

Volume 1 / **Fundamental Algorithms**

**DONALD E. KNUTH** *Stanford University*

**THE ART OF  
COMPUTER PROGRAMMING  
SECOND EDITION**

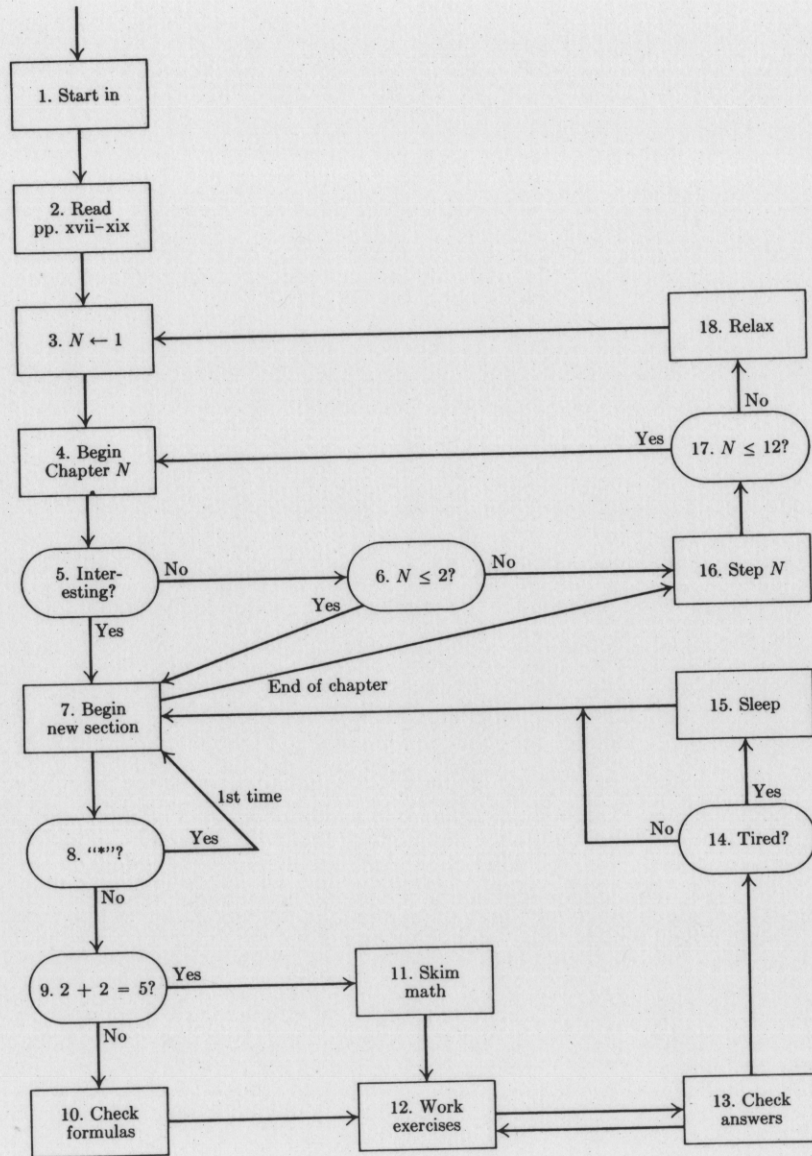


**ADDISON-WESLEY PUBLISHING COMPANY**

© 1973

Reading, Massachusetts  
Menlo Park, California · London · Amsterdam · Don Mills, Ontario · Sydney

## Procedure for Reading This Set of Books



Flow chart for reading this set of books.

1. Begin reading this procedure, unless you have already begun to read it. *Continue to follow the steps faithfully.* (The general form of this procedure and its accompanying flowchart will be used throughout this book.)
2. Read the Notes on the Exercises, pp. xvii-xix.
3. Set  $N$  equal to 1.
4. Begin reading Chapter  $N$ . Do *not* read the quotations which appear at the beginning of the chapter.
5. Is the subject of the chapter interesting to you? If so, go to step 7; if not, go to step 6.
6. Is  $N \leq 2$ ? If not, go to step 16; if so, scan through the chapter anyway. (Chapters 1 and 2 contain important introductory material and also a review of basic programming techniques. You should at least skim over the sections on notation and about MIX.)
7. Begin reading the next section of the chapter; if you have reached the end of the chapter, go to step 16.
8. Is section number marked with "\*" ? If so, you may omit this section on first reading (it covers a rather specialized topic which is interesting but not essential); go back to step 7.
9. Are you mathematically inclined? If math is all Greek to you, go to step 11; otherwise go to step 10.
10. Check the mathematical derivations made in this section (and report errors to the author). Go to step 12.
11. If the current section is full of mathematical computations, you had better omit reading the derivations. However, you should become familiar with the basic results of the section; these are usually stated near the beginning or in italics right at the very end of the hard parts.
12. Work the recommended exercises in this section in accordance with the hints given in the Notes on the Exercises (which you read in step 2).

13. After you have worked on the exercises to your satisfaction, check your answers with the answer printed in the corresponding answer section at the rear of the book (if any answer appears for that problem). Also read the answers to the exercises you did not have time to work. *Note:* In most cases it is reasonable to read the answer to exercise  $n$  before working on exercise  $n + 1$ , so steps 12–13 are usually done simultaneously.
14. Are you tired? If not, go back to step 7.
15. Go to sleep. Then, wake up, and go back to step 7.
16. Increase  $N$  by one. If  $N = 3, 5, 7, 9, 11$ , or  $12$ , begin the next volume of this set of books.
17. If  $N$  is less than or equal to  $12$ , go back to step 4.
18. Congratulations. Now try to get your friends to purchase a copy of volume one and to start reading it. Also, go back to step 3.

*Woe be to him that reads but one book.*

— GEORGE HERBERT, *Jacula Prudentum*, 1144 (1640)

*Le défaut unique de tous les ouvrages  
c'est d'être trop longs.*

— VAUVENARGUES, *Réflexions*, 628 (1746)

*Books are a triviality. Life alone is great.*

— THOMAS CARLYLE, *Journal* (1839)

## Notes on the Exercises

The exercises in this set of books have been designed for self-study as well as classroom study. It is difficult, if not impossible, for anyone to learn a subject purely by reading about it, without applying the information to specific problems and thereby forcing himself to think about what has been read. Furthermore, we all learn best the things that we have discovered for ourselves. Therefore the exercises form a major part of this work; a definite attempt has been made to keep them as informative as possible and to select problems that are enjoyable to solve.

In many books, easy exercises are found mixed randomly among extremely difficult ones. This is sometimes unfortunate because the reader should have some idea about how much time it ought to take him to do a problem before he tackles it (otherwise he may just skip over all the problems). A classic example of this situation is the book *Dynamic Programming* by Richard Bellman; this is an important, pioneering book in which a group of problems is collected together at the end of some chapters under the heading "Exercises and Research Problems," with extremely trivial questions appearing in the midst of deep, unsolved problems. It is rumored that someone once asked Dr. Bellman how to tell the exercises apart from the research problems, and he replied, "If you can solve it, it is an exercise; otherwise it's a research problem."

Good arguments can be made for including both research problems and very easy exercises in a book of this kind; therefore, to save the reader from the possible dilemma of determining which are which, *rating numbers* have been provided to indicate the level of difficulty. These numbers have the following general significance:

### *Rating Interpretation*

- 00 An extremely easy exercise which can be answered immediately if the material of the text has been understood, and which can almost always be worked "in your head."
- 10 A simple problem, which makes a person think over the material just read, but which is by no means difficult. It should be possible to do this in one minute at most; pencil and paper may be useful in obtaining the solution.
- 20 An average problem which tests basic understanding of the text material but which may take about fifteen to twenty minutes to answer completely.

## CHAPTER ONE

### BASIC CONCEPTS

*Many persons who are not conversant with mathematical studies imagine that because the business of [Babbage's Analytical Engine] is to give its results in numerical notation, the nature of its processes must consequently be arithmetical and numerical, rather than algebraical and analytical. This is an error. The engine can arrange and combine its numerical quantities exactly as if they were letters or any other general symbols; and in fact it might bring out its results in algebraical notation, were provisions made accordingly.*

— ADA AUGUSTA, Countess of Lovelace (1844)

*Practise yourself, for heaven's sake, in little things; and thence proceed to greater.*

— EPICTETUS (*Discourses* IV. i)



#### 1.1. ALGORITHMS

The notion of an *algorithm* is basic to all of computer programming, so we should begin with a careful analysis of this concept.

The word “algorithm” itself is quite interesting; at first glance it may look as though someone intended to write “logarithm” but jumbled up the first four letters. The word did not appear in *Webster's New World Dictionary* as late as 1957; we find only the older form “algorism” with its ancient meaning, i.e., the process of doing arithmetic using Arabic numerals. In the middle ages, abacists computed on the abacus and algorists computed by algorism. Following the middle ages, the origin of this word was in doubt, and early linguists attempted to guess at its derivation by making combinations like *algiros* [painful] + *arithmos* [number]; others said no, the word comes from “King Algor of Castile.” Finally, historians of mathematics found the true origin of the word algorism: it comes from the name of a famous Persian textbook author, Abu Ja'far Mohammed ibn Mûsâ al-Khowârizmî (c. 825)—literally, “Father of Ja'far, Mohammed, son of Moses, native of Khowârizm.” Khowârizm is today the small Soviet city of Khiva. Al-Khowârizmî wrote the celebrated book *Kitab al-jabr w'al-muqabala* (“Rules of restoration and reduction”); another word, “algebra,” stems from the title of his book, although the book wasn't really very algebraic.

Gradually the form and meaning of "algorism" became corrupted; as explained by the Oxford English Dictionary, the word was "erroneously refashioned" by "learned confusion" with the word *arithmetic*. The change from "algorism" to "algorithm" is not hard to understand in view of the fact that people had forgotten the original derivation of the word. An early German mathematical dictionary, *Vollständiges Mathematisches Lexicon* (Leipzig, 1747), gives the following definition for the word *Algorithmus*: "Under this designation are combined the notions of the four types of arithmetic calculations, namely addition, multiplication, subtraction, and division." The latin phrase *algorithmus infinitesimalis* was at that time used to denote "ways of calculation with infinitely small quantities, as invented by Leibnitz."

By 1950, the word algorithm was most frequently associated with "Euclid's algorithm," a process for finding the greatest common divisor of two numbers which appears in Euclid's *Elements* (Book 7, Propositions 1 and 2.) It will be instructive to exhibit Euclid's algorithm here:

**Algorithm E** (*Euclid's algorithm*). Given two positive integers  $m$  and  $n$ , find their greatest common divisor, i.e., the largest positive integer which evenly divides both  $m$  and  $n$ .

**E1.** [Find remainder.] Divide  $m$  by  $n$  and let  $r$  be the remainder. (We will have  $0 \leq r < n$ .)

**E2.** [Is it zero?] If  $r = 0$ , the algorithm terminates;  $n$  is the answer.

**E3.** [Interchange.] Set  $m \leftarrow n$ ,  $n \leftarrow r$ , and go back to step E1. ■

Of course, Euclid did not present his algorithm in just this manner. The above format illustrates the style in which all of the algorithms throughout this book will be presented.

Each algorithm we consider has been given an identifying letter (e.g., E in the above) and the steps of the algorithm are identified by this letter followed by a number (e.g., E1, E2, etc.). The chapters are divided into numbered sections; within a section the algorithms are designated by letter only, but when algorithms are referred to in other sections, the appropriate section number is also used. For example, we are now in Section 1.1; within this section Euclid's algorithm is called Algorithm E, while in later sections it is referred to as Algorithm 1.1E.

Each step of an algorithm (e.g., step E1 above) begins with a phrase in brackets which sums up as briefly as possible the principal content of that step. This phrase also usually appears in an accompanying *flow chart* (e.g., Fig. 1), so the reader will be able to picture the algorithm more readily.

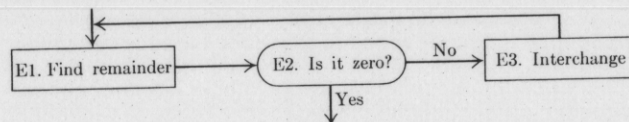


Fig. 1. Flow chart for Algorithm E.

After the summarizing phrase comes a description in words and symbols of some *action* to be performed or some decision to be made. There are also occasionally *parenthesized comments* (e.g., the second sentence in step E1) which are included as explanatory information about that step, often indicating certain characteristics of the variables or the current goals at that step, etc.; the parenthesized remarks do not specify actions which belong to the algorithm, they are only for the reader's benefit as possible aids to comprehension.

The " $\leftarrow$ " arrow in step E3 is the all-important *replacement* operation (sometimes called *assignment* or *substitution*); " $m \leftarrow n$ " means the value of variable  $m$  is to be replaced by the current value of variable  $n$ . When Algorithm E begins, the values of  $m$  and  $n$  are the originally given numbers; but when it ends, these variables will have, in general, different values. An arrow is used to distinguish the replacement operation from the equality relation: We will not say, "Set  $m = n$ ," but we will perhaps ask, "Does  $m = n$ ?" The " $=$ " sign denotes a condition which can be tested, the " $\leftarrow$ " sign denotes an action which can be performed. The operation of *increasing  $n$  by one* is denoted by " $n \leftarrow n + 1$ " (read " $n$  is replaced by  $n + 1$ "); in general, "variable  $\leftarrow$  formula" means the formula is to be computed using the present values of any variables appearing within it, and the result replaces the previous value of the variable at the left of the arrow. Persons untrained in computer work sometimes have a tendency to denote the operation of increasing  $n$  by one by " $n \rightarrow n + 1$ ," saying " $n$  becomes  $n + 1$ "; this can only lead to confusion because of its conflict with the standard conventions, and it should be avoided.

Note that the order of the actions in step E3 is important; "set  $m \leftarrow n$ ,  $n \leftarrow r$ " is quite different from "set  $n \leftarrow r$ ,  $m \leftarrow n$ ," since the latter would imply that the previous value of  $n$  is lost before it can be used to set  $m$ . Thus the latter sequence is equivalent to "set  $n \leftarrow r$ ,  $m \leftarrow r$ ." When several variables are all to be set equal to the same quantity, we use multiple arrows; thus " $n \leftarrow r$ ,  $m \leftarrow r$ " may be written as " $n \leftarrow m \leftarrow r$ ." To interchange the values of two variables, we can write "Exchange  $m \leftrightarrow n$ "; this action may also be specified by using a new variable  $t$  and writing "set  $t \leftarrow m$ ,  $m \leftarrow n$ ,  $n \leftarrow t$ ."

An algorithm starts at the lowest-numbered step, usually step 1, and steps are executed in sequential order, unless otherwise specified. In step E3, the imperative "go back to step E1" specifies the computational order in an obvious fashion. In step E2, the action is prefaced by the condition "if  $r = 0$ "; so if  $r \neq 0$ , the rest of that sentence does not apply and no action is specified. We might have added the redundant sentence, "If  $r \neq 0$ , go on to step E3."

The heavy vertical line, "■", appearing at the end of step E3 is used to indicate the end of an algorithm and the resumption of text.

We have now discussed virtually all the notational conventions used in the algorithms of this book, except for a notation used to denote "subscripted" or "indexed" items which are elements of an ordered array. Suppose we have  $n$  quantities,  $v_1, v_2, \dots, v_n$ ; instead of writing  $v_j$  for the  $j$ 'th element, the notation  $v[j]$  is often used. Similarly,  $a[i, j]$  is sometimes used in preference to a doubly-subscripted notation like  $a_{ij}$ . Sometimes multiple-letter names are used for

variables and are usually set in capital letters, e.g., TEMP might be the name of a variable used for temporarily holding a computed value, PRIME[K] might denote the Kth prime number, etc.

So much for the *form* of algorithms; now let us *perform* one. It should be mentioned immediately that the reader should *not* expect to read an algorithm as he reads a novel; such an attempt would make it pretty difficult to understand what is going on. An algorithm must be seen to be believed, and the best way to learn what an algorithm is all about is to try it. The reader should always take pencil and paper and work through an example of each algorithm immediately upon encountering it in the text. Usually the outline of a worked example will be given, or else the reader can easily conjure one up. This is a simple and painless method for obtaining an understanding of a given algorithm, and all other approaches are generally unsuccessful.

Let us therefore work out an example of Algorithm E. Suppose that we are given  $m = 119$  and  $n = 544$ ; we are ready to begin, at step E1. (The reader should now follow the algorithm as we give a play-by-play account.) Dividing  $m$  by  $n$  in this case is quite simple, almost too simple, since the quotient is zero and the remainder is 119. Thus,  $r \leftarrow 119$ . We proceed to step E2, and since  $r \neq 0$  no action occurs. In step E3 we set  $m \leftarrow 544$ ,  $n \leftarrow 119$ . It is clear that if  $m < n$  originally, the quotient in step E1 will always be zero and the algorithm will always proceed to interchange  $m$  and  $n$  in this rather cumbersome fashion. We could add a new step:

“E0. [Ensure  $m \geq n$ .] If  $m < n$ , exchange  $m \leftrightarrow n$ .”

if desired, without making an essential change in the algorithm except to increase its length as well as to decrease the time required to perform it in about one half of the cases.

Back at step E1, we find that  $\frac{544}{119} = 4\frac{68}{119}$ , so  $r \leftarrow 68$ . Again E2 is inapplicable, and at E3 we set  $m \leftarrow 119$ ,  $n \leftarrow 68$ . The next round sets  $r \leftarrow 51$ , and ultimately  $m \leftarrow 68$ ,  $n \leftarrow 51$ . Next  $r \leftarrow 17$ , and  $m \leftarrow 51$ ,  $n \leftarrow 17$ . Finally, when 51 is divided by 17,  $r \leftarrow 0$ , so at step E2 the algorithm terminates. The greatest common divisor of 119 and 544 is 17.

So this is an algorithm. The modern meaning for algorithm is quite similar to that of *recipe*, *process*, *method*, *technique*, *procedure*, *routine*, except that the word “algorithm” connotes something just a little different. Besides merely being a finite set of rules which gives a sequence of operations for solving a specific type of problem, an algorithm has five important features:

1) **Finiteness.** An algorithm must always terminate after a finite number of steps. Algorithm E satisfies this condition, because after step E1 the value of  $r$  is less than  $n$ , so if  $r \neq 0$ , the value of  $n$  decreases the next time that step E1 is encountered. A decreasing sequence of positive integers must eventually terminate, so step E1 is executed only a finite number of times for any given original value of  $n$ . Note, however, that the number of steps can become arbi-

trarily large; certain huge choices of  $m$  and  $n$  will cause step E1 to be executed over a million times.

(A procedure which has all of the characteristics of an algorithm except that it possibly lacks finiteness may be called a “computational method.” Besides his algorithm for the greatest common divisor of two integers, Euclid also gave a geometrical construction that is essentially equivalent to Algorithm E, except it is a procedure for obtaining the “greatest common measure” of the lengths of two line segments; this is a computational method that does not terminate if the given lengths are “incommensurate.”)

2) **Definiteness.** Each step of an algorithm must be precisely defined; the actions to be carried out must be rigorously and unambiguously specified for each case. The algorithms of this book will hopefully meet this criterion, but since they are specified in the English language, there is a possibility the reader might not understand exactly what the author intended. To get around this difficulty, formally defined *programming languages* or *computer languages* are designed for specifying algorithms, in which every statement has a very definite meaning. Many of the algorithms of this book will be given both in English and in a computer language. An expression of a computational method in a computer language is called a *program*.

In Algorithm E, the criterion of definiteness as applied to step E1 means that the reader is supposed to understand exactly what it means to divide  $m$  by  $n$  and what the remainder is. In actual fact, there is no universal agreement about what this means if  $m$  and  $n$  are not positive integers; what is the remainder of  $-8$  divided by  $-\pi$ ? What is the remainder of  $59/13$  divided by zero? Therefore the criterion of definiteness means we must make sure the values of  $m$  and  $n$  are always positive integers whenever step E1 is to be executed. This is initially true, by hypothesis, and after step E1  $r$  is a nonnegative integer which must be nonzero if we get to step E3; so  $m$  and  $n$  are indeed positive integers as required.

3) **Input.** An algorithm has zero or more inputs, i.e., quantities which are given to it initially before the algorithm begins. These inputs are taken from specified sets of objects. In Algorithm E, for example, there are two inputs, namely  $m$  and  $n$ , which are both taken from the set of *positive integers*.

4) **Output.** An algorithm has one or more outputs, i.e., quantities which have a specified relation to the inputs. Algorithm E has one output, namely  $n$  in step E2, which is the greatest common divisor of the two inputs.

(We can easily *prove* that this number is indeed the greatest common divisor, as follows. After step E1, we have

$$m = qn + r,$$

for some integer  $q$ . If  $r = 0$ , then  $m$  is a multiple of  $n$ , and clearly in such a case  $n$  is the greatest common divisor of  $m$  and  $n$ . If  $r \neq 0$ , note that any number which divides both  $m$  and  $n$  must divide  $m - qn = r$ , and any number which

divides both  $n$  and  $r$  must divide  $qn + r = m$ ; so the set of divisors of  $m, n$  is the same as the set of divisors of  $n, r$  and, in particular, the *greatest* common divisor of  $m, n$  is the same as the greatest common divisor of  $n, r$ . Therefore step E3 does not change the answer to the original problem.)

5) **Effectiveness.** An algorithm is also generally expected to be *effective*. This means that all of the operations to be performed in the algorithm must be sufficiently basic that they can in principle be done exactly and in a finite length of time by a man using pencil and paper. Algorithm E uses only the operations of dividing one positive integer by another, testing if an integer is zero, and setting the value of one variable equal to the value of another. These operations are effective, because integers can be represented on paper in a finite manner and there is at least one method (the "division algorithm") for dividing one by another. But the same operations would *not* be effective if the values involved were arbitrary real numbers specified by an infinite decimal expansion, nor if the values were the lengths of physical line segments, which cannot be specified exactly. Another example of a noneffective step is, "If 2 is the largest integer  $n$  for which there is a solution to the equation  $x^n + y^n = z^n$  in positive integers  $x, y,$  and  $z,$  then go to step E4." Such a statement would not be an effective operation until someone succeeds in showing that there is an algorithm to determine whether 2 is or is not the largest integer with the stated property.

Let us try to compare the concept of an algorithm with that of a cookbook recipe: A recipe presumably has the qualities of finiteness (although it is said that a watched pot never boils), input (eggs, flour, etc.) and output (TV dinner, etc.) but notoriously lacks definiteness. There are frequent cases in which the definiteness is missing, e.g., "Add a dash of salt." A "dash" is defined as "less than  $\frac{1}{2}$  teaspoon"; salt is perhaps well enough defined; but where should the salt be added (on top, side, etc.)? Instructions like "toss lightly until mixture is crumbly," "warm cognac in small saucepan," etc., are quite adequate as explanations to a trained cook, perhaps, but an algorithm must be specified to such a degree that even a computer can follow the directions. Still, a computer programmer can learn much by studying a good recipe book. (In fact, the author has barely resisted the temptation to name the present volume "The Programmer's Cookbook." Perhaps someday he will attempt a book called "Algorithms for the Kitchen.")

We should remark that the "finiteness" restriction is really not strong enough for practical use; a useful algorithm should require not only a finite number of steps, but a very finite number, a reasonable number. For example, there is an algorithm which determines whether or not the game of chess is a forced victory for the White pieces (see exercise 2.2.3-28); here is an algorithm which can solve a problem of intense interest to thousands of people, yet it is a safe bet that we will never in our lifetimes know the answer to this problem, because the algorithm requires fantastically large amounts of time for its execution, even though it is "finite." See also Chapter 8 for a discussion of some finite numbers which are so large as to actually be beyond comprehension.

In practice we not only want algorithms, we want good algorithms in some loosely-defined aesthetic sense. One criterion of goodness is the length of time taken to perform the algorithm; this can be expressed in terms of the number of times each step is executed. Other criteria are the adaptability of the algorithm to computers, its simplicity and elegance, etc.

Occasionally, we will have several algorithms for the same problem, and we must decide which is best. This leads us to the extremely interesting and all-important field of algorithmic analysis: given an algorithm, the problem is to determine its performance characteristics.

For example, we can consider Euclid's algorithm from this point of view. Suppose we ask the question, "Assuming that the value of  $n$  is known but  $m$  is allowed to range over all positive integers, what is the *average* number of times,  $T_n,$  that step E1 of Algorithm E will be performed?" In the first place, we have to check that this question does have a meaningful answer (since we are trying to take an average over infinitely many choices for  $m$ ). But it is evident that after the first execution of step E1 only the remainder of  $m$  after division by  $n$  is relevant. So all we must do to find the average,  $T_n,$  is to try the algorithm for  $m = 1, m = 2, \dots, m = n,$  count the total number of times step E1 has been executed, and divide by  $n.$

Now the important question is to determine the *nature* of  $T_n;$  is it approximately equal to  $\frac{1}{2}n,$  or  $\sqrt{n},$  etc.? As a matter of fact, the answer to this question is an extremely difficult and fascinating mathematical problem, not yet completely resolved, which is examined in more detail in Section 4.5.3. For large values of  $n$  it is possible to prove that  $T_n$  is approximately  $(12 \ln 2 / \pi^2) \ln n,$  that is, proportional to the *natural logarithm* of  $n,$  with a constant of proportionality that might not have been guessed offhand! For further details about Euclid's algorithm, and other ways to calculate the greatest common divisor, see Section 4.5.

"Analysis of algorithms" is the name the author likes to use to describe investigations such as this. The general idea is to take a particular algorithm and to determine its average behavior; occasionally we also study whether or not an algorithm is "optimal" in some sense. The theory of algorithms is another subject entirely, dealing primarily with the existence or nonexistence of effective algorithms to compute particular quantities; such theory is not investigated very deeply in this set of books, although it is considered briefly in Chapter 11.

So far our discussion of algorithms has been rather imprecise, and a mathematically oriented reader is justified in thinking that the preceding commentary makes a very shaky foundation on which to erect any theory about algorithms. We therefore close this section with a brief indication of one method by which the concept of algorithm can be firmly grounded in terms of mathematical set theory. Let us formally define a computational method to be a quadruple  $(Q, I, \Omega, f),$  in which  $Q$  is a set containing subsets  $I$  and  $\Omega,$  and  $f$  is a function from  $Q$  into itself. Furthermore  $f$  should leave  $\Omega$  pointwise fixed; that is,  $f(q)$  should equal  $q$  for all elements  $q$  of  $\Omega.$  The four quantities  $Q, I, \Omega, f$  are intended to represent respectively the states of the computation, the input, the output,

and the computational rule. Each input  $x$  in the set  $I$  defines a computational sequence,  $x_0, x_1, x_2, \dots$ , as follows:

$$x_0 = x \quad \text{and} \quad x_{k+1} = f(x_k) \quad \text{for} \quad k \geq 0. \quad (1)$$

The computational sequence is said to terminate in  $k$  steps if  $k$  is the smallest integer for which  $x_k$  is in  $\Omega$ , and in this case it is said to produce the output  $x_k$  from  $x$ . (Note that if  $x_k$  is in  $\Omega$ , so is  $x_{k+1}$ , because  $x_{k+1} = x_k$  in such a case.) Some computational sequences may never terminate; an *algorithm* is a computational method which terminates in finitely many steps for all  $x$  in  $I$ .

Algorithm E may, for example, be formalized in these terms as follows: Let  $Q$  be the set of all singletons  $(n)$ , all ordered pairs  $(m, n)$ , and all ordered quadruples  $(m, n, r, 1)$ ,  $(m, n, r, 2)$ , and  $(m, n, p, 3)$ , where  $m, n$ , and  $p$  are positive integers and  $r$  is a nonnegative integer. Let  $I$  be the subset of all pairs  $(m, n)$  and let  $\Omega$  be the subset of all singletons  $(n)$ . Let  $f$  be defined as follows:

$$\begin{aligned} f(m, n) &= (m, n, 0, 1); & f(n) &= (n); \\ f(m, n, r, 1) &= (m, n, \text{remainder of } m \text{ divided by } n, 2); \\ f(m, n, r, 2) &= (n) \text{ if } r = 0, (m, n, r, 3) \text{ otherwise}; \\ f(m, n, p, 3) &= (n, p, p, 1). \end{aligned} \quad (2)$$

The correspondence between this notation and Algorithm E is evident.

The above formulation of the concept "algorithm" does not include the restriction of "effectiveness" mentioned earlier; for example,  $Q$  might denote infinite sequences which are not computable by pencil and paper methods, or  $f$  might involve operations that mortal man cannot always perform. If we wish to restrict the notion of algorithm so that only elementary operations are involved, we can place restrictions on  $Q, I, \Omega$ , and  $f$ , for example as follows: Let  $A$  be a finite set of letters, and let  $A^*$  be the set of all strings on  $A$  (i.e., the set of all ordered sequences  $x_1x_2 \dots x_n$ , where  $n \geq 0$  and  $x_j$  is in  $A$  for  $1 \leq j \leq n$ ). The idea is to encode the states of the computation so that they are represented by strings of  $A^*$ . Now let  $N$  be a nonnegative integer and let  $Q$  be the set of all  $(\sigma, j)$ , where  $\sigma$  is in  $A^*$  and  $j$  is an integer,  $0 \leq j \leq N$ ; let  $I$  be the subset of  $Q$  with  $j = 0$  and let  $\Omega$  be the subset with  $j = N$ . If  $\theta$  and  $\sigma$  are strings in  $A^*$ , we say that  $\theta$  occurs in  $\sigma$  if  $\sigma$  has the form  $\alpha\theta\omega$  for strings  $\alpha$  and  $\omega$ . To complete our definition, let  $f$  be a function of the following type, defined by the strings  $\theta_j, \phi_j$  and the integers  $a_j, b_j$  for  $0 \leq j < N$ :

$$\begin{aligned} f(\sigma, j) &= (\sigma, a_j) && \text{if } \theta_j \text{ does not occur in } \sigma; \\ f(\sigma, j) &= (\alpha\phi_j\omega, b_j) && \text{if } \alpha \text{ is the shortest possible string} \\ &&& \text{for which } \sigma = \alpha\theta_j\omega; \\ f(\sigma, N) &= (\sigma, N). \end{aligned} \quad (3)$$

Such a computational method is clearly "effective," and experience shows that it is also powerful enough to do anything we can do by hand. There are many

other essentially equivalent ways to formulate the concept of an effective computational method (for example, using Turing machines). The above formulation is virtually the same as that given by A. A. Markov in 1951, in his book *The Theory of Algorithms* (tr. from the Russian by J. J. Schorr-Kon, U.S. Dept. of Commerce, Office of Technical Services, number OTS 60-51085).

## EXERCISES

1. [10] The text showed how to interchange the values of variables  $m$  and  $n$ , using the replacement notation, by setting  $t \leftarrow m, m \leftarrow n, n \leftarrow t$ . Show how the values  $(a, b, c, d)$  of four variables can be rearranged to  $(b, c, d, a)$  by a sequence of replacements. In other words, the new value of  $a$  is to be the original value of  $b$ , etc. Try to use the minimum number of replacements.
2. [15] Prove that  $m$  is always greater than  $n$  at the beginning of step E1, except possibly the first time this step occurs.
3. [20] Change Algorithm E (for the sake of efficiency) so that at step E3 we do not interchange values but immediately divide  $n$  by  $r$  and let  $m$  be the remainder. Add appropriate new steps so as to avoid all trivial replacement operations. Write this new algorithm in the style of Algorithm E, and call it Algorithm F.
4. [16] What is the greatest common divisor of 2166 and 6099?
- ▶ 5. [12] Show that the "Procedure for Reading This Set of Books" which appears in the preface actually fails to be a genuine algorithm on three of our five counts! Also mention some differences in format between it and Algorithm E.
6. [20] What is  $T_5$ , according to the notation near the end of this section?
- ▶ 7. [M21] Suppose that  $m$  is known and  $n$  is allowed to range over all positive integers; let  $U_m$  be the average number of times that step E1 is executed in Algorithm E. Show that  $U_m$  is well defined. Is  $U_m$  in any way related to  $T_m$ ?
8. [M25] Give an "effective" formal algorithm for computing the greatest common divisor of positive integers  $m$  and  $n$ , by specifying  $\theta_j, \phi_j, a_j, b_j$  as in Eqs. (3). Let the input be represented by the string  $a^m b^n$ , that is,  $m$   $a$ 's followed by  $n$   $b$ 's. Try to make your solution as simple as possible. [Hint: Use Algorithm E, but instead of division in step E1, set  $r \leftarrow |m - n|, n \leftarrow \min(m, n)$ .]
- ▶ 9. [M30] Suppose that  $C_1 = (Q_1, I_1, \Omega_1, f_1)$  and  $C_2 = (Q_2, I_2, \Omega_2, f_2)$  are computational methods. For example,  $C_1$  might stand for Algorithm E as in Eqs. (2), except that  $m, n$  are restricted in magnitude, and  $C_2$  might stand for a computer program implementation of Algorithm E. ( $Q_2$  might be the set of all states of the machine, i.e., all possible configurations of its memory and registers;  $f_2$  might be the definition of single machine actions; and  $I_2$  might be the initial state including the program for determining the greatest common divisor, as well as the values of  $m$  and  $n$ .) Formulate a set-theoretic definition for the concept " $C_2$  is a representation of  $C_1$ ": This is to mean intuitively that any computation sequence of  $C_1$  is mimicked by  $C_2$ , except that  $C_2$  might take more steps in which to do the computation and it might retain more information in its states. (We thereby obtain a rigorous interpretation of the statement, "Program  $X$  is an implementation of Algorithm  $Y$ .")