Javed Aslam
11/18/97

# Average Case Analysis of Quicksort ①

1.  $T(n) = \dfrac{2}{n} \displaystyle\sum_{k=1}^{n-1} T(k) + n$

| First rank | (small) left problem | (large) right problem |
|---|---|---|
| $n$ | $n-1$ | $0$ |
| $n-1$ | $n-2$ | $1$ |
| $n-2$ | $n-3$ | $2$ |
| $\vdots$ | $0$ | $n-2$ |
| | | $n-1$ |

Consider $T(n-1)$ as well:

So, $T(n) = \dfrac{1}{n} \displaystyle\sum_{j=0}^{n-1} \left( T(j) + T(n-j-1) \right) + n$

$= \dfrac{2}{n} \displaystyle\sum_{k=0}^{n-1} T(k) + n$

$= \dfrac{2}{n} \displaystyle\sum_{k=1}^{n-1} T(k) + n$

(no "recursive" call to 0-sized subproblem)

2.  $T(n-1) = \dfrac{2}{n-1} \displaystyle\sum_{k=1}^{n-2} T(k) + n-1$

Multiplying (1) by $n$ and (2) by $n-1$:

$n\,T(n) = 2 \displaystyle\sum_{k=1}^{n-1} T(k) + n^2$

$(n-1)\,T(n-1) = 2 \displaystyle\sum_{k=1}^{n-2} T(k) + (n-1)^2$

---

$n\,T(n) - (n-1)\,T(n-1) = 2\,T(n-1) + 2n-1$

$\Rightarrow \quad n\,T(n) = (n+1)\,T(n-1) + 2n-1$

$\Rightarrow \quad T(n) = \dfrac{n+1}{n}\,T(n-1) + 2 - 1/n$

Ignoring the $1/n$ term, we have

$$T(n) \le \dfrac{n+1}{n}\,T(n-1) + 2$$

Solve by iteration:

$$T(n) \leq 2 + \frac{n+1}{n} T(n-1)$$

$$\leq 2 + \frac{n+1}{n}\left(2 + \frac{n}{n-1} T(n-2)\right)$$

$$= 2 + 2\frac{n+1}{n} + \frac{n+1}{n-1} T(n-2)$$

$$\leq 2 + 2\frac{n+1}{n} + \frac{n+1}{n-1}\left(2 + \frac{n-1}{n-2} T(n-3)\right)$$

$$= 2 + 2\frac{n+1}{n} + 2\frac{n+1}{n-1} + \frac{n+1}{n-2} T(n-3)$$

$$= 2 \cdot \frac{n+1}{n+1} + 2\frac{n+1}{n} + 2\frac{n+1}{n-1} + \frac{n+1}{n-2} T(n-3)$$

$$\vdots$$

$$= 2(n+1) \sum_{i=0}^{k-1} \frac{1}{(n+1)-i} + \frac{n+1}{n-(k-1)} T(n-k)$$

$\boxed{\begin{array}{c} k = n-1 \\ T(n-k) = T(1) \end{array}}$

$$= 2(n+1) \sum_{i=0}^{n-2} \frac{1}{(n+1)-i} + \frac{n+1}{2} T(1)$$

$$= 2(n+1) \sum_{j=3}^{n+1} \frac{1}{j} + \frac{n+1}{2} \Theta(1)$$

$$= 2(n+1) \Theta(\ln n) + \Theta(n)$$

$$= \Theta(n \log n)$$

Alternate page ② :

$$T(n) \leq \frac{n+1}{n} T(n-1) + 2$$

Divide by $n+1$ :

$$\frac{T(n)}{n+1} \leq \frac{T(n-1)}{n} + \frac{2}{n+1}$$

Let $R(n) = \frac{T(n)}{n+1}$ (and thus, $R(n-1) = \frac{T(n-1)}{n}$ )

$\Rightarrow$
$$R(n) \leq R(n-1) + \frac{2}{n+1}$$

Note that $R(1) = \frac{T(1)}{1+1} = \Theta(1)$. Now, simple iteration :

$$R(n) \leq \frac{2}{n+1} + \frac{2}{n} + \frac{2}{n-1} + \cdots + \frac{2}{3} + R(1)$$

$$= 2 \sum_{k=3}^{n+1} \frac{1}{k} + \Theta(1)$$

$$= \Theta(\ln n)$$

But $T(n) = (n+1) R(n)$

$$= (n+1) \cdot \Theta(\ln n)$$

$$= \Theta(n \ln n)$$

In the following discussion, our expression for the running time of INSERTION-SORT will evolve from a messy formula that uses all the statement costs $c_i$ to a much simpler notation that is more concise and more easily manipulated. This simpler notation will also make it easy to determine whether one algorithm is more efficient than another.

We start by presenting the INSERTION-SORT procedure with the time "cost" of each statement and the number of times each statement is executed. For each $j = 2, 3, \ldots, n$, where $n = length[A]$, we let $t_j$ be the number of times the **while** loop test in line 5 is executed for that value of $j$. When **a for or while** loop exits in the usual way (i.e., due to the test in the loop header), the test is executed one time more than the loop body. We assume that comments are not executable statements, and so they take no time.

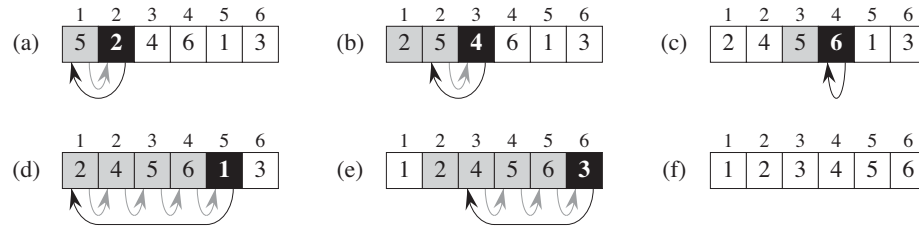| INSERTION-SORT$(A)$ | | *cost* | *times* |
|---|---|---|---|
| 1 | **for** $j \leftarrow 2$ **to** $length[A]$ | $c_1$ | $n$ |
| 2 | **do** $key \leftarrow A[j]$ | $c_2$ | $n - 1$ |
| 3 | $\triangleright$ Insert $A[j]$ into the sorted sequence $A[1 .. j - 1]$. | $0$ | $n - 1$ |
| 4 | $i \leftarrow j - 1$ | $c_4$ | $n - 1$ |
| 5 | **while** $i > 0$ and $A[i] > key$ | $c_5$ | $\sum_{j=2}^{n} t_j$ |
| 6 | **do** $A[i + 1] \leftarrow A[i]$ | $c_6$ | $\sum_{j=2}^{n} (t_j - 1)$ |
| 7 | $i \leftarrow i - 1$ | $c_7$ | $\sum_{j=2}^{n} (t_j - 1)$ |
| 8 | $A[i + 1] \leftarrow key$ | $c_8$ | $n - 1$ |

The running time of the algorithm is the sum of running times for each statement executed; a statement that takes $c_i$ steps to execute and is executed $n$ times will contribute $c_i n$ to the total running time.[5] To compute $T(n)$, the running time of INSERTION-SORT, we sum the products of the *cost* and *times* columns, obtaining

$$T(n) = c_1 n + c_2(n - 1) + c_4(n - 1) + c_5 \sum_{j=2}^{n} t_j + c_6 \sum_{j=2}^{n} (t_j - 1)$$
$$+ c_7 \sum_{j=2}^{n} (t_j - 1) + c_8(n - 1).$$

Even for inputs of a given size, an algorithm's running time may depend on *which* input of that size is given. For example, in INSERTION-SORT, the best

---

[5]This characteristic does not necessarily hold for a resource such as memory. A statement that references $m$ words of memory and is executed $n$ times does not necessarily consume $mn$ words of memory in total.

**Figure 2.2** The operation of INSERTION-SORT on the array $A = \langle 5, 2, 4, 6, 1, 3 \rangle$. Array indices appear above the rectangles, and values stored in the array positions appear within the rectangles. **(a)–(e)** The iterations of the **for** loop of lines 1–8. In each iteration, the black rectangle holds the key taken from $A[j]$, which is compared with the values in shaded rectangles to its left in the test of line 5. Shaded arrows show array values moved one position to the right in line 6, and black arrows indicate where the key moves to in line 8. **(f)** The final sorted array.

INSERTION-SORT$(A)$

```
1   for j = 2 to A.length
2       key = A[j]
3       // Insert A[j] into the sorted sequence A[1 .. j − 1].
4       i = j − 1
5       while i > 0 and A[i] > key
6           A[i + 1] = A[i]
7           i = i − 1
8       A[i + 1] = key
```

**Loop invariants and the correctness of insertion sort**

Figure 2.2 shows how this algorithm works for $A = \langle 5, 2, 4, 6, 1, 3 \rangle$. The index $j$ indicates the "current card" being inserted into the hand. At the beginning of each iteration of the **for** loop, which is indexed by $j$, the subarray consisting of elements $A[1 .. j − 1]$ constitutes the currently sorted hand, and the remaining subarray $A[j + 1 .. n]$ corresponds to the pile of cards still on the table. In fact, elements $A[1 .. j − 1]$ are the elements *originally* in positions 1 through $j − 1$, but now in sorted order. We state these properties of $A[1 .. j − 1]$ formally as a *loop invariant*:

> At the start of each iteration of the **for** loop of lines 1–8, the subarray $A[1 .. j − 1]$ consists of the elements originally in $A[1 .. j − 1]$, but in sorted order.

We use loop invariants to help us understand why an algorithm is correct. We must show three things about a loop invariant:

motif, suppose we have two piles of cards face up on a table. Each pile is sorted, with the smallest cards on top. We wish to merge the two piles into a single sorted output pile, which is to be face down on the table. Our basic step consists of choosing the smaller of the two cards on top of the face-up piles, removing it from its pile (which exposes a new top card), and placing this card face down onto the output pile. We repeat this step until one input pile is empty, at which time we just take the remaining input pile and place it face down onto the output pile. Computationally, each basic step takes constant time, since we are checking just two top cards. Since we perform at most $n$ basic steps, merging takes $\Theta(n)$ time.
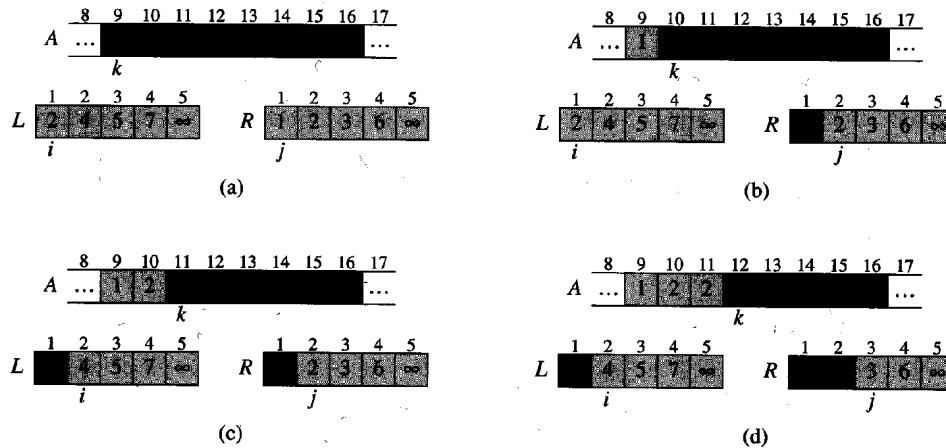
The following pseudocode implements the above idea, but with an additional twist that avoids having to check whether either pile is empty in each basic step. The idea is to put on the bottom of each pile a *sentinel* card, which contains a special value that we use to simplify our code. Here, we use $\infty$ as the sentinel value, so that whenever a card with $\infty$ is exposed, it cannot be the smaller card unless both piles have their sentinel cards exposed. But once that happens, all the nonsentinel cards have already been placed onto the output pile. Since we know in advance that exactly $r - p + 1$ cards will be placed onto the output pile, we can stop once we have performed that many basic steps.

MERGE($A, p, q, r$)
```
 1   n₁ ← q − p + 1
 2   n₂ ← r − q
 3   create arrays L[1 .. n₁ + 1] and R[1 .. n₂ + 1]
 4   for i ← 1 to n₁
 5       do L[i] ← A[p + i − 1]
 6   for j ← 1 to n₂
 7       do R[j] ← A[q + j]
 8   L[n₁ + 1] ← ∞
 9   R[n₂ + 1] ← ∞
10   i ← 1
11   j ← 1
12   for k ← p to r
13       do if L[i] ≤ R[j]
14           then A[k] ← L[i]
15                i ← i + 1
16           else A[k] ← R[j]
17                j ← j + 1
```

In detail, the MERGE procedure works as follows. Line 1 computes the length $n_1$ of the subarray $A[p..q]$, and line 2 computes the length $n_2$ of the subarray $A[q + 1..r]$. We create arrays $L$ and $R$ ("left" and "right"), of lengths $n_1 + 1$ and $n_2 + 1$, respectively, in line 3. The **for** loop of lines 4–5 copies the subar-

**Figure 2.3** The operation of lines 10–17 in the call MERGE($A$, 9, 12, 16), when the subarray $A[9..16]$ contains the sequence ⟨2, 4, 5, 7, 1, 2, 3, 6⟩. After copying and inserting sentinels, the array $L$ contains ⟨2, 4, 5, 7, ∞⟩, and the array $R$ contains ⟨1, 2, 3, 6, ∞⟩. Lightly shaded positions in $A$ contain their final values, and lightly shaded positions in $L$ and $R$ contain values that have yet to be copied back into $A$. Taken together, the lightly shaded positions always comprise the values originally in $A[9..16]$, along with the two sentinels. Heavily shaded positions in $A$ contain values that will be copied over, and heavily shaded positions in $L$ and $R$ contain values that have already been copied back into $A$. (a)–(h) The arrays $A$, $L$, and $R$, and their respective indices $k$, $i$, and $j$ prior to each iteration of the loop of lines 12–17. (i) The arrays and indices at termination. At this point, the subarray in $A[9..16]$ is sorted, and the two sentinels in $L$ and $R$ are the only two elements in these arrays that have not been copied into $A$.

ray $A[p..q]$ into $L[1..n_1]$, and the **for** loop of lines 6–7 copies the subarray $A[q+1..r]$ into $R[1..n_2]$. Lines 8–9 put the sentinels at the ends of the arrays $L$ and $R$. Lines 10–17, illustrated in Figure 2.3, perform the $r - p + 1$ basic steps by maintaining the following loop invariant:

> At the start of each iteration of the **for** loop of lines 12–17, the subarray $A[p..k-1]$ contains the $k - p$ smallest elements of $L[1..n_1+1]$ and $R[1..n_2+1]$, in sorted order. Moreover, $L[i]$ and $R[j]$ are the smallest elements of their arrays that have not been copied back into $A$.

We must show that this loop invariant holds prior to the first iteration of the **for** loop of lines 12–17, 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, we have $k = p$, so that the subarray $A[p..k-1]$ is empty. This empty subarray contains the $k - p = 0$

lines 4–7 take $\Theta(n_1 + n_2) = \Theta(n)$ time,[6] and there are $n$ iterations of the **for** loop of lines 12–17, each of which takes constant time.

We can now use the MERGE procedure as a subroutine in the merge sort algorithm. The procedure MERGE-SORT$(A, p, r)$ sorts the elements in the subarray $A[p..r]$. If $p \geq r$, the subarray has at most one element and is therefore already sorted. Otherwise, the divide step simply computes an index $q$ that partitions $A[p..r]$ into two subarrays: $A[p..q]$, containing $\lceil n/2 \rceil$ elements, and $A[q+1..r]$, containing $\lfloor n/2 \rfloor$ elements.[7]

MERGE-SORT$(A, p, r)$

1  **if** $p < r$
2      **then** $q \leftarrow \lfloor (p+r)/2 \rfloor$
3          MERGE-SORT$(A, p, q)$
4          MERGE-SORT$(A, q+1, r)$
5          MERGE$(A, p, q, r)$

To sort the entire sequence $A = \langle A[1], A[2], \ldots, A[n] \rangle$, we make the initial call MERGE-SORT$(A, 1, length[A])$, where once again $length[A] = n$. Figure 2.4 illustrates the operation of the procedure bottom-up when $n$ is a power of 2. The algorithm consists of merging pairs of 1-item sequences to form sorted sequences of length 2, merging pairs of sequences of length 2 to form sorted sequences of length 4, and so on, until two sequences of length $n/2$ are merged to form the final sorted sequence of length $n$.
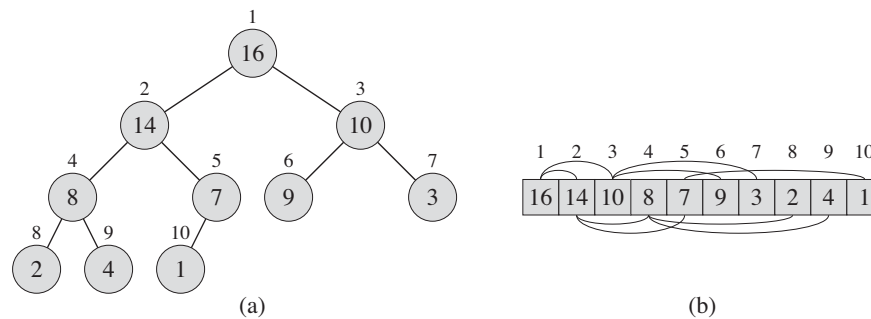
### 2.3.2  Analyzing divide-and-conquer algorithms

When an algorithm contains a recursive call to itself, its running time can often be described by a *recurrence equation* or *recurrence*, which describes the overall running time on a problem of size $n$ in terms of the running time on smaller inputs. We can then use mathematical tools to solve the recurrence and provide bounds on the performance of the algorithm.

A recurrence for the running time of a divide-and-conquer algorithm is based on the three steps of the basic paradigm. As before, we let $T(n)$ be the running time on a problem of size $n$. If the problem size is small enough, say $n \leq c$

---

[6]We shall see in Chapter 3 how to formally interpret equations containing $\Theta$-notation.

[7]The expression $\lceil x \rceil$ denotes the least integer greater than or equal to $x$, and $\lfloor x \rfloor$ denotes the greatest integer less than or equal to $x$. These notations are defined in Chapter 3. The easiest way to verify that setting $q$ to $\lfloor (p+r)/2 \rfloor$ yields subarrays $A[p..q]$ and $A[q+1..r]$ of sizes $\lceil n/2 \rceil$ and $\lfloor n/2 \rfloor$, respectively, is to examine the four cases that arise depending on whether each of $p$ and $r$ is odd or even.

**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 $\lfloor n/2 \rfloor + 1, \lfloor n/2 \rfloor + 2, \ldots, n$.
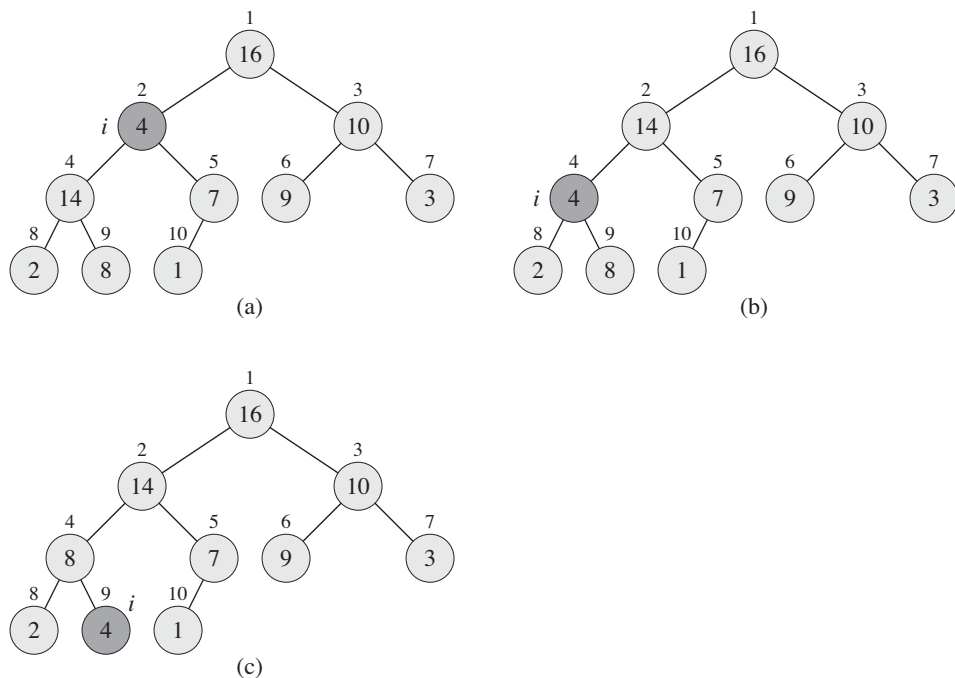
## 6.2   Maintaining the heap property

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 LEFT($i$) and 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 = LEFT(i)
 2  r = RIGHT(i)
 3  if l ≤ A.heap-size and A[l] > A[i]
 4      largest = l
 5  else largest = i
 6  if r ≤ A.heap-size and A[r] > A[largest]
 7      largest = r
 8  if largest ≠ i
 9      exchange A[i] with A[largest]
10      MAX-HEAPIFY(A, 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 *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 $A[largest]$, which causes node $i$ and its

**Figure 6.2** The action of MAX-HEAPIFY($A, 2$), where $A.heap\text{-}size = 10$. **(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 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 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 MAX-HEAPIFY recursively on that subtree.

The running time of 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 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 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 $i > A.heap\text{-}size/2$?

#### 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 $A[1 \mathinner{.\,.} n]$, where $n = A.length$, into a max-heap. By Exercise 6.1-7, the elements in the subarray $A[(\lfloor n/2 \rfloor + 1) \mathinner{.\,.} n]$ 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    $A.heap\text{-}size = A.length$
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 $i + 1$, $i + 2, \ldots, n$ 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 $\lfloor n/2 \rfloor + 1, \lfloor n/2 \rfloor + 2, \ldots, n$ is a leaf and is thus the root of a trivial max-heap.
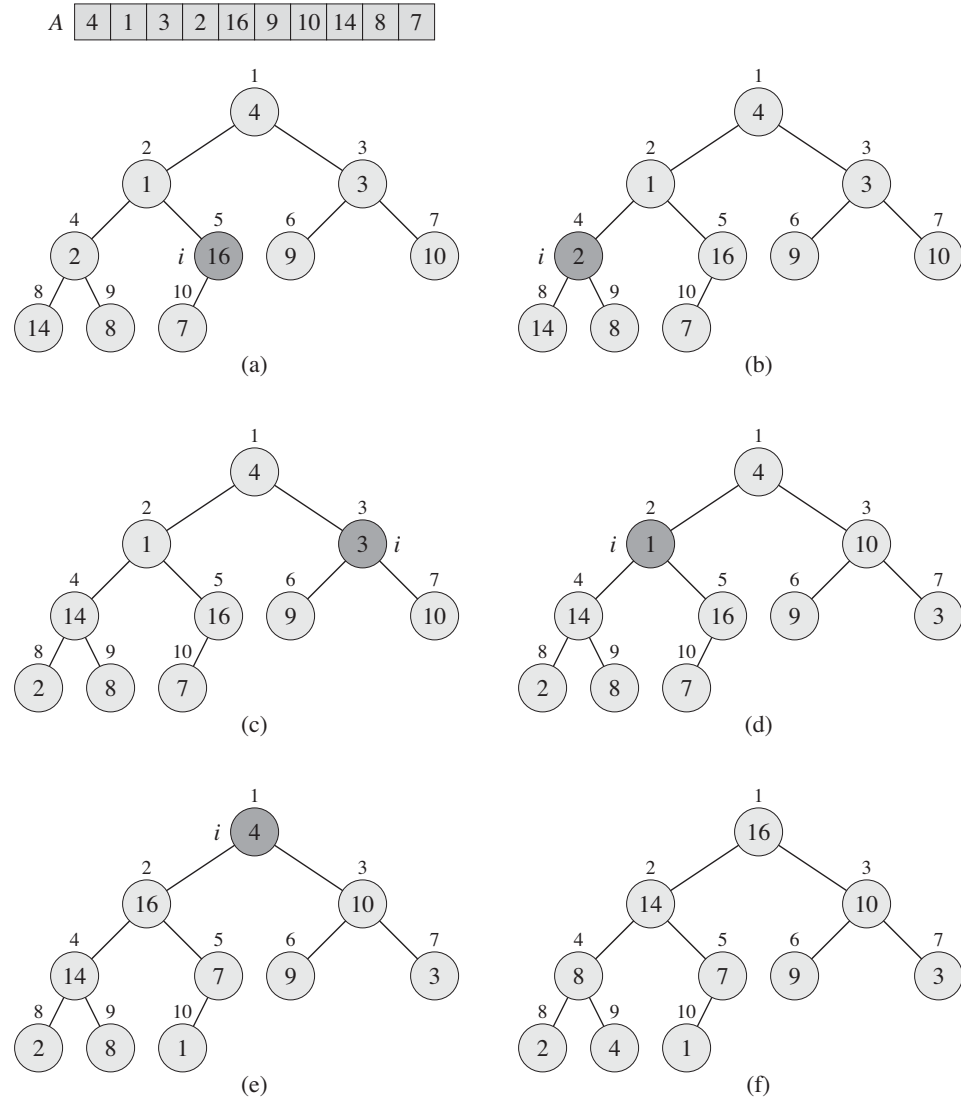
**Maintenance:** To see that each iteration maintains the loop invariant, observe that the children of node $i$ are numbered higher than $i$. By the loop invariant, therefore, they are both roots of max-heaps. This is precisely the condition required for the call MAX-HEAPIFY$(A, i)$ to make node $i$ a max-heap root. Moreover, the MAX-HEAPIFY call preserves the property that nodes $i + 1, i + 2, \ldots, n$ 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 $1, 2, \ldots, n$ 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 $\lfloor \lg n \rfloor$ (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 evalaute the last summation by substituting $x = 1/2$ in the formula (A.8), yielding

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

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

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

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 $A.heap\text{-}size$—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 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.length downto 2
3       exchange A[1] with A[i]
4       A.heap-size = A.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***
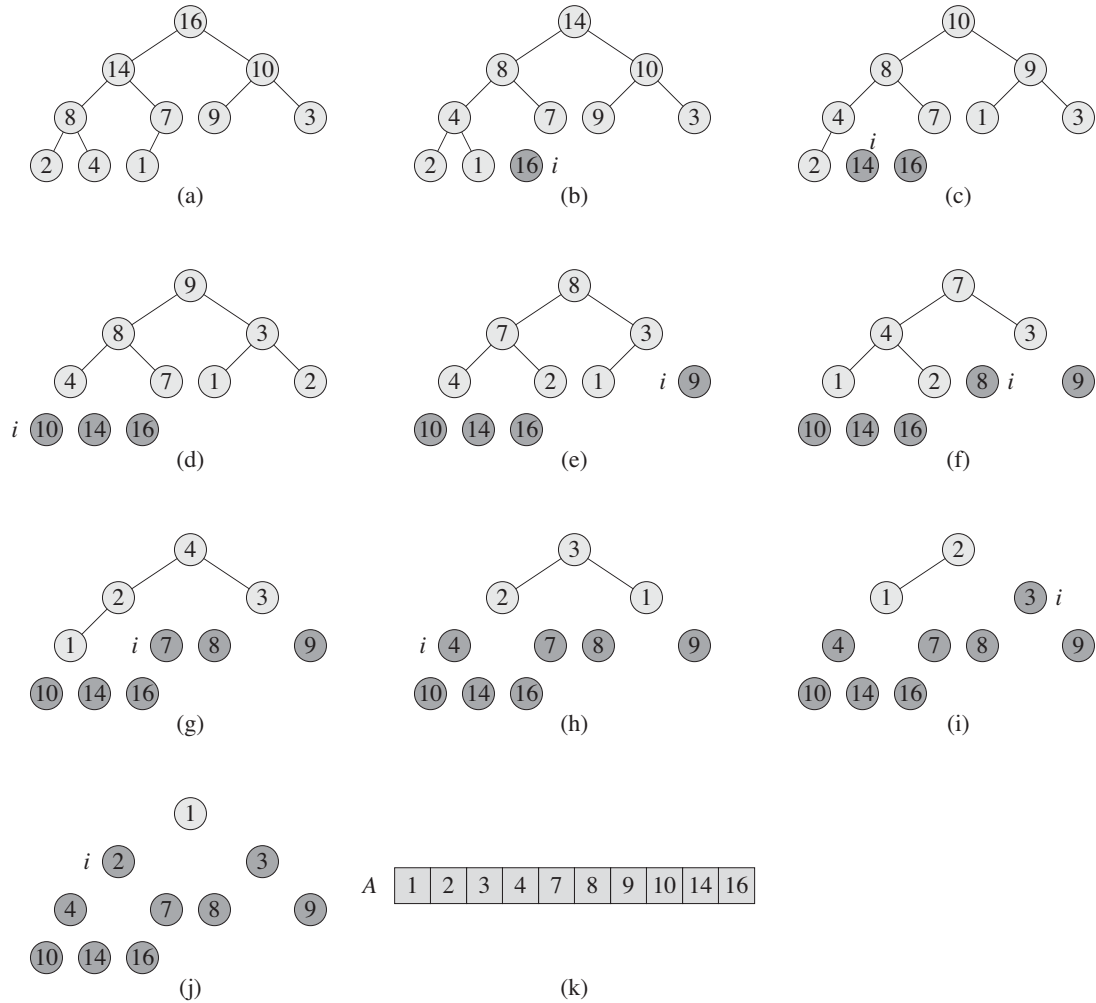Argue the correctness of HEAPSORT using the following loop invariant:

> 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** $A.heap\text{-}size < 1$
2       **error** "heap underflow"
3   $max = A[1]$
4   $A[1] = A[A.heap\text{-}size]$
5   $A.heap\text{-}size = A.heap\text{-}size - 1$
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] = -∞
3   HEAP-INCREASE-KEY(A, A.heap-size, key)
```
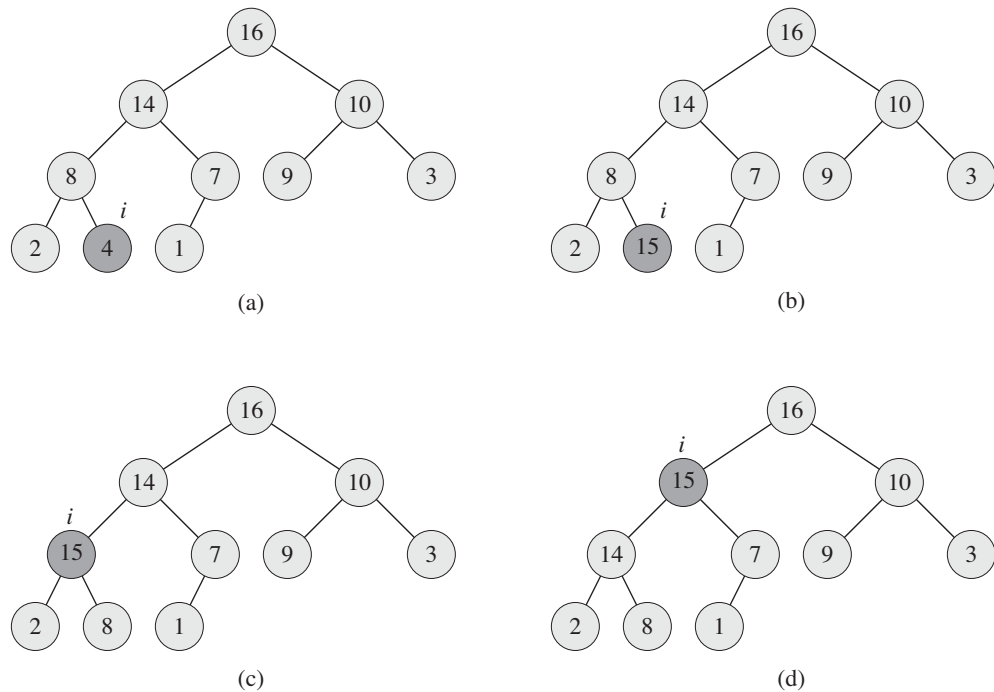
The running time of MAX-HEAP-INSERT on an $n$-element heap is $O(\lg n)$.

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?

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

The following procedure implements quicksort:

QUICKSORT($A, p, r$)

1   **if** $p < r$
2       $q = $ 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 QUICKSORT($A, 1, A.length$).

**Partitioning the array**

The key to the algorithm is the PARTITION procedure, which rearranges the subar-
ray $A[p \mathinner{\ldotp\ldotp} r]$ in place.

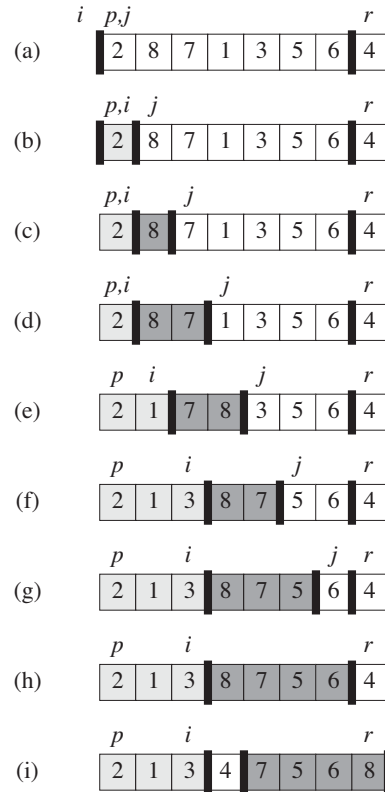PARTITION($A, p, r$)

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 $A[p \mathinner{\ldotp\ldotp} r]$. 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 $A[k] > x$.
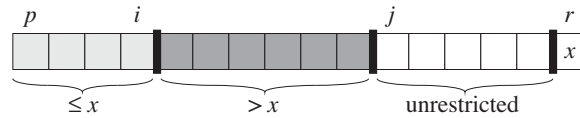> 3. If $k = r$, then $A[k] = x$.

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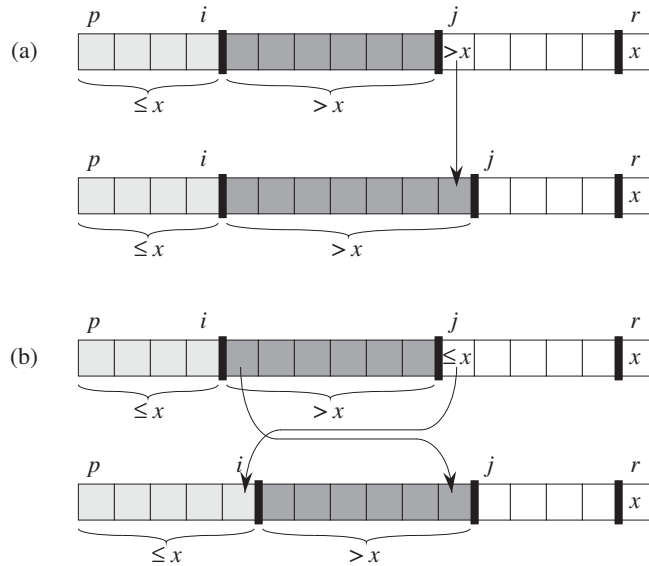
The final two lines of PARTITION finish up by swapping the pivot element with the leftmost element greater than $x$, thereby moving the pivot into its correct place in the partitioned array, and then returning the pivot's new index. The output of PARTITION now satisfies the specifications given for the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

The running time of PARTITION on the subarray $A[p \mathinner{..} r]$ 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 $A[p \mathinner{.\,.} r]$ have the same value?   Modify PARTITION so that $q = \lfloor (p + r)/2 \rfloor$ when all elements in the array $A[p \mathinner{.\,.} r]$ 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.
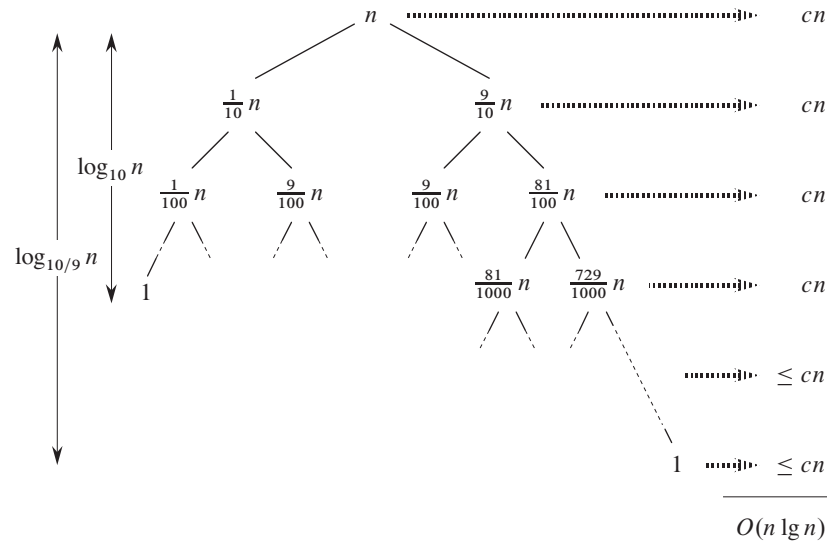
### Best-case partitioning

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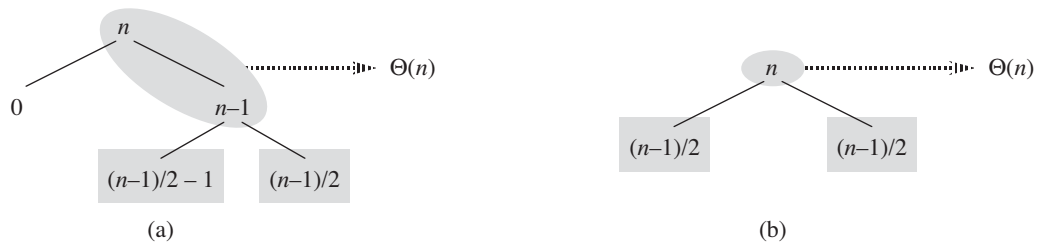
**Figure 7.4** A recursion tree for QUICKSORT in which PARTITION always produces a 9-to-1 split, yielding a running time of $O(n \lg n)$. Nodes show subproblem sizes, with per-level costs on the right. The per-level costs include the constant $c$ implicit in the $\Theta(n)$ term.

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 sub-arrays 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 $0 < \alpha \le 1/2$ 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 $0 < \alpha \le 1/2$, 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 $A[p \mathinner{.\,.} r]$. We do so by first exchanging element $A[r]$ with an element chosen at random from $A[p \mathinner{.\,.} r]$. By randomly sampling the range $p, \ldots, r$, 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 = 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 = 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

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

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 $T(n) \le cn^2$ for some constant $c$. Substituting this guess into recurrence (7.1), we obtain

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

The expression $q^2 + (n-q-1)^2$ achieves a maximum over the parameter's range $0 \le q \le n-1$ 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 $z_1, z_2, \ldots, z_n$, with $z_i$ being the $i$th smallest element. We also define the set $Z_{ij} = \{z_i, z_{i+1}, \ldots, z_j\}$ 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

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

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

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

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

$$= \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} \Pr\{z_i \text{ is compared to } z_j\} \ . \tag{7.2}$$

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

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

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

***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 $q = 0, 1, \ldots, n - 1$ when $q = 0$ or $q = n - 1$.

***7.4-4***
Show that RANDOMIZED-QUICKSORT's expected running time is $\Omega(n \lg n)$.