O(n \lg n)$ . Thus, with a 9-to-1 proportional split at every level of recursion, which intuitively seems quite unbalanced, quicksort runs in  $O(n \lg n)$  time—asymptotically the same as if the split were right down the middle. Indeed, even a 99-to-1 split yields an  $O(n \lg n)$  running time. In fact, any split of *constant* proportionality yields a recursion tree of depth  $\Theta(\lg n)$ , where the cost at each level is  $O(n)$ . The running time is therefore  $O(n \lg n)$  whenever the split has constant proportionality.



**Figure 7.5** (a) Two levels of a recursion tree for quicksort. The partitioning at the root costs  $n$  and produces a “bad” split: two subarrays of sizes  $0$  and  $n - 1$ . The partitioning of the subarray of size  $n - 1$  costs  $n - 1$  and produces a “good” split: subarrays of size  $(n - 1)/2 - 1$  and  $(n - 1)/2$ . (b) A single level of a recursion tree that is very well balanced. In both parts, the partitioning cost for the subproblems shown with elliptical shading is  $\Theta(n)$ . Yet the subproblems remaining to be solved in (a), shown with square shading, are no larger than the corresponding subproblems remaining to be solved in (b).

### Intuition for the average case

To develop a clear notion of the randomized behavior of quicksort, we must make an assumption about how frequently we expect to encounter the various inputs. The behavior of quicksort depends on the relative ordering of the values in the array elements given as the input, and not by the particular values in the array. As in our probabilistic analysis of the hiring problem in Section 5.2, we will assume for now that all permutations of the input numbers are equally likely.

When we run quicksort on a random input array, the partitioning is highly unlikely to happen in the same way at every level, as our informal analysis has assumed. We expect that some of the splits will be reasonably well balanced and that some will be fairly unbalanced. For example, Exercise 7.2-6 asks you to show that about 80 percent of the time `PARTITION` produces a split that is more balanced than 9 to 1, and about 20 percent of the time it produces a split that is less balanced than 9 to 1.

In the average case, `PARTITION` produces a mix of “good” and “bad” splits. In a recursion tree for an average-case execution of `PARTITION`, the good and bad splits are distributed randomly throughout the tree. Suppose, for the sake of intuition, that the good and bad splits alternate levels in the tree, and that the good splits are best-case splits and the bad splits are worst-case splits. Figure 7.5(a) shows the splits at two consecutive levels in the recursion tree. At the root of the tree, the cost is  $n$  for partitioning, and the subarrays produced have sizes  $n - 1$  and  $0$ : the worst case. At the next level, the subarray of size  $n - 1$  undergoes best-case partitioning into subarrays of size  $(n - 1)/2 - 1$  and  $(n - 1)/2$ . Let’s assume that the boundary-condition cost is 1 for the subarray of size 0.



The combination of the bad split followed by the good split produces three subarrays of sizes 0,  $(n - 1)/2 - 1$ , and  $(n - 1)/2$  at a combined partitioning cost of  $\Theta(n) + \Theta(n - 1) = \Theta(n)$ . Certainly, this situation is no worse than that in Figure 7.5(b), namely a single level of partitioning that produces two subarrays of size  $(n - 1)/2$ , at a cost of  $\Theta(n)$ . Yet this latter situation is balanced! Intuitively, the  $\Theta(n - 1)$  cost of the bad split can be absorbed into the  $\Theta(n)$  cost of the good split, and the resulting split is good. Thus, the running time of quicksort, when levels alternate between good and bad splits, is like the running time for good splits alone: still  $O(n \lg n)$ , but with a slightly larger constant hidden by the  $O$ -notation. We shall give a rigorous analysis of the expected running time of a randomized version of quicksort in Section 7.4.2.

## Exercises

### 7.2-1

Use the substitution method to prove that the recurrence  $T(n) = T(n - 1) + \Theta(n)$  has the solution  $T(n) = \Theta(n^2)$ , as claimed at the beginning of Section 7.2.

### 7.2-2

What is the running time of QUICKSORT when all elements of array  $A$  have the same value?

### 7.2-3

Show that the running time of QUICKSORT is  $\Theta(n^2)$  when the array  $A$  contains distinct elements and is sorted in decreasing order.

### 7.2-4

Banks often record transactions on an account in order of the times of the transactions, but many people like to receive their bank statements with checks listed in order by check number. People usually write checks in order by check number, and merchants usually cash them with reasonable dispatch. The problem of converting time-of-transaction ordering to check-number ordering is therefore the problem of sorting almost-sorted input. Argue that the procedure INSERTION-SORT would tend to beat the procedure QUICKSORT on this problem.

### 7.2-5

Suppose that the splits at every level of quicksort are in the proportion  $1 - \alpha$  to  $\alpha$ , where  $0 < \alpha \leq 1/2$  is a constant. Show that the minimum depth of a leaf in the recursion tree is approximately  $-\lg n / \lg \alpha$  and the maximum depth is approximately  $-\lg n / \lg(1 - \alpha)$ . (Don't worry about integer round-off.)

**7.2-6 ★**

Argue that for any constant  $0 < \alpha \leq 1/2$ , the probability is approximately  $1 - 2\alpha$  that on a random input array, PARTITION produces a split more balanced than  $1 - \alpha$  to  $\alpha$ .

---

**7.3 A randomized version of quicksort**

In exploring the average-case behavior of quicksort, we have made an assumption that all permutations of the input numbers are equally likely. In an engineering situation, however, we cannot always expect this assumption to hold. (See Exercise 7.2-4.) As we saw in Section 5.3, we can sometimes add randomization to an algorithm in order to obtain good expected performance over all inputs. Many people regard the resulting randomized version of quicksort as the sorting algorithm of choice for large enough inputs.

In Section 5.3, we randomized our algorithm by explicitly permuting the input. We could do so for quicksort also, but a different randomization technique, called *random sampling*, yields a simpler analysis. Instead of always using  $A[r]$  as the pivot, we will select a randomly chosen element from the subarray  $A[p \dots r]$ . We do so by first exchanging element  $A[r]$  with an element chosen at random from  $A[p \dots r]$ . By randomly sampling the range  $p, \dots, r$ , we ensure that the pivot element  $x = A[r]$  is equally likely to be any of the  $r - p + 1$  elements in the subarray. Because we randomly choose the pivot element, we expect the split of the input array to be reasonably well balanced on average.

The changes to PARTITION and QUICKSORT are small. In the new partition procedure, we simply implement the swap before actually partitioning:

RANDOMIZED-PARTITION( $A, p, r$ )

```

1   $i = \text{RANDOM}(p, r)$ 
2  exchange  $A[r]$  with  $A[i]$ 
3  return PARTITION( $A, p, r$ )
```

The new quicksort calls RANDOMIZED-PARTITION in place of PARTITION:

RANDOMIZED-QUICKSORT( $A, p, r$ )

```

1  if  $p < r$ 
2       $q = \text{RANDOMIZED-PARTITION}(A, p, r)$ 
3      RANDOMIZED-QUICKSORT( $A, p, q - 1$ )
4      RANDOMIZED-QUICKSORT( $A, q + 1, r$ )
```

We analyze this algorithm in the next section.

## Exercises

### 7.3-1

Why do we analyze the expected running time of a randomized algorithm and not its worst-case running time?

### 7.3-2

When RANDOMIZED-QUICKSORT runs, how many calls are made to the random-number generator RANDOM in the worst case? How about in the best case? Give your answer in terms of  $\Theta$ -notation.

---

## 7.4 Analysis of quicksort

Section 7.2 gave some intuition for the worst-case behavior of quicksort and for why we expect it to run quickly. In this section, we analyze the behavior of quicksort more rigorously. We begin with a worst-case analysis, which applies to either QUICKSORT or RANDOMIZED-QUICKSORT, and conclude with an analysis of the expected running time of RANDOMIZED-QUICKSORT.

### 7.4.1 Worst-case analysis

We saw in Section 7.2 that a worst-case split at every level of recursion in quicksort produces a  $\Theta(n^2)$  running time, which, intuitively, is the worst-case running time of the algorithm. We now prove this assertion.

Using the substitution method (see Section 4.3), we can show that the running time of quicksort is  $O(n^2)$ . Let  $T(n)$  be the worst-case time for the procedure QUICKSORT on an input of size  $n$ . We have the recurrence

$$T(n) = \max_{0 \leq q \leq n-1} (T(q) + T(n - q - 1)) + \Theta(n) , \quad (7.1)$$

where the parameter  $q$  ranges from 0 to  $n - 1$  because the procedure PARTITION produces two subproblems with total size  $n - 1$ . We guess that  $T(n) \leq cn^2$  for some constant  $c$ . Substituting this guess into recurrence (7.1), we obtain

$$\begin{aligned} T(n) &\leq \max_{0 \leq q \leq n-1} (cq^2 + c(n - q - 1)^2) + \Theta(n) \\ &= c \cdot \max_{0 \leq q \leq n-1} (q^2 + (n - q - 1)^2) + \Theta(n) . \end{aligned}$$

The expression  $q^2 + (n - q - 1)^2$  achieves a maximum over the parameter's range  $0 \leq q \leq n - 1$  at either endpoint. To verify this claim, note that the second derivative of the expression with respect to  $q$  is positive (see Exercise 7.4-3). This

observation gives us the bound  $\max_{0 \leq q \leq n-1} (q^2 + (n - q - 1)^2) \leq (n - 1)^2 = n^2 - 2n + 1$ . Continuing with our bounding of  $T(n)$ , we obtain

$$\begin{aligned} T(n) &\leq cn^2 - c(2n - 1) + \Theta(n) \\ &\leq cn^2, \end{aligned}$$

since we can pick the constant  $c$  large enough so that the  $c(2n - 1)$  term dominates the  $\Theta(n)$  term. Thus,  $T(n) = O(n^2)$ . We saw in Section 7.2 a specific case in which quicksort takes  $\Omega(n^2)$  time: when partitioning is unbalanced. Alternatively, Exercise 7.4-1 asks you to show that recurrence (7.1) has a solution of  $T(n) = \Omega(n^2)$ . Thus, the (worst-case) running time of quicksort is  $\Theta(n^2)$ .

### 7.4.2 Expected running time

We have already seen the intuition behind why the expected running time of RANDOMIZED-QUICKSORT is  $O(n \lg n)$ : if, in each level of recursion, the split induced by RANDOMIZED-PARTITION puts any constant fraction of the elements on one side of the partition, then the recursion tree has depth  $\Theta(\lg n)$ , and  $O(n)$  work is performed at each level. Even if we add a few new levels with the most unbalanced split possible between these levels, the total time remains  $O(n \lg n)$ . We can analyze the expected running time of RANDOMIZED-QUICKSORT precisely by first understanding how the partitioning procedure operates and then using this understanding to derive an  $O(n \lg n)$  bound on the expected running time. This upper bound on the expected running time, combined with the  $\Theta(n \lg n)$  best-case bound we saw in Section 7.2, yields a  $\Theta(n \lg n)$  expected running time. We assume throughout that the values of the elements being sorted are distinct.

### Running time and comparisons

The QUICKSORT and RANDOMIZED-QUICKSORT procedures differ only in how they select pivot elements; they are the same in all other respects. We can therefore couch our analysis of RANDOMIZED-QUICKSORT by discussing the QUICKSORT and PARTITION procedures, but with the assumption that pivot elements are selected randomly from the subarray passed to RANDOMIZED-PARTITION.

The running time of QUICKSORT is dominated by the time spent in the PARTITION procedure. Each time the PARTITION procedure is called, it selects a pivot element, and this element is never included in any future recursive calls to QUICKSORT and PARTITION. Thus, there can be at most  $n$  calls to PARTITION over the entire execution of the quicksort algorithm. One call to PARTITION takes  $O(1)$  time plus an amount of time that is proportional to the number of iterations of the **for** loop in lines 3–6. Each iteration of this **for** loop performs a comparison in line 4, comparing the pivot element to another element of the array  $A$ . Therefore,

if we can count the total number of times that line 4 is executed, we can bound the total time spent in the **for** loop during the entire execution of QUICKSORT.

**Lemma 7.1**

Let  $X$  be the number of comparisons performed in line 4 of PARTITION over the entire execution of QUICKSORT on an  $n$ -element array. Then the running time of QUICKSORT is  $O(n + X)$ .

**Proof** By the discussion above, the algorithm makes at most  $n$  calls to PARTITION, each of which does a constant amount of work and then executes the **for** loop some number of times. Each iteration of the **for** loop executes line 4. ■

Our goal, therefore, is to compute  $X$ , the total number of comparisons performed in all calls to PARTITION. We will not attempt to analyze how many comparisons are made in *each* call to PARTITION. Rather, we will derive an overall bound on the total number of comparisons. To do so, we must understand when the algorithm compares two elements of the array and when it does not. For ease of analysis, we rename the elements of the array  $A$  as  $z_1, z_2, \dots, z_n$ , with  $z_i$  being the  $i$ th smallest element. We also define the set  $Z_{ij} = \{z_i, z_{i+1}, \dots, z_j\}$  to be the set of elements between  $z_i$  and  $z_j$ , inclusive.

When does the algorithm compare  $z_i$  and  $z_j$ ? To answer this question, we first observe that each pair of elements is compared at most once. Why? Elements are compared only to the pivot element and, after a particular call of PARTITION finishes, the pivot element used in that call is never again compared to any other elements.

Our analysis uses indicator random variables (see Section 5.2). We define

$$X_{ij} = \mathbf{I}\{z_i \text{ is compared to } z_j\} ,$$

where we are considering whether the comparison takes place at any time during the execution of the algorithm, not just during one iteration or one call of PARTITION. Since each pair is compared at most once, we can easily characterize the total number of comparisons performed by the algorithm:

$$X = \sum_{i=1}^{n-1} \sum_{j=i+1}^n X_{ij} .$$

Taking expectations of both sides, and then using linearity of expectation and Lemma 5.1, we obtain

$$\mathbf{E}[X] = \mathbf{E}\left[\sum_{i=1}^{n-1} \sum_{j=i+1}^n X_{ij}\right]$$

$$\begin{aligned}
&= \sum_{i=1}^{n-1} \sum_{j=i+1}^n E[X_{ij}] \\
&= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr\{z_i \text{ is compared to } z_j\} .
\end{aligned} \tag{7.2}$$

It remains to compute  $\Pr\{z_i \text{ is compared to } z_j\}$ . Our analysis assumes that the RANDOMIZED-PARTITION procedure chooses each pivot randomly and independently.

Let us think about when two items are *not* compared. Consider an input to quicksort of the numbers 1 through 10 (in any order), and suppose that the first pivot element is 7. Then the first call to PARTITION separates the numbers into two sets:  $\{1, 2, 3, 4, 5, 6\}$  and  $\{8, 9, 10\}$ . In doing so, the pivot element 7 is compared to all other elements, but no number from the first set (e.g., 2) is or ever will be compared to any number from the second set (e.g., 9).

In general, because we assume that element values are distinct, once a pivot  $x$  is chosen with  $z_i < x < z_j$ , we know that  $z_i$  and  $z_j$  cannot be compared at any subsequent time. If, on the other hand,  $z_i$  is chosen as a pivot before any other item in  $Z_{ij}$ , then  $z_i$  will be compared to each item in  $Z_{ij}$ , except for itself. Similarly, if  $z_j$  is chosen as a pivot before any other item in  $Z_{ij}$ , then  $z_j$  will be compared to each item in  $Z_{ij}$ , except for itself. In our example, the values 7 and 9 are compared because 7 is the first item from  $Z_{7,9}$  to be chosen as a pivot. In contrast, 2 and 9 will never be compared because the first pivot element chosen from  $Z_{2,9}$  is 7. Thus,  $z_i$  and  $z_j$  are compared if and only if the first element to be chosen as a pivot from  $Z_{ij}$  is either  $z_i$  or  $z_j$ .

We now compute the probability that this event occurs. Prior to the point at which an element from  $Z_{ij}$  has been chosen as a pivot, the whole set  $Z_{ij}$  is together in the same partition. Therefore, any element of  $Z_{ij}$  is equally likely to be the first one chosen as a pivot. Because the set  $Z_{ij}$  has  $j-i+1$  elements, and because pivots are chosen randomly and independently, the probability that any given element is the first one chosen as a pivot is  $1/(j-i+1)$ . Thus, we have

$$\begin{aligned}
\Pr\{z_i \text{ is compared to } z_j\} &= \Pr\{z_i \text{ or } z_j \text{ is first pivot chosen from } Z_{ij}\} \\
&= \Pr\{z_i \text{ is first pivot chosen from } Z_{ij}\} \\
&\quad + \Pr\{z_j \text{ is first pivot chosen from } Z_{ij}\} \\
&= \frac{1}{j-i+1} + \frac{1}{j-i+1} \\
&= \frac{2}{j-i+1} .
\end{aligned} \tag{7.3}$$

The second line follows because the two events are mutually exclusive. Combining equations (7.2) and (7.3), we get that

$$E[X] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1}.$$

We can evaluate this sum using a change of variables ( $k = j - i$ ) and the bound on the harmonic series in equation (A.7):

$$\begin{aligned} E[X] &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1} \\ &= \sum_{i=1}^{n-1} \sum_{k=1}^{n-i} \frac{2}{k+1} \\ &< \sum_{i=1}^{n-1} \sum_{k=1}^n \frac{2}{k} \\ &= \sum_{i=1}^{n-1} O(\lg n) \\ &= O(n \lg n). \end{aligned} \tag{7.4}$$

Thus we conclude that, using RANDOMIZED-PARTITION, the expected running time of quicksort is  $O(n \lg n)$  when element values are distinct.

## Exercises

### 7.4-1

Show that in the recurrence

$$T(n) = \max_{0 \leq q \leq n-1} (T(q) + T(n-q-1)) + \Theta(n),$$

$$T(n) = \Omega(n^2).$$

### 7.4-2

Show that quicksort's best-case running time is  $\Omega(n \lg n)$ .

### 7.4-3

Show that the expression  $q^2 + (n - q - 1)^2$  achieves a maximum over  $q = 0, 1, \dots, n - 1$  when  $q = 0$  or  $q = n - 1$ .

### 7.4-4

Show that RANDOMIZED-QUICKSORT's expected running time is  $\Omega(n \lg n)$ .

**7.4-5**

We can improve the running time of quicksort in practice by taking advantage of the fast running time of insertion sort when its input is “nearly” sorted. Upon calling quicksort on a subarray with fewer than  $k$  elements, let it simply return without sorting the subarray. After the top-level call to quicksort returns, run insertion sort on the entire array to finish the sorting process. Argue that this sorting algorithm runs in  $O(nk + n \lg(n/k))$  expected time. How should we pick  $k$ , both in theory and in practice?

**7.4-6 ★**

Consider modifying the PARTITION procedure by randomly picking three elements from array  $A$  and partitioning about their median (the middle value of the three elements). Approximate the probability of getting at worst an  $\alpha$ -to- $(1 - \alpha)$  split, as a function of  $\alpha$  in the range  $0 < \alpha < 1$ .

---

**Problems**
**7-1 Hoare partition correctness**

The version of PARTITION given in this chapter is not the original partitioning algorithm. Here is the original partition algorithm, which is due to C. A. R. Hoare:

HOARE-PARTITION( $A, p, r$ )

```

1   $x = A[p]$ 
2   $i = p - 1$ 
3   $j = r + 1$ 
4  while TRUE
5      repeat
6           $j = j - 1$ 
7      until  $A[j] \leq x$ 
8      repeat
9           $i = i + 1$ 
10     until  $A[i] \geq x$ 
11     if  $i < j$ 
12         exchange  $A[i]$  with  $A[j]$ 
13     else return  $j$ 
```

- a. Demonstrate the operation of HOARE-PARTITION on the array  $A = \langle 13, 19, 9, 5, 12, 8, 7, 4, 11, 2, 6, 21 \rangle$ , showing the values of the array and auxiliary values after each iteration of the **while** loop in lines 4–13.



The next three questions ask you to give a careful argument that the procedure HOARE-PARTITION is correct. Assuming that the subarray  $A[p..r]$  contains at least two elements, prove the following:

- b.* The indices  $i$  and  $j$  are such that we never access an element of  $A$  outside the subarray  $A[p..r]$ .
- c.* When HOARE-PARTITION terminates, it returns a value  $j$  such that  $p \leq j < r$ .
- d.* Every element of  $A[p..j]$  is less than or equal to every element of  $A[j+1..r]$  when HOARE-PARTITION terminates.

The PARTITION procedure in Section 7.1 separates the pivot value (originally in  $A[r]$ ) from the two partitions it forms. The HOARE-PARTITION procedure, on the other hand, always places the pivot value (originally in  $A[p]$ ) into one of the two partitions  $A[p..j]$  and  $A[j+1..r]$ . Since  $p \leq j < r$ , this split is always nontrivial.

- e.* Rewrite the QUICKSORT procedure to use HOARE-PARTITION.

## 7-2 Quicksort with equal element values

The analysis of the expected running time of randomized quicksort in Section 7.4.2 assumes that all element values are distinct. In this problem, we examine what happens when they are not.

- a.* Suppose that all element values are equal. What would be randomized quicksort's running time in this case?
- b.* The PARTITION procedure returns an index  $q$  such that each element of  $A[p..q-1]$  is less than or equal to  $A[q]$  and each element of  $A[q+1..r]$  is greater than  $A[q]$ . Modify the PARTITION procedure to produce a procedure  $\text{PARTITION}'(A, p, r)$ , which permutes the elements of  $A[p..r]$  and returns two indices  $q$  and  $t$ , where  $p \leq q \leq t \leq r$ , such that
  - all elements of  $A[q..t]$  are equal,
  - each element of  $A[p..q-1]$  is less than  $A[q]$ , and
  - each element of  $A[t+1..r]$  is greater than  $A[q]$ .

Like PARTITION, your  $\text{PARTITION}'$  procedure should take  $\Theta(r-p)$  time.

- c.* Modify the RANDOMIZED-QUICKSORT procedure to call  $\text{PARTITION}'$ , and name the new procedure  $\text{RANDOMIZED-QUICKSORT}'$ . Then modify the QUICKSORT procedure to produce a procedure  $\text{QUICKSORT}'(p, r)$  that calls

RANDOMIZED-PARTITION' and recurses only on partitions of elements not known to be equal to each other.

- d.* Using QUICKSORT', how would you adjust the analysis in Section 7.4.2 to avoid the assumption that all elements are distinct?

### 7-3 Alternative quicksort analysis

An alternative analysis of the running time of randomized quicksort focuses on the expected running time of each individual recursive call to RANDOMIZED-QUICKSORT, rather than on the number of comparisons performed.

- a.* Argue that, given an array of size  $n$ , the probability that any particular element is chosen as the pivot is  $1/n$ . Use this to define indicator random variables  $X_i = I\{i\text{th smallest element is chosen as the pivot}\}$ . What is  $E[X_i]$ ?
- b.* Let  $T(n)$  be a random variable denoting the running time of quicksort on an array of size  $n$ . Argue that

$$E[T(n)] = E\left[\sum_{q=1}^n X_q (T(q-1) + T(n-q) + \Theta(n))\right]. \quad (7.5)$$

- c.* Show that we can rewrite equation (7.5) as

$$E[T(n)] = \frac{2}{n} \sum_{q=2}^{n-1} E[T(q)] + \Theta(n). \quad (7.6)$$

- d.* Show that

$$\sum_{k=2}^{n-1} k \lg k \leq \frac{1}{2} n^2 \lg n - \frac{1}{8} n^2. \quad (7.7)$$

(Hint: Split the summation into two parts, one for  $k = 2, 3, \dots, \lceil n/2 \rceil - 1$  and one for  $k = \lceil n/2 \rceil, \dots, n-1$ .)

- e.* Using the bound from equation (7.7), show that the recurrence in equation (7.6) has the solution  $E[T(n)] = \Theta(n \lg n)$ . (Hint: Show, by substitution, that  $E[T(n)] \leq an \lg n$  for sufficiently large  $n$  and for some positive constant  $a$ .)

#### 7-4 Stack depth for quicksort

The QUICKSORT algorithm of Section 7.1 contains two recursive calls to itself. After QUICKSORT calls PARTITION, it recursively sorts the left subarray and then it recursively sorts the right subarray. The second recursive call in QUICKSORT is not really necessary; we can avoid it by using an iterative control structure. This technique, called *tail recursion*, is provided automatically by good compilers. Consider the following version of quicksort, which simulates tail recursion:

TAIL-RECURSIVE-QUICKSORT( $A, p, r$ )

```

1  while  $p < r$ 
2      // Partition and sort left subarray.
3       $q = \text{PARTITION}(A, p, r)$ 
4      TAIL-RECURSIVE-QUICKSORT( $A, p, q - 1$ )
5       $p = q + 1$ 
```

- a. Argue that TAIL-RECURSIVE-QUICKSORT( $A, 1, A.length$ ) correctly sorts the array  $A$ .

Compilers usually execute recursive procedures by using a *stack* that contains pertinent information, including the parameter values, for each recursive call. The information for the most recent call is at the top of the stack, and the information for the initial call is at the bottom. Upon calling a procedure, its information is *pushed* onto the stack; when it terminates, its information is *popped*. Since we assume that array parameters are represented by pointers, the information for each procedure call on the stack requires  $O(1)$  stack space. The *stack depth* is the maximum amount of stack space used at any time during a computation.

- b. Describe a scenario in which TAIL-RECURSIVE-QUICKSORT's stack depth is  $\Theta(n)$  on an  $n$ -element input array.
- c. Modify the code for TAIL-RECURSIVE-QUICKSORT so that the worst-case stack depth is  $\Theta(\lg n)$ . Maintain the  $O(n \lg n)$  expected running time of the algorithm.

#### 7-5 Median-of-3 partition

One way to improve the RANDOMIZED-QUICKSORT procedure is to partition around a pivot that is chosen more carefully than by picking a random element from the subarray. One common approach is the *median-of-3* method: choose the pivot as the median (middle element) of a set of 3 elements randomly selected from the subarray. (See Exercise 7.4-6.) For this problem, let us assume that the elements in the input array  $A[1..n]$  are distinct and that  $n \geq 3$ . We denote the

sorted output array by  $A'[1..n]$ . Using the median-of-3 method to choose the pivot element  $x$ , define  $p_i = \Pr\{x = A'[i]\}$ .

- a. Give an exact formula for  $p_i$  as a function of  $n$  and  $i$  for  $i = 2, 3, \dots, n-1$ . (Note that  $p_1 = p_n = 0$ .)
- b. By what amount have we increased the likelihood of choosing the pivot as  $x = A'[(n+1)/2]$ , the median of  $A[1..n]$ , compared with the ordinary implementation? Assume that  $n \rightarrow \infty$ , and give the limiting ratio of these probabilities.
- c. If we define a “good” split to mean choosing the pivot as  $x = A'[i]$ , where  $n/3 \leq i \leq 2n/3$ , by what amount have we increased the likelihood of getting a good split compared with the ordinary implementation? (*Hint*: Approximate the sum by an integral.)
- d. Argue that in the  $\Omega(n \lg n)$  running time of quicksort, the median-of-3 method affects only the constant factor.

### 7-6 Fuzzy sorting of intervals

Consider a sorting problem in which we do not know the numbers exactly. Instead, for each number, we know an interval on the real line to which it belongs. That is, we are given  $n$  closed intervals of the form  $[a_i, b_i]$ , where  $a_i \leq b_i$ . We wish to **fuzzy-sort** these intervals, i.e., to produce a permutation  $\langle i_1, i_2, \dots, i_n \rangle$  of the intervals such that for  $j = 1, 2, \dots, n$ , there exist  $c_j \in [a_{i_j}, b_{i_j}]$  satisfying  $c_1 \leq c_2 \leq \dots \leq c_n$ .

- a. Design a randomized algorithm for fuzzy-sorting  $n$  intervals. Your algorithm should have the general structure of an algorithm that quicksorts the left endpoints (the  $a_i$  values), but it should take advantage of overlapping intervals to improve the running time. (As the intervals overlap more and more, the problem of fuzzy-sorting the intervals becomes progressively easier. Your algorithm should take advantage of such overlapping, to the extent that it exists.)
- b. Argue that your algorithm runs in expected time  $\Theta(n \lg n)$  in general, but runs in expected time  $\Theta(n)$  when all of the intervals overlap (i.e., when there exists a value  $x$  such that  $x \in [a_i, b_i]$  for all  $i$ ). Your algorithm should not be checking for this case explicitly; rather, its performance should naturally improve as the amount of overlap increases.

---

**Chapter notes**

The quicksort procedure was invented by Hoare [170]; Hoare's version appears in Problem 7-1. The PARTITION procedure given in Section 7.1 is due to N. Lomuto. The analysis in Section 7.4 is due to Avrim Blum. Sedgewick [305] and Bentley [43] provide a good reference on the details of implementation and how they matter.

McIlroy [248] showed how to engineer a “killer adversary” that produces an array on which virtually any implementation of quicksort takes  $\Theta(n^2)$  time. If the implementation is randomized, the adversary produces the array after seeing the random choices of the quicksort algorithm.

---

## 8      Sorting in Linear Time

We have now introduced several algorithms that can sort  $n$  numbers in  $O(n \lg n)$  time. Merge sort and heapsort achieve this upper bound in the worst case; quicksort achieves it on average. Moreover, for each of these algorithms, we can produce a sequence of  $n$  input numbers that causes the algorithm to run in  $\Omega(n \lg n)$  time.

These algorithms share an interesting property: *the sorted order they determine is based only on comparisons between the input elements*. We call such sorting algorithms **comparison sorts**. All the sorting algorithms introduced thus far are comparison sorts.

In Section 8.1, we shall prove that any comparison sort must make  $\Omega(n \lg n)$  comparisons in the worst case to sort  $n$  elements. Thus, merge sort and heapsort are asymptotically optimal, and no comparison sort exists that is faster by more than a constant factor.

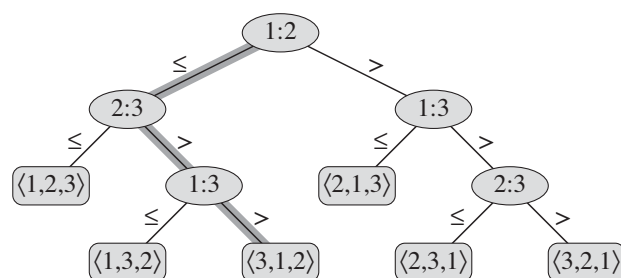
Sections 8.2, 8.3, and 8.4 examine three sorting algorithms—counting sort, radix sort, and bucket sort—that run in linear time. Of course, these algorithms use operations other than comparisons to determine the sorted order. Consequently, the  $\Omega(n \lg n)$  lower bound does not apply to them.

---

### 8.1 Lower bounds for sorting

In a comparison sort, we use only comparisons between elements to gain order information about an input sequence  $\langle a_1, a_2, \dots, a_n \rangle$ . That is, given two elements  $a_i$  and  $a_j$ , we perform one of the tests  $a_i < a_j$ ,  $a_i \leq a_j$ ,  $a_i = a_j$ ,  $a_i \geq a_j$ , or  $a_i > a_j$  to determine their relative order. We may not inspect the values of the elements or gain order information about them in any other way.

In this section, we assume without loss of generality that all the input elements are distinct. Given this assumption, comparisons of the form  $a_i = a_j$  are useless, so we can assume that no comparisons of this form are made. We also note that the comparisons  $a_i \leq a_j$ ,  $a_i \geq a_j$ ,  $a_i > a_j$ , and  $a_i < a_j$  are all equivalent in that



**Figure 8.1** The decision tree for insertion sort operating on three elements. An internal node annotated by  $i:j$  indicates a comparison between  $a_i$  and  $a_j$ . A leaf annotated by the permutation  $\langle \pi(1), \pi(2), \dots, \pi(n) \rangle$  indicates the ordering  $a_{\pi(1)} \leq a_{\pi(2)} \leq \dots \leq a_{\pi(n)}$ . The shaded path indicates the decisions made when sorting the input sequence  $\langle a_1 = 6, a_2 = 8, a_3 = 5 \rangle$ ; the permutation  $\langle 3, 1, 2 \rangle$  at the leaf indicates that the sorted ordering is  $a_3 = 5 \leq a_1 = 6 \leq a_2 = 8$ . There are  $3! = 6$  possible permutations of the input elements, and so the decision tree must have at least 6 leaves.

they yield identical information about the relative order of  $a_i$  and  $a_j$ . We therefore assume that all comparisons have the form  $a_i \leq a_j$ .

### The decision-tree model

We can view comparison sorts abstractly in terms of decision trees. A **decision tree** is a full binary tree that represents the comparisons between elements that are performed by a particular sorting algorithm operating on an input of a given size. Control, data movement, and all other aspects of the algorithm are ignored. Figure 8.1 shows the decision tree corresponding to the insertion sort algorithm from Section 2.1 operating on an input sequence of three elements.

In a decision tree, we annotate each internal node by  $i:j$  for some  $i$  and  $j$  in the range  $1 \leq i, j \leq n$ , where  $n$  is the number of elements in the input sequence. We also annotate each leaf by a permutation  $\langle \pi(1), \pi(2), \dots, \pi(n) \rangle$ . (See Section C.1 for background on permutations.) The execution of the sorting algorithm corresponds to tracing a simple path from the root of the decision tree down to a leaf. Each internal node indicates a comparison  $a_i \leq a_j$ . The left subtree then dictates subsequent comparisons once we know that  $a_i \leq a_j$ , and the right subtree dictates subsequent comparisons knowing that  $a_i > a_j$ . When we come to a leaf, the sorting algorithm has established the ordering  $a_{\pi(1)} \leq a_{\pi(2)} \leq \dots \leq a_{\pi(n)}$ . Because any correct sorting algorithm must be able to produce each permutation of its input, each of the  $n!$  permutations on  $n$  elements must appear as one of the leaves of the decision tree for a comparison sort to be correct. Furthermore, each of these leaves must be reachable from the root by a downward path corresponding to an actual

execution of the comparison sort. (We shall refer to such leaves as “reachable.”) Thus, we shall consider only decision trees in which each permutation appears as a reachable leaf.

### A lower bound for the worst case

The length of the longest simple path from the root of a decision tree to any of its reachable leaves represents the worst-case number of comparisons that the corresponding sorting algorithm performs. Consequently, the worst-case number of comparisons for a given comparison sort algorithm equals the height of its decision tree. A lower bound on the heights of all decision trees in which each permutation appears as a reachable leaf is therefore a lower bound on the running time of any comparison sort algorithm. The following theorem establishes such a lower bound.

#### **Theorem 8.1**

Any comparison sort algorithm requires  $\Omega(n \lg n)$  comparisons in the worst case.

**Proof** From the preceding discussion, it suffices to determine the height of a decision tree in which each permutation appears as a reachable leaf. Consider a decision tree of height  $h$  with  $l$  reachable leaves corresponding to a comparison sort on  $n$  elements. Because each of the  $n!$  permutations of the input appears as some leaf, we have  $n! \leq l$ . Since a binary tree of height  $h$  has no more than  $2^h$  leaves, we have

$$n! \leq l \leq 2^h,$$

which, by taking logarithms, implies

$$\begin{aligned} h &\geq \lg(n!) && \text{(since the } \lg \text{ function is monotonically increasing)} \\ &= \Omega(n \lg n) && \text{(by equation (3.19))} . \end{aligned}$$

■

#### **Corollary 8.2**

Heapsort and merge sort are asymptotically optimal comparison sorts.

**Proof** The  $O(n \lg n)$  upper bounds on the running times for heapsort and merge sort match the  $\Omega(n \lg n)$  worst-case lower bound from Theorem 8.1. ■

### Exercises

#### **8.1-1**

What is the smallest possible depth of a leaf in a decision tree for a comparison sort?



**8.1-2**

Obtain asymptotically tight bounds on  $\lg(n!)$  without using Stirling's approximation. Instead, evaluate the summation  $\sum_{k=1}^n \lg k$  using techniques from Section A.2.

**8.1-3**

Show that there is no comparison sort whose running time is linear for at least half of the  $n!$  inputs of length  $n$ . What about a fraction of  $1/n$  of the inputs of length  $n$ ? What about a fraction  $1/2^n$ ?

**8.1-4**

Suppose that you are given a sequence of  $n$  elements to sort. The input sequence consists of  $n/k$  subsequences, each containing  $k$  elements. The elements in a given subsequence are all smaller than the elements in the succeeding subsequence and larger than the elements in the preceding subsequence. Thus, all that is needed to sort the whole sequence of length  $n$  is to sort the  $k$  elements in each of the  $n/k$  subsequences. Show an  $\Omega(n \lg k)$  lower bound on the number of comparisons needed to solve this variant of the sorting problem. (*Hint:* It is not rigorous to simply combine the lower bounds for the individual subsequences.)

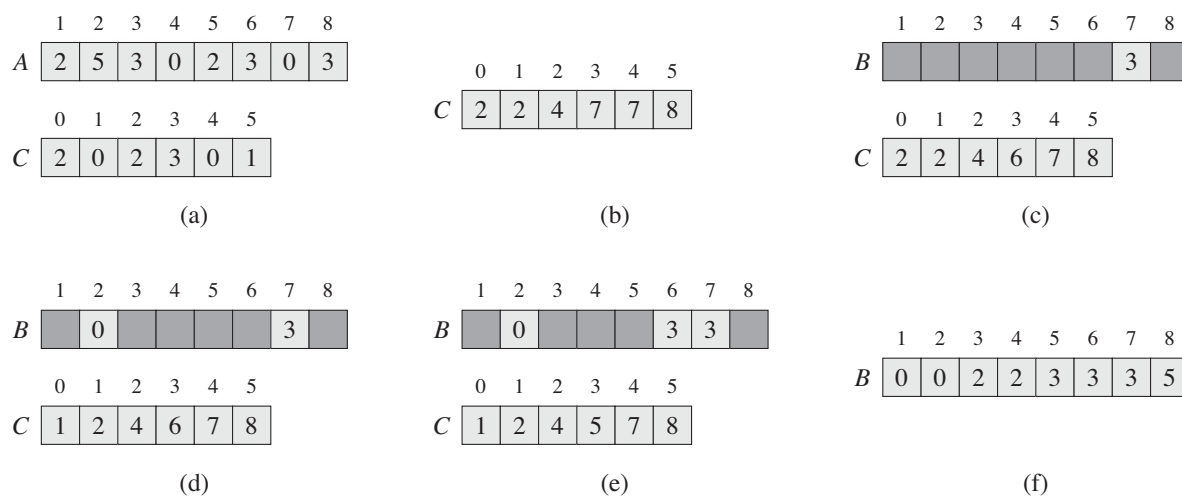
---

## 8.2 Counting sort

**Counting sort** assumes that each of the  $n$  input elements is an integer in the range 0 to  $k$ , for some integer  $k$ . When  $k = O(n)$ , the sort runs in  $\Theta(n)$  time.

Counting sort determines, for each input element  $x$ , the number of elements less than  $x$ . It uses this information to place element  $x$  directly into its position in the output array. For example, if 17 elements are less than  $x$ , then  $x$  belongs in output position 18. We must modify this scheme slightly to handle the situation in which several elements have the same value, since we do not want to put them all in the same position.

In the code for counting sort, we assume that the input is an array  $A[1..n]$ , and thus  $A.length = n$ . We require two other arrays: the array  $B[1..n]$  holds the sorted output, and the array  $C[0..k]$  provides temporary working storage.



**Figure 8.2** The operation of COUNTING-SORT on an input array  $A[1..8]$ , where each element of  $A$  is a nonnegative integer no larger than  $k = 5$ . (a) The array  $A$  and the auxiliary array  $C$  after line 5. (b) The array  $C$  after line 8. (c)–(e) The output array  $B$  and the auxiliary array  $C$  after one, two, and three iterations of the loop in lines 10–12, respectively. Only the lightly shaded elements of array  $B$  have been filled in. (f) The final sorted output array  $B$ .

COUNTING-SORT( $A, B, k$ )

```

1  let  $C[0..k]$  be a new array
2  for  $i = 0$  to  $k$ 
3       $C[i] = 0$ 
4  for  $j = 1$  to  $A.length$ 
5       $C[A[j]] = C[A[j]] + 1$ 
6  //  $C[i]$  now contains the number of elements equal to  $i$ .
7  for  $i = 1$  to  $k$ 
8       $C[i] = C[i] + C[i - 1]$ 
9  //  $C[i]$  now contains the number of elements less than or equal to  $i$ .
10 for  $j = A.length$  downto 1
11      $B[C[A[j]]] = A[j]$ 
12      $C[A[j]] = C[A[j]] - 1$ 

```

Figure 8.2 illustrates counting sort. After the **for** loop of lines 2–3 initializes the array  $C$  to all zeros, the **for** loop of lines 4–5 inspects each input element. If the value of an input element is  $i$ , we increment  $C[i]$ . Thus, after line 5,  $C[i]$  holds the number of input elements equal to  $i$  for each integer  $i = 0, 1, \dots, k$ . Lines 7–8 determine for each  $i = 0, 1, \dots, k$  how many input elements are less than or equal to  $i$  by keeping a running sum of the array  $C$ .

Finally, the **for** loop of lines 10–12 places each element  $A[j]$  into its correct sorted position in the output array  $B$ . If all  $n$  elements are distinct, then when we first enter line 10, for each  $A[j]$ , the value  $C[A[j]]$  is the correct final position of  $A[j]$  in the output array, since there are  $C[A[j]]$  elements less than or equal to  $A[j]$ . Because the elements might not be distinct, we decrement  $C[A[j]]$  each time we place a value  $A[j]$  into the  $B$  array. Decrementing  $C[A[j]]$  causes the next input element with a value equal to  $A[j]$ , if one exists, to go to the position immediately before  $A[j]$  in the output array.

How much time does counting sort require? The **for** loop of lines 2–3 takes time  $\Theta(k)$ , the **for** loop of lines 4–5 takes time  $\Theta(n)$ , the **for** loop of lines 7–8 takes time  $\Theta(k)$ , and the **for** loop of lines 10–12 takes time  $\Theta(n)$ . Thus, the overall time is  $\Theta(k + n)$ . In practice, we usually use counting sort when we have  $k = O(n)$ , in which case the running time is  $\Theta(n)$ .

Counting sort beats the lower bound of  $\Omega(n \lg n)$  proved in Section 8.1 because it is not a comparison sort. In fact, no comparisons between input elements occur anywhere in the code. Instead, counting sort uses the actual values of the elements to index into an array. The  $\Omega(n \lg n)$  lower bound for sorting does not apply when we depart from the comparison sort model.

An important property of counting sort is that it is *stable*: numbers with the same value appear in the output array in the same order as they do in the input array. That is, it breaks ties between two numbers by the rule that whichever number appears first in the input array appears first in the output array. Normally, the property of stability is important only when satellite data are carried around with the element being sorted. Counting sort's stability is important for another reason: counting sort is often used as a subroutine in radix sort. As we shall see in the next section, in order for radix sort to work correctly, counting sort must be stable.

## Exercises

### 8.2-1

Using Figure 8.2 as a model, illustrate the operation of COUNTING-SORT on the array  $A = \langle 6, 0, 2, 0, 1, 3, 4, 6, 1, 3, 2 \rangle$ .

### 8.2-2

Prove that COUNTING-SORT is stable.

### 8.2-3

Suppose that we were to rewrite the **for** loop header in line 10 of the COUNTING-SORT as

```
10  for  $j = 1$  to  $A.length$ 
```

Show that the algorithm still works properly. Is the modified algorithm stable?

**8.2-4**

Describe an algorithm that, given  $n$  integers in the range 0 to  $k$ , preprocesses its input and then answers any query about how many of the  $n$  integers fall into a range  $[a..b]$  in  $O(1)$  time. Your algorithm should use  $\Theta(n + k)$  preprocessing time.

---

**8.3 Radix sort**

*Radix sort* is the algorithm used by the card-sorting machines you now find only in computer museums. The cards have 80 columns, and in each column a machine can punch a hole in one of 12 places. The sorter can be mechanically “programmed” to examine a given column of each card in a deck and distribute the card into one of 12 bins depending on which place has been punched. An operator can then gather the cards bin by bin, so that cards with the first place punched are on top of cards with the second place punched, and so on.

For decimal digits, each column uses only 10 places. (The other two places are reserved for encoding nonnumeric characters.) A  $d$ -digit number would then occupy a field of  $d$  columns. Since the card sorter can look at only one column at a time, the problem of sorting  $n$  cards on a  $d$ -digit number requires a sorting algorithm.

Intuitively, you might sort numbers on their *most significant* digit, sort each of the resulting bins recursively, and then combine the decks in order. Unfortunately, since the cards in 9 of the 10 bins must be put aside to sort each of the bins, this procedure generates many intermediate piles of cards that you would have to keep track of. (See Exercise 8.3-5.)

Radix sort solves the problem of card sorting—counterintuitively—by sorting on the *least significant* digit first. The algorithm then combines the cards into a single deck, with the cards in the 0 bin preceding the cards in the 1 bin preceding the cards in the 2 bin, and so on. Then it sorts the entire deck again on the second-least significant digit and recombines the deck in a like manner. The process continues until the cards have been sorted on all  $d$  digits. Remarkably, at that point the cards are fully sorted on the  $d$ -digit number. Thus, only  $d$  passes through the deck are required to sort. Figure 8.3 shows how radix sort operates on a “deck” of seven 3-digit numbers.

In order for radix sort to work correctly, the digit sorts must be stable. The sort performed by a card sorter is stable, but the operator has to be wary about not changing the order of the cards as they come out of a bin, even though all the cards in a bin have the same digit in the chosen column.

329	720	720	329
457	355	329	355
657	436	436	436
839	457	839	457
436	657	355	657
720	329	457	720
355	839	657	839

**Figure 8.3** The operation of radix sort on a list of seven 3-digit numbers. The leftmost column is the input. The remaining columns show the list after successive sorts on increasingly significant digit positions. Shading indicates the digit position sorted on to produce each list from the previous one.

In a typical computer, which is a sequential random-access machine, we sometimes use radix sort to sort records of information that are keyed by multiple fields. For example, we might wish to sort dates by three keys: year, month, and day. We could run a sorting algorithm with a comparison function that, given two dates, compares years, and if there is a tie, compares months, and if another tie occurs, compares days. Alternatively, we could sort the information three times with a stable sort: first on day, next on month, and finally on year.

The code for radix sort is straightforward. The following procedure assumes that each element in the  $n$ -element array  $A$  has  $d$  digits, where digit 1 is the lowest-order digit and digit  $d$  is the highest-order digit.

```

RADIX-SORT( $A, d$ )
1  for  $i = 1$  to  $d$ 
2      use a stable sort to sort array  $A$  on digit  $i$ 

```

### Lemma 8.3

Given  $n$   $d$ -digit numbers in which each digit can take on up to  $k$  possible values, RADIX-SORT correctly sorts these numbers in  $\Theta(d(n + k))$  time if the stable sort it uses takes  $\Theta(n + k)$  time.

**Proof** The correctness of radix sort follows by induction on the column being sorted (see Exercise 8.3-3). The analysis of the running time depends on the stable sort used as the intermediate sorting algorithm. When each digit is in the range  $0$  to  $k - 1$  (so that it can take on  $k$  possible values), and  $k$  is not too large, counting sort is the obvious choice. Each pass over  $n$   $d$ -digit numbers then takes time  $\Theta(n + k)$ . There are  $d$  passes, and so the total time for radix sort is  $\Theta(d(n + k))$ . ■

When  $d$  is constant and  $k = O(n)$ , we can make radix sort run in linear time. More generally, we have some flexibility in how to break each key into digits.

**Lemma 8.4**

Given  $n$   $b$ -bit numbers and any positive integer  $r \leq b$ , RADIX-SORT correctly sorts these numbers in  $\Theta((b/r)(n + 2^r))$  time if the stable sort it uses takes  $\Theta(n + k)$  time for inputs in the range 0 to  $k$ .

**Proof** For a value  $r \leq b$ , we view each key as having  $d = \lceil b/r \rceil$  digits of  $r$  bits each. Each digit is an integer in the range 0 to  $2^r - 1$ , so that we can use counting sort with  $k = 2^r - 1$ . (For example, we can view a 32-bit word as having four 8-bit digits, so that  $b = 32$ ,  $r = 8$ ,  $k = 2^8 - 1 = 255$ , and  $d = b/r = 4$ .) Each pass of counting sort takes time  $\Theta(n + k) = \Theta(n + 2^r)$  and there are  $d$  passes, for a total running time of  $\Theta(d(n + 2^r)) = \Theta((b/r)(n + 2^r))$ . ■

For given values of  $n$  and  $b$ , we wish to choose the value of  $r$ , with  $r \leq b$ , that minimizes the expression  $(b/r)(n + 2^r)$ . If  $b < \lfloor \lg n \rfloor$ , then for any value of  $r \leq b$ , we have that  $(n + 2^r) = \Theta(n)$ . Thus, choosing  $r = b$  yields a running time of  $(b/b)(n + 2^b) = \Theta(n)$ , which is asymptotically optimal. If  $b \geq \lfloor \lg n \rfloor$ , then choosing  $r = \lfloor \lg n \rfloor$  gives the best time to within a constant factor, which we can see as follows. Choosing  $r = \lfloor \lg n \rfloor$  yields a running time of  $\Theta(bn/\lg n)$ . As we increase  $r$  above  $\lfloor \lg n \rfloor$ , the  $2^r$  term in the numerator increases faster than the  $r$  term in the denominator, and so increasing  $r$  above  $\lfloor \lg n \rfloor$  yields a running time of  $\Omega(bn/\lg n)$ . If instead we were to decrease  $r$  below  $\lfloor \lg n \rfloor$ , then the  $b/r$  term increases and the  $n + 2^r$  term remains at  $\Theta(n)$ .

Is radix sort preferable to a comparison-based sorting algorithm, such as quicksort? If  $b = O(\lg n)$ , as is often the case, and we choose  $r \approx \lg n$ , then radix sort's running time is  $\Theta(n)$ , which appears to be better than quicksort's expected running time of  $\Theta(n \lg n)$ . The constant factors hidden in the  $\Theta$ -notation differ, however. Although radix sort may make fewer passes than quicksort over the  $n$  keys, each pass of radix sort may take significantly longer. Which sorting algorithm we prefer depends on the characteristics of the implementations, of the underlying machine (e.g., quicksort often uses hardware caches more effectively than radix sort), and of the input data. Moreover, the version of radix sort that uses counting sort as the intermediate stable sort does not sort in place, which many of the  $\Theta(n \lg n)$ -time comparison sorts do. Thus, when primary memory storage is at a premium, we might prefer an in-place algorithm such as quicksort.

**Exercises****8.3-1**

Using Figure 8.3 as a model, illustrate the operation of RADIX-SORT on the following list of English words: COW, DOG, SEA, RUG, ROW, MOB, BOX, TAB, BAR, EAR, TAR, DIG, BIG, TEA, NOW, FOX.

**8.3-2**

Which of the following sorting algorithms are stable: insertion sort, merge sort, heapsort, and quicksort? Give a simple scheme that makes any sorting algorithm stable. How much additional time and space does your scheme entail?

**8.3-3**

Use induction to prove that radix sort works. Where does your proof need the assumption that the intermediate sort is stable?

**8.3-4**

Show how to sort  $n$  integers in the range 0 to  $n^3 - 1$  in  $O(n)$  time.

**8.3-5 ★**

In the first card-sorting algorithm in this section, exactly how many sorting passes are needed to sort  $d$ -digit decimal numbers in the worst case? How many piles of cards would an operator need to keep track of in the worst case?

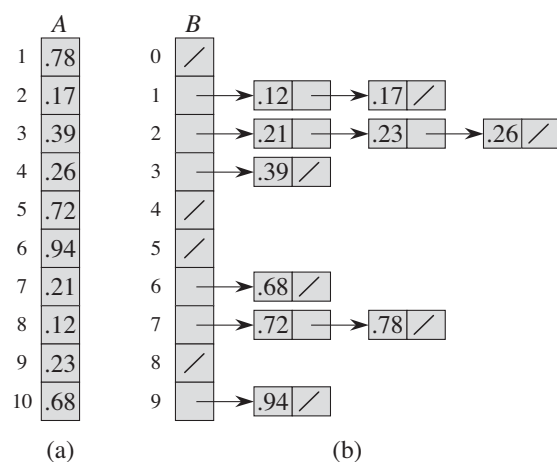
---

**8.4 Bucket sort**

**Bucket sort** assumes that the input is drawn from a uniform distribution and has an average-case running time of  $O(n)$ . Like counting sort, bucket sort is fast because it assumes something about the input. Whereas counting sort assumes that the input consists of integers in a small range, bucket sort assumes that the input is generated by a random process that distributes elements uniformly and independently over the interval  $[0, 1)$ . (See Section C.2 for a definition of uniform distribution.)

Bucket sort divides the interval  $[0, 1)$  into  $n$  equal-sized subintervals, or **buckets**, and then distributes the  $n$  input numbers into the buckets. Since the inputs are uniformly and independently distributed over  $[0, 1)$ , we do not expect many numbers to fall into each bucket. To produce the output, we simply sort the numbers in each bucket and then go through the buckets in order, listing the elements in each.

Our code for bucket sort assumes that the input is an  $n$ -element array  $A$  and that each element  $A[i]$  in the array satisfies  $0 \leq A[i] < 1$ . The code requires an auxiliary array  $B[0..n-1]$  of linked lists (buckets) and assumes that there is a mechanism for maintaining such lists. (Section 10.2 describes how to implement basic operations on linked lists.)



**Figure 8.4** The operation of BUCKET-SORT for  $n = 10$ . (a) The input array  $A[1 \dots 10]$ . (b) The array  $B[0 \dots 9]$  of sorted lists (buckets) after line 8 of the algorithm. Bucket  $i$  holds values in the half-open interval  $[i/10, (i + 1)/10)$ . The sorted output consists of a concatenation in order of the lists  $B[0], B[1], \dots, B[9]$ .

BUCKET-SORT( $A$ )

```

1  let  $B[0 \dots n - 1]$  be a new array
2   $n = A.length$ 
3  for  $i = 0$  to  $n - 1$ 
4      make  $B[i]$  an empty list
5  for  $i = 1$  to  $n$ 
6      insert  $A[i]$  into list  $B[\lfloor nA[i] \rfloor]$ 
7  for  $i = 0$  to  $n - 1$ 
8      sort list  $B[i]$  with insertion sort
9  concatenate the lists  $B[0], B[1], \dots, B[n - 1]$  together in order
```

Figure 8.4 shows the operation of bucket sort on an input array of 10 numbers.

To see that this algorithm works, consider two elements  $A[i]$  and  $A[j]$ . Assume without loss of generality that  $A[i] \leq A[j]$ . Since  $\lfloor nA[i] \rfloor \leq \lfloor nA[j] \rfloor$ , either element  $A[i]$  goes into the same bucket as  $A[j]$  or it goes into a bucket with a lower index. If  $A[i]$  and  $A[j]$  go into the same bucket, then the **for** loop of lines 7–8 puts them into the proper order. If  $A[i]$  and  $A[j]$  go into different buckets, then line 9 puts them into the proper order. Therefore, bucket sort works correctly.

To analyze the running time, observe that all lines except line 8 take  $O(n)$  time in the worst case. We need to analyze the total time taken by the  $n$  calls to insertion sort in line 8.



To analyze the cost of the calls to insertion sort, let  $n_i$  be the random variable denoting the number of elements placed in bucket  $B[i]$ . Since insertion sort runs in quadratic time (see Section 2.2), the running time of bucket sort is

$$T(n) = \Theta(n) + \sum_{i=0}^{n-1} O(n_i^2) .$$

We now analyze the average-case running time of bucket sort, by computing the expected value of the running time, where we take the expectation over the input distribution. Taking expectations of both sides and using linearity of expectation, we have

$$\begin{aligned} E[T(n)] &= E\left[\Theta(n) + \sum_{i=0}^{n-1} O(n_i^2)\right] \\ &= \Theta(n) + \sum_{i=0}^{n-1} E[O(n_i^2)] \quad (\text{by linearity of expectation}) \\ &= \Theta(n) + \sum_{i=0}^{n-1} O(E[n_i^2]) \quad (\text{by equation (C.22)}) . \end{aligned} \tag{8.1}$$

We claim that

$$E[n_i^2] = 2 - 1/n \tag{8.2}$$

for  $i = 0, 1, \dots, n-1$ . It is no surprise that each bucket  $i$  has the same value of  $E[n_i^2]$ , since each value in the input array  $A$  is equally likely to fall in any bucket. To prove equation (8.2), we define indicator random variables

$$X_{ij} = I\{A[j] \text{ falls in bucket } i\}$$

for  $i = 0, 1, \dots, n-1$  and  $j = 1, 2, \dots, n$ . Thus,

$$n_i = \sum_{j=1}^n X_{ij} .$$

To compute  $E[n_i^2]$ , we expand the square and regroup terms:

$$\begin{aligned}
E[n_i^2] &= E\left[\left(\sum_{j=1}^n X_{ij}\right)^2\right] \\
&= E\left[\sum_{j=1}^n \sum_{k=1}^n X_{ij} X_{ik}\right] \\
&= E\left[\sum_{j=1}^n X_{ij}^2 + \sum_{1 \leq j \leq n} \sum_{\substack{1 \leq k \leq n \\ k \neq j}} X_{ij} X_{ik}\right] \\
&= \sum_{j=1}^n E[X_{ij}^2] + \sum_{1 \leq j \leq n} \sum_{\substack{1 \leq k \leq n \\ k \neq j}} E[X_{ij} X_{ik}] , \tag{8.3}
\end{aligned}$$

where the last line follows by linearity of expectation. We evaluate the two summations separately. Indicator random variable  $X_{ij}$  is 1 with probability  $1/n$  and 0 otherwise, and therefore

$$\begin{aligned}
E[X_{ij}^2] &= 1^2 \cdot \frac{1}{n} + 0^2 \cdot \left(1 - \frac{1}{n}\right) \\
&= \frac{1}{n} .
\end{aligned}$$

When  $k \neq j$ , the variables  $X_{ij}$  and  $X_{ik}$  are independent, and hence

$$\begin{aligned}
E[X_{ij} X_{ik}] &= E[X_{ij}] E[X_{ik}] \\
&= \frac{1}{n} \cdot \frac{1}{n} \\
&= \frac{1}{n^2} .
\end{aligned}$$

Substituting these two expected values in equation (8.3), we obtain

$$\begin{aligned}
E[n_i^2] &= \sum_{j=1}^n \frac{1}{n} + \sum_{1 \leq j \leq n} \sum_{\substack{1 \leq k \leq n \\ k \neq j}} \frac{1}{n^2} \\
&= n \cdot \frac{1}{n} + n(n-1) \cdot \frac{1}{n^2} \\
&= 1 + \frac{n-1}{n} \\
&= 2 - \frac{1}{n} ,
\end{aligned}$$

which proves equation (8.2).

Using this expected value in equation (8.1), we conclude that the average-case running time for bucket sort is  $\Theta(n) + n \cdot O(2 - 1/n) = \Theta(n)$ .

Even if the input is not drawn from a uniform distribution, bucket sort may still run in linear time. As long as the input has the property that the sum of the squares of the bucket sizes is linear in the total number of elements, equation (8.1) tells us that bucket sort will run in linear time.

## Exercises

### 8.4-1

Using Figure 8.4 as a model, illustrate the operation of BUCKET-SORT on the array  $A = \langle .79, .13, .16, .64, .39, .20, .89, .53, .71, .42 \rangle$ .

### 8.4-2

Explain why the worst-case running time for bucket sort is  $\Theta(n^2)$ . What simple change to the algorithm preserves its linear average-case running time and makes its worst-case running time  $O(n \lg n)$ ?

### 8.4-3

Let  $X$  be a random variable that is equal to the number of heads in two flips of a fair coin. What is  $E[X^2]$ ? What is  $E^2[X]$ ?

### 8.4-4 ★

We are given  $n$  points in the unit circle,  $p_i = (x_i, y_i)$ , such that  $0 < x_i^2 + y_i^2 \leq 1$  for  $i = 1, 2, \dots, n$ . Suppose that the points are uniformly distributed; that is, the probability of finding a point in any region of the circle is proportional to the area of that region. Design an algorithm with an average-case running time of  $\Theta(n)$  to sort the  $n$  points by their distances  $d_i = \sqrt{x_i^2 + y_i^2}$  from the origin. (*Hint:* Design the bucket sizes in BUCKET-SORT to reflect the uniform distribution of the points in the unit circle.)

### 8.4-5 ★

A **probability distribution function**  $P(x)$  for a random variable  $X$  is defined by  $P(x) = \Pr\{X \leq x\}$ . Suppose that we draw a list of  $n$  random variables  $X_1, X_2, \dots, X_n$  from a continuous probability distribution function  $P$  that is computable in  $O(1)$  time. Give an algorithm that sorts these numbers in linear average-case time.

---

**Problems**
**8-1 Probabilistic lower bounds on comparison sorting**

In this problem, we prove a probabilistic  $\Omega(n \lg n)$  lower bound on the running time of any deterministic or randomized comparison sort on  $n$  distinct input elements. We begin by examining a deterministic comparison sort  $A$  with decision tree  $T_A$ . We assume that every permutation of  $A$ 's inputs is equally likely.

- a. Suppose that each leaf of  $T_A$  is labeled with the probability that it is reached given a random input. Prove that exactly  $n!$  leaves are labeled  $1/n!$  and that the rest are labeled 0.
- b. Let  $D(T)$  denote the external path length of a decision tree  $T$ ; that is,  $D(T)$  is the sum of the depths of all the leaves of  $T$ . Let  $T$  be a decision tree with  $k > 1$  leaves, and let  $LT$  and  $RT$  be the left and right subtrees of  $T$ . Show that  $D(T) = D(LT) + D(RT) + k$ .
- c. Let  $d(k)$  be the minimum value of  $D(T)$  over all decision trees  $T$  with  $k > 1$  leaves. Show that  $d(k) = \min_{1 \leq i \leq k-1} \{d(i) + d(k-i) + k\}$ . (Hint: Consider a decision tree  $T$  with  $k$  leaves that achieves the minimum. Let  $i_0$  be the number of leaves in  $LT$  and  $k - i_0$  the number of leaves in  $RT$ .)
- d. Prove that for a given value of  $k > 1$  and  $i$  in the range  $1 \leq i \leq k - 1$ , the function  $i \lg i + (k - i) \lg(k - i)$  is minimized at  $i = k/2$ . Conclude that  $d(k) = \Omega(k \lg k)$ .
- e. Prove that  $D(T_A) = \Omega(n! \lg(n!))$ , and conclude that the average-case time to sort  $n$  elements is  $\Omega(n \lg n)$ .

Now, consider a *randomized* comparison sort  $B$ . We can extend the decision-tree model to handle randomization by incorporating two kinds of nodes: ordinary comparison nodes and “randomization” nodes. A randomization node models a random choice of the form  $\text{RANDOM}(1, r)$  made by algorithm  $B$ ; the node has  $r$  children, each of which is equally likely to be chosen during an execution of the algorithm.

- f. Show that for any randomized comparison sort  $B$ , there exists a deterministic comparison sort  $A$  whose expected number of comparisons is no more than those made by  $B$ .

**8-2 Sorting in place in linear time**

Suppose that we have an array of  $n$  data records to sort and that the key of each record has the value 0 or 1. An algorithm for sorting such a set of records might possess some subset of the following three desirable characteristics:

1. The algorithm runs in  $O(n)$  time.
  2. The algorithm is stable.
  3. The algorithm sorts in place, using no more than a constant amount of storage space in addition to the original array.
- a.* Give an algorithm that satisfies criteria 1 and 2 above.
  - b.* Give an algorithm that satisfies criteria 1 and 3 above.
  - c.* Give an algorithm that satisfies criteria 2 and 3 above.
  - d.* Can you use any of your sorting algorithms from parts (a)–(c) as the sorting method used in line 2 of RADIX-SORT, so that RADIX-SORT sorts  $n$  records with  $b$ -bit keys in  $O(bn)$  time? Explain how or why not.
  - e.* Suppose that the  $n$  records have keys in the range from 1 to  $k$ . Show how to modify counting sort so that it sorts the records in place in  $O(n + k)$  time. You may use  $O(k)$  storage outside the input array. Is your algorithm stable? (*Hint:* How would you do it for  $k = 3$ ?)

**8-3 Sorting variable-length items**

- a.* You are given an array of integers, where different integers may have different numbers of digits, but the total number of digits over *all* the integers in the array is  $n$ . Show how to sort the array in  $O(n)$  time.
- b.* You are given an array of strings, where different strings may have different numbers of characters, but the total number of characters over all the strings is  $n$ . Show how to sort the strings in  $O(n)$  time.

(Note that the desired order here is the standard alphabetical order; for example,  $a < ab < b$ .)

**8-4 Water jugs**

Suppose that you are given  $n$  red and  $n$  blue water jugs, all of different shapes and sizes. All red jugs hold different amounts of water, as do the blue ones. Moreover, for every red jug, there is a blue jug that holds the same amount of water, and vice versa.

Your task is to find a grouping of the jugs into pairs of red and blue jugs that hold the same amount of water. To do so, you may perform the following operation: pick a pair of jugs in which one is red and one is blue, fill the red jug with water, and then pour the water into the blue jug. This operation will tell you whether the red or the blue jug can hold more water, or that they have the same volume. Assume that such a comparison takes one time unit. Your goal is to find an algorithm that makes a minimum number of comparisons to determine the grouping. Remember that you may not directly compare two red jugs or two blue jugs.

- a. Describe a deterministic algorithm that uses  $\Theta(n^2)$  comparisons to group the jugs into pairs.
- b. Prove a lower bound of  $\Omega(n \lg n)$  for the number of comparisons that an algorithm solving this problem must make.
- c. Give a randomized algorithm whose expected number of comparisons is  $O(n \lg n)$ , and prove that this bound is correct. What is the worst-case number of comparisons for your algorithm?

### 8-5 Average sorting

Suppose that, instead of sorting an array, we just require that the elements increase on average. More precisely, we call an  $n$ -element array  $A$   **$k$ -sorted** if, for all  $i = 1, 2, \dots, n - k$ , the following holds:

$$\frac{\sum_{j=i}^{i+k-1} A[j]}{k} \leq \frac{\sum_{j=i+1}^{i+k} A[j]}{k}.$$

- a. What does it mean for an array to be 1-sorted?
- b. Give a permutation of the numbers  $1, 2, \dots, 10$  that is 2-sorted, but not sorted.
- c. Prove that an  $n$ -element array is  $k$ -sorted if and only if  $A[i] \leq A[i + k]$  for all  $i = 1, 2, \dots, n - k$ .
- d. Give an algorithm that  $k$ -sorts an  $n$ -element array in  $O(n \lg(n/k))$  time.

We can also show a lower bound on the time to produce a  $k$ -sorted array, when  $k$  is a constant.

- e. Show that we can sort a  $k$ -sorted array of length  $n$  in  $O(n \lg k)$  time. (*Hint:* Use the solution to Exercise 6.5-9. )
- f. Show that when  $k$  is a constant,  $k$ -sorting an  $n$ -element array requires  $\Omega(n \lg n)$  time. (*Hint:* Use the solution to the previous part along with the lower bound on comparison sorts.)

**8-6 Lower bound on merging sorted lists**

The problem of merging two sorted lists arises frequently. We have seen a procedure for it as the subroutine MERGE in Section 2.3.1. In this problem, we will prove a lower bound of  $2n - 1$  on the worst-case number of comparisons required to merge two sorted lists, each containing  $n$  items.

First we will show a lower bound of  $2n - o(n)$  comparisons by using a decision tree.

- a. Given  $2n$  numbers, compute the number of possible ways to divide them into two sorted lists, each with  $n$  numbers.
- b. Using a decision tree and your answer to part (a), show that any algorithm that correctly merges two sorted lists must perform at least  $2n - o(n)$  comparisons.

Now we will show a slightly tighter  $2n - 1$  bound.

- c. Show that if two elements are consecutive in the sorted order and from different lists, then they must be compared.
- d. Use your answer to the previous part to show a lower bound of  $2n - 1$  comparisons for merging two sorted lists.

**8-7 The 0-1 sorting lemma and columnsort**

A **compare-exchange** operation on two array elements  $A[i]$  and  $A[j]$ , where  $i < j$ , has the form

```
COMPARE-EXCHANGE( $A, i, j$ )
1  if  $A[i] > A[j]$ 
2      exchange  $A[i]$  with  $A[j]$ 
```

After the compare-exchange operation, we know that  $A[i] \leq A[j]$ .

An **oblivious compare-exchange algorithm** operates solely by a sequence of prespecified compare-exchange operations. The indices of the positions compared in the sequence must be determined in advance, and although they can depend on the number of elements being sorted, they cannot depend on the values being sorted, nor can they depend on the result of any prior compare-exchange operation. For example, here is insertion sort expressed as an oblivious compare-exchange algorithm:

```
INSERTION-SORT( $A$ )
1  for  $j = 2$  to  $A.length$ 
2      for  $i = j - 1$  downto 1
3          COMPARE-EXCHANGE( $A, i, i + 1$ )
```

The **0-1 sorting lemma** provides a powerful way to prove that an oblivious compare-exchange algorithm produces a sorted result. It states that if an oblivious compare-exchange algorithm correctly sorts all input sequences consisting of only 0s and 1s, then it correctly sorts all inputs containing arbitrary values.

You will prove the 0-1 sorting lemma by proving its contrapositive: if an oblivious compare-exchange algorithm fails to sort an input containing arbitrary values, then it fails to sort some 0-1 input. Assume that an oblivious compare-exchange algorithm  $X$  fails to correctly sort the array  $A[1..n]$ . Let  $A[p]$  be the smallest value in  $A$  that algorithm  $X$  puts into the wrong location, and let  $A[q]$  be the value that algorithm  $X$  moves to the location into which  $A[p]$  should have gone. Define an array  $B[1..n]$  of 0s and 1s as follows:

$$B[i] = \begin{cases} 0 & \text{if } A[i] \leq A[p], \\ 1 & \text{if } A[i] > A[p]. \end{cases}$$

- a. Argue that  $A[q] > A[p]$ , so that  $B[p] = 0$  and  $B[q] = 1$ .
- b. To complete the proof of the 0-1 sorting lemma, prove that algorithm  $X$  fails to sort array  $B$  correctly.

Now you will use the 0-1 sorting lemma to prove that a particular sorting algorithm works correctly. The algorithm, **columnsort**, works on a rectangular array of  $n$  elements. The array has  $r$  rows and  $s$  columns (so that  $n = rs$ ), subject to three restrictions:

- $r$  must be even,
- $s$  must be a divisor of  $r$ , and
- $r \geq 2s^2$ .

When columnsort completes, the array is sorted in **column-major order**: reading down the columns, from left to right, the elements monotonically increase.

Columnsort operates in eight steps, regardless of the value of  $n$ . The odd steps are all the same: sort each column individually. Each even step is a fixed permutation. Here are the steps:

1. Sort each column.
2. Transpose the array, but reshape it back to  $r$  rows and  $s$  columns. In other words, turn the leftmost column into the top  $r/s$  rows, in order; turn the next column into the next  $r/s$  rows, in order; and so on.
3. Sort each column.
4. Perform the inverse of the permutation performed in step 2.



10	14	5	4	1	2	4	8	10	1	3	6	1	4	11
8	7	17	8	3	5	12	16	18	2	5	7	3	8	14
12	1	6	10	7	6	1	3	7	4	8	10	6	10	17
16	9	11	12	9	11	9	14	15	9	13	15	2	9	12
4	15	2	16	14	13	2	5	6	11	14	17	5	13	16
18	3	13	18	15	17	11	13	17	12	16	18	7	15	18
(a)			(b)			(c)			(d)			(e)		
1	4	11	5	10	16	4	10	16	1	7	13			
2	8	12	6	13	17	5	11	17	2	8	14			
3	9	14	7	15	18	6	12	18	3	9	15			
5	10	16	1	4	11	1	7	13	4	10	16			
6	13	17	2	8	12	2	8	14	5	11	17			
7	15	18	3	9	14	3	9	15	6	12	18			
(f)			(g)			(h)			(i)					

**Figure 8.5** The steps of columnsort. **(a)** The input array with 6 rows and 3 columns. **(b)** After sorting each column in step 1. **(c)** After transposing and reshaping in step 2. **(d)** After sorting each column in step 3. **(e)** After performing step 4, which inverts the permutation from step 2. **(f)** After sorting each column in step 5. **(g)** After shifting by half a column in step 6. **(h)** After sorting each column in step 7. **(i)** After performing step 8, which inverts the permutation from step 6. The array is now sorted in column-major order.

5. Sort each column.
6. Shift the top half of each column into the bottom half of the same column, and shift the bottom half of each column into the top half of the next column to the right. Leave the top half of the leftmost column empty. Shift the bottom half of the last column into the top half of a new rightmost column, and leave the bottom half of this new column empty.
7. Sort each column.
8. Perform the inverse of the permutation performed in step 6.

Figure 8.5 shows an example of the steps of columnsort with  $r = 6$  and  $s = 3$ . (Even though this example violates the requirement that  $r \geq 2s^2$ , it happens to work.)

- c. Argue that we can treat columnsort as an oblivious compare-exchange algorithm, even if we do not know what sorting method the odd steps use.

Although it might seem hard to believe that columnsort actually sorts, you will use the 0-1 sorting lemma to prove that it does. The 0-1 sorting lemma applies because we can treat columnsort as an oblivious compare-exchange algorithm. A

couple of definitions will help you apply the 0-1 sorting lemma. We say that an area of an array is *clean* if we know that it contains either all 0s or all 1s. Otherwise, the area might contain mixed 0s and 1s, and it is *dirty*. From here on, assume that the input array contains only 0s and 1s, and that we can treat it as an array with  $r$  rows and  $s$  columns.

- d. Prove that after steps 1–3, the array consists of some clean rows of 0s at the top, some clean rows of 1s at the bottom, and at most  $s$  dirty rows between them.
- e. Prove that after step 4, the array, read in column-major order, starts with a clean area of 0s, ends with a clean area of 1s, and has a dirty area of at most  $s^2$  elements in the middle.
- f. Prove that steps 5–8 produce a fully sorted 0-1 output. Conclude that column-sort correctly sorts all inputs containing arbitrary values.
- g. Now suppose that  $s$  does not divide  $r$ . Prove that after steps 1–3, the array consists of some clean rows of 0s at the top, some clean rows of 1s at the bottom, and at most  $2s - 1$  dirty rows between them. How large must  $r$  be, compared with  $s$ , for column-sort to correctly sort when  $s$  does not divide  $r$ ?
- h. Suggest a simple change to step 1 that allows us to maintain the requirement that  $r \geq 2s^2$  even when  $s$  does not divide  $r$ , and prove that with your change, column-sort correctly sorts.

---

## Chapter notes

The decision-tree model for studying comparison sorts was introduced by Ford and Johnson [110]. Knuth's comprehensive treatise on sorting [211] covers many variations on the sorting problem, including the information-theoretic lower bound on the complexity of sorting given here. Ben-Or [39] studied lower bounds for sorting using generalizations of the decision-tree model.

Knuth credits H. H. Seward with inventing counting sort in 1954, as well as with the idea of combining counting sort with radix sort. Radix sorting starting with the least significant digit appears to be a folk algorithm widely used by operators of mechanical card-sorting machines. According to Knuth, the first published reference to the method is a 1929 document by L. J. Comrie describing punched-card equipment. Bucket sorting has been in use since 1956, when the basic idea was proposed by E. J. Isaac and R. C. Singleton [188].

Munro and Raman [263] give a stable sorting algorithm that performs  $O(n^{1+\epsilon})$  comparisons in the worst case, where  $0 < \epsilon \leq 1$  is any fixed constant. Although

any of the  $O(n \lg n)$ -time algorithms make fewer comparisons, the algorithm by Munro and Raman moves data only  $O(n)$  times and operates in place.

The case of sorting  $n$   $b$ -bit integers in  $o(n \lg n)$  time has been considered by many researchers. Several positive results have been obtained, each under slightly different assumptions about the model of computation and the restrictions placed on the algorithm. All the results assume that the computer memory is divided into addressable  $b$ -bit words. Fredman and Willard [115] introduced the fusion tree data structure and used it to sort  $n$  integers in  $O(n \lg n / \lg \lg n)$  time. This bound was later improved to  $O(n \sqrt{\lg n})$  time by Andersson [16]. These algorithms require the use of multiplication and several precomputed constants. Andersson, Hagerup, Nilsson, and Raman [17] have shown how to sort  $n$  integers in  $O(n \lg \lg n)$  time without using multiplication, but their method requires storage that can be unbounded in terms of  $n$ . Using multiplicative hashing, we can reduce the storage needed to  $O(n)$ , but then the  $O(n \lg \lg n)$  worst-case bound on the running time becomes an expected-time bound. Generalizing the exponential search trees of Andersson [16], Thorup [335] gave an  $O(n(\lg \lg n)^2)$ -time sorting algorithm that does not use multiplication or randomization, and it uses linear space. Combining these techniques with some new ideas, Han [158] improved the bound for sorting to  $O(n \lg \lg n \lg \lg \lg n)$  time. Although these algorithms are important theoretical breakthroughs, they are all fairly complicated and at the present time seem unlikely to compete with existing sorting algorithms in practice.

The columnsort algorithm in Problem 8-7 is by Leighton [227].

The  $i$ th *order statistic* of a set of  $n$  elements is the  $i$ th smallest element. For example, the *minimum* of a set of elements is the first order statistic ( $i = 1$ ), and the *maximum* is the  $n$ th order statistic ( $i = n$ ). A *median*, informally, is the “halfway point” of the set. When  $n$  is odd, the median is unique, occurring at  $i = (n + 1)/2$ . When  $n$  is even, there are two medians, occurring at  $i = n/2$  and  $i = n/2 + 1$ . Thus, regardless of the parity of  $n$ , medians occur at  $i = \lfloor (n + 1)/2 \rfloor$  (the *lower median*) and  $i = \lceil (n + 1)/2 \rceil$  (the *upper median*). For simplicity in this text, however, we consistently use the phrase “the median” to refer to the lower median.

This chapter addresses the problem of selecting the  $i$ th order statistic from a set of  $n$  distinct numbers. We assume for convenience that the set contains distinct numbers, although virtually everything that we do extends to the situation in which a set contains repeated values. We formally specify the *selection problem* as follows:

**Input:** A set  $A$  of  $n$  (distinct) numbers and an integer  $i$ , with  $1 \leq i \leq n$ .

**Output:** The element  $x \in A$  that is larger than exactly  $i - 1$  other elements of  $A$ .

We can solve the selection problem in  $O(n \lg n)$  time, since we can sort the numbers using heapsort or merge sort and then simply index the  $i$ th element in the output array. This chapter presents faster algorithms.

In Section 9.1, we examine the problem of selecting the minimum and maximum of a set of elements. More interesting is the general selection problem, which we investigate in the subsequent two sections. Section 9.2 analyzes a practical randomized algorithm that achieves an  $O(n)$  expected running time, assuming distinct elements. Section 9.3 contains an algorithm of more theoretical interest that achieves the  $O(n)$  running time in the worst case.

---

## 9.1 Minimum and maximum

How many comparisons are necessary to determine the minimum of a set of  $n$  elements? We can easily obtain an upper bound of  $n - 1$  comparisons: examine each element of the set in turn and keep track of the smallest element seen so far. In the following procedure, we assume that the set resides in array  $A$ , where  $A.length = n$ .

```
MINIMUM( $A$ )
1   $min = A[1]$ 
2  for  $i = 2$  to  $A.length$ 
3      if  $min > A[i]$ 
4           $min = A[i]$ 
5  return  $min$ 
```

We can, of course, find the maximum with  $n - 1$  comparisons as well.

Is this the best we can do? Yes, since we can obtain a lower bound of  $n - 1$  comparisons for the problem of determining the minimum. Think of any algorithm that determines the minimum as a tournament among the elements. Each comparison is a match in the tournament in which the smaller of the two elements wins. Observing that every element except the winner must lose at least one match, we conclude that  $n - 1$  comparisons are necessary to determine the minimum. Hence, the algorithm MINIMUM is optimal with respect to the number of comparisons performed.

### Simultaneous minimum and maximum

In some applications, we must find both the minimum and the maximum of a set of  $n$  elements. For example, a graphics program may need to scale a set of  $(x, y)$  data to fit onto a rectangular display screen or other graphical output device. To do so, the program must first determine the minimum and maximum value of each coordinate.

At this point, it should be obvious how to determine both the minimum and the maximum of  $n$  elements using  $\Theta(n)$  comparisons, which is asymptotically optimal: simply find the minimum and maximum independently, using  $n - 1$  comparisons for each, for a total of  $2n - 2$  comparisons.

In fact, we can find both the minimum and the maximum using at most  $3 \lfloor n/2 \rfloor$  comparisons. We do so by maintaining both the minimum and maximum elements seen thus far. Rather than processing each element of the input by comparing it against the current minimum and maximum, at a cost of 2 comparisons per element,

we process elements in pairs. We compare pairs of elements from the input first *with each other*, and then we compare the smaller with the current minimum and the larger to the current maximum, at a cost of 3 comparisons for every 2 elements.

How we set up initial values for the current minimum and maximum depends on whether  $n$  is odd or even. If  $n$  is odd, we set both the minimum and maximum to the value of the first element, and then we process the rest of the elements in pairs. If  $n$  is even, we perform 1 comparison on the first 2 elements to determine the initial values of the minimum and maximum, and then process the rest of the elements in pairs as in the case for odd  $n$ .

Let us analyze the total number of comparisons. If  $n$  is odd, then we perform  $3 \lfloor n/2 \rfloor$  comparisons. If  $n$  is even, we perform 1 initial comparison followed by  $3(n-2)/2$  comparisons, for a total of  $3n/2 - 2$ . Thus, in either case, the total number of comparisons is at most  $3 \lfloor n/2 \rfloor$ .

### Exercises

#### 9.1-1

Show that the second smallest of  $n$  elements can be found with  $n + \lceil \lg n \rceil - 2$  comparisons in the worst case. (*Hint:* Also find the smallest element.)

#### 9.1-2 ★

Prove the lower bound of  $\lceil 3n/2 \rceil - 2$  comparisons in the worst case to find both the maximum and minimum of  $n$  numbers. (*Hint:* Consider how many numbers are potentially either the maximum or minimum, and investigate how a comparison affects these counts.)

---

## 9.2 Selection in expected linear time

The general selection problem appears more difficult than the simple problem of finding a minimum. Yet, surprisingly, the asymptotic running time for both problems is the same:  $\Theta(n)$ . In this section, we present a divide-and-conquer algorithm for the selection problem. The algorithm RANDOMIZED-SELECT is modeled after the quicksort algorithm of Chapter 7. As in quicksort, we partition the input array recursively. But unlike quicksort, which recursively processes both sides of the partition, RANDOMIZED-SELECT works on only one side of the partition. This difference shows up in the analysis: whereas quicksort has an expected running time of  $\Theta(n \lg n)$ , the expected running time of RANDOMIZED-SELECT is  $\Theta(n)$ , assuming that the elements are distinct.

RANDOMIZED-SELECT uses the procedure RANDOMIZED-PARTITION introduced in Section 7.3. Thus, like RANDOMIZED-QUICKSORT, it is a randomized algorithm, since its behavior is determined in part by the output of a random-number generator. The following code for RANDOMIZED-SELECT returns the  $i$ th smallest element of the array  $A[p..r]$ .

```

RANDOMIZED-SELECT( $A, p, r, i$ )
1  if  $p == r$ 
2      return  $A[p]$ 
3   $q = \text{RANDOMIZED-PARTITION}(A, p, r)$ 
4   $k = q - p + 1$ 
5  if  $i == k$            // the pivot value is the answer
6      return  $A[q]$ 
7  elseif  $i < k$ 
8      return RANDOMIZED-SELECT( $A, p, q - 1, i$ )
9  else return RANDOMIZED-SELECT( $A, q + 1, r, i - k$ )

```

The RANDOMIZED-SELECT procedure works as follows. Line 1 checks for the base case of the recursion, in which the subarray  $A[p..r]$  consists of just one element. In this case,  $i$  must equal 1, and we simply return  $A[p]$  in line 2 as the  $i$ th smallest element. Otherwise, the call to RANDOMIZED-PARTITION in line 3 partitions the array  $A[p..r]$  into two (possibly empty) subarrays  $A[p..q-1]$  and  $A[q+1..r]$  such that each element of  $A[p..q-1]$  is less than or equal to  $A[q]$ , which in turn is less than each element of  $A[q+1..r]$ . As in quicksort, we will refer to  $A[q]$  as the *pivot* element. Line 4 computes the number  $k$  of elements in the subarray  $A[p..q]$ , that is, the number of elements in the low side of the partition, plus one for the pivot element. Line 5 then checks whether  $A[q]$  is the  $i$ th smallest element. If it is, then line 6 returns  $A[q]$ . Otherwise, the algorithm determines in which of the two subarrays  $A[p..q-1]$  and  $A[q+1..r]$  the  $i$ th smallest element lies. If  $i < k$ , then the desired element lies on the low side of the partition, and line 8 recursively selects it from the subarray. If  $i > k$ , however, then the desired element lies on the high side of the partition. Since we already know  $k$  values that are smaller than the  $i$ th smallest element of  $A[p..r]$ —namely, the elements of  $A[p..q]$ —the desired element is the  $(i - k)$ th smallest element of  $A[q+1..r]$ , which line 9 finds recursively. The code appears to allow recursive calls to subarrays with 0 elements, but Exercise 9.2-1 asks you to show that this situation cannot happen.

The worst-case running time for RANDOMIZED-SELECT is  $\Theta(n^2)$ , even to find the minimum, because we could be extremely unlucky and always partition around the largest remaining element, and partitioning takes  $\Theta(n)$  time. We will see that

the algorithm has a linear expected running time, though, and because it is randomized, no particular input elicits the worst-case behavior.

To analyze the expected running time of RANDOMIZED-SELECT, we let the running time on an input array  $A[p..r]$  of  $n$  elements be a random variable that we denote by  $T(n)$ , and we obtain an upper bound on  $E[T(n)]$  as follows. The procedure RANDOMIZED-PARTITION is equally likely to return any element as the pivot. Therefore, for each  $k$  such that  $1 \leq k \leq n$ , the subarray  $A[p..q]$  has  $k$  elements (all less than or equal to the pivot) with probability  $1/n$ . For  $k = 1, 2, \dots, n$ , we define indicator random variables  $X_k$  where

$$X_k = I \{ \text{the subarray } A[p..q] \text{ has exactly } k \text{ elements} \} ,$$

and so, assuming that the elements are distinct, we have

$$E[X_k] = 1/n . \tag{9.1}$$

When we call RANDOMIZED-SELECT and choose  $A[q]$  as the pivot element, we do not know, a priori, if we will terminate immediately with the correct answer, recurse on the subarray  $A[p..q-1]$ , or recurse on the subarray  $A[q+1..r]$ . This decision depends on where the  $i$ th smallest element falls relative to  $A[q]$ . Assuming that  $T(n)$  is monotonically increasing, we can upper-bound the time needed for the recursive call by the time needed for the recursive call on the largest possible input. In other words, to obtain an upper bound, we assume that the  $i$ th element is always on the side of the partition with the greater number of elements. For a given call of RANDOMIZED-SELECT, the indicator random variable  $X_k$  has the value 1 for exactly one value of  $k$ , and it is 0 for all other  $k$ . When  $X_k = 1$ , the two subarrays on which we might recurse have sizes  $k-1$  and  $n-k$ . Hence, we have the recurrence

$$\begin{aligned} T(n) &\leq \sum_{k=1}^n X_k \cdot (T(\max(k-1, n-k)) + O(n)) \\ &= \sum_{k=1}^n X_k \cdot T(\max(k-1, n-k)) + O(n) . \end{aligned}$$



Taking expected values, we have

$$\begin{aligned}
& E[T(n)] \\
& \leq E \left[ \sum_{k=1}^n X_k \cdot T(\max(k-1, n-k)) + O(n) \right] \\
& = \sum_{k=1}^n E[X_k \cdot T(\max(k-1, n-k))] + O(n) \quad (\text{by linearity of expectation}) \\
& = \sum_{k=1}^n E[X_k] \cdot E[T(\max(k-1, n-k))] + O(n) \quad (\text{by equation (C.24)}) \\
& = \sum_{k=1}^n \frac{1}{n} \cdot E[T(\max(k-1, n-k))] + O(n) \quad (\text{by equation (9.1)}) .
\end{aligned}$$

In order to apply equation (C.24), we rely on  $X_k$  and  $T(\max(k-1, n-k))$  being independent random variables. Exercise 9.2-2 asks you to justify this assertion.

Let us consider the expression  $\max(k-1, n-k)$ . We have

$$\max(k-1, n-k) = \begin{cases} k-1 & \text{if } k > \lceil n/2 \rceil, \\ n-k & \text{if } k \leq \lceil n/2 \rceil. \end{cases}$$

If  $n$  is even, each term from  $T(\lceil n/2 \rceil)$  up to  $T(n-1)$  appears exactly twice in the summation, and if  $n$  is odd, all these terms appear twice and  $T(\lfloor n/2 \rfloor)$  appears once. Thus, we have

$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} E[T(k)] + O(n) .$$

We show that  $E[T(n)] = O(n)$  by substitution. Assume that  $E[T(n)] \leq cn$  for some constant  $c$  that satisfies the initial conditions of the recurrence. We assume that  $T(n) = O(1)$  for  $n$  less than some constant; we shall pick this constant later. We also pick a constant  $a$  such that the function described by the  $O(n)$  term above (which describes the non-recursive component of the running time of the algorithm) is bounded from above by  $an$  for all  $n > 0$ . Using this inductive hypothesis, we have

$$\begin{aligned}
E[T(n)] & \leq \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} ck + an \\
& = \frac{2c}{n} \left( \sum_{k=1}^{n-1} k - \sum_{k=1}^{\lceil n/2 \rceil - 1} k \right) + an
\end{aligned}$$

$$\begin{aligned}
&= \frac{2c}{n} \left( \frac{(n-1)n}{2} - \frac{(\lfloor n/2 \rfloor - 1) \lfloor n/2 \rfloor}{2} \right) + an \\
&\leq \frac{2c}{n} \left( \frac{(n-1)n}{2} - \frac{(n/2 - 2)(n/2 - 1)}{2} \right) + an \\
&= \frac{2c}{n} \left( \frac{n^2 - n}{2} - \frac{n^2/4 - 3n/2 + 2}{2} \right) + an \\
&= \frac{c}{n} \left( \frac{3n^2}{4} + \frac{n}{2} - 2 \right) + an \\
&= c \left( \frac{3n}{4} + \frac{1}{2} - \frac{2}{n} \right) + an \\
&\leq \frac{3cn}{4} + \frac{c}{2} + an \\
&= cn - \left( \frac{cn}{4} - \frac{c}{2} - an \right) .
\end{aligned}$$

In order to complete the proof, we need to show that for sufficiently large  $n$ , this last expression is at most  $cn$  or, equivalently, that  $cn/4 - c/2 - an \geq 0$ . If we add  $c/2$  to both sides and factor out  $n$ , we get  $n(c/4 - a) \geq c/2$ . As long as we choose the constant  $c$  so that  $c/4 - a > 0$ , i.e.,  $c > 4a$ , we can divide both sides by  $c/4 - a$ , giving

$$n \geq \frac{c/2}{c/4 - a} = \frac{2c}{c - 4a} .$$

Thus, if we assume that  $T(n) = O(1)$  for  $n < 2c/(c - 4a)$ , then  $E[T(n)] = O(n)$ . We conclude that we can find any order statistic, and in particular the median, in expected linear time, assuming that the elements are distinct.

## Exercises

### 9.2-1

Show that RANDOMIZED-SELECT never makes a recursive call to a 0-length array.

### 9.2-2

Argue that the indicator random variable  $X_k$  and the value  $T(\max(k - 1, n - k))$  are independent.

### 9.2-3

Write an iterative version of RANDOMIZED-SELECT.

**9.2-4**

Suppose we use RANDOMIZED-SELECT to select the minimum element of the array  $A = \langle 3, 2, 9, 0, 7, 5, 4, 8, 6, 1 \rangle$ . Describe a sequence of partitions that results in a worst-case performance of RANDOMIZED-SELECT.

---

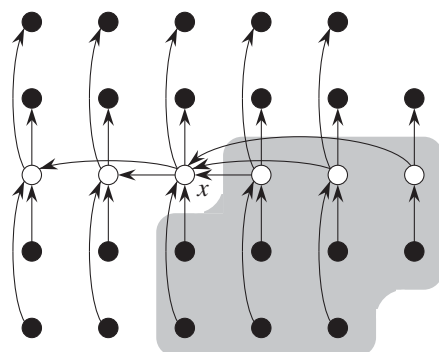
**9.3 Selection in worst-case linear time**

We now examine a selection algorithm whose running time is  $O(n)$  in the worst case. Like RANDOMIZED-SELECT, the algorithm SELECT finds the desired element by recursively partitioning the input array. Here, however, we *guarantee* a good split upon partitioning the array. SELECT uses the deterministic partitioning algorithm PARTITION from quicksort (see Section 7.1), but modified to take the element to partition around as an input parameter.

The SELECT algorithm determines the  $i$ th smallest of an input array of  $n > 1$  distinct elements by executing the following steps. (If  $n = 1$ , then SELECT merely returns its only input value as the  $i$ th smallest.)

1. Divide the  $n$  elements of the input array into  $\lfloor n/5 \rfloor$  groups of 5 elements each and at most one group made up of the remaining  $n \bmod 5$  elements.
2. Find the median of each of the  $\lfloor n/5 \rfloor$  groups by first insertion-sorting the elements of each group (of which there are at most 5) and then picking the median from the sorted list of group elements.
3. Use SELECT recursively to find the median  $x$  of the  $\lfloor n/5 \rfloor$  medians found in step 2. (If there are an even number of medians, then by our convention,  $x$  is the lower median.)
4. Partition the input array around the median-of-medians  $x$  using the modified version of PARTITION. Let  $k$  be one more than the number of elements on the low side of the partition, so that  $x$  is the  $k$ th smallest element and there are  $n - k$  elements on the high side of the partition.
5. If  $i = k$ , then return  $x$ . Otherwise, use SELECT recursively to find the  $i$ th smallest element on the low side if  $i < k$ , or the  $(i - k)$ th smallest element on the high side if  $i > k$ .

To analyze the running time of SELECT, we first determine a lower bound on the number of elements that are greater than the partitioning element  $x$ . Figure 9.1 helps us to visualize this bookkeeping. At least half of the medians found in



**Figure 9.1** Analysis of the algorithm SELECT. The  $n$  elements are represented by small circles, and each group of 5 elements occupies a column. The medians of the groups are whitened, and the median-of-medians  $x$  is labeled. (When finding the median of an even number of elements, we use the lower median.) Arrows go from larger elements to smaller, from which we can see that 3 out of every full group of 5 elements to the right of  $x$  are greater than  $x$ , and 3 out of every group of 5 elements to the left of  $x$  are less than  $x$ . The elements known to be greater than  $x$  appear on a shaded background.

step 2 are greater than or equal to the median-of-medians  $x$ .<sup>1</sup> Thus, at least half of the  $\lceil n/5 \rceil$  groups contribute at least 3 elements that are greater than  $x$ , except for the one group that has fewer than 5 elements if 5 does not divide  $n$  exactly, and the one group containing  $x$  itself. Discounting these two groups, it follows that the number of elements greater than  $x$  is at least

$$3 \left( \left\lceil \frac{1}{2} \left\lceil \frac{n}{5} \right\rceil \right\rceil - 2 \right) \geq \frac{3n}{10} - 6.$$

Similarly, at least  $3n/10 - 6$  elements are less than  $x$ . Thus, in the worst case, step 5 calls SELECT recursively on at most  $7n/10 + 6$  elements.

We can now develop a recurrence for the worst-case running time  $T(n)$  of the algorithm SELECT. Steps 1, 2, and 4 take  $O(n)$  time. (Step 2 consists of  $O(n)$  calls of insertion sort on sets of size  $O(1)$ .) Step 3 takes time  $T(\lceil n/5 \rceil)$ , and step 5 takes time at most  $T(7n/10 + 6)$ , assuming that  $T$  is monotonically increasing. We make the assumption, which seems unmotivated at first, that any input of fewer than 140 elements requires  $O(1)$  time; the origin of the magic constant 140 will be clear shortly. We can therefore obtain the recurrence

<sup>1</sup>Because of our assumption that the numbers are distinct, all medians except  $x$  are either greater than or less than  $x$ .

$$T(n) \leq \begin{cases} O(1) & \text{if } n < 140, \\ T(\lceil n/5 \rceil) + T(7n/10 + 6) + O(n) & \text{if } n \geq 140. \end{cases}$$

We show that the running time is linear by substitution. More specifically, we will show that  $T(n) \leq cn$  for some suitably large constant  $c$  and all  $n > 0$ . We begin by assuming that  $T(n) \leq cn$  for some suitably large constant  $c$  and all  $n < 140$ ; this assumption holds if  $c$  is large enough. We also pick a constant  $a$  such that the function described by the  $O(n)$  term above (which describes the non-recursive component of the running time of the algorithm) is bounded above by  $an$  for all  $n > 0$ . Substituting this inductive hypothesis into the right-hand side of the recurrence yields

$$\begin{aligned} T(n) &\leq c \lceil n/5 \rceil + c(7n/10 + 6) + an \\ &\leq cn/5 + c + 7cn/10 + 6c + an \\ &= 9cn/10 + 7c + an \\ &= cn + (-cn/10 + 7c + an), \end{aligned}$$

which is at most  $cn$  if

$$-cn/10 + 7c + an \leq 0. \tag{9.2}$$

Inequality (9.2) is equivalent to the inequality  $c \geq 10a(n/(n-70))$  when  $n > 70$ . Because we assume that  $n \geq 140$ , we have  $n/(n-70) \leq 2$ , and so choosing  $c \geq 20a$  will satisfy inequality (9.2). (Note that there is nothing special about the constant 140; we could replace it by any integer strictly greater than 70 and then choose  $c$  accordingly.) The worst-case running time of SELECT is therefore linear.

As in a comparison sort (see Section 8.1), SELECT and RANDOMIZED-SELECT determine information about the relative order of elements only by comparing elements. Recall from Chapter 8 that sorting requires  $\Omega(n \lg n)$  time in the comparison model, even on average (see Problem 8-1). The linear-time sorting algorithms in Chapter 8 make assumptions about the input. In contrast, the linear-time selection algorithms in this chapter do not require any assumptions about the input. They are not subject to the  $\Omega(n \lg n)$  lower bound because they manage to solve the selection problem without sorting. Thus, solving the selection problem by sorting and indexing, as presented in the introduction to this chapter, is asymptotically inefficient.

## Exercises

### 9.3-1

In the algorithm SELECT, the input elements are divided into groups of 5. Will the algorithm work in linear time if they are divided into groups of 7? Argue that SELECT does not run in linear time if groups of 3 are used.

### 9.3-2

Analyze SELECT to show that if  $n \geq 140$ , then at least  $\lceil n/4 \rceil$  elements are greater than the median-of-medians  $x$  and at least  $\lceil n/4 \rceil$  elements are less than  $x$ .

### 9.3-3

Show how quicksort can be made to run in  $O(n \lg n)$  time in the worst case, assuming that all elements are distinct.

### 9.3-4 ★

Suppose that an algorithm uses only comparisons to find the  $i$ th smallest element in a set of  $n$  elements. Show that it can also find the  $i - 1$  smaller elements and the  $n - i$  larger elements without performing any additional comparisons.

### 9.3-5

Suppose that you have a “black-box” worst-case linear-time median subroutine. Give a simple, linear-time algorithm that solves the selection problem for an arbitrary order statistic.

### 9.3-6

The  $k$ th *quantiles* of an  $n$ -element set are the  $k - 1$  order statistics that divide the sorted set into  $k$  equal-sized sets (to within 1). Give an  $O(n \lg k)$ -time algorithm to list the  $k$ th quantiles of a set.

### 9.3-7

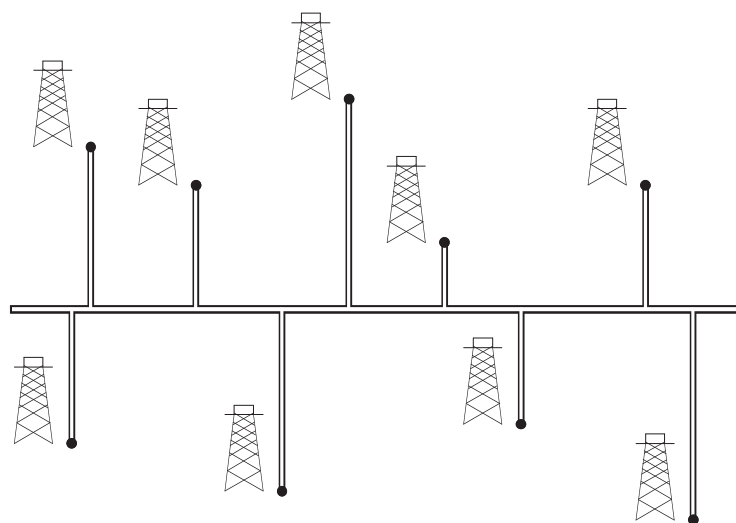
Describe an  $O(n)$ -time algorithm that, given a set  $S$  of  $n$  distinct numbers and a positive integer  $k \leq n$ , determines the  $k$  numbers in  $S$  that are closest to the median of  $S$ .

### 9.3-8

Let  $X[1..n]$  and  $Y[1..n]$  be two arrays, each containing  $n$  numbers already in sorted order. Give an  $O(\lg n)$ -time algorithm to find the median of all  $2n$  elements in arrays  $X$  and  $Y$ .

### 9.3-9

Professor Olay is consulting for an oil company, which is planning a large pipeline running east to west through an oil field of  $n$  wells. The company wants to connect



**Figure 9.2** Professor Olay needs to determine the position of the east-west oil pipeline that minimizes the total length of the north-south spurs.

a spur pipeline from each well directly to the main pipeline along a shortest route (either north or south), as shown in Figure 9.2. Given the  $x$ - and  $y$ -coordinates of the wells, how should the professor pick the optimal location of the main pipeline, which would be the one that minimizes the total length of the spurs? Show how to determine the optimal location in linear time.

---

## Problems

### 9-1 Largest $i$ numbers in sorted order

Given a set of  $n$  numbers, we wish to find the  $i$  largest in sorted order using a comparison-based algorithm. Find the algorithm that implements each of the following methods with the best asymptotic worst-case running time, and analyze the running times of the algorithms in terms of  $n$  and  $i$ .

- Sort the numbers, and list the  $i$  largest.
- Build a max-priority queue from the numbers, and call EXTRACT-MAX  $i$  times.
- Use an order-statistic algorithm to find the  $i$ th largest number, partition around that number, and sort the  $i$  largest numbers.

### 9-2 Weighted median

For  $n$  distinct elements  $x_1, x_2, \dots, x_n$  with positive weights  $w_1, w_2, \dots, w_n$  such that  $\sum_{i=1}^n w_i = 1$ , the **weighted (lower) median** is the element  $x_k$  satisfying

$$\sum_{x_i < x_k} w_i < \frac{1}{2}$$

and

$$\sum_{x_i > x_k} w_i \leq \frac{1}{2}.$$

For example, if the elements are 0.1, 0.35, 0.05, 0.1, 0.15, 0.05, 0.2 and each element equals its weight (that is,  $w_i = x_i$  for  $i = 1, 2, \dots, 7$ ), then the median is 0.1, but the weighted median is 0.2.

- a. Argue that the median of  $x_1, x_2, \dots, x_n$  is the weighted median of the  $x_i$  with weights  $w_i = 1/n$  for  $i = 1, 2, \dots, n$ .
- b. Show how to compute the weighted median of  $n$  elements in  $O(n \lg n)$  worst-case time using sorting.
- c. Show how to compute the weighted median in  $\Theta(n)$  worst-case time using a linear-time median algorithm such as SELECT from Section 9.3.

The **post-office location problem** is defined as follows. We are given  $n$  points  $p_1, p_2, \dots, p_n$  with associated weights  $w_1, w_2, \dots, w_n$ . We wish to find a point  $p$  (not necessarily one of the input points) that minimizes the sum  $\sum_{i=1}^n w_i d(p, p_i)$ , where  $d(a, b)$  is the distance between points  $a$  and  $b$ .

- d. Argue that the weighted median is a best solution for the 1-dimensional post-office location problem, in which points are simply real numbers and the distance between points  $a$  and  $b$  is  $d(a, b) = |a - b|$ .
- e. Find the best solution for the 2-dimensional post-office location problem, in which the points are  $(x, y)$  coordinate pairs and the distance between points  $a = (x_1, y_1)$  and  $b = (x_2, y_2)$  is the **Manhattan distance** given by  $d(a, b) = |x_1 - x_2| + |y_1 - y_2|$ .

### 9-3 Small order statistics

We showed that the worst-case number  $T(n)$  of comparisons used by SELECT to select the  $i$ th order statistic from  $n$  numbers satisfies  $T(n) = \Theta(n)$ , but the constant hidden by the  $\Theta$ -notation is rather large. When  $i$  is small relative to  $n$ , we can implement a different procedure that uses SELECT as a subroutine but makes fewer comparisons in the worst case.



- a. Describe an algorithm that uses  $U_i(n)$  comparisons to find the  $i$ th smallest of  $n$  elements, where

$$U_i(n) = \begin{cases} T(n) & \text{if } i \geq n/2, \\ \lfloor n/2 \rfloor + U_i(\lceil n/2 \rceil) + T(2i) & \text{otherwise.} \end{cases}$$

(Hint: Begin with  $\lfloor n/2 \rfloor$  disjoint pairwise comparisons, and recurse on the set containing the smaller element from each pair.)

- b. Show that, if  $i < n/2$ , then  $U_i(n) = n + O(T(2i) \lg(n/i))$ .
- c. Show that if  $i$  is a constant less than  $n/2$ , then  $U_i(n) = n + O(\lg n)$ .
- d. Show that if  $i = n/k$  for  $k \geq 2$ , then  $U_i(n) = n + O(T(2n/k) \lg k)$ .

#### 9-4 Alternative analysis of randomized selection

In this problem, we use indicator random variables to analyze the RANDOMIZED-SELECT procedure in a manner akin to our analysis of RANDOMIZED-QUICKSORT in Section 7.4.2.

As in the quicksort analysis, we assume that all elements are distinct, and we rename the elements of the input array  $A$  as  $z_1, z_2, \dots, z_n$ , where  $z_i$  is the  $i$ th smallest element. Thus, the call RANDOMIZED-SELECT( $A, 1, n, k$ ) returns  $z_k$ .

For  $1 \leq i < j \leq n$ , let

$$X_{ijk} = \mathbf{I}\{z_i \text{ is compared with } z_j \text{ sometime during the execution of the algorithm to find } z_k\}.$$

- a. Give an exact expression for  $E[X_{ijk}]$ . (Hint: Your expression may have different values, depending on the values of  $i, j$ , and  $k$ .)
- b. Let  $X_k$  denote the total number of comparisons between elements of array  $A$  when finding  $z_k$ . Show that

$$E[X_k] \leq 2 \left( \sum_{i=1}^k \sum_{j=k}^n \frac{1}{j-i+1} + \sum_{j=k+1}^n \frac{j-k-1}{j-k+1} + \sum_{i=1}^{k-2} \frac{k-i-1}{k-i+1} \right).$$

- c. Show that  $E[X_k] \leq 4n$ .
- d. Conclude that, assuming all elements of array  $A$  are distinct, RANDOMIZED-SELECT runs in expected time  $O(n)$ .

---

**Chapter notes**

The worst-case linear-time median-finding algorithm was devised by Blum, Floyd, Pratt, Rivest, and Tarjan [50]. The fast randomized version is due to Hoare [169]. Floyd and Rivest [108] have developed an improved randomized version that partitions around an element recursively selected from a small sample of the elements.

It is still unknown exactly how many comparisons are needed to determine the median. Bent and John [41] gave a lower bound of  $2n$  comparisons for median finding, and Schönhage, Paterson, and Pippenger [302] gave an upper bound of  $3n$ . Dor and Zwick have improved on both of these bounds. Their upper bound [93] is slightly less than  $2.95n$ , and their lower bound [94] is  $(2 + \epsilon)n$ , for a small positive constant  $\epsilon$ , thereby improving slightly on related work by Dor et al. [92]. Paterson [272] describes some of these results along with other related work.



---

### *III Data Structures*

---

## Introduction

Sets are as fundamental to computer science as they are to mathematics. Whereas mathematical sets are unchanging, the sets manipulated by algorithms can grow, shrink, or otherwise change over time. We call such sets *dynamic*. The next five chapters present some basic techniques for representing finite dynamic sets and manipulating them on a computer.

Algorithms may require several different types of operations to be performed on sets. For example, many algorithms need only the ability to insert elements into, delete elements from, and test membership in a set. We call a dynamic set that supports these operations a *dictionary*. Other algorithms require more complicated operations. For example, min-priority queues, which Chapter 6 introduced in the context of the heap data structure, support the operations of inserting an element into and extracting the smallest element from a set. The best way to implement a dynamic set depends upon the operations that must be supported.

### Elements of a dynamic set

In a typical implementation of a dynamic set, each element is represented by an object whose attributes can be examined and manipulated if we have a pointer to the object. (Section 10.3 discusses the implementation of objects and pointers in programming environments that do not contain them as basic data types.) Some kinds of dynamic sets assume that one of the object's attributes is an identifying *key*. If the keys are all different, we can think of the dynamic set as being a set of key values. The object may contain *satellite data*, which are carried around in other object attributes but are otherwise unused by the set implementation. It may

also have attributes that are manipulated by the set operations; these attributes may contain data or pointers to other objects in the set.

Some dynamic sets presuppose that the keys are drawn from a totally ordered set, such as the real numbers, or the set of all words under the usual alphabetic ordering. A total ordering allows us to define the minimum element of the set, for example, or to speak of the next element larger than a given element in a set.

### Operations on dynamic sets

Operations on a dynamic set can be grouped into two categories: *queries*, which simply return information about the set, and *modifying operations*, which change the set. Here is a list of typical operations. Any specific application will usually require only a few of these to be implemented.

#### SEARCH( $S, k$ )

A query that, given a set  $S$  and a key value  $k$ , returns a pointer  $x$  to an element in  $S$  such that  $x.key = k$ , or NIL if no such element belongs to  $S$ .

#### INSERT( $S, x$ )

A modifying operation that augments the set  $S$  with the element pointed to by  $x$ . We usually assume that any attributes in element  $x$  needed by the set implementation have already been initialized.

#### DELETE( $S, x$ )

A modifying operation that, given a pointer  $x$  to an element in the set  $S$ , removes  $x$  from  $S$ . (Note that this operation takes a pointer to an element  $x$ , not a key value.)

#### MINIMUM( $S$ )

A query on a totally ordered set  $S$  that returns a pointer to the element of  $S$  with the smallest key.

#### MAXIMUM( $S$ )

A query on a totally ordered set  $S$  that returns a pointer to the element of  $S$  with the largest key.

#### SUCCESSOR( $S, x$ )

A query that, given an element  $x$  whose key is from a totally ordered set  $S$ , returns a pointer to the next larger element in  $S$ , or NIL if  $x$  is the maximum element.

#### PREDECESSOR( $S, x$ )

A query that, given an element  $x$  whose key is from a totally ordered set  $S$ , returns a pointer to the next smaller element in  $S$ , or NIL if  $x$  is the minimum element.

In some situations, we can extend the queries `SUCCESSOR` and `PREDECESSOR` so that they apply to sets with nondistinct keys. For a set on  $n$  keys, the normal presumption is that a call to `MINIMUM` followed by  $n - 1$  calls to `SUCCESSOR` enumerates the elements in the set in sorted order.

We usually measure the time taken to execute a set operation in terms of the size of the set. For example, Chapter 13 describes a data structure that can support any of the operations listed above on a set of size  $n$  in time  $O(\lg n)$ .

### Overview of Part III

Chapters 10–14 describe several data structures that we can use to implement dynamic sets; we shall use many of these later to construct efficient algorithms for a variety of problems. We already saw another important data structure—the heap—in Chapter 6.

Chapter 10 presents the essentials of working with simple data structures such as stacks, queues, linked lists, and rooted trees. It also shows how to implement objects and pointers in programming environments that do not support them as primitives. If you have taken an introductory programming course, then much of this material should be familiar to you.

Chapter 11 introduces hash tables, which support the dictionary operations `INSERT`, `DELETE`, and `SEARCH`. In the worst case, hashing requires  $\Theta(n)$  time to perform a `SEARCH` operation, but the expected time for hash-table operations is  $O(1)$ . The analysis of hashing relies on probability, but most of the chapter requires no background in the subject.

Binary search trees, which are covered in Chapter 12, support all the dynamic-set operations listed above. In the worst case, each operation takes  $\Theta(n)$  time on a tree with  $n$  elements, but on a randomly built binary search tree, the expected time for each operation is  $O(\lg n)$ . Binary search trees serve as the basis for many other data structures.

Chapter 13 introduces red-black trees, which are a variant of binary search trees. Unlike ordinary binary search trees, red-black trees are guaranteed to perform well: operations take  $O(\lg n)$  time in the worst case. A red-black tree is a balanced search tree; Chapter 18 in Part V presents another kind of balanced search tree, called a B-tree. Although the mechanics of red-black trees are somewhat intricate, you can glean most of their properties from the chapter without studying the mechanics in detail. Nevertheless, you probably will find walking through the code to be quite instructive.

In Chapter 14, we show how to augment red-black trees to support operations other than the basic ones listed above. First, we augment them so that we can dynamically maintain order statistics for a set of keys. Then, we augment them in a different way to maintain intervals of real numbers.

---

## 10 Elementary Data Structures

In this chapter, we examine the representation of dynamic sets by simple data structures that use pointers. Although we can construct many complex data structures using pointers, we present only the rudimentary ones: stacks, queues, linked lists, and rooted trees. We also show ways to synthesize objects and pointers from arrays.

---

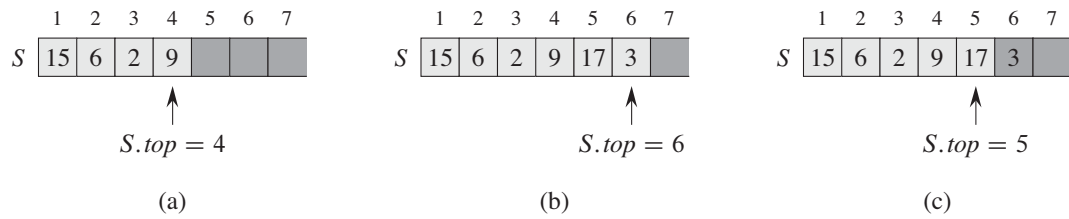
### 10.1 Stacks and queues

Stacks and queues are dynamic sets in which the element removed from the set by the DELETE operation is prespecified. In a *stack*, the element deleted from the set is the one most recently inserted: the stack implements a *last-in, first-out*, or *LIFO*, policy. Similarly, in a *queue*, the element deleted is always the one that has been in the set for the longest time: the queue implements a *first-in, first-out*, or *FIFO*, policy. There are several efficient ways to implement stacks and queues on a computer. In this section we show how to use a simple array to implement each.

#### Stacks

The INSERT operation on a stack is often called PUSH, and the DELETE operation, which does not take an element argument, is often called POP. These names are allusions to physical stacks, such as the spring-loaded stacks of plates used in cafeterias. The order in which plates are popped from the stack is the reverse of the order in which they were pushed onto the stack, since only the top plate is accessible.

As Figure 10.1 shows, we can implement a stack of at most  $n$  elements with an array  $S[1..n]$ . The array has an attribute  $S.top$  that indexes the most recently



**Figure 10.1** An array implementation of a stack  $S$ . Stack elements appear only in the lightly shaded positions. **(a)** Stack  $S$  has 4 elements. The top element is 9. **(b)** Stack  $S$  after the calls  $\text{PUSH}(S, 17)$  and  $\text{PUSH}(S, 3)$ . **(c)** Stack  $S$  after the call  $\text{POP}(S)$  has returned the element 3, which is the one most recently pushed. Although element 3 still appears in the array, it is no longer in the stack; the top is element 17.

inserted element. The stack consists of elements  $S[1 \dots S.top]$ , where  $S[1]$  is the element at the bottom of the stack and  $S[S.top]$  is the element at the top.

When  $S.top = 0$ , the stack contains no elements and is *empty*. We can test to see whether the stack is empty by query operation  $\text{STACK-EMPTY}$ . If we attempt to pop an empty stack, we say the stack *underflows*, which is normally an error. If  $S.top$  exceeds  $n$ , the stack *overflows*. (In our pseudocode implementation, we don't worry about stack overflow.)

We can implement each of the stack operations with just a few lines of code:

$\text{STACK-EMPTY}(S)$

```

1  if  $S.top == 0$ 
2      return TRUE
3  else return FALSE

```

$\text{PUSH}(S, x)$

```

1   $S.top = S.top + 1$ 
2   $S[S.top] = x$ 

```

$\text{POP}(S)$

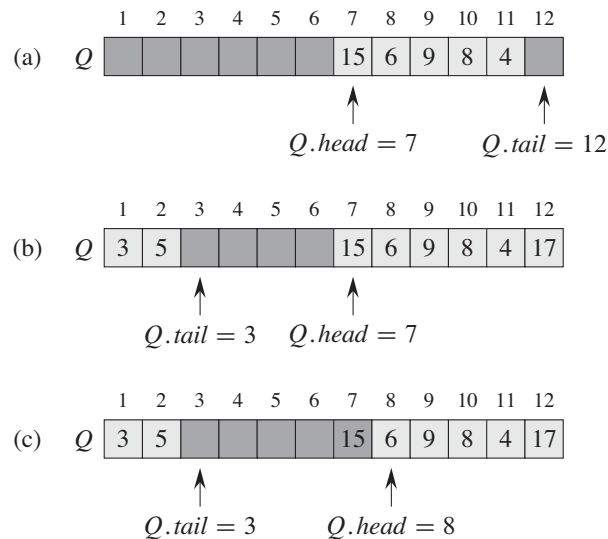
```

1  if  $\text{STACK-EMPTY}(S)$ 
2      error "underflow"
3  else  $S.top = S.top - 1$ 
4      return  $S[S.top + 1]$ 

```

Figure 10.1 shows the effects of the modifying operations  $\text{PUSH}$  and  $\text{POP}$ . Each of the three stack operations takes  $O(1)$  time.





**Figure 10.2** A queue implemented using an array  $Q[1..12]$ . Queue elements appear only in the lightly shaded positions. (a) The queue has 5 elements, in locations  $Q[7..11]$ . (b) The configuration of the queue after the calls  $ENQUEUE(Q, 17)$ ,  $ENQUEUE(Q, 3)$ , and  $ENQUEUE(Q, 5)$ . (c) The configuration of the queue after the call  $DEQUEUE(Q)$  returns the key value 15 formerly at the head of the queue. The new head has key 6.

## Queues

We call the INSERT operation on a queue **ENQUEUE**, and we call the DELETE operation **DEQUEUE**; like the stack operation **POP**, **DEQUEUE** takes no element argument. The FIFO property of a queue causes it to operate like a line of customers waiting to pay a cashier. The queue has a **head** and a **tail**. When an element is enqueued, it takes its place at the tail of the queue, just as a newly arriving customer takes a place at the end of the line. The element dequeued is always the one at the head of the queue, like the customer at the head of the line who has waited the longest.

Figure 10.2 shows one way to implement a queue of at most  $n - 1$  elements using an array  $Q[1..n]$ . The queue has an attribute  $Q.head$  that indexes, or points to, its head. The attribute  $Q.tail$  indexes the next location at which a newly arriving element will be inserted into the queue. The elements in the queue reside in locations  $Q.head, Q.head + 1, \dots, Q.tail - 1$ , where we “wrap around” in the sense that location 1 immediately follows location  $n$  in a circular order. When  $Q.head = Q.tail$ , the queue is empty. Initially, we have  $Q.head = Q.tail = 1$ . If we attempt to dequeue an element from an empty queue, the queue underflows.

When  $Q.head = Q.tail + 1$ , the queue is full, and if we attempt to enqueue an element, then the queue overflows.

In our procedures ENQUEUE and DEQUEUE, we have omitted the error checking for underflow and overflow. (Exercise 10.1-4 asks you to supply code that checks for these two error conditions.) The pseudocode assumes that  $n = Q.length$ .

```

ENQUEUE( $Q, x$ )
1   $Q[Q.tail] = x$ 
2  if  $Q.tail == Q.length$ 
3       $Q.tail = 1$ 
4  else  $Q.tail = Q.tail + 1$ 

DEQUEUE( $Q$ )
1   $x = Q[Q.head]$ 
2  if  $Q.head == Q.length$ 
3       $Q.head = 1$ 
4  else  $Q.head = Q.head + 1$ 
5  return  $x$ 

```

Figure 10.2 shows the effects of the ENQUEUE and DEQUEUE operations. Each operation takes  $O(1)$  time.

## Exercises

### 10.1-1

Using Figure 10.1 as a model, illustrate the result of each operation in the sequence PUSH( $S, 4$ ), PUSH( $S, 1$ ), PUSH( $S, 3$ ), POP( $S$ ), PUSH( $S, 8$ ), and POP( $S$ ) on an initially empty stack  $S$  stored in array  $S[1..6]$ .

### 10.1-2

Explain how to implement two stacks in one array  $A[1..n]$  in such a way that neither stack overflows unless the total number of elements in both stacks together is  $n$ . The PUSH and POP operations should run in  $O(1)$  time.

### 10.1-3

Using Figure 10.2 as a model, illustrate the result of each operation in the sequence ENQUEUE( $Q, 4$ ), ENQUEUE( $Q, 1$ ), ENQUEUE( $Q, 3$ ), DEQUEUE( $Q$ ), ENQUEUE( $Q, 8$ ), and DEQUEUE( $Q$ ) on an initially empty queue  $Q$  stored in array  $Q[1..6]$ .

### 10.1-4

Rewrite ENQUEUE and DEQUEUE to detect underflow and overflow of a queue.

**10.1-5**

Whereas a stack allows insertion and deletion of elements at only one end, and a queue allows insertion at one end and deletion at the other end, a **deque** (double-ended queue) allows insertion and deletion at both ends. Write four  $O(1)$ -time procedures to insert elements into and delete elements from both ends of a deque implemented by an array.

**10.1-6**

Show how to implement a queue using two stacks. Analyze the running time of the queue operations.

**10.1-7**

Show how to implement a stack using two queues. Analyze the running time of the stack operations.

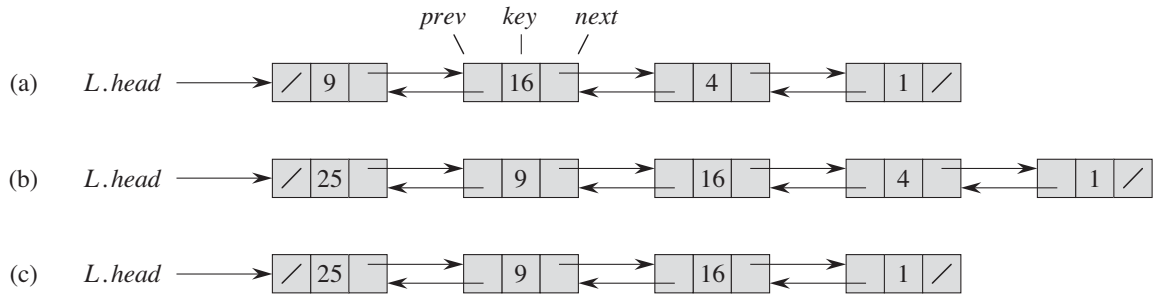
---

**10.2 Linked lists**

A **linked list** is a data structure in which the objects are arranged in a linear order. Unlike an array, however, in which the linear order is determined by the array indices, the order in a linked list is determined by a pointer in each object. Linked lists provide a simple, flexible representation for dynamic sets, supporting (though not necessarily efficiently) all the operations listed on page 230.

As shown in Figure 10.3, each element of a **doubly linked list**  $L$  is an object with an attribute *key* and two other pointer attributes: *next* and *prev*. The object may also contain other satellite data. Given an element  $x$  in the list,  $x.next$  points to its successor in the linked list, and  $x.prev$  points to its predecessor. If  $x.prev = \text{NIL}$ , the element  $x$  has no predecessor and is therefore the first element, or **head**, of the list. If  $x.next = \text{NIL}$ , the element  $x$  has no successor and is therefore the last element, or **tail**, of the list. An attribute  $L.head$  points to the first element of the list. If  $L.head = \text{NIL}$ , the list is empty.

A list may have one of several forms. It may be either singly linked or doubly linked, it may be sorted or not, and it may be circular or not. If a list is **singly linked**, we omit the *prev* pointer in each element. If a list is **sorted**, the linear order of the list corresponds to the linear order of keys stored in elements of the list; the minimum element is then the head of the list, and the maximum element is the tail. If the list is **unsorted**, the elements can appear in any order. In a **circular list**, the *prev* pointer of the head of the list points to the tail, and the *next* pointer of the tail of the list points to the head. We can think of a circular list as a ring of



**Figure 10.3** (a) A doubly linked list  $L$  representing the dynamic set  $\{1, 4, 9, 16\}$ . Each element in the list is an object with attributes for the key and pointers (shown by arrows) to the next and previous objects. The  $next$  attribute of the tail and the  $prev$  attribute of the head are NIL, indicated by a diagonal slash. The attribute  $L.head$  points to the head. (b) Following the execution of  $LIST-INSERT(L, x)$ , where  $x.key = 25$ , the linked list has a new object with key 25 as the new head. This new object points to the old head with key 9. (c) The result of the subsequent call  $LIST-DELETE(L, x)$ , where  $x$  points to the object with key 4.

elements. In the remainder of this section, we assume that the lists with which we are working are unsorted and doubly linked.

### Searching a linked list

The procedure  $LIST-SEARCH(L, k)$  finds the first element with key  $k$  in list  $L$  by a simple linear search, returning a pointer to this element. If no object with key  $k$  appears in the list, then the procedure returns NIL. For the linked list in Figure 10.3(a), the call  $LIST-SEARCH(L, 4)$  returns a pointer to the third element, and the call  $LIST-SEARCH(L, 7)$  returns NIL.

$LIST-SEARCH(L, k)$

```

1   $x = L.head$ 
2  while  $x \neq NIL$  and  $x.key \neq k$ 
3       $x = x.next$ 
4  return  $x$ 
```

To search a list of  $n$  objects, the  $LIST-SEARCH$  procedure takes  $\Theta(n)$  time in the worst case, since it may have to search the entire list.

### Inserting into a linked list

Given an element  $x$  whose  $key$  attribute has already been set, the  $LIST-INSERT$  procedure “splices”  $x$  onto the front of the linked list, as shown in Figure 10.3(b).

LIST-INSERT( $L, x$ )

```

1   $x.next = L.head$ 
2  if  $L.head \neq \text{NIL}$ 
3       $L.head.prev = x$ 
4   $L.head = x$ 
5   $x.prev = \text{NIL}$ 

```

(Recall that our attribute notation can cascade, so that  $L.head.prev$  denotes the *prev* attribute of the object that  $L.head$  points to.) The running time for LIST-INSERT on a list of  $n$  elements is  $O(1)$ .

### Deleting from a linked list

The procedure LIST-DELETE removes an element  $x$  from a linked list  $L$ . It must be given a pointer to  $x$ , and it then “splices”  $x$  out of the list by updating pointers. If we wish to delete an element with a given key, we must first call LIST-SEARCH to retrieve a pointer to the element.

LIST-DELETE( $L, x$ )

```

1  if  $x.prev \neq \text{NIL}$ 
2       $x.prev.next = x.next$ 
3  else  $L.head = x.next$ 
4  if  $x.next \neq \text{NIL}$ 
5       $x.next.prev = x.prev$ 

```

Figure 10.3(c) shows how an element is deleted from a linked list. LIST-DELETE runs in  $O(1)$  time, but if we wish to delete an element with a given key,  $\Theta(n)$  time is required in the worst case because we must first call LIST-SEARCH to find the element.

### Sentinels

The code for LIST-DELETE would be simpler if we could ignore the boundary conditions at the head and tail of the list:

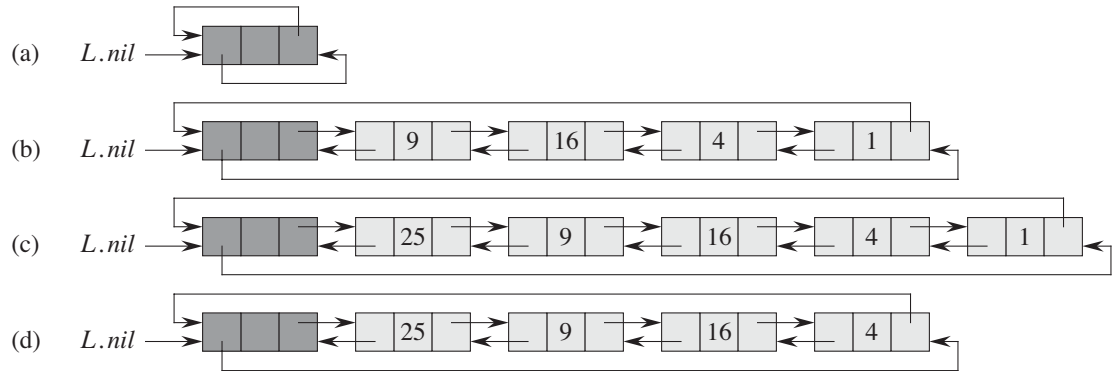
LIST-DELETE'( $L, x$ )

```

1   $x.prev.next = x.next$ 
2   $x.next.prev = x.prev$ 

```

A *sentinel* is a dummy object that allows us to simplify boundary conditions. For example, suppose that we provide with list  $L$  an object  $L.nil$  that represents NIL



**Figure 10.4** A circular, doubly linked list with a sentinel. The sentinel  $L.nil$  appears between the head and tail. The attribute  $L.head$  is no longer needed, since we can access the head of the list by  $L.nil.next$ . (a) An empty list. (b) The linked list from Figure 10.3(a), with key 9 at the head and key 1 at the tail. (c) The list after executing  $LIST-INSERT'(L, x)$ , where  $x.key = 25$ . The new object becomes the head of the list. (d) The list after deleting the object with key 1. The new tail is the object with key 4.

but has all the attributes of the other objects in the list. Wherever we have a reference to NIL in list code, we replace it by a reference to the sentinel  $L.nil$ . As shown in Figure 10.4, this change turns a regular doubly linked list into a **circular, doubly linked list with a sentinel**, in which the sentinel  $L.nil$  lies between the head and tail. The attribute  $L.nil.next$  points to the head of the list, and  $L.nil.prev$  points to the tail. Similarly, both the *next* attribute of the tail and the *prev* attribute of the head point to  $L.nil$ . Since  $L.nil.next$  points to the head, we can eliminate the attribute  $L.head$  altogether, replacing references to it by references to  $L.nil.next$ . Figure 10.4(a) shows that an empty list consists of just the sentinel, and both  $L.nil.next$  and  $L.nil.prev$  point to  $L.nil$ .

The code for  $LIST-SEARCH$  remains the same as before, but with the references to NIL and  $L.head$  changed as specified above:

```
LIST-SEARCH'(L, k)
1   $x = L.nil.next$ 
2  while  $x \neq L.nil$  and  $x.key \neq k$ 
3       $x = x.next$ 
4  return  $x$ 
```

We use the two-line procedure  $LIST-DELETE'$  from before to delete an element from the list. The following procedure inserts an element into the list:

LIST-INSERT'( $L, x$ )

```
1   $x.next = L.nil.next$   
2   $L.nil.next.prev = x$   
3   $L.nil.next = x$   
4   $x.prev = L.nil$ 
```

Figure 10.4 shows the effects of LIST-INSERT' and LIST-DELETE' on a sample list.

Sentinels rarely reduce the asymptotic time bounds of data structure operations, but they can reduce constant factors. The gain from using sentinels within loops is usually a matter of clarity of code rather than speed; the linked list code, for example, becomes simpler when we use sentinels, but we save only  $O(1)$  time in the LIST-INSERT' and LIST-DELETE' procedures. In other situations, however, the use of sentinels helps to tighten the code in a loop, thus reducing the coefficient of, say,  $n$  or  $n^2$  in the running time.

We should use sentinels judiciously. When there are many small lists, the extra storage used by their sentinels can represent significant wasted memory. In this book, we use sentinels only when they truly simplify the code.

## Exercises

### 10.2-1

Can you implement the dynamic-set operation INSERT on a singly linked list in  $O(1)$  time? How about DELETE?

### 10.2-2

Implement a stack using a singly linked list  $L$ . The operations PUSH and POP should still take  $O(1)$  time.

### 10.2-3

Implement a queue by a singly linked list  $L$ . The operations ENQUEUE and DEQUEUE should still take  $O(1)$  time.

### 10.2-4

As written, each loop iteration in the LIST-SEARCH' procedure requires two tests: one for  $x \neq L.nil$  and one for  $x.key \neq k$ . Show how to eliminate the test for  $x \neq L.nil$  in each iteration.

### 10.2-5

Implement the dictionary operations INSERT, DELETE, and SEARCH using singly linked, circular lists. What are the running times of your procedures?

**10.2-6**

The dynamic-set operation UNION takes two disjoint sets  $S_1$  and  $S_2$  as input, and it returns a set  $S = S_1 \cup S_2$  consisting of all the elements of  $S_1$  and  $S_2$ . The sets  $S_1$  and  $S_2$  are usually destroyed by the operation. Show how to support UNION in  $O(1)$  time using a suitable list data structure.

**10.2-7**

Give a  $\Theta(n)$ -time nonrecursive procedure that reverses a singly linked list of  $n$  elements. The procedure should use no more than constant storage beyond that needed for the list itself.

**10.2-8 ★**

Explain how to implement doubly linked lists using only one pointer value  $x.np$  per item instead of the usual two ( $next$  and  $prev$ ). Assume that all pointer values can be interpreted as  $k$ -bit integers, and define  $x.np$  to be  $x.np = x.next \text{ XOR } x.prev$ , the  $k$ -bit “exclusive-or” of  $x.next$  and  $x.prev$ . (The value NIL is represented by 0.) Be sure to describe what information you need to access the head of the list. Show how to implement the SEARCH, INSERT, and DELETE operations on such a list. Also show how to reverse such a list in  $O(1)$  time.

---

## 10.3 Implementing pointers and objects

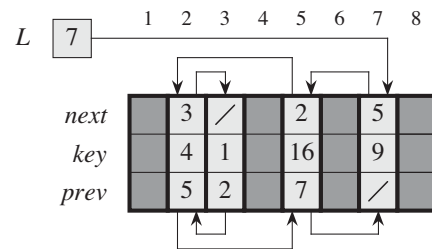
How do we implement pointers and objects in languages that do not provide them? In this section, we shall see two ways of implementing linked data structures without an explicit pointer data type. We shall synthesize objects and pointers from arrays and array indices.

**A multiple-array representation of objects**

We can represent a collection of objects that have the same attributes by using an array for each attribute. As an example, Figure 10.5 shows how we can implement the linked list of Figure 10.3(a) with three arrays. The array *key* holds the values of the keys currently in the dynamic set, and the pointers reside in the arrays *next* and *prev*. For a given array index  $x$ , the array entries  $key[x]$ ,  $next[x]$ , and  $prev[x]$  represent an object in the linked list. Under this interpretation, a pointer  $x$  is simply a common index into the *key*, *next*, and *prev* arrays.

In Figure 10.3(a), the object with key 4 follows the object with key 16 in the linked list. In Figure 10.5, key 4 appears in  $key[2]$ , and key 16 appears in  $key[5]$ , and so  $next[5] = 2$  and  $prev[2] = 5$ . Although the constant NIL appears in the *next*





**Figure 10.5** The linked list of Figure 10.3(a) represented by the arrays *key*, *next*, and *prev*. Each vertical slice of the arrays represents a single object. Stored pointers correspond to the array indices shown at the top; the arrows show how to interpret them. Lightly shaded object positions contain list elements. The variable *L* keeps the index of the head.

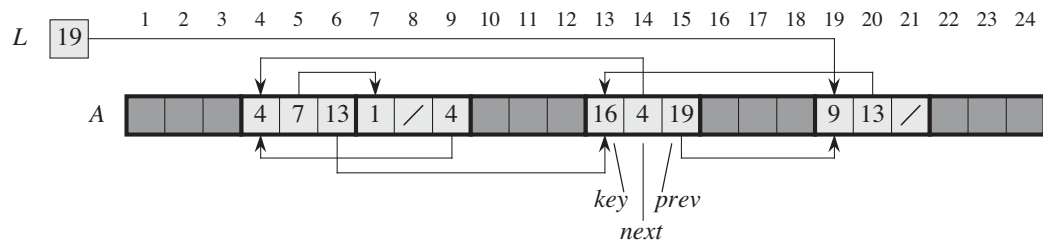
attribute of the tail and the *prev* attribute of the head, we usually use an integer (such as 0 or  $-1$ ) that cannot possibly represent an actual index into the arrays. A variable *L* holds the index of the head of the list.

### A single-array representation of objects

The words in a computer memory are typically addressed by integers from 0 to  $M - 1$ , where  $M$  is a suitably large integer. In many programming languages, an object occupies a contiguous set of locations in the computer memory. A pointer is simply the address of the first memory location of the object, and we can address other memory locations within the object by adding an offset to the pointer.

We can use the same strategy for implementing objects in programming environments that do not provide explicit pointer data types. For example, Figure 10.6 shows how to use a single array *A* to store the linked list from Figures 10.3(a) and 10.5. An object occupies a contiguous subarray  $A[j \dots k]$ . Each attribute of the object corresponds to an offset in the range from 0 to  $k - j$ , and a pointer to the object is the index *j*. In Figure 10.6, the offsets corresponding to *key*, *next*, and *prev* are 0, 1, and 2, respectively. To read the value of  $i.\text{prev}$ , given a pointer *i*, we add the value *i* of the pointer to the offset 2, thus reading  $A[i + 2]$ .

The single-array representation is flexible in that it permits objects of different lengths to be stored in the same array. The problem of managing such a heterogeneous collection of objects is more difficult than the problem of managing a homogeneous collection, where all objects have the same attributes. Since most of the data structures we shall consider are composed of homogeneous elements, it will be sufficient for our purposes to use the multiple-array representation of objects.



**Figure 10.6** The linked list of Figures 10.3(a) and 10.5 represented in a single array  $A$ . Each list element is an object that occupies a contiguous subarray of length 3 within the array. The three attributes *key*, *next*, and *prev* correspond to the offsets 0, 1, and 2, respectively, within each object. A pointer to an object is the index of the first element of the object. Objects containing list elements are lightly shaded, and arrows show the list ordering.

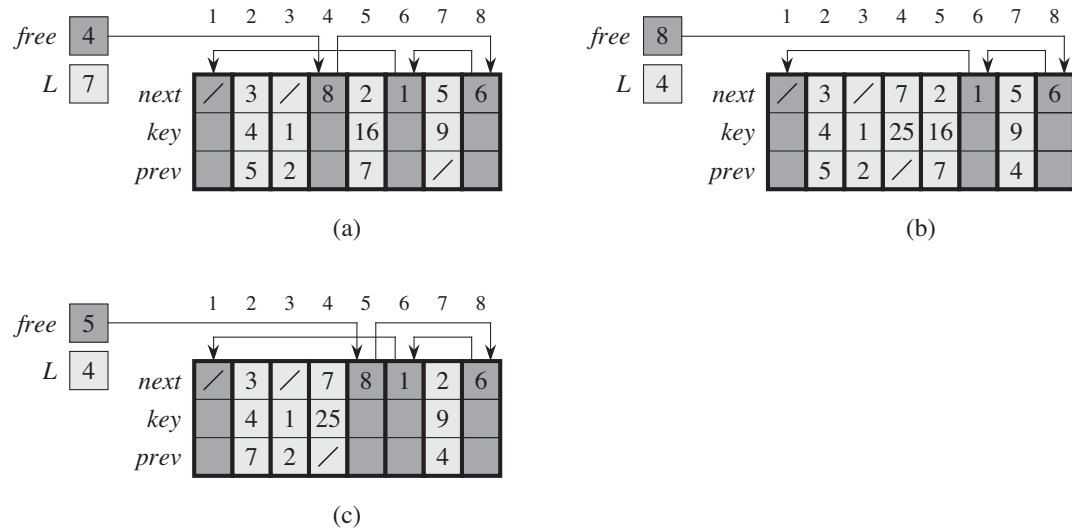
## Allocating and freeing objects

To insert a key into a dynamic set represented by a doubly linked list, we must allocate a pointer to a currently unused object in the linked-list representation. Thus, it is useful to manage the storage of objects not currently used in the linked-list representation so that one can be allocated. In some systems, a ***garbage collector*** is responsible for determining which objects are unused. Many applications, however, are simple enough that they can bear responsibility for returning an unused object to a storage manager. We shall now explore the problem of allocating and freeing (or deallocating) homogeneous objects using the example of a doubly linked list represented by multiple arrays.

Suppose that the arrays in the multiple-array representation have length  $m$  and that at some moment the dynamic set contains  $n \leq m$  elements. Then  $n$  objects represent elements currently in the dynamic set, and the remaining  $m - n$  objects are **free**; the free objects are available to represent elements inserted into the dynamic set in the future.

We keep the free objects in a singly linked list, which we call the *free list*. The free list uses only the *next* array, which stores the *next* pointers within the list. The head of the free list is held in the global variable *free*. When the dynamic set represented by linked list *L* is nonempty, the free list may be intertwined with list *L*, as shown in Figure 10.7. Note that each object in the representation is either in list *L* or in the free list, but not in both.

The free list acts like a stack: the next object allocated is the last one freed. We can use a list implementation of the stack operations PUSH and POP to implement the procedures for allocating and freeing objects, respectively. We assume that the global variable *free* used in the following procedures points to the first element of the free list.



**Figure 10.7** The effect of the `ALLOCATE-OBJECT` and `FREE-OBJECT` procedures. (a) The list of Figure 10.5 (lightly shaded) and a free list (heavily shaded). Arrows show the free-list structure. (b) The result of calling `ALLOCATE-OBJECT()` (which returns index 4), setting `key[4]` to 25, and calling `LIST-INSERT(L, 4)`. The new free-list head is object 8, which had been `next[4]` on the free list. (c) After executing `LIST-DELETE(L, 5)`, we call `FREE-OBJECT(5)`. Object 5 becomes the new free-list head, with object 8 following it on the free list.

#### `ALLOCATE-OBJECT()`

```

1  if free == NIL
2      error "out of space"
3  else x = free
4      free = x.next
5      return x

```

#### `FREE-OBJECT(x)`

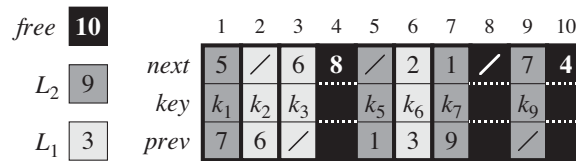
```

1  x.next = free
2  free = x

```

The free list initially contains all  $n$  unallocated objects. Once the free list has been exhausted, running the `ALLOCATE-OBJECT` procedure signals an error. We can even service several linked lists with just a single free list. Figure 10.8 shows two linked lists and a free list intertwined through `key`, `next`, and `prev` arrays.

The two procedures run in  $O(1)$  time, which makes them quite practical. We can modify them to work for any homogeneous collection of objects by letting any one of the attributes in the object act like a `next` attribute in the free list.



**Figure 10.8** Two linked lists,  $L_1$  (lightly shaded) and  $L_2$  (heavily shaded), and a free list (darkened) intertwined.

## Exercises

### 10.3-1

Draw a picture of the sequence  $\langle 13, 4, 8, 19, 5, 11 \rangle$  stored as a doubly linked list using the multiple-array representation. Do the same for the single-array representation.

### 10.3-2

Write the procedures `ALLOCATE-OBJECT` and `FREE-OBJECT` for a homogeneous collection of objects implemented by the single-array representation.

### 10.3-3

Why don't we need to set or reset the *prev* attributes of objects in the implementation of the `ALLOCATE-OBJECT` and `FREE-OBJECT` procedures?

### 10.3-4

It is often desirable to keep all elements of a doubly linked list compact in storage, using, for example, the first  $m$  index locations in the multiple-array representation. (This is the case in a paged, virtual-memory computing environment.) Explain how to implement the procedures `ALLOCATE-OBJECT` and `FREE-OBJECT` so that the representation is compact. Assume that there are no pointers to elements of the linked list outside the list itself. (*Hint*: Use the array implementation of a stack.)

### 10.3-5

Let  $L$  be a doubly linked list of length  $n$  stored in arrays *key*, *prev*, and *next* of length  $m$ . Suppose that these arrays are managed by `ALLOCATE-OBJECT` and `FREE-OBJECT` procedures that keep a doubly linked free list  $F$ . Suppose further that of the  $m$  items, exactly  $n$  are on list  $L$  and  $m - n$  are on the free list. Write a procedure `COMPACTIFY-LIST( $L, F$ )` that, given the list  $L$  and the free list  $F$ , moves the items in  $L$  so that they occupy array positions  $1, 2, \dots, n$  and adjusts the free list  $F$  so that it remains correct, occupying array positions  $n + 1, n + 2, \dots, m$ . The running time of your procedure should be  $\Theta(n)$ , and it should use only a constant amount of extra space. Argue that your procedure is correct.

## 10.4 Representing rooted trees

The methods for representing lists given in the previous section extend to any homogeneous data structure. In this section, we look specifically at the problem of representing rooted trees by linked data structures. We first look at binary trees, and then we present a method for rooted trees in which nodes can have an arbitrary number of children.

We represent each node of a tree by an object. As with linked lists, we assume that each node contains a *key* attribute. The remaining attributes of interest are pointers to other nodes, and they vary according to the type of tree.

### Binary trees

Figure 10.9 shows how we use the attributes *p*, *left*, and *right* to store pointers to the parent, left child, and right child of each node in a binary tree *T*. If  $x.p = \text{NIL}$ , then *x* is the root. If node *x* has no left child, then  $x.\text{left} = \text{NIL}$ , and similarly for the right child. The root of the entire tree *T* is pointed to by the attribute *T.root*. If  $T.\text{root} = \text{NIL}$ , then the tree is empty.

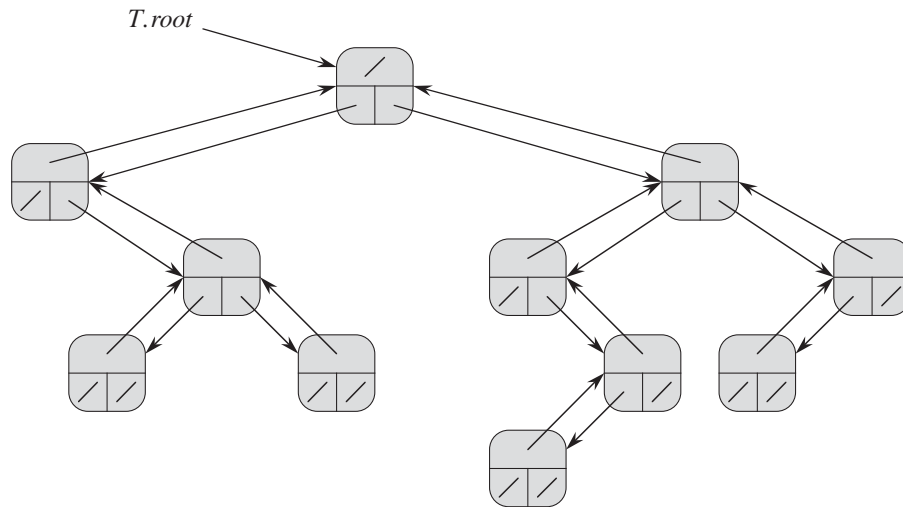
### Rooted trees with unbounded branching

We can extend the scheme for representing a binary tree to any class of trees in which the number of children of each node is at most some constant *k*: we replace the *left* and *right* attributes by  $\text{child}_1, \text{child}_2, \dots, \text{child}_k$ . This scheme no longer works when the number of children of a node is unbounded, since we do not know how many attributes (arrays in the multiple-array representation) to allocate in advance. Moreover, even if the number of children *k* is bounded by a large constant but most nodes have a small number of children, we may waste a lot of memory.

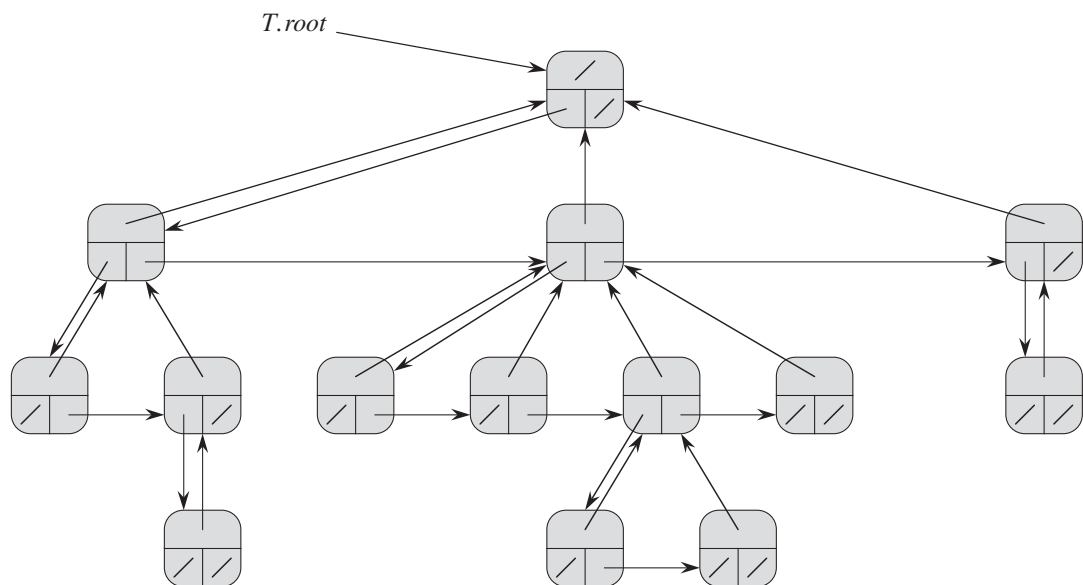
Fortunately, there is a clever scheme to represent trees with arbitrary numbers of children. It has the advantage of using only  $O(n)$  space for any *n*-node rooted tree. The **left-child, right-sibling representation** appears in Figure 10.10. As before, each node contains a parent pointer *p*, and *T.root* points to the root of tree *T*. Instead of having a pointer to each of its children, however, each node *x* has only two pointers:

1. *x.left-child* points to the leftmost child of node *x*, and
2. *x.right-sibling* points to the sibling of *x* immediately to its right.

If node *x* has no children, then  $x.\text{left-child} = \text{NIL}$ , and if node *x* is the rightmost child of its parent, then  $x.\text{right-sibling} = \text{NIL}$ .



**Figure 10.9** The representation of a binary tree  $T$ . Each node  $x$  has the attributes  $x.p$  (top),  $x.left$  (lower left), and  $x.right$  (lower right). The key attributes are not shown.



**Figure 10.10** The left-child, right-sibling representation of a tree  $T$ . Each node  $x$  has attributes  $x.p$  (top),  $x.left-child$  (lower left), and  $x.right-sibling$  (lower right). The key attributes are not shown.

### Other tree representations

We sometimes represent rooted trees in other ways. In Chapter 6, for example, we represented a heap, which is based on a complete binary tree, by a single array plus the index of the last node in the heap. The trees that appear in Chapter 21 are traversed only toward the root, and so only the parent pointers are present; there are no pointers to children. Many other schemes are possible. Which scheme is best depends on the application.

### Exercises

#### 10.4-1

Draw the binary tree rooted at index 6 that is represented by the following attributes:

index	key	left	right
1	12	7	3
2	15	8	NIL
3	4	10	NIL
4	10	5	9
5	2	NIL	NIL
6	18	1	4
7	7	NIL	NIL
8	14	6	2
9	21	NIL	NIL
10	5	NIL	NIL

#### 10.4-2

Write an  $O(n)$ -time recursive procedure that, given an  $n$ -node binary tree, prints out the key of each node in the tree.

#### 10.4-3

Write an  $O(n)$ -time nonrecursive procedure that, given an  $n$ -node binary tree, prints out the key of each node in the tree. Use a stack as an auxiliary data structure.

#### 10.4-4

Write an  $O(n)$ -time procedure that prints all the keys of an arbitrary rooted tree with  $n$  nodes, where the tree is stored using the left-child, right-sibling representation.

#### 10.4-5 ★

Write an  $O(n)$ -time nonrecursive procedure that, given an  $n$ -node binary tree, prints out the key of each node. Use no more than constant extra space outside

of the tree itself and do not modify the tree, even temporarily, during the procedure.

#### 10.4-6 ★

The left-child, right-sibling representation of an arbitrary rooted tree uses three pointers in each node: *left-child*, *right-sibling*, and *parent*. From any node, its parent can be reached and identified in constant time and all its children can be reached and identified in time linear in the number of children. Show how to use only two pointers and one boolean value in each node so that the parent of a node or all of its children can be reached and identified in time linear in the number of children.

---

## Problems

### 10-1 Comparisons among lists

For each of the four types of lists in the following table, what is the asymptotic worst-case running time for each dynamic-set operation listed?

	unsorted, singly linked	sorted, singly linked	unsorted, doubly linked	sorted, doubly linked
SEARCH( $L, k$ )				
INSERT( $L, x$ )				
DELETE( $L, x$ )				
SUCCESSOR( $L, x$ )				
PREDECESSOR( $L, x$ )				
MINIMUM( $L$ )				
MAXIMUM( $L$ )				



### 10-2 Mergeable heaps using linked lists

A *mergeable heap* supports the following operations: MAKE-HEAP (which creates an empty mergeable heap), INSERT, MINIMUM, EXTRACT-MIN, and UNION.<sup>1</sup> Show how to implement mergeable heaps using linked lists in each of the following cases. Try to make each operation as efficient as possible. Analyze the running time of each operation in terms of the size of the dynamic set(s) being operated on.

- a. Lists are sorted.
- b. Lists are unsorted.
- c. Lists are unsorted, and dynamic sets to be merged are disjoint.

### 10-3 Searching a sorted compact list

Exercise 10.3-4 asked how we might maintain an  $n$ -element list compactly in the first  $n$  positions of an array. We shall assume that all keys are distinct and that the compact list is also sorted, that is,  $key[i] < key[next[i]]$  for all  $i = 1, 2, \dots, n$  such that  $next[i] \neq \text{NIL}$ . We will also assume that we have a variable  $L$  that contains the index of the first element on the list. Under these assumptions, you will show that we can use the following randomized algorithm to search the list in  $O(\sqrt{n})$  expected time.

```

COMPACT-LIST-SEARCH( $L, n, k$ )
1   $i = L$ 
2  while  $i \neq \text{NIL}$  and  $key[i] < k$ 
3       $j = \text{RANDOM}(1, n)$ 
4      if  $key[i] < key[j]$  and  $key[j] \leq k$ 
5           $i = j$ 
6          if  $key[i] == k$ 
7              return  $i$ 
8       $i = next[i]$ 
9  if  $i == \text{NIL}$  or  $key[i] > k$ 
10     return NIL
11 else return  $i$ 

```

If we ignore lines 3–7 of the procedure, we have an ordinary algorithm for searching a sorted linked list, in which index  $i$  points to each position of the list in

---

<sup>1</sup>Because we have defined a mergeable heap to support MINIMUM and EXTRACT-MIN, 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*.

turn. The search terminates once the index  $i$  “falls off” the end of the list or once  $\text{key}[i] \geq k$ . In the latter case, if  $\text{key}[i] = k$ , clearly we have found a key with the value  $k$ . If, however,  $\text{key}[i] > k$ , then we will never find a key with the value  $k$ , and so terminating the search was the right thing to do.

Lines 3–7 attempt to skip ahead to a randomly chosen position  $j$ . Such a skip benefits us if  $\text{key}[j]$  is larger than  $\text{key}[i]$  and no larger than  $k$ ; in such a case,  $j$  marks a position in the list that  $i$  would have to reach during an ordinary list search. Because the list is compact, we know that any choice of  $j$  between 1 and  $n$  indexes some object in the list rather than a slot on the free list.

Instead of analyzing the performance of COMPACT-LIST-SEARCH directly, we shall analyze a related algorithm, COMPACT-LIST-SEARCH', which executes two separate loops. This algorithm takes an additional parameter  $t$  which determines an upper bound on the number of iterations of the first loop.

COMPACT-LIST-SEARCH'( $L, n, k, t$ )

```

1   $i = L$ 
2  for  $q = 1$  to  $t$ 
3       $j = \text{RANDOM}(1, n)$ 
4      if  $\text{key}[i] < \text{key}[j]$  and  $\text{key}[j] \leq k$ 
5           $i = j$ 
6          if  $\text{key}[i] == k$ 
7              return  $i$ 
8  while  $i \neq \text{NIL}$  and  $\text{key}[i] < k$ 
9       $i = \text{next}[i]$ 
10 if  $i == \text{NIL}$  or  $\text{key}[i] > k$ 
11     return  $\text{NIL}$ 
12 else return  $i$ 
```

To compare the execution of the algorithms COMPACT-LIST-SEARCH( $L, n, k$ ) and COMPACT-LIST-SEARCH'( $L, n, k, t$ ), assume that the sequence of integers returned by the calls of RANDOM( $1, n$ ) is the same for both algorithms.

- a. Suppose that COMPACT-LIST-SEARCH( $L, n, k$ ) takes  $t$  iterations of the **while** loop of lines 2–8. Argue that COMPACT-LIST-SEARCH'( $L, n, k, t$ ) returns the same answer and that the total number of iterations of both the **for** and **while** loops within COMPACT-LIST-SEARCH' is at least  $t$ .

In the call COMPACT-LIST-SEARCH'( $L, n, k, t$ ), let  $X_t$  be the random variable that describes the distance in the linked list (that is, through the chain of *next* pointers) from position  $i$  to the desired key  $k$  after  $t$  iterations of the **for** loop of lines 2–7 have occurred.

- b.* Argue that the expected running time of  $\text{COMPACT-LIST-SEARCH}'(L, n, k, t)$  is  $O(t + E[X_t])$ .
- c.* Show that  $E[X_t] \leq \sum_{r=1}^n (1 - r/n)^t$ . (*Hint:* Use equation (C.25).)
- d.* Show that  $\sum_{r=0}^{n-1} r^t \leq n^{t+1}/(t+1)$ .
- e.* Prove that  $E[X_t] \leq n/(t+1)$ .
- f.* Show that  $\text{COMPACT-LIST-SEARCH}'(L, n, k, t)$  runs in  $O(t + n/t)$  expected time.
- g.* Conclude that  $\text{COMPACT-LIST-SEARCH}$  runs in  $O(\sqrt{n})$  expected time.
- h.* Why do we assume that all keys are distinct in  $\text{COMPACT-LIST-SEARCH}$ ? Argue that random skips do not necessarily help asymptotically when the list contains repeated key values.

---

## Chapter notes

Aho, Hopcroft, and Ullman [6] and Knuth [209] are excellent references for elementary data structures. Many other texts cover both basic data structures and their implementation in a particular programming language. Examples of these types of textbooks include Goodrich and Tamassia [147], Main [241], Shaffer [311], and Weiss [352, 353, 354]. Gonnet [145] provides experimental data on the performance of many data-structure operations.

The origin of stacks and queues as data structures in computer science is unclear, since corresponding notions already existed in mathematics and paper-based business practices before the introduction of digital computers. Knuth [209] cites A. M. Turing for the development of stacks for subroutine linkage in 1947.

Pointer-based data structures also seem to be a folk invention. According to Knuth, pointers were apparently used in early computers with drum memories. The A-1 language developed by G. M. Hopper in 1951 represented algebraic formulas as binary trees. Knuth credits the IPL-II language, developed in 1956 by A. Newell, J. C. Shaw, and H. A. Simon, for recognizing the importance and promoting the use of pointers. Their IPL-III language, developed in 1957, included explicit stack operations.