# Chapter 5

## (Divide-and-Conquer)

# 5

# Divide-and-Conquer

*Whatever man prays for, he prays for a miracle. Every prayer reduces itself to this—Great God, grant that twice two be not four.*
—Ivan Turgenev (1818–1883), Russian novelist and short-story writer

**D**ivide-and-conquer is probably the best-known general algorithm design technique. Though its fame may have something to do with its catchy name, it is well deserved: quite a few very efficient algorithms are specific implementations of this general strategy. Divide-and-conquer algorithms work according to the following general plan:

1. A problem is divided into several subproblems of the same type, ideally of about equal size.
2. The subproblems are solved (typically recursively, though sometimes a different algorithm is employed, especially when subproblems become small enough).
3. If necessary, the solutions to the subproblems are combined to get a solution to the original problem.

The divide-and-conquer technique is diagrammed in Figure 5.1, which depicts the case of dividing a problem into two smaller subproblems, by far the most widely occurring case (at least for divide-and-conquer algorithms designed to be executed on a single-processor computer).

As an example, let us consider the problem of computing the sum of $n$ numbers $a_0, \ldots, a_{n-1}$. If $n > 1$, we can divide the problem into two instances of the same problem: to compute the sum of the first $\lfloor n/2 \rfloor$ numbers and to compute the sum of the remaining $\lceil n/2 \rceil$ numbers. (Of course, if $n = 1$, we simply return $a_0$ as the answer.) Once each of these two sums is computed by applying the same method recursively, we can add their values to get the sum in question:

$$a_0 + \cdots + a_{n-1} = (a_0 + \cdots + a_{\lfloor n/2 \rfloor - 1}) + (a_{\lfloor n/2 \rfloor} + \cdots + a_{n-1}).$$

Is this an efficient way to compute the sum of $n$ numbers? A moment of reflection (why could it be more efficient than the brute-force summation?), a
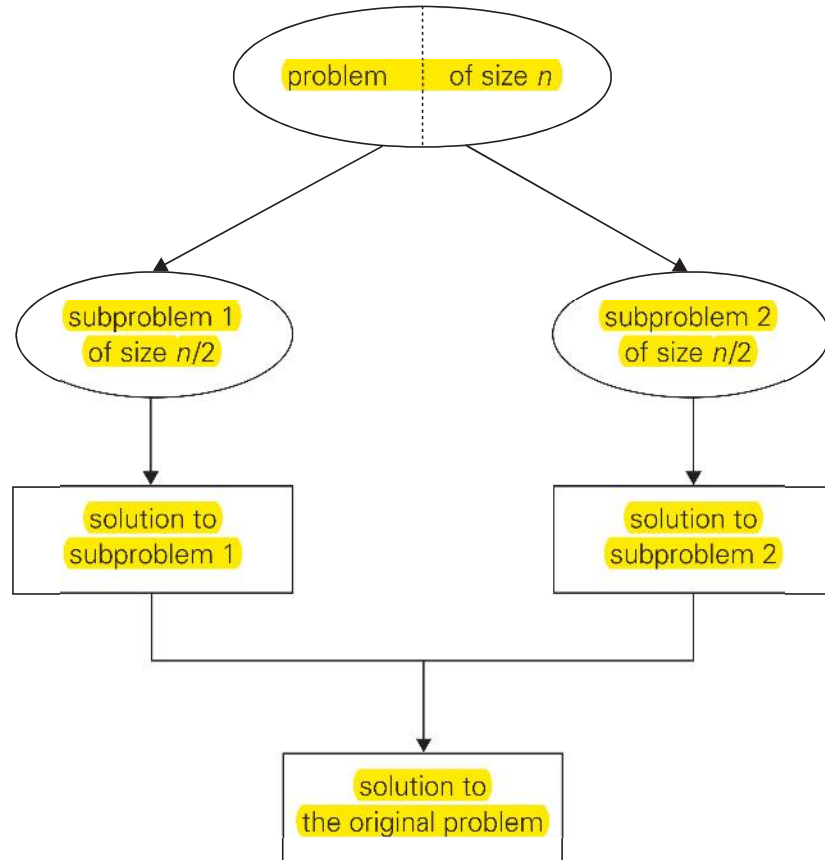
**FIGURE 5.1** Divide-and-conquer technique (typical case).

small example of summing, say, four numbers by this algorithm, a formal analysis (which follows), and common sense (we do not normally compute sums this way, do we?) all lead to a negative answer to this question.[1]

Thus, not every divide-and-conquer algorithm is necessarily more efficient than even a brute-force solution. But often our prayers to the Goddess of Algorithmics—see the chapter's epigraph—*are* answered, and the time spent on executing the divide-and-conquer plan turns out to be significantly smaller than solving a problem by a different method. In fact, the divide-and-conquer approach yields some of the most important and efficient algorithms in computer science. We discuss a few classic examples of such algorithms in this chapter. Though we consider only sequential algorithms here, it is worth keeping in mind that the divide-and-conquer technique is ideally suited for parallel computations, in which each subproblem can be solved simultaneously by its own processor.

---

1.   Actually, the divide-and-conquer algorithm, called the ***pairwise summation***, may substantially reduce the accumulated round-off error of the sum of numbers that can be represented only approximately in a digital computer [Hig93].

As mentioned above, in the most typical case of divide-and-conquer a problem's instance of size $n$ is divided into two instances of size $n/2$. More generally, an instance of size $n$ can be divided into $b$ instances of size $n/b$, with $a$ of them needing to be solved. (Here, $a$ and $b$ are constants; $a \geq 1$ and $b > 1$.) Assuming that size $n$ is a power of $b$ to simplify our analysis, we get the following recurrence for the running time $T(n)$:

$$T(n) = aT(n/b) + f(n), \tag{5.1}$$

where $f(n)$ is a function that accounts for the time spent on dividing an instance of size $n$ into instances of size $n/b$ and combining their solutions. (For the sum example above, $a = b = 2$ and $f(n) = 1$.) Recurrence (5.1) is called the **general divide-and-conquer recurrence**. Obviously, the order of growth of its solution $T(n)$ depends on the values of the constants $a$ and $b$ and the order of growth of the function $f(n)$. The efficiency analysis of many divide-and-conquer algorithms is greatly simplified by the following theorem (see Appendix B).

**Master Theorem**   If $f(n) \in \Theta(n^d)$ where $d \geq 0$ in recurrence (5.1), then

$$T(n) \in \begin{cases} \Theta(n^d) & \text{if } a < b^d, \\ \Theta(n^d \log n) & \text{if } a = b^d, \\ \Theta(n^{\log_b a}) & \text{if } a > b^d. \end{cases}$$

Analogous results hold for the $O$ and $\Omega$ notations, too.

For example, the recurrence for the number of additions $A(n)$ made by the divide-and-conquer sum-computation algorithm (see above) on inputs of size $n = 2^k$ is

$$A(n) = 2A(n/2) + 1.$$

Thus, for this example, $a = 2$, $b = 2$, and $d = 0$; hence, since $a > b^d$,

$$A(n) \in \Theta(n^{\log_b a}) = \Theta(n^{\log_2 2}) = \Theta(n).$$

Note that we were able to find the solution's efficiency class without going through the drudgery of solving the recurrence. But, of course, this approach can only establish a solution's order of growth to within an unknown multiplicative constant, whereas solving a recurrence equation with a specific initial condition yields an exact answer (at least for $n$'s that are powers of $b$).

It is also worth pointing out that if $a = 1$, recurrence (5.1) covers decrease-by-a-constant-factor algorithms discussed in the previous chapter. In fact, some people consider such algorithms as binary search degenerate cases of divide-and-conquer, where just one of two subproblems of half the size needs to be solved. It is better not to do this and consider decrease-by-a-constant-factor and divide-and-conquer as different design paradigms.

## 5.1 Mergesort

Mergesort is a perfect example of a successful application of the divide-and-conquer technique. It sorts a given array $A[0..n-1]$ by dividing it into two halves $A[0..\lfloor n/2 \rfloor - 1]$ and $A[\lfloor n/2 \rfloor..n-1]$, sorting each of them recursively, and then merging the two smaller sorted arrays into a single sorted one.

**ALGORITHM** *Mergesort($A[0..n-1]$)*

    //Sorts array $A[0..n-1]$ by recursive mergesort
    //Input: An array $A[0..n-1]$ of orderable elements
    //Output: Array $A[0..n-1]$ sorted in nondecreasing order
    **if** $n > 1$
        copy $A[0..\lfloor n/2 \rfloor - 1]$ to $B[0..\lfloor n/2 \rfloor - 1]$
        copy $A[\lfloor n/2 \rfloor..n-1]$ to $C[0..\lceil n/2 \rceil - 1]$
        *Mergesort($B[0..\lfloor n/2 \rfloor - 1]$)*
        *Mergesort($C[0..\lceil n/2 \rceil - 1]$)*
        *Merge(B, C, A)*    //see below

The ***merging*** of two sorted arrays can be done as follows. Two pointers (array indices) are initialized to point to the first elements of the arrays being merged. The elements pointed to are compared, and the smaller of them is added to a new array being constructed; after that, the index of the smaller element is incremented to point to its immediate successor in the array it was copied from. This operation is repeated until one of the two given arrays is exhausted, and then the remaining elements of the other array are copied to the end of the new array.

**ALGORITHM** *Merge($B[0..p-1]$, $C[0..q-1]$, $A[0..p+q-1]$)*

    //Merges two sorted arrays into one sorted array
    //Input: Arrays $B[0..p-1]$ and $C[0..q-1]$ both sorted
    //Output: Sorted array $A[0..p+q-1]$ of the elements of $B$ and $C$
    $i \leftarrow 0$;  $j \leftarrow 0$;  $k \leftarrow 0$
    **while** $i < p$ **and** $j < q$ **do**
        **if** $B[i] \leq C[j]$
            $A[k] \leftarrow B[i]$;  $i \leftarrow i + 1$
        **else** $A[k] \leftarrow C[j]$;  $j \leftarrow j + 1$
        $k \leftarrow k + 1$
    **if** $i = p$
        copy $C[j..q-1]$ to $A[k..p+q-1]$
    **else** copy $B[i..p-1]$ to $A[k..p+q-1]$

The operation of the algorithm on the list 8, 3, 2, 9, 7, 1, 5, 4 is illustrated in Figure 5.2.
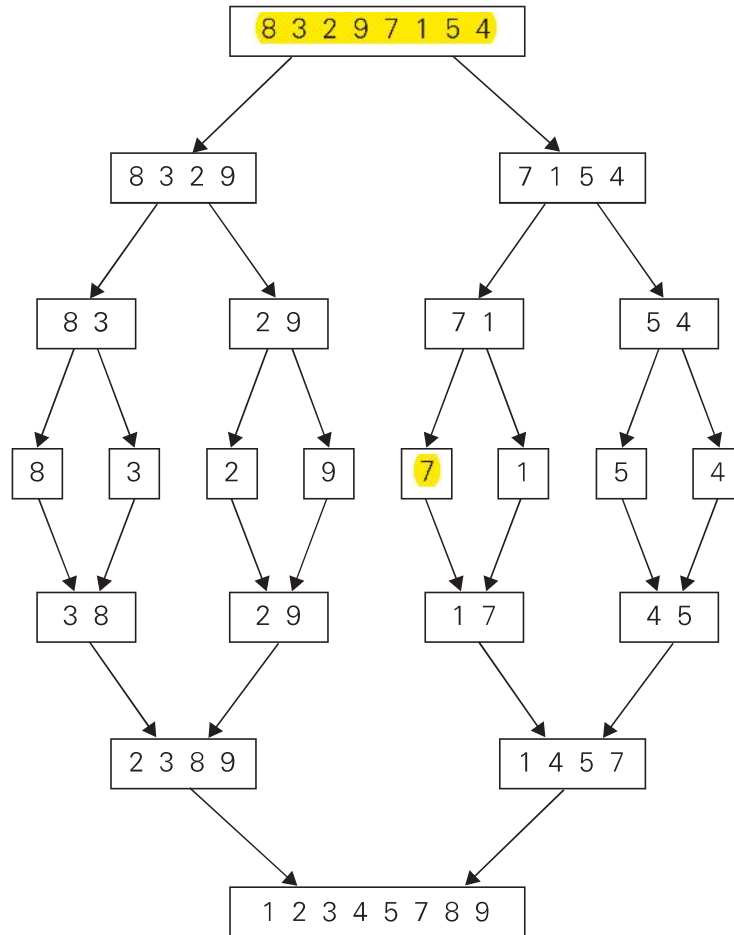
**FIGURE 5.2** Example of mergesort operation.

How efficient is mergesort? Assuming for simplicity that $n$ is a power of 2, the recurrence relation for the number of key comparisons $C(n)$ is

$$C(n) = 2C(n/2) + C_{merge}(n) \quad \text{for } n > 1, \quad C(1) = 0.$$

Let us analyze $C_{merge}(n)$, the number of key comparisons performed during the merging stage. At each step, exactly one comparison is made, after which the total number of elements in the two arrays still needing to be processed is reduced by 1. In the worst case, neither of the two arrays becomes empty before the other one contains just one element (e.g., smaller elements may come from the alternating arrays). Therefore, for the worst case, $C_{merge}(n) = n - 1$, and we have the recurrence

$$C_{worst}(n) = 2C_{worst}(n/2) + n - 1 \quad \text{for } n > 1, \quad C_{worst}(1) = 0.$$

Hence, according to the Master Theorem, $C_{worst}(n) \in \Theta(n \log n)$ (why?). In fact, it is easy to find the exact solution to the worst-case recurrence for $n = 2^k$:

$$C_{worst}(n) = n \log_2 n - n + 1.$$

The number of key comparisons made by mergesort in the worst case comes very close to the theoretical minimum[2] that any general comparison-based sorting algorithm can have. For large $n$, the number of comparisons made by this algorithm in the average case turns out to be about $0.25n$ less (see [Gon91, p. 173]) and hence is also in $\Theta(n \log n)$. A noteworthy advantage of mergesort over quicksort and heapsort—the two important advanced sorting algorithms to be discussed later—is its stability (see Problem 7 in this section's exercises). The principal shortcoming of mergesort is the linear amount of extra storage the algorithm requires. Though merging can be done in-place, the resulting algorithm is quite complicated and of theoretical interest only.
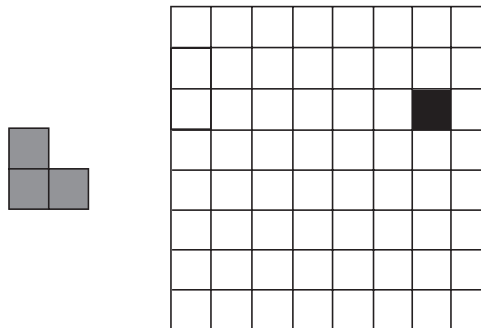
There are two main ideas leading to several variations of mergesort. First, the algorithm can be implemented bottom up by merging pairs of the array's elements, then merging the sorted pairs, and so on. (If $n$ is not a power of 2, only slight bookkeeping complications arise.) This avoids the time and space overhead of using a stack to handle recursive calls. Second, we can divide a list to be sorted in more than two parts, sort each recursively, and then merge them together. This scheme, which is particularly useful for sorting files residing on secondary memory devices, is called ***multiway mergesort***.

─────────── **Exercises 5.1** ───────────

**1. a.** Write pseudocode for a divide-and-conquer algorithm for finding the position of the largest element in an array of $n$ numbers.

   **b.** What will be your algorithm's output for arrays with several elements of the largest value?

   **c.** Set up and solve a recurrence relation for the number of key comparisons made by your algorithm.

   **d.** How does this algorithm compare with the brute-force algorithm for this problem?

**2. a.** Write pseudocode for a divide-and-conquer algorithm for finding values of both the largest and smallest elements in an array of $n$ numbers.

   **b.** Set up and solve (for $n = 2^k$) a recurrence relation for the number of key comparisons made by your algorithm.

   **c.** How does this algorithm compare with the brute-force algorithm for this problem?

**3. a.** Write pseudocode for a divide-and-conquer algorithm for the exponentiation problem of computing $a^n$ where $n$ is a positive integer.

   **b.** Set up and solve a recurrence relation for the number of multiplications made by this algorithm.

───────────

2.   As we shall see in Section 11.2, this theoretical minimum is $\lceil \log_2 n! \rceil \approx \lceil n \log_2 n - 1.44n \rceil$.

    **c.** How does this algorithm compare with the brute-force algorithm for this problem?

**4.** As mentioned in Chapter 2, logarithm bases are irrelevant in most contexts arising in analyzing an algorithm's efficiency class. Is this true for both assertions of the Master Theorem that include logarithms?

**5.** Find the order of growth for solutions of the following recurrences.
    **a.** $T(n) = 4T(n/2) + n$, $T(1) = 1$
    **b.** $T(n) = 4T(n/2) + n^2$, $T(1) = 1$
    **c.** $T(n) = 4T(n/2) + n^3$, $T(1) = 1$

**6.** Apply mergesort to sort the list $E, X, A, M, P, L, E$ in alphabetical order.

**7.** Is mergesort a stable sorting algorithm?

**8. a.** Solve the recurrence relation for the number of key comparisons made by mergesort in the worst case. You may assume that $n = 2^k$.

    **b.** Set up a recurrence relation for the number of key comparisons made by mergesort on best-case inputs and solve it for $n = 2^k$.

    **c.** Set up a recurrence relation for the number of key moves made by the version of mergesort given in Section 5.1. Does taking the number of key moves into account change the algorithm's efficiency class?

**9.** Let $A[0..n-1]$ be an array of $n$ real numbers. A pair $(A[i], A[j])$ is said to be an ***inversion*** if these numbers are out of order, i.e., $i < j$ but $A[i] > A[j]$. Design an $O(n \log n)$ algorithm for counting the number of inversions.

**10.** Implement the bottom-up version of mergesort in the language of your choice.

**11.** *Tromino puzzle* A tromino (more accurately, a right tromino) is an L-shaped tile formed by three $1 \times 1$ squares. The problem is to cover any $2^n \times 2^n$ chessboard with a missing square with trominoes. Trominoes can be oriented in an arbitrary way, but they should cover all the squares of the board except the missing one exactly and with no overlaps. [Gol94]

Design a divide-and-conquer algorithm for this problem.

## 5.2 Quicksort

Quicksort is the other important sorting algorithm that is based on the divide-and-conquer approach. Unlike mergesort, which divides its input elements according to their position in the array, quicksort divides them according to their value. We already encountered this idea of an array partition in Section 4.5, where we discussed the selection problem. A partition is an arrangement of the array's elements so that all the elements to the left of some element $A[s]$ are less than or equal to $A[s]$, and all the elements to the right of $A[s]$ are greater than or equal to it:

$$\underbrace{A[0] \ldots A[s-1]}_{\text{all are } \leq A[s]} \; A[s] \; \underbrace{A[s+1] \ldots A[n-1]}_{\text{all are } \geq A[s]}$$

Obviously, after a partition is achieved, $A[s]$ will be in its final position in the sorted array, and we can continue sorting the two subarrays to the left and to the right of $A[s]$ independently (e.g., by the same method). Note the difference with mergesort: there, the division of the problem into two subproblems is immediate and the entire work happens in combining their solutions; here, the entire work happens in the division stage, with no work required to combine the solutions to the subproblems.

Here is pseudocode of quicksort: call $Quicksort(A[0..n-1])$ where

**ALGORITHM**  *Quicksort(A[l..r])*

    //Sorts a subarray by quicksort
    //Input: Subarray of array $A[0..n-1]$, defined by its left and right
    //        indices $l$ and $r$
    //Output: Subarray $A[l..r]$ sorted in nondecreasing order
    **if** $l < r$
        $s \leftarrow Partition(A[l..r])$  //$s$ is a split position
        $Quicksort(A[l..s-1])$
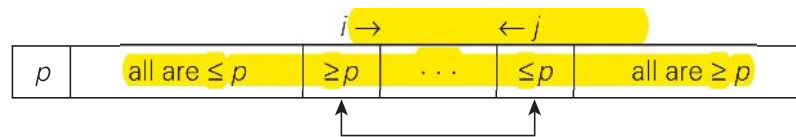        $Quicksort(A[s+1..r])$

As a partition algorithm, we can certainly use the Lomuto partition discussed in Section 4.5. Alternatively, we can partition $A[0..n-1]$ and, more generally, its subarray $A[l..r]$ $(0 \leq l < r \leq n-1)$ by the more sophisticated method suggested by C.A.R. Hoare, the prominent British computer scientist who invented quicksort.[3]

---

3.    C.A.R. Hoare, at age 26, invented his algorithm in 1960 while trying to sort words for a machine translation project from Russian to English. Says Hoare, "My first thought on how to do this was bubblesort and, by an amazing stroke of luck, my second thought was Quicksort." It is hard to disagree with his overall assessment: "I have been very lucky. What a wonderful way to start a career in Computing, by discovering a new sorting algorithm!" [Hoa96]. Twenty years later, he received the Turing Award for "fundamental contributions to the definition and design of programming languages"; in 1980, he was also knighted for services to education and computer science.
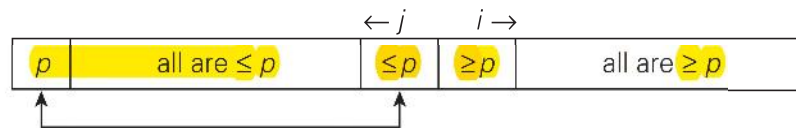
As before, we start by selecting a pivot—an element with respect to whose value we are going to divide the subarray. There are several different strategies for selecting a pivot; we will return to this issue when we analyze the algorithm's efficiency. For now, we use the simplest strategy of selecting the subarray's first element: $p = A[l]$.

Unlike the Lomuto algorithm, we will now scan the subarray from both ends, comparing the subarray's elements to the pivot. The left-to-right scan, denoted below by index pointer $i$, starts with the second element. Since we want elements smaller than the pivot to be in the left part of the subarray, this scan skips over elements that are smaller than the pivot and stops upon encountering the first element greater than or equal to the pivot. The right-to-left scan, denoted below by index pointer $j$, starts with the last element of the subarray. Since we want elements larger than the pivot to be in the right part of the subarray, this scan skips over elements that are larger than the pivot and stops on encountering the first element smaller than or equal to the pivot. (Why is it worth stopping the scans after encountering an element equal to the pivot? Because doing this tends to yield more even splits for arrays with a lot of duplicates, which makes the algorithm run faster. For example, if we did otherwise for an array of $n$ equal elements, we would have gotten a split into subarrays of sizes $n - 1$ and 0, reducing the problem size just by 1 after scanning the entire array.)
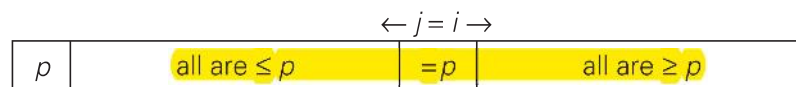
After both scans stop, three situations may arise, depending on whether or not the scanning indices have crossed. If scanning indices $i$ and $j$ have not crossed, i.e., $i < j$, we simply exchange $A[i]$ and $A[j]$ and resume the scans by incrementing $i$ and decrementing $j$, respectively:

| | | $i \rightarrow$ | | $\leftarrow j$ | | |
|---|---|---|---|---|---|---|
| $p$ | all are $\leq p$ | $\geq p$ | $\cdots$ | $\leq p$ | all are $\geq p$ |

If the scanning indices have crossed over, i.e., $i > j$, we will have partitioned the subarray after exchanging the pivot with $A[j]$:

| | | | $\leftarrow j$ | $i \rightarrow$ | |
|---|---|---|---|---|---|
| $p$ | all are $\leq p$ | $\leq p$ | $\geq p$ | all are $\geq p$ |

Finally, if the scanning indices stop while pointing to the same element, i.e., $i = j$, the value they are pointing to must be equal to $p$ (why?). Thus, we have the subarray partitioned, with the split position $s = i = j$:

| | | $\leftarrow j = i \rightarrow$ | |
|---|---|---|---|
| $p$ | all are $\leq p$ | $= p$ | all are $\geq p$ |

We can combine the last case with the case of crossed-over indices ($i > j$) by exchanging the pivot with $A[j]$ whenever $i \geq j$.

Here is pseudocode implementing this partitioning procedure.

**ALGORITHM** *HoarePartition(A[l..r])*

    //Partitions a subarray by Hoare's algorithm, using the first element
    //        as a pivot
    //Input: Subarray of array $A[0..n-1]$, defined by its left and right
    //        indices $l$ and $r$ $(l < r)$
    //Output: Partition of $A[l..r]$, with the split position returned as
    //        this function's value
    $p \leftarrow A[l]$
    $i \leftarrow l;\ j \leftarrow r + 1$
    **repeat**
        **repeat** $i \leftarrow i + 1$ **until** $A[i] \geq p$
        **repeat** $j \leftarrow j - 1$ **until** $A[j] \leq p$
        swap($A[i], A[j]$)
    **until** $i \geq j$
    swap($A[i], A[j]$)    //undo last swap when $i \geq j$
    swap($A[l], A[j]$)
    **return** $j$

Note that index $i$ can go out of the subarray's bounds in this pseudocode. Rather than checking for this possibility every time index $i$ is incremented, we can append to array $A[0..n-1]$ a "sentinel" that would prevent index $i$ from advancing beyond position $n$. Note that the more sophisticated method of pivot selection mentioned at the end of the section makes such a sentinel unnecessary.
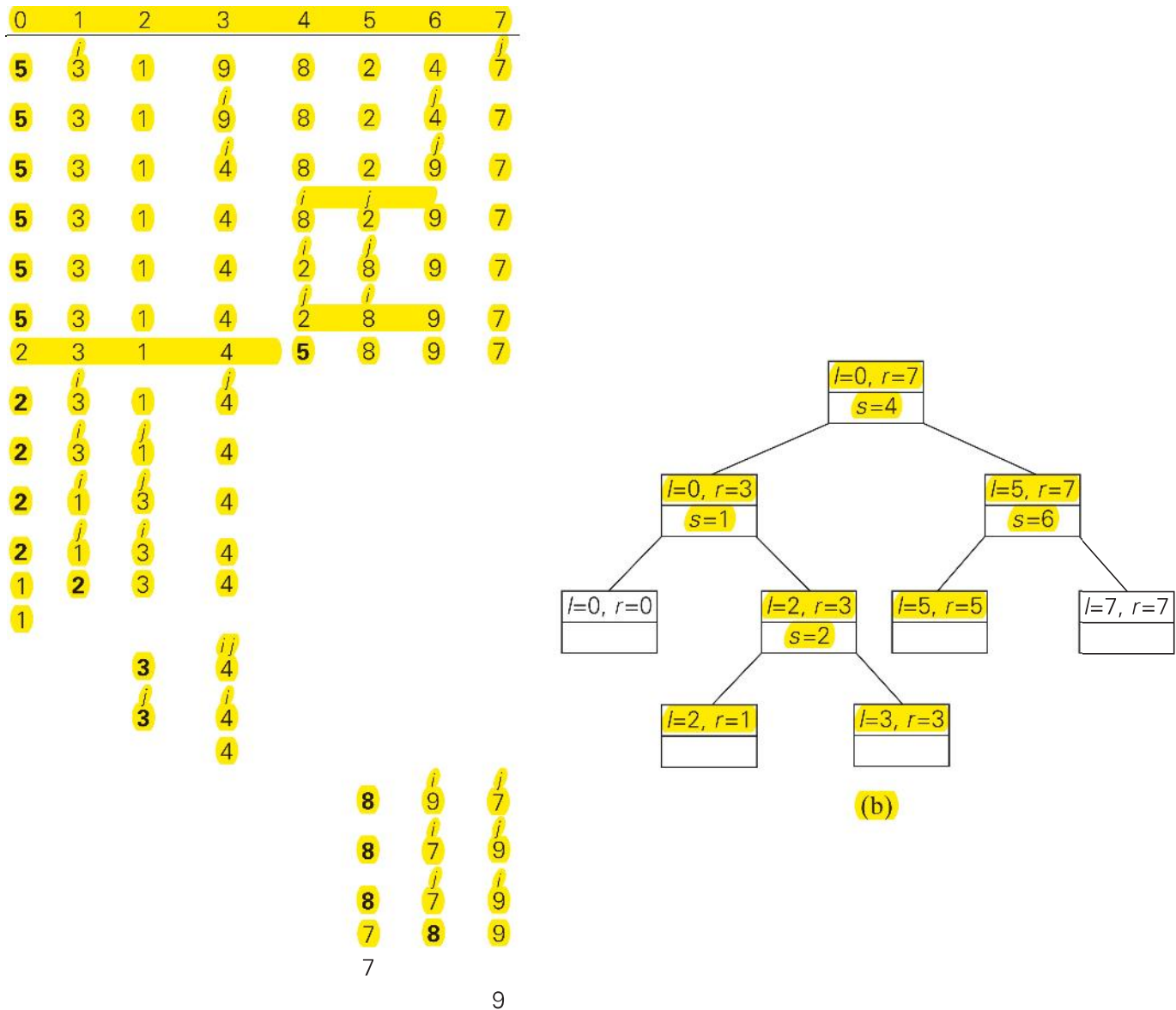
An example of sorting an array by quicksort is given in Figure 5.3.

We start our discussion of quicksort's efficiency by noting that the number of key comparisons made before a partition is achieved is $n + 1$ if the scanning indices cross over and $n$ if they coincide (why?). If all the splits happen in the middle of corresponding subarrays, we will have the best case. The number of key comparisons in the best case satisfies the recurrence

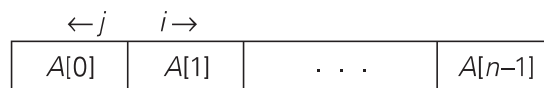$$C_{best}(n) = 2C_{best}(n/2) + n \quad \text{for } n > 1, \quad C_{best}(1) = 0.$$

According to the Master Theorem, $C_{best}(n) \in \Theta(n \log_2 n)$; solving it exactly for $n = 2^k$ yields $C_{best}(n) = n \log_2 n$.

In the worst case, all the splits will be skewed to the extreme: one of the two subarrays will be empty, and the size of the other will be just 1 less than the size of the subarray being partitioned. This unfortunate situation will happen, in particular, for increasing arrays, i.e., for inputs for which the problem is already solved! Indeed, if $A[0..n-1]$ is a strictly increasing array and we use $A[0]$ as the pivot, the left-to-right scan will stop on $A[1]$ while the right-to-left scan will go all the way to reach $A[0]$, indicating the split at position 0:

**FIGURE 5.3** Example of quicksort operation. (a) Array's transformations with pivots shown in bold. (b) Tree of recursive calls to *Quicksort* with input values *l* and *r* of subarray bounds and split position *s* of a partition obtained.



So, after making $n + 1$ comparisons to get to this partition and exchanging the pivot $A[0]$ with itself, the algorithm will be left with the strictly increasing array $A[1..n - 1]$ to sort. This sorting of strictly increasing arrays of diminishing sizes will

the size of this stack can be made to be in $O(\log n)$ by always sorting first the smaller of two subarrays obtained by partitioning, it is worse than the $O(1)$ space efficiency of heapsort. Although more sophisticated ways of choosing a pivot make the quadratic running time of the worst case very unlikely, they do not eliminate it completely. And even the performance on randomly ordered arrays is known to be sensitive not only to implementation details of the algorithm but also to both computer architecture and data type. Still, the January/February 2000 issue of *Computing in Science & Engineering,* a joint publication of the American Institute of Physics and the IEEE Computer Society, selected quicksort as one of the 10 algorithms "with the greatest influence on the development and practice of science and engineering in the 20th century."

---
## Exercises 5.2
---

1. Apply quicksort to sort the list $E, X, A, M, P, L, E$ in alphabetical order. Draw the tree of the recursive calls made.

2. For the partitioning procedure outlined in this section:
   a. Prove that if the scanning indices stop while pointing to the same element, i.e., $i = j$, the value they are pointing to must be equal to $p$.
   b. Prove that when the scanning indices stop, $j$ cannot point to an element more than one position to the left of the one pointed to by $i$.

3. Give an example showing that quicksort is not a stable sorting algorithm.

4. Give an example of an array of $n$ elements for which the sentinel mentioned in the text is actually needed. What should be its value? Also explain why a single sentinel suffices for any input.

5. For the version of quicksort given in this section:
   a. Are arrays made up of all equal elements the worst-case input, the best-case input, or neither?
   b. Are strictly decreasing arrays the worst-case input, the best-case input, or neither?

6. a. For quicksort with the median-of-three pivot selection, are strictly increasing arrays the worst-case input, the best-case input, or neither?
   b. Answer the same question for strictly decreasing arrays.

7. a. Estimate how many times faster quicksort will sort an array of one million random numbers than insertion sort.
   b. True or false: For every $n > 1$, there are $n$-element arrays that are sorted faster by insertion sort than by quicksort?

8. Design an algorithm to rearrange elements of a given array of $n$ real numbers so that all its negative elements precede all its positive elements. Your algorithm should be both time efficient and space efficient.