

Chapter 1

Introduction and Overview

1.1 INTRODUCTION

This chapter introduces the subject of data structures and presents an overview of the content of the text. Basic terminology and concepts will be defined and relevant examples provided. An overview of data organization and certain data structures will be covered along with a discussion of the different operations which are applied to these data structures. Last, we will introduce the notion of an algorithm and its complexity, and we will discuss the time-space tradeoff that may occur in choosing a particular algorithm and data structure for a given problem.

1.2 BASIC TERMINOLOGY; ELEMENTARY DATA ORGANIZATION

Data are simply values or sets of values. A *data item* refers to a single unit of values. Data items that are divided into subitems are called *group items*; those that are not are called *elementary items*. For example, an employee's name may be divided into three subitems—first name, middle initial and last name—but the social security number would normally be treated as a single item.

Collections of data are frequently organized into a hierarchy of *fields*, *records* and *files*. In order to make these terms more precise, we introduce some additional terminology.

An entity is something that has certain attributes or properties which may be assigned values. The values themselves may be either numeric or nonnumeric. For example, the following are possible attributes and their corresponding values for an entity, an employee of a given organization:

Attributes:	Name	Age	Sex	Social Security Number
Values:	ROHLAND, GAIL	34	F	134-24-5533

Entities with similar attributes (e.g., all the employees in an organization) form an *entity set*. Each attribute of an entity set has a *range* of values, the set of all possible values that could be assigned to the particular attribute.

The term "information" is sometimes used for data with given attributes, or, in other words, meaningful or processed data.

The way that data are organized into the hierarchy of fields, records and files reflects the relationship between attributes, entities and entity sets. That is, a *field* is a single elementary unit of information representing an attribute of an entity, a *record* is the collection of field values of a given entity and a *file* is the collection of records of the entities in a given entity set.

Each record in a file may contain many field items, but the value in a certain field may uniquely determine the record in the file. Such a field K is called a *primary key*, and the values k_1, k_2, \dots in such a field are called *keys* or *key values*.

EXAMPLE 1.1

- (a) Suppose an automobile dealership maintains an inventory file where each record contains the following data:

Serial Number, Type, Year, Price, Accessories

The Serial Number field can serve as a primary key for the file, since each automobile has a unique serial number.

- (b) Suppose an organization maintains a membership file where each record contains the following data:

Name, Address, Telephone Number, Dues Owed

Although there are four data items, Name and Address may be group items. Here the Name field is a

primary key. Note that the Address and Telephone Number fields may not serve as primary keys, since some members may belong to the same family and have the same address and telephone number.

Records may also be classified according to length. A file can have fixed-length records or variable-length records. In *fixed-length records*, all the records contain the same data items with the same amount of space assigned to each data item. In *variable-length records*, file records may contain different lengths. For example, student records usually have variable lengths, since different students take different numbers of courses. Usually, variable-length records have a minimum and a maximum length.

The above organization of data into fields, records and files may not be complex enough to maintain and efficiently process certain collections of data. For this reason, data are also organized into more complex types of structures. The study of such data structures, which forms the subject matter of this text, includes the following three steps:

- (1) Logical or mathematical description of the structure
- (2) Implementation of the structure on a computer
- (3) Quantitative analysis of the structure, which includes determining the amount of memory needed to store the structure and the time required to process the structure

The next section introduces us to some of these data structures.

Remark: The second and third of the steps in the study of data structures depend on whether the data are stored (a) in the main (primary) memory of the computer or (b) in a secondary (external) storage unit. This text will mainly cover the first case. This means that, given the address of a memory location, the time required to access the content of the memory cell does not depend on the particular cell or upon the previous cell accessed. The second case, called *file management* or *data base management*, is a subject unto itself and lies beyond the scope of this text.

1.3 DATA STRUCTURES

Data may be organized in many different ways; the logical or mathematical model of a particular organization of data is called a *data structure*. The choice of a particular data model depends on two considerations. First, it must be rich enough in structure to mirror the actual relationships of the data in the real world. On the other hand, the structure should be simple enough that one can effectively process the data when necessary. This section will introduce us to some of the data structures which will be discussed in detail later in the text.

Arrays

The simplest type of data structure is a *linear* (or *one-dimensional*) *array*. By a linear array, we mean a list of a finite number n of similar data elements referenced respectively by a set of n consecutive numbers, usually $1, 2, 3, \dots, n$. If we choose the name A for the array, then the elements of A are denoted by subscript notation

$$a_1, a_2, a_3, \dots, a_n$$

or by the parenthesis notation

$$A(1), A(2), A(3), \dots, A(N)$$

or by the bracket notation

$$A[1], A[2], A[3], \dots, A[N]$$

Regardless of the notation, the number K in $A[K]$ is called a *subscript* and $A[K]$ is called a *subscripted variable*.

Remark: The parentheses notation and the bracket notation are frequently used when the array name consists of more than one letter or when the array name appears in an algorithm. When using this

notation we will use ordinary uppercase letters for the name and subscripts as indicated above by the *A* and *N*. Otherwise, we may use the usual subscript notation of italics for the name and subscripts and lowercase letters for the subscripts as indicated above by the *a* and *n*. The former notation follows the practice of computer-oriented texts whereas the latter notation follows the practice of mathematics in print.

EXAMPLE 1.2

A linear array *STUDENT* consisting of the names of six students is pictured in Fig. 1-1. Here *STUDENT*[1] denotes John Brown, *STUDENT*[2] denotes Sandra Gold, and so on.

STUDENT	
1.	John Brown
2	Sandra Gold
3	Tom Jones
4	June Kelly
5	Mary Reed
6	Alan Smith

Fig. 1-1

Linear arrays are called one-dimensional arrays because each element in such an array is referenced by one subscript. A *two-dimensional array* is a collection of similar data elements where each element is referenced by two subscripts. (Such arrays are called *matrices* in mathematics, and *tables* in business applications.) Multidimensional arrays are defined analogously. Arrays will be covered in detail in Chap. 4.

EXAMPLE 1.3

A chain of 28 stores, each store having 4 departments, may list its weekly sales (to the nearest dollar) as in Fig. 1-2. Such data can be stored in the computer using a two-dimensional array in which the first subscript denotes the store and the second subscript the department. If *SALES* is the name given to the array, then

$$SALES[1, 1] = 2872, \quad SALES[1, 2] = 805, \quad SALES[1, 3] = 3211, \dots, \quad SALES[28, 4] = 982$$

The size of this array is denoted by 28×4 (read 28 by 4), since it contains 28 *rows* (the horizontal lines of numbers) and 4 *columns* (the vertical lines of numbers).

Dept. Store	1	2	3	4
1	2872	805	3211	1560
2	2196	1223	2525	1744
3	3257	1017	3686	1951
...
28	2618	931	2333	982

Fig. 1-2

Linked Lists

Linked lists will be introduced by means of an example. Suppose a brokerage firm maintains a file where each record contains a customer's name and his or her salesperson, and suppose the file contains the data appearing in Fig. 1-3. Clearly the file could be stored in the computer by such a table, i.e., by two columns of nine names. However, this may not be the most useful way to store the data, as the following discussion shows.

	Customer	Salesperson
1	Adams	Smith
2	Brown	Ray
3	Clark	Jones
4	Drew	Ray
5	Evans	Smith
6	Farmer	Jones
7	Geller	Ray
8	Hill	Smith
9	Infeld	Ray

Fig. 1-3

Another way of storing the data in Fig. 1-3 is to have a separate array for the salespeople and an entry (called a *pointer*) in the customer file which gives the location of each customer's salesperson. This is done in Fig. 1-4, where some of the pointers are pictured by an arrow from the location of the pointer to the location of the corresponding salesperson. Practically speaking, an integer used as a pointer requires less space than a name; hence this representation saves space, especially if there are hundreds of customers for each salesperson.

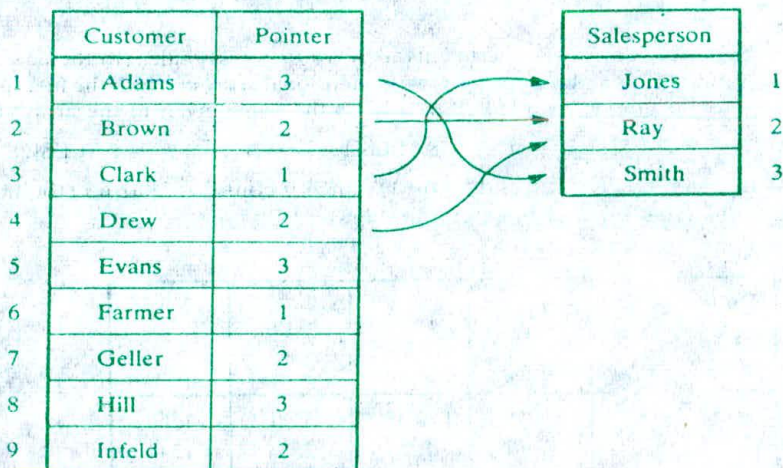


Fig. 1-4

Suppose the firm wants the list of customers for a given salesperson. Using the data representation in Fig. 1-4 the firm would have to search through the entire customer file. One way to simplify such a

search is to have the arrows in Fig. 1-4 point the other way; each salesperson would now have a set of pointers giving the positions of his or her customers, as in Fig. 1-5. The main disadvantage of this representation is that each salesperson may have many pointers and the set of pointers will change as customers are added and deleted.

	Salesperson	Pointer
1	Jones	3, 6
2	Ray	2, 4, 7, 9
3	Smith	1, 5, 8

Fig. 1-5

Another very popular way to store the type of data in Fig. 1-3 is shown in Fig. 1-6. Here each salesperson has one pointer which points to his or her first customer, whose pointer in turn points to the second customer, and so on, with the salesperson's last customer indicated by a 0. This is pictured with arrows in Fig. 1-6 for the salesperson Ray. Using this representation one can easily obtain the entire list of customers for a given salesperson and, as we will see in Chap. 5, one can easily insert and delete customers.

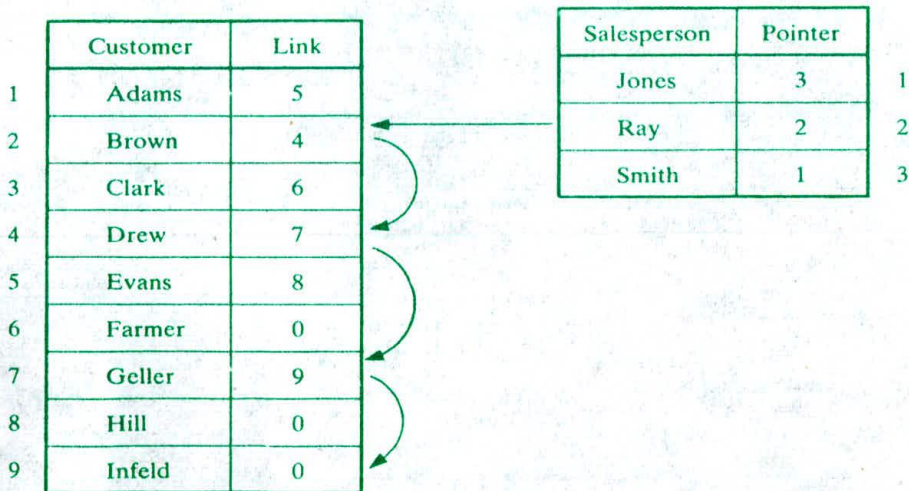


Fig. 1-6

The representation of the data in Fig. 1-6 is an example of linked lists. Although the terms "pointer" and "link" are usually used synonymously, we will try to use the term "pointer" when an element in one list points to an element in a different list, and to reserve the term "link" for the case when an element in a list points to an element in that same list.

Trees

Data frequently contain a hierarchical relationship between various elements. The data structure which reflects this relationship is called a *rooted tree graph* or, simply, a *tree*. Trees will be defined and discussed in detail in Chap. 7. Here we indicate some of their basic properties by means of two examples.

EXAMPLE 1.4 Record Structure

Although a file may be maintained by means of one or more arrays, a record, where one indicates both the group items and the elementary items, can best be described by means of a tree structure. For example, an employee personnel record may contain the following data items:

Social Security Number, Name, Address, Age, Salary, Dependents

However, Name may be a group item with the subitems Last, First and MI (middle initial). Also, Address may be a group item with the subitems Street address and Area address, where Area itself may be a group item having subitems City, State and ZIP code number. This hierarchical structure is pictured in Fig. 1-7(a). Another way of picturing such a tree structure is in terms of levels, as in Fig. 1-7(b).

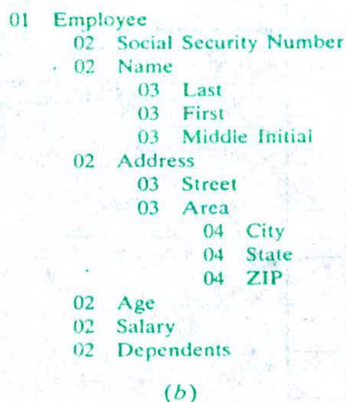
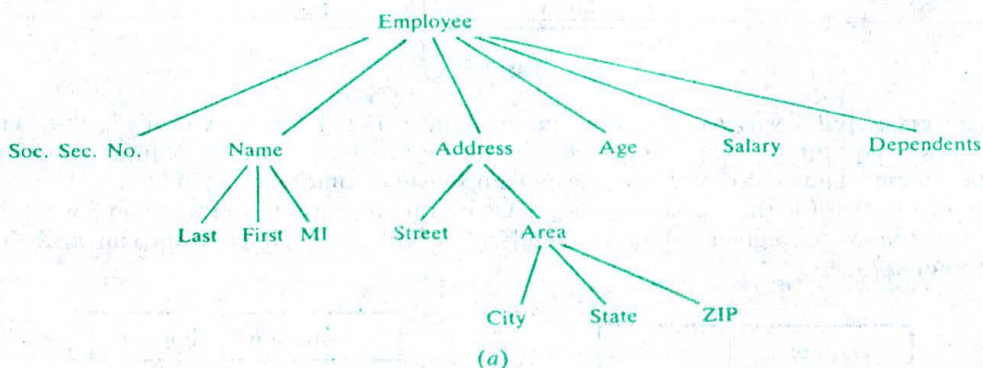


Fig. 1-7

EXAMPLE 1.5 Algebraic Expressions

Consider the algebraic expression

$$(2x + y)(a - 7b)^3$$

Using a vertical arrow (\uparrow) for exponentiation and an asterisk ($*$) for multiplication, we can represent the expression by the tree in Fig. 1-8. Observe that the order in which the operations will be performed is reflected in the diagram, the exponentiation must take place after the subtraction, and the multiplication at the top of the tree must be executed last.

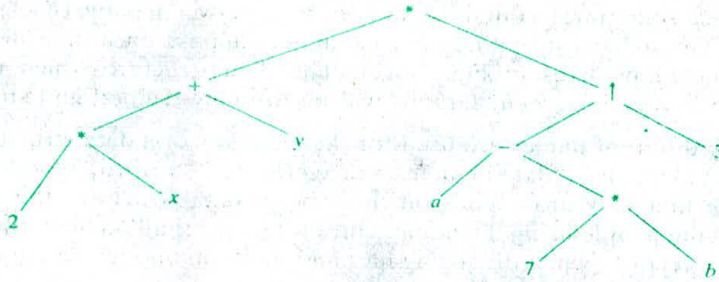
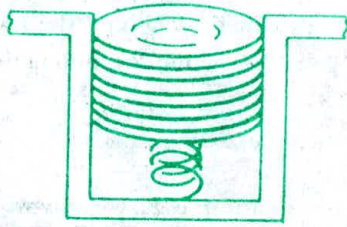


Fig. 1-8

There are data structures other than arrays, linked lists and trees which we shall study. Some of these structures are briefly described below.

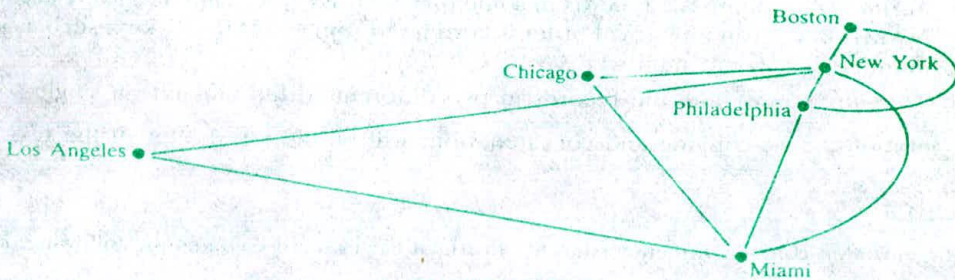
- (a) **Stack** A stack, also called a last-in first-out (LIFO) system, is a linear list in which insertions and deletions can take place only at one end, called the *top*. This structure is similar in its operation to a stack of dishes on a spring system, as pictured in Fig. 1-9(a). Note that new dishes are inserted only at the top of the stack and dishes can be deleted only from the top of the stack.



(a) Stack of dishes.



(b) Queue waiting for a bus.



(c) Airline flights.

Fig. 1-9

- (b) **Queue**. A queue, also called a first-in first-out (FIFO) system, is a linear list in which deletions can take place only at one end of the list, the "front" of the list, and insertions can take place only at the other end of the list, the "rear" of the list. This structure operates in much the same way as a line of people waiting at a bus stop, as pictured in Fig. 1-9(b): the first person in line is the first person to board the bus. Another analogy is with automobiles waiting to pass through an intersection—the first car in line is the first car through.

- (c) *Graph*. Data sometimes contain a relationship between pairs of elements which is not necessarily hierarchical in nature. For example, suppose an airline flies only between the cities connected by lines in Fig. 1-9(c). The data structure which reflects this type of relationship is called a *graph*. Graphs will be formally defined and studied in Chap. 8.

Remark: Many different names are used for the elements of a data structure. Some commonly used names are "data element," "data item," "item aggregate," "record," "node" and "data object." The particular name that is used depends on the type of data structure, the context in which the structure is used and the people using the name. Our preference shall be the term "data element," but we will use the term "record" when discussing files and the term "node" when discussing linked lists, trees and graphs.

1.4 DATA STRUCTURE OPERATIONS

The data appearing in our data structures are processed by means of certain operations. In fact, the particular data structure that one chooses for a given situation depends largely on the frequency with which specific operations are performed. This section introduces the reader to some of the most frequently used of these operations.

The following four operations play a major role in this text:

- B {
- (1) *Traversing*: Accessing each record exactly once so that certain items in the record may be processed. (This accessing and processing is sometimes called "visiting" the record.)
 - (2) *Searching*: Finding the location of the record with a given key value, or finding the locations of all records which satisfy one or more conditions.
 - (3) *Inserting*: Adding a new record to the structure.
 - (4) *Deleting*: Removing a record from the structure.

Sometimes two or more of the operations may be used in a given situation; e.g., we may want to delete the record with a given key, which may mean we first need to search for the location of the record.

The following two operations, which are used in special situations, will also be considered:

- (1) *Sorting*: Arranging the records in some logical order (e.g., alphabetically according to some NAME key, or in numerical order according to some NUMBER key, such as social security number or account number)
- (2) *Merging*: Combining the records in two different sorted files into a single sorted file

Other operations, e.g., copying and concatenation, will be discussed later in the text.

EXAMPLE 1.6

An organization contains a membership file in which each record contains the following data for a given member:

Name,	Address,	Telephone Number,	Age,	Sex
-------	----------	-------------------	------	-----

- (a) Suppose the organization wants to announce a meeting through a mailing. Then one would traverse the file to obtain Name and Address for each member.
- (b) Suppose one wants to find the names of all members living in a certain area. Again one would traverse the file to obtain the data.
- (c) Suppose one wants to obtain Address for a given Name. Then one would search the file for the record containing Name.
- (d) Suppose a new person joins the organization. Then one would insert his or her record into the file.
- (e) Suppose a member dies. Then one would delete his or her record from the file.

- (f) Suppose a member has moved and has a new address and telephone number. Given the name of the member, one would first need to search for the record in the file. Then one would perform the "update"—i.e., change items in the record with the new data.
- (g) Suppose one wants to find the number of members 65 or older. Again one would traverse the file, counting such members.

1.5 ALGORITHMS: COMPLEXITY, TIME-SPACE TRADEOFF

An algorithm is a well-defined list of steps for solving a particular problem. One major purpose of this text is to develop efficient algorithms for the processing of our data. The time and space it uses are two major measures of the efficiency of an algorithm. The complexity of an algorithm is the function which gives the running time and/or space in terms of the input size. (The notion of complexity will be treated in Chap. 2.)

Each of our algorithms will involve a particular data structure. Accordingly, we may not always be able to use the most efficient algorithm, since the choice of data structure depends on many things, including the type of data and the frequency with which various data operations are applied. Sometimes the choice of data structure involves a time-space tradeoff: by increasing the amount of space for storing the data, one may be able to reduce the time needed for processing the data, or vice versa. We illustrate these ideas with two examples.

Searching Algorithms

Consider a membership file, as in Example 1.6, in which each record contains, among other data, the name and telephone number of its member. Suppose we are given the name of a member and we want to find his or her telephone number. One way to do this is to linearly search through the file, i.e., to apply the following algorithm:

Linear Search: Search each record of the file, one at a time, until finding the given Name and hence the corresponding telephone number.

First of all, it is clear that the time required to execute the algorithm is proportional to the number of comparisons. Also, assuming that each name in the file is equally likely to be picked, it is intuitively clear that the average number of comparisons for a file with n records is equal to $n/2$; that is, the complexity of the linear search algorithm is given by $C(n) = n/2$.

The above algorithm would be impossible in practice if we were searching through a list consisting of thousands of names, as in a telephone book. However, if the names are sorted alphabetically, as in telephone books, then we can use an efficient algorithm called binary search. This algorithm is discussed in detail in Chap. 4, but we briefly describe its general idea below.

Binary Search: Compare the given Name with the name in the middle of the list; this tells which half of the list contains Name. Then compare Name with the name in the middle of the correct half to determine which quarter of the list contains Name. Continue the process until finding Name in the list.

One can show that the complexity of the binary search algorithm is given by

$$C(n) = \log_2 n$$

Thus, for example, one will not require more than 15 comparisons to find a given Name in a list containing 25 000 names.

Although the binary search algorithm is a very efficient algorithm, it has some major drawbacks. Specifically, the algorithm assumes that one has direct access to the middle name in the list or a sublist. This means that the list must be stored in some type of array. Unfortunately, inserting an element in an array requires elements to be moved down the list, and deleting an element from an array requires elements to be moved up the list.

The telephone company solves the above problem by printing a new directory every year while keeping a separate temporary file for new telephone customers. That is, the telephone company updates its files every year. On the other hand, a bank may want to insert a new customer in its file almost instantaneously. Accordingly, a linearly sorted list may not be the best data structure for a bank.

An Example of Time-Space Tradeoff

Suppose a file of records contains names, social security numbers and much additional information among its fields. Sorting the file alphabetically and using a binary search is a very efficient way to find the record for a given name. On the other hand, suppose we are given only the social security number of the person. Then we would have to do a linear search for the record, which is extremely time-consuming for a very large number of records. How can we solve such a problem? One way is to have another file which is sorted numerically according to social security number. This, however, would double the space required for storing the data. Another way, pictured in Fig. 1-10, is to have the main file sorted numerically by social security number and to have an auxiliary array with only two columns, the first column containing an alphabetized list of the names and the second column containing pointers which give the locations of the corresponding records in the main file. This is one way of solving the problem that is used frequently, since the additional space, containing only two columns, is minimal for the amount of extra information it provides.

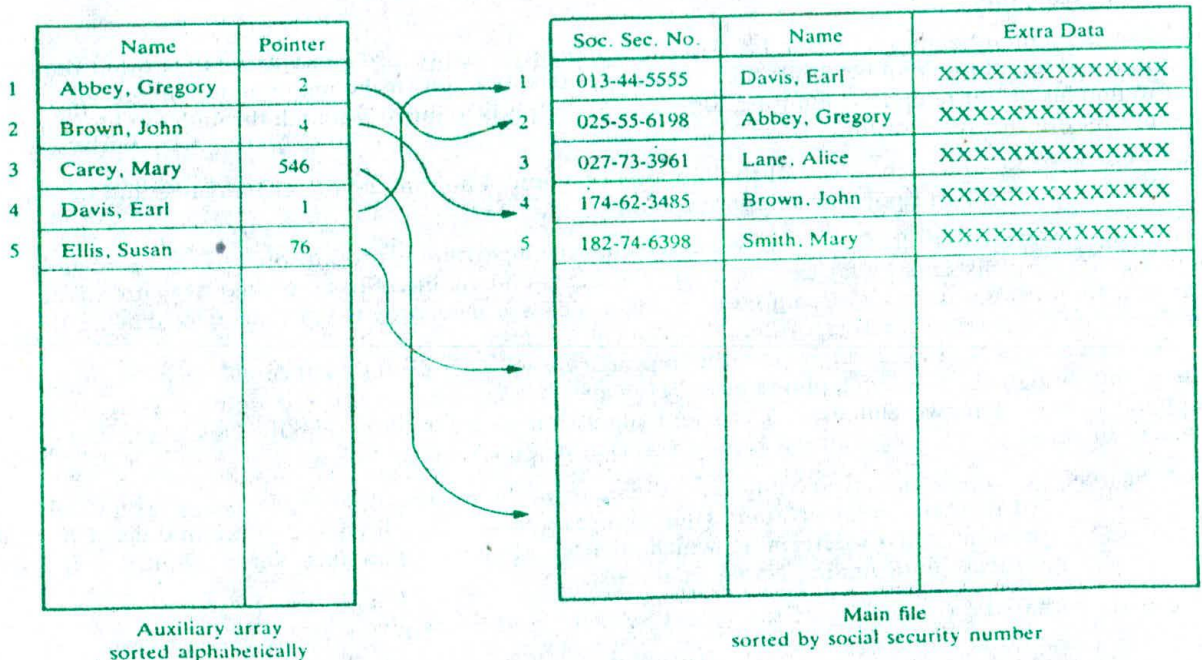


Fig. 1-10

Remark: Suppose a file is sorted numerically by social security number. As new records are inserted into the file, data must be constantly moved to new locations in order to maintain the sorted order. One simple way to minimize the movement of data is to have the social security number serve as the address of each record. Not only would there be no movement of data when records are inserted,

but there would be instant access to any record. However, this method of storing data would require one billion (10^9) memory locations for only hundreds or possibly thousands of records. Clearly, this tradeoff of space for time is not worth the expense. An alternative method is to define a function H from the set K of key values—social security numbers—into the set L of addresses of memory cells. Such a function H is called a *hashing function*. Hashing functions and their properties will be covered in Chap. 9.

Solved Problems

BASIC TERMINOLOGY

1.1. A professor keeps a class list containing the following data for each student:

Name, Major, Student Number, Test Scores, Final Grade

- State the entities, attributes and entity set of the list.
 - Describe the field values, records and file.
 - Which attributes can serve as primary keys for the list?
- Each student is an entity, and the collection of students is the entity set. The properties, name, major, and so on, of the students are the attributes.
 - The field values are the values assigned to the attributes, i.e., the actual names, test scores, and so on. The field values for each student constitute a record, and the collection of all the student records is the file.
 - Either Name or Student Number can serve as a primary key, since each uniquely determines the student's record. Normally the professor uses Name as the primary key, but the registrar may use Student Number.

1.2. A hospital maintains a patient file in which each record contains the following data:

Name, Admission Date, Social Security Number, Room, Bed Number, Doctor

- Which items can serve as primary keys?
 - Which pair of items can serve as a primary key?
 - Which items can be group items?
- Name and Social Security Number can serve as primary keys. (We assume that no two patients have the same name.)
 - Room and Bed Number in combination also uniquely determine a given patient.
 - Name, Admission Date and Doctor may be group items.

1.3. Which of the following data items may lead to variable-length records when included as items in the record: (a) age, (b) sex, (c) name of spouse, (d) names of children, (e) education, (f) previous employers?

Since (d) and (f) may contain a few or many items, they may lead to variable-length records. Also, (e) may contain many items, unless it asks only for the highest level obtained.

1.4 Data base systems will be only briefly covered in this text. Why?

"Data base systems" refers to data stored in the secondary memory of the computer. The implementation and analysis of data structures in the secondary memory are very different from those in the main memory of the computer. This text is primarily concerned with data structures in main memory, not secondary memory.

DATA STRUCTURES AND OPERATIONS

1.5 Give a brief description of (a) traversing, (b) sorting and (c) searching.

- (a) Accessing and processing each record exactly once
- (b) Arranging the data in some given order
- (c) Finding the location of the record with a given key or keys

1.6 Give a brief description of (a) inserting and (b) deleting.

- (a) Adding a new record to the data structure, usually keeping a particular ordering
- (b) Removing a particular record from the data structure

1.7 Consider the linear array NAME in Fig. 1-11, which is sorted alphabetically.

- (a) Find NAME[2], NAME[4] and NAME[7].
- (b) Suppose Davis is to be inserted into the array. How many names must be moved to new locations?
- (c) Suppose Gupta is to be deleted from the array. How many names must be moved to new locations?

(a) Here NAME[K] is the k th name in the list. Hence,

$$\text{NAME}[2] = \text{Clark}, \quad \text{NAME}[4] = \text{Gupta}, \quad \text{NAME}[7] = \text{Pace}$$

- (b) Since Davis will be assigned to NAME[3], the names Evans through Smith must be moved. Hence six names are moved.
- (c) The names Jones through Smith must be moved up the array. Hence four names must be moved.

NAME	
1	Adams
2	Clark
3	Evans
4	Gupta
5	Jones
6	Lane
7	Pace
8	Smith

Fig. 1-11

1.8 Consider the linear array NAME in Fig. 1-12. The values of FIRST and LINK[K] in the figure determine a linear ordering of the names as follows. FIRST gives the location of the first name in the list, and LINK[K] gives the location of the name following NAME[K], with 0 denoting the end of the list. Find the linear ordering of the names.

The ordering is obtained as follows:

FIRST = 5, so the first name in the list is NAME[5], which is Brooks.

LINK[5] = 2, so the next name is NAME[2], which is Clark.

LINK[2] = 8, so the next name is NAME[8], which is Fisher.

LINK[8] = 4, so the next name is NAME[4], which is Hansen.

LINK[4] = 10, so the next name is NAME[10], which is Leary.

LINK[10] = 6, so the next name is NAME[6], which is Pitt.

LINK[6] = 1, so the next name is NAME[1], which is Rogers.

LINK[1] = 7, so the next name is NAME[7], which is Walker.

LINK[7] = 0, which indicates the end of the list.

Thus the linear ordering of the names is Brooks, Clark, Fisher, Hansen, Leary, Pitt, Rogers, Walker. Note that this is the alphabetical ordering of the names.

	FIRST		NAME	LINK
	5	1	Rogers	7
		2	Clark	8
		3		
		4	Hansen	10
		5	Brooks	2
		6	Pitt	1
		7	Walker	0
		8	Fisher	4
		9		
		10	Leary	6

Fig. 1-12

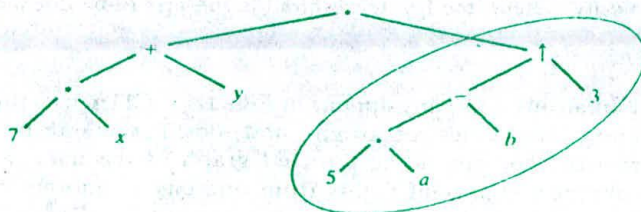


Fig. 1-13

1.9 Consider the algebraic expression $(7x + y)(5a - b)^3$. (a) Draw the corresponding tree diagram as in Example 1.5. (b) Find the scope of the exponential operation. (The scope of a node v in a tree is the subtree consisting of v and the nodes following v .)

- (a) Use a vertical arrow (\uparrow) for exponentiation and an asterisk (*) for multiplication to obtain the tree in Fig. 1-13.
- (b) The scope of the exponentiation operation \uparrow is the subtree circled in the diagram. It corresponds to the expression $(5a - b)^3$.

1.10 The following is a tree structure given by means of level numbers as discussed in Example 1.4:
 01 Employee 02 Name 02 Number 02 Hours 03 Regular 03 Overtime 02 Rate
 Draw the corresponding tree diagram.

The tree diagram appears in Fig. 1-14. Here each node v is the successor of the node which precedes v and has a lower level number than v .

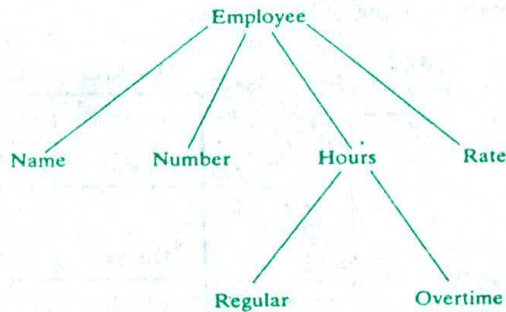


Fig. 1-14

1.11 Discuss whether a stack or a queue is the appropriate structure for determining the order in which elements are processed in each of the following situations.

- (a) Batch computer programs are submitted to the computer center.
- (b) Program A calls subprogram B, which calls subprogram C, and so on.
- (c) Employees have a contract which calls for a seniority system for hiring and firing.
- (a) Queue. Excluding priority cases, programs are executed on a first come, first served basis.
- (b) Stack. The last subprogram is executed first, and its results are transferred to the next-to-last program, which is then executed, and so on, until the original calling program is executed.
- (c) Stack. In a seniority system, the last to be hired is the first to be discharged.

1.12 The daily flights of an airline company appear in Fig. 1-15. CITY lists the cities, and ORIG[K] and DEST[K] denote the cities of origin and destination, respectively, of the flight NUMBER[K]. Draw the corresponding directed graph of the data. (The graph is directed because the flight numbers represent flights from one city to another but not returning.)

The nodes of the graph are the five cities. Draw an arrow from city A to city B if there is a flight from A to B, and label the arrow with the flight number. The directed graph appears in Fig. 1-16.

	CITY
1	Atlanta
2	Boston
3	Chicago
4	Miami
5	Philadelphia

(a)

	NUMBER	ORIG	DEST
1	701	2	3
2	702	3	2
3	705	5	3
4	708	3	4
5	711	2	5
6	712	5	2
7	713	5	1
8	715	1	4
9	717	5	4
10	718	4	5

(b)

Fig. 1-15

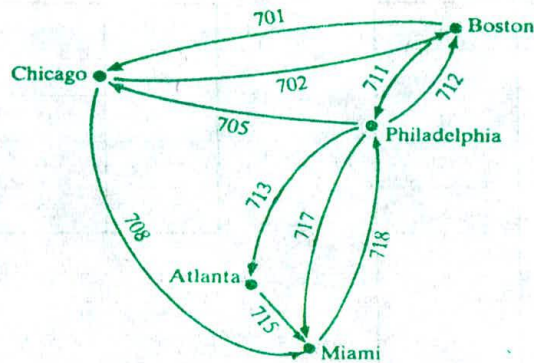


Fig. 1-16

COMPLEXITY; SPACE-TIME TRADEOFFS

1.13 Briefly describe the notions of (a) the complexity of an algorithm and (b) the space-time tradeoff of algorithms.

- (a) The complexity of an algorithm is a function $f(n)$ which measures the time and/or space used by an algorithm in terms of the input size n .
- (b) The space-time tradeoff refers to a choice between algorithmic solutions of a data processing problem that allows one to decrease the running time of an algorithmic solution by increasing the space to store the data and vice versa.

1.14 Suppose a data set S contains n elements.

- (a) Compare the running time T_1 of the linear search algorithm with the running time T_2 of the binary search algorithm when (i) $n = 1000$ and (ii) $n = 10\,000$.

- (b) Discuss searching for a given item in S when S is stored as a linked list.
- (a) Recall (Sec. 1.5) that the expected running of the linear search algorithm is $f(n) = n/2$ and that the binary search algorithm is $f(n) = \log_2 n$. Accordingly, (i) for $n = 1000$, $T_1 = 500$ but $T_2 = \log_2 1000 \approx 10$; and (ii) for $n = 10\,000$, $T_1 = 5000$ but $T_2 = \log_2 10\,000 \approx 14$.
- (b) The binary search algorithm assumes that one can directly access the middle element in the set S . But one cannot directly access the middle element in a linked list. Hence one may have to use a linear search algorithm when S is stored as a linked list.

1.15 Consider the data in Fig. 1-15, which gives the different flights of an airline. Discuss different ways of storing the data so as to decrease the time in executing the following:

- (a) Find the origin and destination of a flight, given the flight number.
- (b) Given city A and city B, find whether there is a flight from A to B, and if there is, find its flight number.
- (a) Store the data of Fig. 1-15(b) in arrays ORIG and DEST where the subscript is the flight number, as pictured in Fig. 1-17(a)
- (b) Store the data of Fig. 1-15(b) in a two-dimensional array FLIGHT where FLIGHT[J, K] contains the flight number of the flight from CITY[J] to CITY[K], or contains 0 when there is no such flight, as pictured in Fig. 1-17(b).

	ORIG	DEST
701	2	3
702	3	2
703	0	0
704	0	0
705	5	3
706	0	0
⋮	⋮	⋮
715	1	4
716	0	0
717	5	4
718	4	5

(a)

FLIGHT	1	2	3	4	5
1	0	0	0	715	0
2	0	0	701	0	711
3	0	702	0	708	0
4	0	0	0	0	718
5	713	712	705	717	0

(b)

Fig. 1-17

1.16 Suppose an airline serves n cities with s flights. Discuss drawbacks to the data representations used in Fig. 1-17(a) and Fig. 1-17(b).

- (a) Suppose the flight numbers are spaced very far apart; i.e., suppose the ratio of the number s of flights to the number of memory locations is very small, e.g., approximately 0.05. Then the extra storage space may not be worth the expense.
- (b) Suppose the ratio of the number s of flights to the number n of memory locations in the array FLIGHT is very small, i.e., that the array FLIGHT is one that contains a large number of zeros (such an array is called a sparse matrix). Then the extra storage space may not be worth the expense.

Chapter 2

Preliminaries

2.1 INTRODUCTION

The development of algorithms for the creation and processing of data structures is a major feature of this text. This chapter describes, by means of simple examples, the format that will be used to present our algorithms. The format we have selected is similar to the format used by Knuth in his well-known text *Fundamental Algorithms*. Although our format is language-free, the algorithms will be sufficiently well structured and detailed that they can be easily translated into some programming language such as Pascal, FORTRAN, PL/1 or BASIC. In fact, some of our algorithms will be translated into such languages in the problems sections.

Algorithms may be quite complex. The computer programs implementing the more complex algorithms can be more easily understood if these programs are organized into hierarchies of modules similar to the one in Fig. 2-1. In such an organization, each program contains first a main module, which gives a general description of the algorithm; this main module refers to certain submodules, which contain more detailed information than the main module; each of the submodules may refer to more detailed submodules; and so on. The organization of a program into such a hierarchy of modules normally requires the use of certain basic flow patterns and logical structures which are usually associated with the notion of structured programming. These flow patterns and logical structures will be reviewed in this chapter.

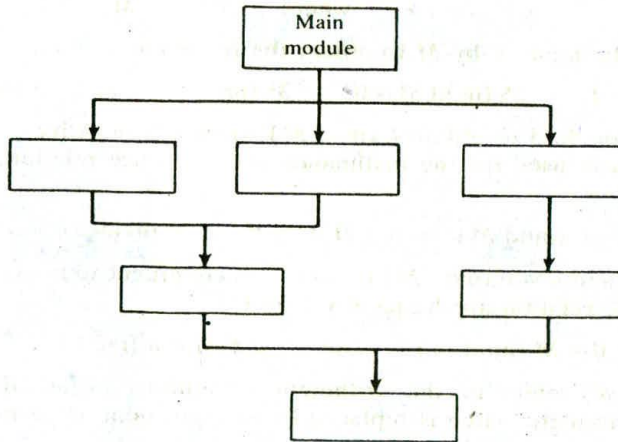


Fig. 2-1 A hierarchy of modules.

The chapter begins with a brief outline and discussion of various mathematical functions which occur in the study of algorithms and in computer science in general, and the chapter ends with a discussion of the different kinds of variables that can appear in our algorithms and programs.

The notion of the complexity of an algorithm is also covered in this chapter. This important measurement of algorithms gives us a tool to compare different algorithmic solutions to a particular problem such as searching or sorting. The concept of an algorithm and its complexity is fundamental not only to data structures but also to almost all areas of computer science.

2.2 MATHEMATICAL NOTATION AND FUNCTIONS

This section gives various mathematical functions which appear very often in the analysis of algorithms and in computer science in general, together with their notation.

Floor and Ceiling Functions

Let x be any real number. Then x lies between two integers called the floor and the ceiling of x . Specifically,

$\lfloor x \rfloor$, called the *floor* of x , denotes the greatest integer that does not exceed x .

$\lceil x \rceil$, called the *ceiling* of x , denotes the least integer that is not less than x .

If x is itself an integer, then $\lfloor x \rfloor = \lceil x \rceil = x$; otherwise $\lfloor x \rfloor + 1 = \lceil x \rceil$.

EXAMPLE 2.1

$$\begin{array}{llll} \lfloor 3.14 \rfloor = 3, & \lfloor \sqrt{5} \rfloor = 2, & \lfloor -8.5 \rfloor = -9, & \lceil 7 \rceil = 7 \\ \lceil 3.14 \rceil = 4, & \lceil \sqrt{5} \rceil = 3, & \lceil -8.5 \rceil = -8, & \lfloor 7 \rfloor = 7 \end{array}$$

Remainder Function; Modular Arithmetic

Let k be any integer and let M be a positive integer. Then

$$k \pmod{M}$$

(read k modulo M) will denote the integer remainder when k is divided by M . More exactly, $k \pmod{M}$ is the unique integer r such that

$$k = Mq + r \quad \text{where} \quad 0 \leq r < M$$

When k is positive, simply divide k by M to obtain the remainder r . Thus

$$25 \pmod{7} = 4, \quad 25 \pmod{5} = 0, \quad 35 \pmod{11} = 2, \quad 3 \pmod{8} = 3$$

Problem 2.2(b) shows a method to obtain $k \pmod{M}$ when k is negative.

The term "mod" is also used for the mathematical congruence relation, which is denoted and defined as follows:

$$a \equiv b \pmod{M} \quad \text{if and only if} \quad M \text{ divides } b - a$$

M is called the *modulus*, and $a \equiv b \pmod{M}$ is read " a is congruent to b modulo M ." The following aspects of the congruence relation are frequently useful:

$$0 \equiv M \pmod{M} \quad \text{and} \quad a \pm M \equiv a \pmod{M}$$

Arithmetic modulo M refers to the arithmetic operations of addition, multiplication and subtraction where the arithmetic value is replaced by its equivalent value in the set

$$\{0, 1, 2, \dots, M-1\}$$

or in the set

$$\{1, 2, 3, \dots, M\}$$

For example, in arithmetic modulo 12, sometimes called "clock" arithmetic,

$$6 + 9 \equiv 3, \quad 7 \times 5 \equiv 11, \quad 1 - 5 \equiv 8, \quad 2 + 10 \equiv 0 \equiv 12$$

(The use of 0 or M depends on the application.)

Integer and Absolute Value Functions

Let x be any real number. The *integer value* of x , written $\text{INT}(x)$, converts x into an integer by deleting (truncating) the fractional part of the number. Thus

$$\text{INT}(3.14) = 3, \quad \text{INT}(\sqrt{5}) = 2, \quad \text{INT}(-8.5) = -8, \quad \text{INT}(7) = 7$$

Observe that $\text{INT}(x) = [x]$ or $\text{INT}(x) = \lceil x \rceil$ according to whether x is positive or negative.

The *absolute value* of the real number x , written $\text{ABS}(x)$ or $|x|$, is defined as the greater of x or $-x$. Hence $\text{ABS}(0) = 0$, and, for $x \neq 0$, $\text{ABS}(x) = x$ or $\text{ABS}(x) = -x$, depending on whether x is positive or negative. Thus

$$|-15| = 15, \quad |7| = 7, \quad |-3.33| = 3.33, \quad |4.44| = 4.44, \quad |-0.075| = 0.075$$

We note that $|x| = |-x|$ and, for $x \neq 0$, $|x|$ is positive.

Summation Symbol; Sums

Here we introduce the summation symbol Σ (the Greek letter sigma). Consider a sequence a_1, a_2, a_3, \dots . Then the sums

$$a_1 + a_2 + \dots + a_n \quad \text{and} \quad a_m + a_{m+1} + \dots + a_n$$

will be denoted, respectively, by

$$\sum_{j=1}^n a_j \quad \text{and} \quad \sum_{j=m}^n a_j$$

The letter j in the above expressions is called a *dummy index* or *dummy variable*. Other letters frequently used as dummy variables are i , k , s and t .

EXAMPLE 2.2

$$\begin{aligned} \sum_{i=1}^n a_i b_i &= a_1 b_1 + a_2 b_2 + \dots + a_n b_n \\ \sum_{j=2}^5 j^2 &= 2^2 + 3^2 + 4^2 + 5^2 = 4 + 9 + 16 + 25 = 54 \\ \sum_{j=1}^n j &= 1 + 2 + \dots + n \end{aligned}$$

The last sum in Example 2.2 will appear very often. It has the value $n(n+1)/2$. That is,

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

Thus, for example,

$$1 + 2 + \dots + 50 = \frac{50(51)}{2} = 1275$$

Factorial Function

The product of the positive integers from 1 to n , inclusive, is denoted by $n!$ (read " n factorial"). That is,

$$n! = 1 \cdot 2 \cdot 3 \cdot \dots \cdot (n-2)(n-1)n$$

It is also convenient to define $0! = 1$.

EXAMPLE 2.3

(a) $2! = 1 \cdot 2 = 2; \quad 3! = 1 \cdot 2 \cdot 3 = 6; \quad 4! = 1 \cdot 2 \cdot 3 \cdot 4 = 24$

(b) For $n > 1$, we have $n! = n \cdot (n-1)!$. Hence

$$5! = 5 \cdot 4! = 5 \cdot 24 = 120; \quad 6! = 6 \cdot 5! = 6 \cdot 120 = 720$$

Permutations

A *permutation* of a set of n elements is an arrangement of the elements in a given order. For example, the permutations of the set consisting of the elements a, b, c are as follows:

$$abc, \quad acb, \quad bac, \quad bca, \quad cab, \quad cba$$

One can prove: *There are $n!$ permutations of a set of n elements.* Accordingly, there are $4! = 24$ permutations of a set with 4 elements, $5! = 120$ permutations of a set with 5 elements, and so on.

Exponents and Logarithms

Recall the following definitions for integer exponents (where m is a positive integer):

$$a^m = a \cdot a \cdots a \text{ (} m \text{ times)}, \quad a^0 = 1, \quad a^{-m} = \frac{1}{a^m}$$

Exponents are extended to include all rational numbers by defining, for any rational number m/n ,

$$a^{m/n} = \sqrt[n]{a^m} = (\sqrt[n]{a})^m$$

For example,

$$2^4 = 16, \quad 2^{-4} = \frac{1}{2^4} = \frac{1}{16}, \quad 125^{2/3} = 5^2 = 25$$

In fact, exponents are extended to include all real numbers by defining, for any real number x ,

$$a^x = \lim_{r \rightarrow x} a^r \quad \text{where } r \text{ is a rational number}$$

Accordingly, the exponential function $f(x) = a^x$ is defined for all real numbers.

Logarithms are related to exponents as follows. Let b be a positive number. The logarithm of any positive number x to the base b , written

$$\log_b x$$

represents the exponent to which b must be raised to obtain x . That is,

$$y = \log_b x \quad \text{and} \quad b^y = x$$

are equivalent statements. Accordingly,

$$\begin{array}{llll} \log_2 8 = 3 & \text{since} & 2^3 = 8; & \log_{10} 100 = 2 \quad \text{since} \quad 10^2 = 100 \\ \log_2 64 = 6 & \text{since} & 2^6 = 64; & \log_{10} 0.001 = -3 \quad \text{since} \quad 10^{-3} = 0.001 \end{array}$$

Furthermore, for any base b ,

$$\begin{array}{ll} \log_b 1 = 0 & \text{since} \quad b^0 = 1 \\ \log_b b = 1 & \text{since} \quad b^1 = b \end{array}$$

The logarithm of a negative number and the logarithm of 0 are not defined.

One may also view the exponential and logarithmic functions

$$f(x) = b^x \quad \text{and} \quad g(x) = \log_b x$$

as inverse functions of each other. Accordingly, the graphs of these two functions are related. (See Prob. 2.5.)

Frequently, logarithms are expressed using approximate values. For example, using tables or calculators, one obtains

$$\log_{10} 300 = 2.4771 \quad \text{and} \quad \log_e 40 = 3.6889$$

as approximate answers. (Here $e = 2.718281 \cdots$)

Logarithms to the base 10 (called *common logarithms*), logarithms to the base e (called *natural logarithms*) and logarithms to the base 2 (called *binary logarithms*) are of special importance. Some texts write:

$$\begin{array}{lll} \ln x & \text{instead of} & \log_e x \\ \lg x \text{ or } \text{Log } x & \text{instead of} & \log_2 x \end{array}$$

This text on data structures is mainly concerned with binary logarithms. Accordingly,

The term $\log x$ shall mean $\log_2 x$ unless otherwise specified.

Frequently, we will require only the floor or the ceiling of a binary logarithm. This can be obtained by looking at the powers of 2. For example,

$$\begin{array}{llll} \lceil \log_2 100 \rceil = 6 & \text{since} & 2^6 = 64 & 2^7 = 128 \\ \lceil \log_2 1000 \rceil = 9 & \text{since} & 2^8 = 512 & \text{and } 2^9 = 1024 \end{array}$$

and so on.

2.3 ALGORITHMIC NOTATION

An algorithm, intuitively speaking, is a finite step-by-step list of well-defined instructions for solving a particular problem. The formal definition of an algorithm, which uses the notion of a Turing machine or its equivalent, is very sophisticated and lies beyond the scope of this text. This section describes the format that is used to present algorithms throughout the text. This algorithmic notation is best described by means of examples.

EXAMPLE 2.4

An array DATA of numerical values is in memory. We want to find the location LOC and the value MAX of the largest element of DATA. Given no other information about DATA, one way to solve the problem is as follows:

Initially begin with $\text{LOC} = 1$ and $\text{MAX} = \text{DATA}[1]$. Then compare MAX with each successive element $\text{DATA}[K]$ of DATA. If $\text{DATA}[K]$ exceeds MAX, then update LOC and MAX so that $\text{LOC} = K$ and $\text{MAX} = \text{DATA}[K]$. The final values appearing in LOC and MAX give the location and value of the largest element of DATA.

A formal presentation of this algorithm, whose flowchart appears in Fig. 2-2, follows.

Algorithm 2.1: (Largest Element in Array) A nonempty array DATA with N numerical values is given. This algorithm finds the location LOC and the value MAX of the largest element of DATA. The variable K is used as a counter.

- Step 1. [Initialize.] Set $K := 1$, $\text{LOC} := 1$ and $\text{MAX} := \text{DATA}[1]$.
- Step 2. [Increment counter.] Set $K := K + 1$.
- Step 3. [Test counter.] If $K > N$, then:
Write: LOC, MAX, and Exit.
- Step 4. [Compare and update.] If $\text{MAX} < \text{DATA}[K]$, then:
Set $\text{LOC} := K$ and $\text{MAX} := \text{DATA}[K]$.
- Step 5. [Repeat loop.] Go to Step 2.

The format for the formal presentation of an algorithm consists of two parts. The first part is a paragraph which tells the purpose of the algorithm, identifies the variables which occur in the algorithm and lists the input data. The second part of the algorithm consists of the list of steps that is to be executed.

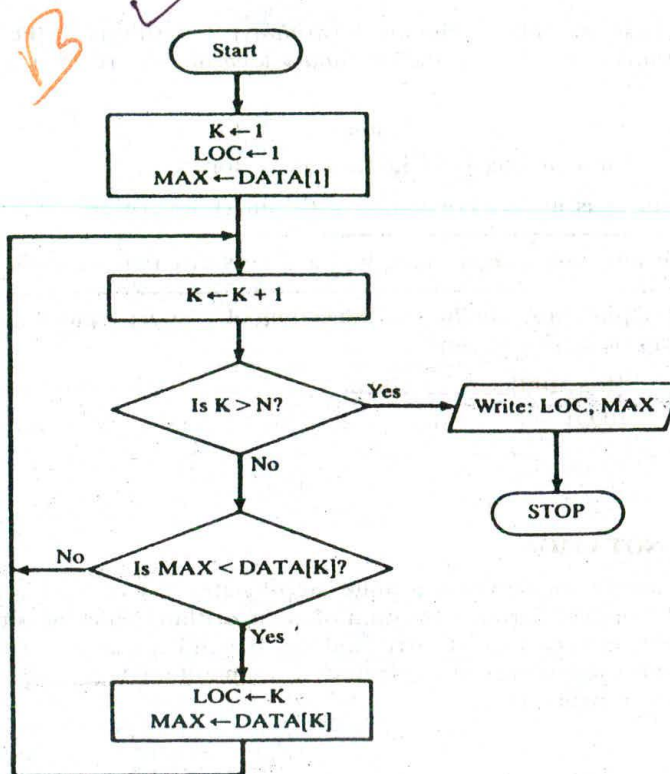


Fig. 2-2

The following summarizes certain conventions that we will use in presenting our algorithms. Some control structures will be covered in the next section.

Identifying Number

Each algorithm is assigned an identifying number as follows: Algorithm 4.3 refers to the third algorithm in Chap. 4; Algorithm P5.3 refers to the algorithm in Prob. 5.3 in Chap. 5. Note that the letter "P" indicates that the algorithm appears in a problem.

Steps, Control, Exit

The steps of the algorithm are executed one after the other, beginning with Step 1, unless indicated otherwise. Control may be transferred to Step n of the algorithm by the statement "Go to Step n ." For example, Step 5 transfers control back to Step 2 in Algorithm 2.1. Generally speaking, these Go to statements may be practically eliminated by using certain control structures discussed in the next section.

If several statements appear in the same step, e.g.,

Set $K := 1$, $LOC := 1$ and $MAX := DATA[1]$.

then they are executed from left to right.

The algorithm is completed when the statement

Exit.

is encountered. This statement is similar to the STOP statement used in FORTRAN and in flowcharts.

Comments

Each step may contain a comment in brackets which indicates the main purpose of the step. The comment will usually appear at the beginning or the end of the step.

Variable Names

Variable names will use capital letters, as in MAX and DATA. Single-letter names of variables used as counters or subscripts will also be capitalized in the algorithms (K and N, for example), even though lowercase may be used for these same variables (*k* and *n*) in the accompanying mathematical description and analysis. (Recall the discussion of italic and lowercase symbols in Sec. 1.3 of Chap. 1, under "Arrays.")

Assignment Statement

Our assignment statements will use the dots-equal notation := that is used in Pascal. For example,

$$\text{Max} := \text{DATA}[1]$$

assigns the value in DATA[1] to MAX. Some texts use the backward arrow ← or the equal sign = for this operation.

Input and Output

Data may be input and assigned to variables by means of a Read statement with the following form:

Read: Variables names.

Similarly, messages, placed in quotation marks, and data in variables may be output by means of a Write or Print statement with the following form:

Write: Messages and/or variable names.

Procedures

The term "procedure" will be used for an independent algorithmic module which solves a particular problem. The use of the word "procedure" or "module" rather than "algorithm" for a given problem is simply a matter of taste. Generally speaking, the word "algorithm" will be reserved for the solution of general problems. The term "procedure" will also be used to describe a certain type of subalgorithm which is discussed in Sec. 2.6.

2.4 CONTROL STRUCTURES

Algorithms and their equivalent computer programs are more easily understood if they mainly use self-contained modules and three types of logic, or flow of control, called

- (1) Sequence logic, or sequential flow
- (2) Selection logic, or conditional flow
- (3) Iteration logic, or repetitive flow

These three types of logic are discussed below, and in each case we show the equivalent flowchart.

Sequence Logic (Sequential Flow)

Sequence logic has already been discussed. Unless instructions are given to the contrary, the modules are executed in the obvious sequence. The sequence may be presented explicitly, by means of numbered steps, or implicitly, by the order in which the modules are written. (See Fig. 2-3.) Most processing, even of complex problems, will generally follow this elementary flow pattern.

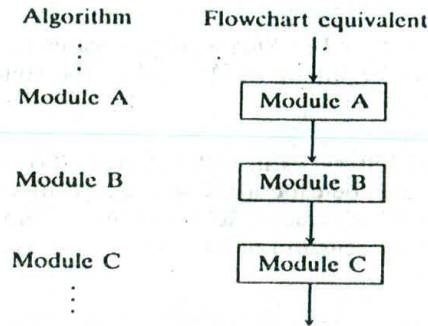


Fig. 2-3 Sequence logic.

Selection Logic (Conditional Flow)

Selection logic employs a number of conditions which lead to a selection of one out of several alternative modules. The structures which implement this logic are called conditional structures or If structures. For clarity, we will frequently indicate the end of such a structure by the statement

[End of If structure.]

or some equivalent.

These conditional structures fall into three types, which are discussed separately.

(1) *Single alternative.* This structure has the form

If condition, then;
[Module A]
[End of If structure.]

The logic of this structure is pictured in Fig. 2-4(a). If the condition holds, then Module A, which may consist of one or more statements, is executed; otherwise Module A is skipped and control transfers to the next step of the algorithm.

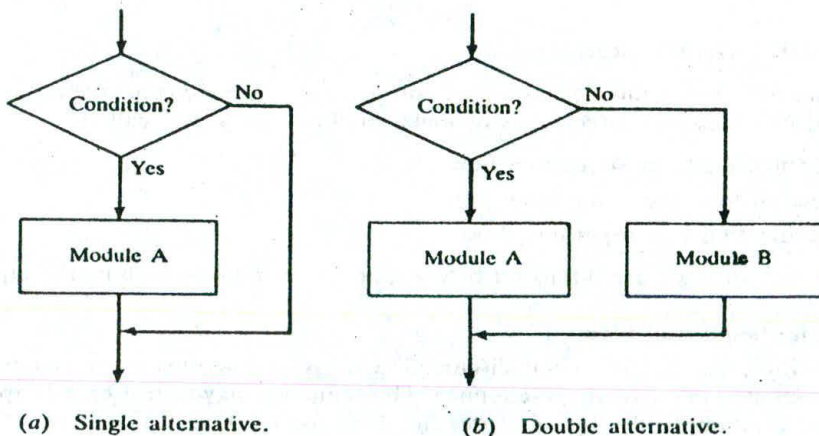


Fig. 2-4

(2) *Double alternative.* This structure has the form

If condition, then:

[Module A]

Else:

[Module B]

[End of If structure.]

The logic of this structure is pictured in Fig. 2-4(b). As indicated by the flowchart, if the condition holds, then Module A is executed; otherwise Module B is executed.

(3) *Multiple alternatives.* This structure has the form:

If condition(1), then:

[Module A₁]

Else if condition(2), then:

[Module A₂]

⋮

Else if condition(M), then:

[Module A_M]

Else:

[Module B]

[End of If structure.]

The logic of this structure allows only one of the modules to be executed. Specifically, either the module which follows the first condition which holds is executed, or the module which follows the final Else statement is executed. In practice, there will rarely be more than three alternatives.

EXAMPLE 2.5

The solutions of the quadratic equation

$$ax^2 + bx + c = 0.$$

where $a \neq 0$, are given by the quadratic formula

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

The quantity $D = b^2 - 4ac$ is called the *discriminant* of the equation. If D is negative, then there are no real solutions. If $D = 0$, then there is only one (double) real solution, $x = -b/2a$. If D is positive, the formula gives the two distinct real solutions. The following algorithm finds the solutions of a quadratic equation.

Algorithm 2.2: (Quadratic Equation) This algorithm inputs the coefficients A, B, C of a quadratic equation and outputs the real solutions, if any.

Step 1. Read: A, B, C.

Step 2. Set $D := B^2 - 4AC$.

Step 3. If $D > 0$, then:

(a) Set $X1 := (-B + \sqrt{D})/2A$ and $X2 := (-B - \sqrt{D})/2A$.

(b) Write: X1, X2.

Else if $D = 0$, then:

(a) Set $X := -B/2A$.

(b) Write: 'UNIQUE SOLUTION', X.

Else:

Write: 'NO REAL SOLUTIONS'.

[End of If structure.]

Step 4. Exit.

Remark: Observe that there are three mutually exclusive conditions in Step 3 of Algorithm 2.2 that depend on whether D is positive, zero or negative. In such a situation, we may alternatively list the different cases as follows:

- Step 3. (1) If $D > 0$, then:

 (2) If $D = 0$, then:

 (3) If $D < 0$, then:

This is similar to the use of the CASE statement in Pascal.

Iteration Logic (Repetitive Flow)

The third kind of logic refers to either of two types of structures involving loops. Each type begins with a Repeat statement and is followed by a module, called the *body of the loop*. For clarity, we will indicate the end of the structure by the statement

[End of loop.]

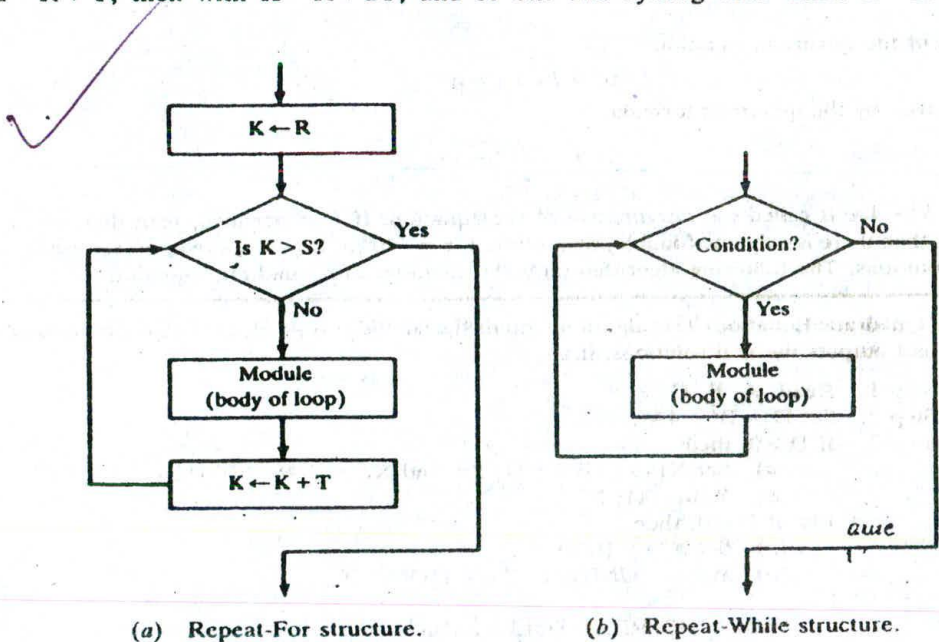
or some equivalent.

Each type of loop structure is discussed separately.

The *repeat-for loop* uses an index variable, such as K , to control the loop. The loop will usually have the form:

Repeat for $K = R$ to S by T :
 [Module]
 [End of loop.]

The logic of this structure is pictured in Fig. 2-5(a). Here R is called the *initial value*, S the *end value* or *test value*, and T the *increment*. Observe that the body of the loop is executed first with $K = R$, then with $K = R + T$, then with $K = R + 2T$, and so on. The cycling ends when $K > S$. The flowchart



(a) Repeat-For structure.

(b) Repeat-While structure.

Fig. 2-5

assumes that the increment T is positive; if T is negative, so that K decreases in value, then the cycling ends when $K < S$.

The *repeat-while loop* uses a condition to control the loop. The loop will usually have the form

```
Repeat while condition:
    [Module]
[End of loop.]
```

The logic of this structure is pictured in Fig. 2-5(b). Observe that the cycling continues until the condition is false. We emphasize that there must be a statement before the structure that initializes the condition controlling the loop, and in order that the looping may eventually cease, there must be a statement in the body of the loop that changes the condition.

EXAMPLE 2.6

Algorithm 2.1 is rewritten using a repeat-while loop rather than a Go to statement:

Algorithm 2.3: (Largest Element in Array) Given a nonempty array $DATA$ with N numerical values, this algorithm finds the location LOC and the value MAX of the largest element of $DATA$.

1. [Initialize.] Set $K := 1$, $LOC := 1$ and $MAX := DATA[1]$.
2. Repeat Steps 3 and 4 while $K \leq N$:
3. If $MAX < DATA[K]$, then:
 - Set $LOC := K$ and $MAX := DATA[K]$.
 - [End of If structure.]
4. Set $K := K + 1$.
 - [End of Step 2 loop.]
5. Write: LOC , MAX .
6. Exit.

Algorithm 2.3 indicates some other properties of our algorithms. Usually we will omit the word "Step." We will try to use repeat structures instead of Go to statements. The repeat statement may explicitly indicate the steps that form the body of the loop. The "End of loop" statement may explicitly indicate the step where the loop begins. The modules contained in our logic structures will normally be indented for easier reading. This conforms to the usual format in structured programming.

Any other new notation or convention either will be self-explanatory or will be explained when it occurs.

2.5 COMPLEXITY OF ALGORITHMS

The analysis of algorithms is a major task in computer science. In order to compare algorithms, we must have some criteria to measure the efficiency of our algorithms. This section discusses this important topic.

Suppose M is an algorithm, and suppose n is the size of the input data. The time and space used by the algorithm M are the two main measures for the efficiency of M . The time is measured by counting the number of key operations—in sorting and searching algorithms, for example, the number of comparisons. That is because key operations are so defined that the time for the other operations is much less than or at most proportional to the time for the key operations. The space is measured by counting the maximum of memory needed by the algorithm.

The *complexity* of an algorithm M is the function $f(n)$ which gives the running time and/or storage space requirement of the algorithm in terms of the size n of the input data. Frequently, the storage

space required by an algorithm is simply a multiple of the data size n . Accordingly, unless otherwise stated or implied, the term "complexity" shall refer to the running time of the algorithm.

The following example illustrates that the function $f(n)$, which gives the running time of an algorithm, depends not only on the size n of the input data but also on the particular data.

EXAMPLE 2.7

Suppose we are given an English short story TEXT, and suppose we want to search through TEXT for the first occurrence of a given 3-letter word W. If W is the 3-letter word "the," then it is likely that W occurs near the beginning of TEXT, so $f(n)$ will be small. On the other hand, if W is the 3-letter word "zoo," then W may not appear in TEXT at all, so $f(n)$ will be large.

The above discussion leads us to the question of finding the complexity function $f(n)$ for certain cases. The two cases one usually investigates in complexity theory are as follows:

- (1) *Worst case*: the maximum value of $f(n)$ for any possible input
- (2) *Average case*: the expected value of $f(n)$

Sometimes we also consider the minimum possible value of $f(n)$, called the *best case*.

The analysis of the average case assumes a certain probabilistic distribution for the input data; one such assumption might be that all possible permutations of an input data set are equally likely. The average case also uses the following concept in probability theory. Suppose the numbers n_1, n_2, \dots, n_k occur with respective probabilities p_1, p_2, \dots, p_k . Then the *expectation* or *average value* E is given by

$$E = n_1p_1 + n_2p_2 + \dots + n_kp_k$$

These ideas are illustrated in the following example.

EXAMPLE 2.8 Linear Search

Suppose a linear array DATA contains n elements, and suppose a specific ITEM of information is given. We want either to find the location LOC of ITEM in the array DATA, or to send some message, such as LOC = 0, to indicate that ITEM does not appear in DATA. The linear search algorithm solves this problem by comparing ITEM, one by one, with each element in DATA. That is, we compare ITEM with DATA[1], then DATA[2], and so on, until we find LOC such that ITEM = DATA[LOC]. A formal presentation of this algorithm follows.

Algorithm 2.4: (Linear Search) A linear array DATA with N elements and a specific ITEM of information are given. This algorithm finds the location LOC of ITEM in the array DATA or sets LOC = 0.

1. [Initialize] Set $K := 1$ and $LOC := 0$.
2. Repeat Steps 3 and 4 while $LOC = 0$ and $K \leq N$.
3. If $ITEM = DATA[K]$, then: Set $LOC := K$.
4. Set $K := K + 1$. [Increments counter.]
[End of Step 2 loop.]
5. [Successful?]
If $LOC = 0$, then:
Write: ITEM is not in the array DATA.
Else:
Write: LOC is the location of ITEM.
[End of If structure.]
6. Exit.

The complexity of the search algorithm is given by the number C of comparisons between ITEM and DATA[K]. We seek $C(n)$ for the worst case and the average case.

Worst Case

Clearly the worst case occurs when ITEM is the last element in the array DATA or is not there at all. In either situation, we have

$$C(n) = n$$

Accordingly, $C(n) = n$ is the worst-case complexity of the linear search algorithm.

Average Case

Here we assume that ITEM does appear in DATA, and that it is equally likely to occur at any position in the array. Accordingly, the number of comparisons can be any of the numbers 1, 2, 3, ..., n , and each number occurs with probability $p = 1/n$. Then

$$\begin{aligned} C(n) &= 1 \cdot \frac{1}{n} + 2 \cdot \frac{1}{n} + \cdots + n \cdot \frac{1}{n} \\ &= (1 + 2 + \cdots + n) \cdot \frac{1}{n} \\ &= \frac{n(n+1)}{2} \cdot \frac{1}{n} = \frac{n+1}{2} \end{aligned}$$

This agrees with our intuitive feeling that the average number of comparisons needed to find the location of ITEM is approximately equal to half the number of elements in the DATA list.

Remark: The complexity of the average case of an algorithm is usually much more complicated to analyze than that of the worst case. Moreover, the probabilistic distribution that one assumes for the average case may not actually apply to real situations. Accordingly, unless otherwise stated or implied, the complexity of an algorithm shall mean the function which gives the running time of the worst case in terms of the input size. This is not too strong an assumption, since the complexity of the average case for many algorithms is proportional to the worst case.

Rate of Growth; Big O Notation

Suppose M is an algorithm, and suppose n is the size of the input data. Clearly the complexity $f(n)$ of M increases as n increases. It is usually the rate of increase of $f(n)$ that we want to examine. This is usually done by comparing $f(n)$ with some standard function, such as

$$\log_2 n, \quad n, \quad n \log_2 n, \quad n^2, \quad n^3, \quad 2^n$$

The rates of growth for these standard functions are indicated in Fig. 2-6, which gives their approximate values for certain values of n . Observe that the functions are listed in the order of their rates of growth: the logarithmic function $\log_2 n$ grows most slowly, the exponential function 2^n grows most rapidly, and the polynomial functions n^c grow according to the exponent c . One way to compare the function $f(n)$ with these standard functions is to use the functional O notation defined as follows:

$n \backslash g(n)$	$\log n$	n	$n \log n$	n^2	n^3	2^n
5	3	5	15	25	125	32
10	4	10	40	100	10^3	10^3
100	7	100	700	10^4	10^6	10^{30}
1000	10	10^3	10^4	10^6	10^9	10^{300}

Fig. 2-6 Rate of growth of standard functions.

Suppose $f(n)$ and $g(n)$ are functions defined on the positive integers with the property that $f(n)$ is bounded by some multiple of $g(n)$ for almost all n . That is, suppose there exist a positive integer n_0 and a positive number M such that, for all $n > n_0$, we have

$$|f(n)| \leq M|g(n)|$$

Then we may write

$$f(n) = O(g(n))$$

which is read " $f(n)$ is of order $g(n)$." For any polynomial $P(n)$ of degree m , we show in Prob. 2.10 that $P(n) = O(n^m)$; e.g.,

$$8n^3 - 576n^2 + 832n - 248 = O(n^3)$$

We can also write

$$f(n) = h(n) + O(g(n)) \quad \text{when} \quad f(n) - h(n) = O(g(n))$$

(This is called the "big O " notation since $f(n) = o(g(n))$ has an entirely different meaning.)

To indicate the convenience of this notation, we give the complexity of certain well-known searching and sorting algorithms:

- (a) Linear search: $O(n)$
- (b) Binary search: $O(\log n)$
- (c) Bubble sort: $O(n^2)$
- (d) Merge-sort: $O(n \log n)$

These results are discussed in Chap. 9, on sorting and searching.

2.6 SUBALGORITHMS

A *subalgorithm* is a complete and independently defined algorithmic module which is used (or *invoked* or *called*) by some main algorithm or by some other subalgorithm. A subalgorithm receives values, called *arguments*, from an originating (calling) algorithm; performs computations; and then sends back the result to the calling algorithm. The subalgorithm is defined independently so that it may be called by many different algorithms or called at different times in the same algorithm. The relationship between an algorithm and a subalgorithm is similar to the relationship between a main program and a subprogram in a programming language.

The main difference between the format of a subalgorithm and that of an algorithm is that the subalgorithm will usually have a heading of the form

$$\text{NAME}(\text{PAR}_1, \text{PAR}_2, \dots, \text{PAR}_K)$$

Here NAME refers to the name of the subalgorithm which is used when the subalgorithm is called, and $\text{PAR}_1, \text{PAR}_2, \dots, \text{PAR}_K$ refer to parameters which are used to transmit data between the subalgorithm and the calling algorithm:

Another difference is that the subalgorithm will have a Return statement rather than an Exit statement; this emphasizes that control is transferred back to the calling program when the execution of the subalgorithm is completed.

Subalgorithms fall into two basic categories: *function* subalgorithms and *procedure* subalgorithms. The similarities and differences between these two types of subalgorithms will be examined below by means of examples. One major difference between the subalgorithms is that the function subalgorithm returns only a single value to the calling algorithm, whereas the procedure subalgorithm may send back more than one value.

EXAMPLE 2.9

The following function subalgorithm MEAN finds the average AVE of three numbers A, B and C.

Function 2.5: MEAN(A, B, C)

1. Set AVE := (A + B + C)/3.
2. Return(AVE).

Note that MEAN is the name of the subalgorithm and A, B and C are the parameters. The Return statement includes, in parentheses, the variable AVE, whose value is returned to the calling program.

The subalgorithm MEAN is invoked by an algorithm in the same way as a function subprogram is invoked by a calling program. For example, suppose an algorithm contains the statement

$$\text{Set TEST} := \text{MEAN}(T_1, T_2, T_3)$$

where T_1 , T_2 and T_3 are test scores. The argument values T_1 , T_2 and T_3 are transferred to the parameters A, B, C in the subalgorithm, the subalgorithm MEAN is executed, and then the value of AVE is returned to the program and replaces MEAN(T_1 , T_2 , T_3) in the statement. Hence the average of T_1 , T_2 and T_3 is assigned to TEST.

EXAMPLE 2.10

The following procedure SWITCH interchanges the values of AAA and BBB.

Procedure 2.6: SWITCH(AAA, BBB)

1. Set TEMP := AAA, AAA := BBB and BBB := TEMP.
2. Return.

The procedure is invoked by means of a Call statement. For example the Call statement

$$\text{Call SWITCH}(BEG, AUX)$$

has the net effect of interchanging the values of BEG and AUX. Specifically, when the procedure SWITCH is invoked, the argument of BEG and AUX are transferred to the parameters AAA and BBB, respectively; the procedure is executed, which interchanges the values of AAA and BBB; and then the new values of AAA and BBB are transferred back to BEG and AUX, respectively.

Remark: Any function subalgorithm can be easily translated into an equivalent procedure by simply adjoining an extra parameter which is used to return the computed value to the calling algorithm. For example, Function 2.5 may be translated into a procedure

$$\text{MEAN}(A, B, C, \text{AVE})$$

where the parameter AVE is assigned the average of A, B, C. Then the statement

$$\text{Call MEAN}(T_1, T_2, T_3, \text{TEST})$$

also has the effect of assigning the average of T_1 , T_2 and T_3 to TEST. Generally speaking, we will use procedures rather than function subalgorithms.

2.7 VARIABLES, DATA TYPES

Each variable in any of our algorithms or programs has a data type which determines the code that is used for storing its value. Four such data types follow:

- (1) *Character.* Here data are coded using some character code such as EBCDIC or ASCII. The 8-bit EBCDIC code of some characters appears in Fig. 2-7. A single character is normally stored in a byte.

Char.	Zone	Numeric	Hex	Char.	Zone	Numeric	Hex	Char.	Zone	Numeric	Hex
A	1100	0001	C1	S	1110	0010	E2	blank	0100	0000	40
B	↓	0010	C2	T	↓	0011	E3	<	↓	1011	4B
C	↓	0011	C3	U	↓	0100	E4	(↓	1100	4C
D	↓	0100	C4	V	↓	0101	E5	+	↓	1101	4D
E	↓	0101	C5	W	↓	0110	E6	&	0100	1110	4E
F	↓	0110	C6	X	↓	0111	E7	\$	0101	0000	50
G	↓	0111	C7	Y	↓	1000	E8	*	↓	1011	5B
H	↓	1000	C8	Z	1110	1001	E9)	↓	1100	5C
I	1100	1001	C9	0	1111	0000	F0	:	0101	1110	5E
J	1101	0001	D1	1	↓	0001	F1	-	0110	0000	60
K	↓	0010	D2	2	↓	0010	F2	/	↓	0001	61
L	↓	0011	D3	3	↓	0011	F3	%	↓	1011	6B
M	↓	0100	D4	4	↓	0100	F4	>	↓	1100	6C
N	↓	0101	D5	5	↓	0101	F5	?	0110	1110	6E
O	↓	0110	D6	6	↓	0110	F6	#	0111	1111	6F
P	↓	0111	D7	7	↓	0111	F7	@	↓	1010	7A
Q	↓	1000	D8	8	↓	1000	F8	=	↓	1011	7B
R	1101	1001	D9	9	1111	1001	F9		0111	1100	7C
										1110	7E

Fig. 2-7 Part of the EBCDIC code.

- (2) *Real (or floating point)*. Here numerical data are coded using the exponential form of the data.
- (3) *Integer (or fixed point)*. Here positive integers are coded using binary representation, and negative integers by some binary variation such as 2's complement.
- (4) *Logical*. Here the variable can have only the value true or false; hence it may be coded using only one bit, 1 for true and 0 for false. (Sometimes the bytes 1111 1111 and 0000 0000 may be used for true and false, respectively.)

The data types of variables in our algorithms will not be explicitly stated as with computer programs but will usually be implied by the context.

EXAMPLE 2.11

Suppose a 32-bit memory location X contains the following sequence of bits:

0110 1100 1100 0111 1101 0110 0110 1100

There is no way to know the content of the cell unless the data type of X is known.

- (a) Suppose X is declared to be of character type and EBCDIC is used. Then the four characters %GO% are stored in X .
- (b) Suppose X is declared to be of some other type, such as integer or real. Then an integer or real number is stored in X .

Local and Global Variables

The organization of a computer program into a main program and various subprograms has led to the notion of local and global variables. Normally, each program module contains its own list of variables, called *local variables*, which can be accessed only by the given program module. Also,

subprogram modules may contain parameters, variables which transfer data between a subprogram and its calling program.

EXAMPLE 2.12

Consider the procedure SWITCH(AAA, BBB) in Example 2.10. The variables AAA and BBB are parameters; they are used to transfer data between the procedure and a calling algorithm. On the other hand, the variable TEMP in the procedure is a local variable. It "lives" only in the procedure; i.e., its value can be accessed and changed only during the execution of the procedure. In fact, the name TEMP may be used for a variable in any other module and the use of the name will not interfere with the execution of the procedure SWITCH.

Language designers realized that it would be convenient to have certain variables which can be accessed by some or even all the program modules in a computer program. Variables that can be accessed by all program modules are called *global* variables, and variables that can be accessed by some program modules are called *nonlocal* variables. Each programming language has its own syntax for declaring such variables. For example, FORTRAN uses a COMMON statement to declare global variables, and Pascal uses scope rules to declare global and nonlocal variables.

Accordingly, there are two basic ways for modules to communicate with each other:

- (1) *Directly*, by means of well-defined parameters
- (2) *Indirectly*, by means of nonlocal and global variables

The indirect change of the value of a variable in one module by another module is called a *side effect*. Readers should be very careful when using nonlocal and global variables, since errors caused by side effects may be difficult to detect.

Solved Problems

MATHEMATICAL NOTATION AND FUNCTIONS

2.1 Find (a) $[7.5]$, $[-7.5]$, $[-18]$, $[\sqrt{30}]$, $[\sqrt[3]{30}]$, $[\pi]$; and (b) $\lceil 7.5 \rceil$, $\lceil -7.5 \rceil$, $\lceil -18 \rceil$, $\lceil \sqrt{30} \rceil$, $\lceil \sqrt[3]{30} \rceil$, $\lceil \pi \rceil$.

(a) By definition, $[x]$ denotes the greatest integer that does not exceed x , called the floor of x . Hence,

$$\begin{array}{lll} [7.5] = 7 & [-7.5] = -8 & [-18] = -18 \\ [\sqrt{30}] = 5 & [\sqrt[3]{30}] = 3 & [\pi] = 3 \end{array}$$

(b) By definition, $\lceil x \rceil$ denotes the least integer that is not less than x , called the ceiling of x . Hence,

$$\begin{array}{lll} \lceil 7.5 \rceil = 8 & \lceil -7.5 \rceil = -7 & \lceil -18 \rceil = -18 \\ \lceil \sqrt{30} \rceil = 6 & \lceil \sqrt[3]{30} \rceil = 4 & \lceil \pi \rceil = 4 \end{array}$$

2.2 (a) Find $26 \pmod{7}$, $34 \pmod{8}$, $2345 \pmod{6}$, $495 \pmod{11}$.

(b) Find $-26 \pmod{7}$, $-2345 \pmod{6}$, $-371 \pmod{8}$, $-39 \pmod{3}$.

(c) Using arithmetic modulo 15, evaluate $9 + 13$, $7 + 11$, $4 - 9$, $2 - 10$.

(a) Since k is positive, simply divide k by the modulus M to obtain the remainder r . Then $r = k \pmod{M}$. Thus

$$5 = 26 \pmod{7} \quad 2 = 34 \pmod{8} \quad 5 = 2345 \pmod{6} \quad 0 = 495 \pmod{11}$$

- (b) When k is negative, divide $|k|$ by the modulus to obtain the remainder r' . Then $k \equiv -r' \pmod{M}$. Hence $k \pmod{M} = M - r'$ when $r' \neq 0$. Thus

$$\begin{aligned} -26 \pmod{7} &= 7 - 5 = 2 & -371 \pmod{8} &= 8 - 3 = 5 \\ -2345 \pmod{6} &= 6 - 5 = 1 & -39 \pmod{3} &= 0 \end{aligned}$$

- (c) Use $a \pm M \equiv a \pmod{M}$:

$$\begin{aligned} 9 + 13 &= 22 \equiv 22 - 15 = 7 & 7 + 11 &= 18 \equiv 18 - 15 = 3 \\ 4 - 9 &= -5 \equiv -5 + 15 = 10 & 2 - 10 &= -8 \equiv -8 + 15 = 7 \end{aligned}$$

2.3 List all the permutations of the numbers 1, 2, 3, 4.

Note first that there are $4! = 24$ such permutations:

1234	1243	1324	1342	1423	1432
2134	2143	2314	2341	2413	2431
3124	3142	3214	3241	3412	3421
4123	4132	4213	4231	4312	4321

Observe that the first row contains the six permutations beginning with 1, the second row those beginning with 2, and so on.

2.4 Find: (a) 2^{-5} , $8^{2/3}$, $25^{-3/2}$; (b) $\log_2 32$, $\log_{10} 1000$, $\log_2 (1/16)$; (c) $\lceil \log_2 1000 \rceil$, $\lfloor \log_2 0.01 \rfloor$.

(a) $2^{-5} = 1/2^5 = 1/32$; $8^{2/3} = (\sqrt[3]{8})^2 = 2^2 = 4$; $25^{-3/2} = 1/25^{3/2} = 1/5^3 = 1/125$.

(b) $\log_2 32 = 5$ since $2^5 = 32$; $\log_{10} 1000 = 3$ since $10^3 = 1000$; $\log_2 (1/16) = -4$ since $2^{-4} = 1/2^4 = 1/16$.

(c) $\lceil \log_2 1000 \rceil = 9$ since $2^9 = 512$ but $2^{10} = 1024$;
 $\lfloor \log_2 0.01 \rfloor = -7$ since $2^{-7} = 1/128 < 0.01 < 2^{-6} = 1/64$.

2.5 Plot the graphs of the exponential function $f(x) = 2^x$, the logarithmic function $g(x) = \log_2 x$ and the linear function $h(x) = x$ on the same coordinate axis. (a) Describe a geometric property of the graphs $f(x)$ and $g(x)$. (b) For any positive number c , how are $f(c)$, $g(c)$ and $h(c)$ related?

Figure 2-8 pictures the three functions.

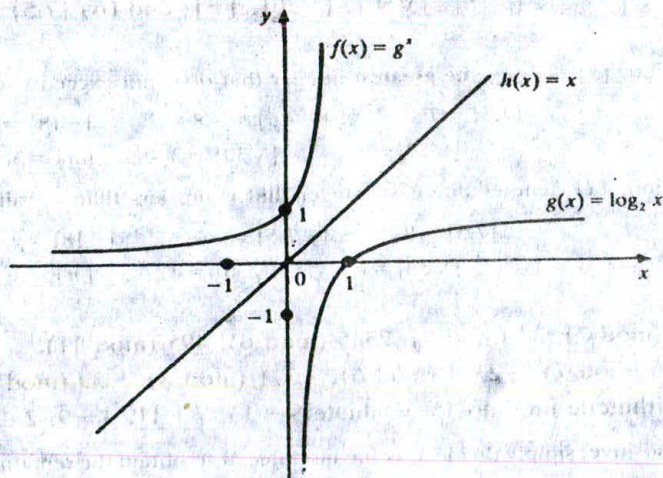


Fig. 2-8

- (a) Since $f(x) = 2^x$ and $g(x) = \log_2 x$ are inverse functions, they are symmetric with respect to the line $y = x$.
- (b) For any positive number c , we have

$$g(c) < h(c) < f(c)$$

In fact, as c increases in value, the vertical distances between the functions,

$$h(c) - g(c) \quad \text{and} \quad f(c) - h(c),$$

increase in value. Moreover, the logarithmic function $g(x)$ grows very slowly compared with the linear function $h(x)$, and the exponential function $f(x)$ grows very quickly compared with $h(x)$.

ALGORITHMS, COMPLEXITY

2.6 Consider Algorithm 2.3, which finds the location LOC and the value MAX of the largest element in an array DATA with n elements. Consider the complexity function $C(n)$, which measures the number of times LOC and MAX are updated in Step 3. (The number of comparisons is independent of the order of the elements in DATA.)

- (a) Describe and find $C(n)$ for the worst case.
- (b) Describe and find $C(n)$ for the best case.
- (c) Find $C(n)$ for the average case when $n = 3$, assuming all arrangements of the elements in DATA are equally likely.
- (a) The worst case occurs when the elements of DATA are in increasing order, where each comparison of MAX with DATA[K] forces LOC and MAX to be updated. In this case, $C(n) = n - 1$.
- (b) The best case occurs when the largest element appears first and so when the comparison of MAX with DATA[K] never forces LOC and MAX to be updated. Accordingly, in this case, $C(n) = 0$.
- (c) Let 1, 2 and 3 denote, respectively, the largest, second largest and smallest elements of DATA. There are six possible ways the elements can appear in DATA, which correspond to the $3! = 6$ permutations of 1, 2, 3. For each permutation p , let n_p denote the number of times LOC and MAX are updated when the algorithm is executed with input p . The six permutations p and the corresponding values n_p follow:

Permutation p :	123	132	213	231	312	321
Value of n_p :	0	0	1	1	1	2

Assuming all permutations p are equally likely,

$$C(3) = \frac{0+0+1+1+1+2}{6} = \frac{5}{6}$$

(The evaluation of the average value of $C(n)$ for arbitrary n lies beyond the scope of this text. One purpose of this problem is to illustrate the difficulty that may occur in finding the complexity of the average case of an algorithm.)

2.7 Suppose Module A requires M units of time to be executed, where M is a constant. Find the complexity $C(n)$ of each algorithm, where n is the size of the input data and b is a positive integer greater than 1.

- (a) **Algorithm P2.7A:**
1. Repeat for $I = 1$ to N :
 2. Repeat for $J = 1$ to N :
 3. Repeat for $K = 1$ to N :
 4. Module A.
 - [End of Step 3 loop.]
 - [End of Step 2 loop.]
 - [End of Step 1 loop.]
 5. Exit.

- (b) **Algorithm P2.7B:**
1. Set $J := 1$.
 2. Repeat Steps 3 and 4 while $J \leq N$:
 3. Module A.
 4. Set $J := B \times J$.
 - [End of Step 2 loop.]
 5. Exit.

Observe that the algorithms use N for n and B for b .)

(a) Here
$$C(n) = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n M$$

The number of times M occurs in the sum is equal to the number of triplets (i, j, k) , where i, j, k are integers from 1 to n inclusive. There are n^3 such triplets. Hence

$$C(n) = Mn^3 = O(n^3)$$

- (b) Observe that the values of the loop index J are the powers of b :

$$1, b, b^2, b^3, b^4, \dots$$

Therefore, Module A will be repeated exactly T times, where T is the first exponent such that

$$b^r > n$$

Hence,

$$T = \lfloor \log_b n \rfloor + 1$$

Accordingly,

$$C(n) = MT = O(\log_b n)$$

- 2.8 (a) Write a procedure FIND(DATA, N, LOC1, LOC2) which finds the location LOC1 of the largest element and the location LOC2 of the second largest element in an array DATA with $n > 1$ elements.

- (b) Why not let FIND also find the values of the largest and second largest elements?

- (a) The elements of DATA are examined one by one. During the execution of the procedure, FIRST and SECOND will denote, respectively, the values of the largest and second largest elements that have already been examined. Each new element DATA[K] is tested as follows. If

$$\text{SECOND} \leq \text{FIRST} < \text{DATA}[K]$$

then FIRST becomes the new SECOND element and DATA[K] becomes the new FIRST element. On the other hand, if

$$\text{SECOND} < \text{DATA}[K] \leq \text{FIRST}$$

then DATA[K] becomes the new SECOND element. Initially, set $\text{FIRST} := \text{DATA}[1]$ and $\text{SECOND} := \text{DATA}[2]$, and check whether or not they are in the right order. A formal presentation of the procedure follows:

Procedure P2.8: FIND(DATA, N, LOC1, LOC2)

1. Set $\text{FIRST} := \text{DATA}[1]$, $\text{SECOND} := \text{DATA}[2]$, $\text{LOC1} := 1$, $\text{LOC2} := 2$.
2. [Are FIRST and SECOND initially correct?]

If $\text{FIRST} < \text{SECOND}$, then:

 - (a) Interchange FIRST and SECOND,
 - (b) Set $\text{LOC1} := 2$ and $\text{LOC2} := 1$.

[End of If structure.]
3. Repeat for $K = 3$ to N :

If $\text{FIRST} < \text{DATA}[K]$, then:

 - (a) Set $\text{SECOND} := \text{FIRST}$ and $\text{FIRST} := \text{DATA}[K]$.
 - (b) Set $\text{LOC2} := \text{LOC1}$ and $\text{LOC1} := K$.

Else if $\text{SECOND} < \text{DATA}[K]$, then:

Set $\text{SECOND} := \text{DATA}[K]$ and $\text{LOC2} := K$.

[End of If structure.]

[End of loop.]
4. Return.

(b) Using additional parameters FIRST and SECOND would be redundant, since LOC1 and LOC2 automatically tell the calling program that DATA[LOC1] and DATA[LOC2] are, respectively, the values of the largest and second largest elements of DATA.

2.9 An integer $n > 1$ is called a *prime* number if its only positive divisors are 1 and n ; otherwise, n is called a *composite* number. For example, the following are the prime numbers less than 20:

2, 3, 5, 7, 11, 13, 17, 19

If $n > 1$ is not prime, i.e., if n is composite, then n must have a divisor $k \neq 1$ such that $k \leq \sqrt{n}$ or, in other words, $k^2 \leq n$.

Suppose we want to find all the prime numbers less than a given number m , such as 30. This can be done by the "sieve method," which consists of the following steps. First list the 30 numbers:

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15
16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30

Cross out 1 and the multiples of 2 from the list as follows:

~~1~~, 2, 3, ~~4~~, 5, ~~6~~, 7, ~~8~~, 9, ~~10~~, 11, ~~12~~, 13, ~~14~~, 15
~~16~~, 17, ~~18~~, 19, ~~20~~, 21, ~~22~~, 23, ~~24~~, 25, ~~26~~, 27, ~~28~~, 29, ~~30~~

Since 3 is the first number following 2 that has not been eliminated, cross out the multiples of 3 from the list as follows:

~~1~~, 2, 3, ~~4~~, 5, ~~6~~, 7, ~~8~~, ~~9~~, ~~10~~, 11, ~~12~~, 13, ~~14~~, ~~15~~
~~16~~, 17, ~~18~~, 19, ~~20~~, ~~21~~, ~~22~~, 23, ~~24~~, 25, ~~26~~, ~~27~~, ~~28~~, 29, ~~30~~

Since 5 is the first number following 3 that has not been eliminated, cross out the multiples of 5 from the list as follows:

~~1~~, 2, 3, ~~4~~, 5, ~~6~~, 7, ~~8~~, ~~9~~, ~~10~~, 11, ~~12~~, 13, ~~14~~, ~~15~~
~~16~~, 17, ~~18~~, 19, ~~20~~, ~~21~~, ~~22~~, 23, ~~24~~, ~~25~~, ~~26~~, ~~27~~, ~~28~~, 29, ~~30~~

Now 7 is the first number following 5 that has not been eliminated, but $7^2 > 30$. This means the algorithm is finished and the numbers left in the list are the primes less than 30:

2, 3, 5, 7, 11, 13, 17, 19, 23, 29

Translate the sieve method into an algorithm to find all prime numbers less than a given number n .

First define an array A such that

$A[1] = 1$, $A[2] = 2$, $A[3] = 3$, $A[4] = 4$, ..., $A[N-1] = N-1$, $A[N] =$

We cross out an integer L from the list by assigning $A[L] = 1$. The following procedure CROSSOUT tests whether $A[K] = 1$, and if not, it sets

$A[2K] = 1$, $A[3K] = 1$, $A[4K] = 1$, ...

That is, it eliminates the multiples of K from the list

Procedure P2.9A: CROSSOUT(A, N, K)

1. If $A[K] = 1$, then: Return.
2. Repeat for $L = 2K$ to N by K :
 Set $A[L] := 1$.
 [End of loop.]
3. Return.

The sieve method can now be simply written:

Algorithm P2.9B: This algorithm prints the prime numbers less than N .

1. [Initialize array A.] Repeat for $K = 1$ to N :
Set $A[K] := K$.
2. [Eliminate multiples of K .] Repeat for $K = 2$ to \sqrt{N} :
Call CROSSOUT(A, N, K).
3. [Print the primes.] Repeat for $K = 2$ to N :
If $A[K] \neq 1$, then: Write: $A[K]$.
4. Exit.

2.10 Suppose $P(n) = a_0 + a_1n + a_2n^2 + \dots + a_mn^m$; that is, suppose degree $P(n) = m$. Prove that $P(n) = O(n^m)$.

Let $b_0 = |a_0|, b_1 = |a_1|, \dots, b_m = |a_m|$. Then, for $n \geq 1$,

$$\begin{aligned} P(n) &\leq b_0 + b_1n + b_2n^2 + \dots + b_mn^m = \left(\frac{b_0}{n^m} + \frac{b_1}{n^{m-1}} + \dots + b_m \right) n^m \\ &\leq (b_0 + b_1 + \dots + b_m) n^m = Mn^m \end{aligned}$$

where $M = |a_0| + |a_1| + \dots + |a_m|$. Hence $P(n) = O(n^m)$.

For example, $x^3 + 3x = O(x^3)$ and $x^5 - 4000000x^2 = O(x^5)$.

VARIABLES, DATA TYPES

2.11 Describe briefly the difference between local variables, parameters and global variables.

Local variables are variables which can be accessed only within a particular program or subprogram. Parameters are variables which are used to transfer data between a subprogram and its calling program. Global variables are variables which can be accessed by all of the program modules in a computer program. Each programming language which allows global variables has its own syntax for declaring them.

2.12 Suppose NUM denotes the number of records in a file. Describe the advantages in defining NUM to be a global variable. Describe the disadvantages in using global variables in general.

Many of the procedures will process all the records in the file using some type of loop. Since NUM will be the same for all these procedures, it would be advantageous to have NUM declared a global variable. Generally speaking, global and nonlocal variables may lead to errors caused by side effects, which may be difficult to detect.

2.13 Suppose a 32-bit memory location AAA contains the following sequence of bits:

0100 1101 1100 0001 1110 1001 0101 1101

Determine the data stored in AAA.

There is no way of knowing the data stored in AAA unless one knows the data type of AAA. If AAA is a character variable and the EBCDIC code is used for storing data, then (AZ) is stored in AAA. If AAA is an integer variable, then the integer with the above binary representation is stored in AAA.

2.14 Mathematically speaking, integers may also be viewed as real numbers. Give some reasons for having two different data types.

The arithmetic for integers, which are stored using some type of binary representation, is much simpler than the arithmetic for real numbers, which are stored using some type of exponential form. Also, certain round-off errors occurring in real arithmetic do not occur in integer arithmetic.

Supplementary Problems

MATHEMATICAL NOTATION AND FUNCTIONS

- 2.15 Find (a) $[3.4]$, $[-3.4]$, $[-7]$, $[\sqrt{75}]$, $[\sqrt[3]{75}]$, $[e]$; (b) $[3.4]$, $[-3.4]$, $[-7]$, $[\sqrt{75}]$, $[\sqrt[3]{75}]$, $[e]$.
- 2.16 (a) Find $48 \pmod{5}$, $48 \pmod{7}$, $1397 \pmod{11}$, $2468 \pmod{9}$.
 (b) Find $-48 \pmod{5}$, $-152 \pmod{7}$, $-358 \pmod{11}$, $-1326 \pmod{13}$.
 (c) Using arithmetic modulo 13, evaluate
 $9 + 10$, $8 + 12$, $3 + 4$, $3 - 4$, $2 - 7$, $5 - 8$
- 2.17 Find (a) $|3 + 8|$, $|3 - 8|$, $|-3 + 8|$, $|-3 - 8|$; (b) $7!$, $8!$, $14!/12!$, $15!/16!$
- 2.18 Find (a) 3^{-4} , $4^{7/2}$, $27^{-2/3}$; (b) $\log_2 64$, $\log_{10} 0.001$, $\log_2 (1/8)$; (c) $\lceil \lg 1\,000\,000 \rceil$, $\lfloor \lg 0.001 \rfloor$.

ALGORITHMS, COMPLEXITY

- 2.19 Consider the complexity function $C(n)$ which measures the number of times LOC is updated in Step 3 of Algorithm 2.3. Find $C(n)$ for the average case when $n = 4$, assuming all arrangements of the given four elements are equally likely. (Compare with Prob. 2.6.)
- 2.20 Consider Procedure P2.8, which finds the location LOC1 of the largest element and the location LOC2 of the second largest element in an array DATA with $n > 1$ elements. Let $C(n)$ denote the number of comparisons during the execution of the procedure.
 (a) Find $C(n)$ for the best case.
 (b) Find $C(n)$ for the worst case.
 (c) Find $C(n)$ for the average case for $n = 4$, assuming all arrangements of the given elements in DATA are equally likely.
- 2.21 Repeat Prob. 2.20, except now let $C(n)$ denote the number of times the values of FIRST and SECOND (or LOC1 and LOC2) must be updated.
- 2.22 Suppose the running time of a Module A is a constant M . Find the order of magnitude of the complexity function $C(n)$ which measures the execution time of each of the following algorithms, where n is the size of the input data (denoted by N in the algorithms).
- (a) Procedure P2.22A:
1. Repeat for $I = 1$ to N :
 2. Repeat for $J = 1$ to I :
 3. Repeat for $K = 1$ to J :
 4. Module A.
 - [End of Step 3 loop.]
 - [End of Step 2 loop.]
 - [End of Step 1 loop.]
 5. Exit.
- (b) Procedure P2.22B:
1. Set $J := N$.
 2. Repeat Steps 3 and 4 while $J > 1$.
 3. Module A.
 4. Set $J := J/2$.
 - [End of Step 2 loop.]
 5. Return.

Programming Problems

- 2.23 Write a function subprogram $\text{DIV}(J, K)$, where J and K are positive integers such that $\text{DIV}(J, K) = 1$ if J divides K but otherwise $\text{DIV}(J, K) = 0$. (For example, $\text{DIV}(3, 15) = 1$ but $\text{DIV}(3, 16) = 0$.)
- 2.24 Write a program using $\text{DIV}(J, K)$ which reads a positive integer $N > 10$ and determines whether or not N is a prime number. (*Hint*: N is prime if (i) $\text{DIV}(2, N) = 0$ (i.e., N is odd) and (ii) $\text{DIV}(K, N) = 0$ for all odd integers K where $1 < K^2 \leq N$.)
- 2.25 Translate Procedure P2.8 into a computer program; i.e., write a program which finds the location LOC1 of the largest element and the location LOC2 of the second largest element in an array DATA with $N > 1$ elements. Test the program using 70, 30, 25, 80, 60, 50, 30, 75, 25, and 60.
- 2.26 Translate the sieve method for finding prime numbers, described in Prob. 2.9, into a program to find the prime numbers less than N . Test the program using (a) $N = 1000$ and (b) $N = 10000$.
- 2.27 Let C denote the number of times LOC is updated using Algorithm 2.3 to find the largest element in an array A with N elements.
- Write a subprogram $\text{COUNT}(A, N, C)$ which finds C .
 - Write a Procedure P2.27 which (i) reads N random numbers between 0 and 1 into an array A and (ii) uses $\text{COUNT}(A, N, C)$ to find the value of C .
 - Write a program which repeats Procedure P2.27 1000 times and finds the average of the 1000 C 's.
 - Test the program for $N = 3$ and compare the result with the value obtained in Prob. 2.6.
 - Test the program for $N = 4$ and compare the result with the value in Prob. 2.19.

String Processing

3.1 INTRODUCTION

Historically, computers were first used for processing numerical data. Today, computers are frequently used for processing nonnumerical data, called *character data*. This chapter discusses how such data are stored and processed by the computer.

One of the primary applications of computers today is in the field of word processing. Such processing usually involves some type of pattern matching, as in checking to see if a particular word S appears in a given text T . We discuss this pattern matching problem in detail and, moreover, present two different pattern matching algorithms. The complexity of these algorithms is also investigated.

Computer terminology usually uses the term "string" for a sequence of characters rather than the term "word," since "word" has another meaning in computer science. For this reason, many texts sometimes use the expression "string processing," "string manipulation" or "text editing" instead of the expression "word processing."

The material in this chapter is essentially tangential and independent of the rest of the text. Accordingly, the reader or instructor may choose to omit this chapter on a first reading or cover this chapter at a later time.

3.2 BASIC TERMINOLOGY

Each programming language contains a *character set* that is used to communicate with the computer. This set usually includes the following:

Alphabet:	A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
Digits:	0 1 2 3 4 5 6 7 8 9
Special characters:	+ - / * () , . \$ = ' □

The set of special characters, which includes the blank space, frequently denoted by \square , varies somewhat from one language to another.

A finite sequence S of zero or more characters is called a *string*. The number of characters in a string is called its *length*. The string with zero characters is called the *empty string* or the *null string*. Specific strings will be denoted by enclosing their characters in single quotation marks. The quotation marks will also serve as string delimiters. Hence

'THE END' 'TO BE OR NOT TO BE' □□

are strings with lengths 7, 18, 0 and 2, respectively. We emphasize that the blank space is a character and hence contributes to the length of the string. Sometimes the quotation marks may be omitted when the context indicates that the expression is a string.

Let S_1 and S_2 be strings. The string consisting of the characters of S_1 followed by the characters of S_2 is called the *concatenation* of S_1 and S_2 ; it will be denoted by $S_1 // S_2$. For example,

'THE' // 'END' = 'THEEND' but 'THE' // '□' // 'END' = 'THE END'

Clearly the length of $S_1 // S_2$ is equal to the sum of the lengths of the strings S_1 and S_2 .

A string Y is called a *substring* of a string S if there exist strings X and Z such that

$$S = X // Y // Z$$

If X is an empty string, then Y is called an *initial substring* of S , and if Z is an empty string then Y is called a *terminal substring* of S . For example,

'BE OR NOT' is a substring of 'TO BE OR NOT TO BE'
 'THE' is an initial substring of 'THE END'

Clearly, if Y is a substring of S, then the length of Y cannot exceed the length of S.

Remark: Characters are stored in the computer using either a 6-bit, a 7-bit or an 8-bit code. The unit equal to the number of bits needed to represent a character is called a *byte*. However, unless otherwise stated or implied, a byte usually means 8 bits. A computer which can access an individual byte is called a *byte-addressable machine*.

3.3 STORING STRINGS

Generally speaking, strings are stored in three types of structures: (1) fixed-length structures, (2) variable-length structures with fixed maximums and (3) linked structures. We discuss each type of structure separately, giving its advantages and disadvantages.

Record-Oriented, Fixed-Length Storage

In fixed-length storage each line of print is viewed as a record, where all records have the same length, i.e., where each record accommodates the same number of characters. Since data are frequently input on terminals with 80-column images or using 80-column cards, we will assume our records have length 80 unless otherwise stated or implied.

EXAMPLE 3.1

Suppose the input consists of the FORTRAN program in Fig. 3-1. Using a record-oriented, fixed-length storage medium, the input data will appear in memory as pictured in Fig. 3-2, where we assume that 200 is the address of the first character of the program.

The main advantages of the above way of storing strings are:

- (1) The ease of accessing data from any given record
- (2) The ease of updating data in any given record (as long as the length of the new data does not exceed the record length)

The main disadvantages are:

- (1) Time is wasted reading an entire record if most of the storage consists of inessential blank spaces.
- (2) Certain records may require more space than available.
- (3) When the correction consists of more or fewer characters than the original text, changing a misspelled word requires the entire record to be changed.

```

C   PROGRAM PRINTING TWO INTEGERS IN INCREASING ORDER
    READ *, J, K
    IF(J.LE.K) THEN
      PRINT *, J, K
    ELSE
      PRINT *, K, J
    ENDIF
    STOP
  END

```

Fig. 3-1 Input data.

B

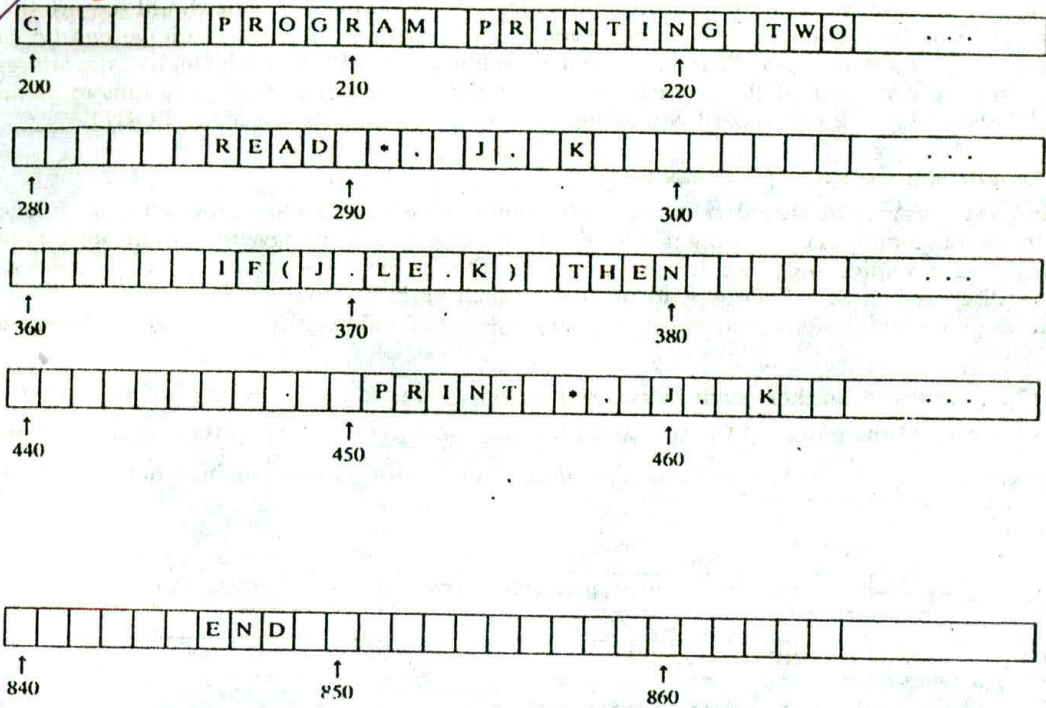


Fig. 3-2 Records stored sequentially in the computer.

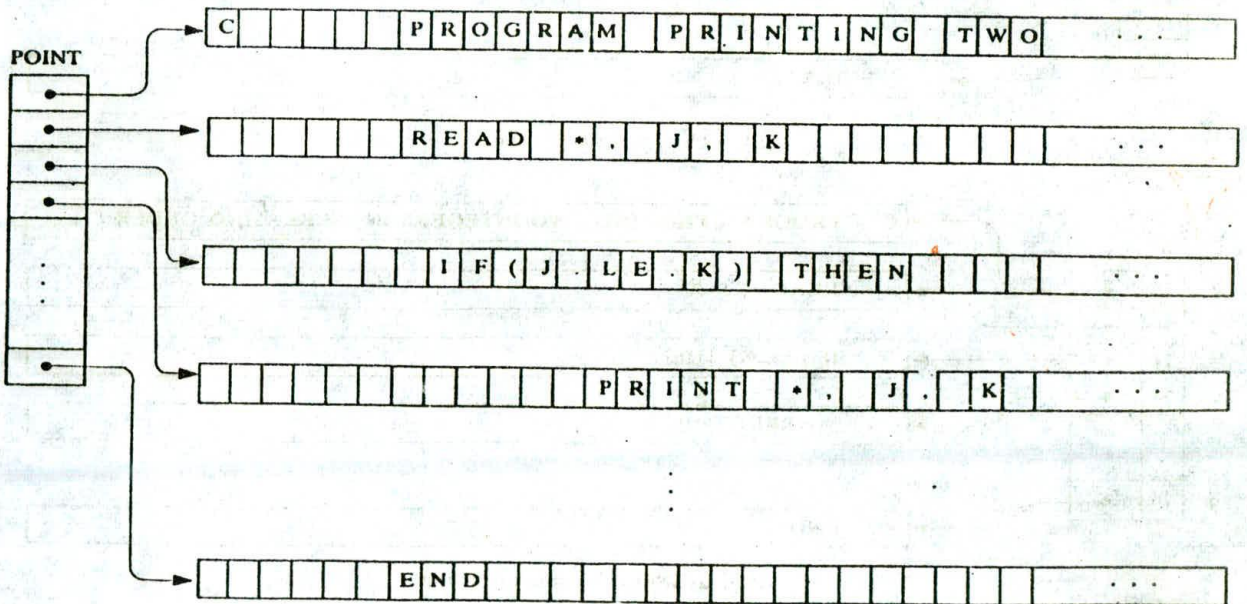


Fig. 3-3 Records stored using pointers.

Remark: Suppose we wanted to insert a new record in Example 3.1. This would require that all succeeding records be moved to new memory locations. However, this disadvantage can be easily remedied as indicated in Fig. 3-3. That is, one can use a linear array POINT which gives the address of each successive record, so that the records need not be stored in consecutive locations in memory. Accordingly, inserting a new record will require only an updating of the array POINT.

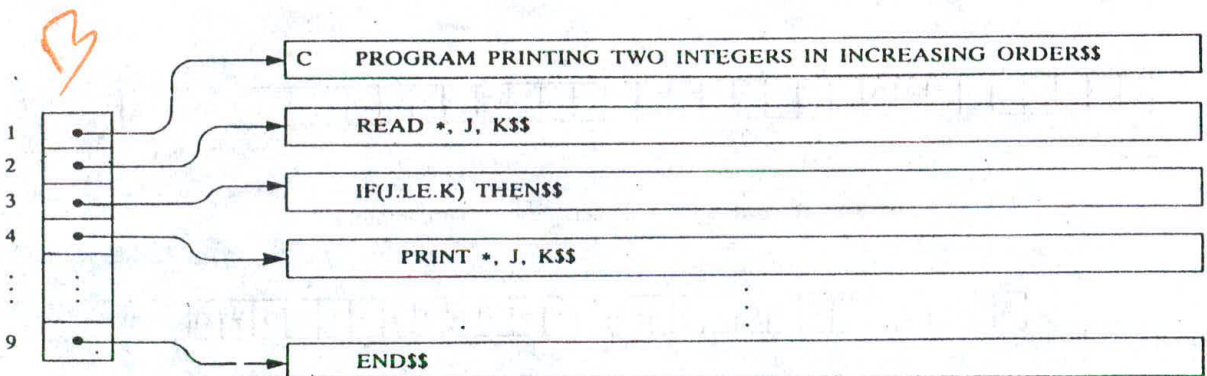
Variable-Length Storage with Fixed Maximum

Although strings may be stored in fixed-length memory locations as above, there are advantages in knowing the actual length of each string. For example, one then does not have to read the entire record when the string occupies only the beginning part of the memory location. Also, certain string operations (discussed in Sec. 3.4) depend on having such variable-length strings.

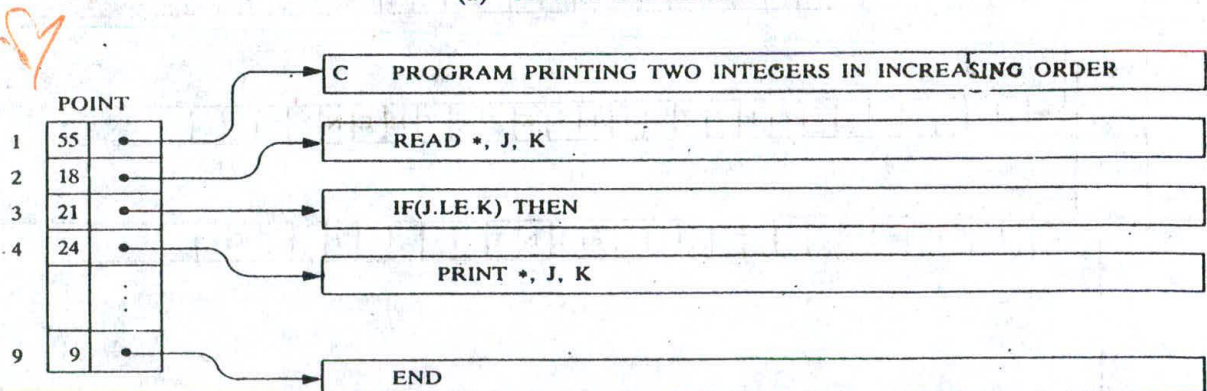
The storage of variable-length strings in memory cells with fixed lengths can be done in two general ways:

- (1) One can use a marker, such as two dollar signs (\$\$), to signal the end of the string.
- (2) One can list the length of the string—as an additional item in the pointer array, for example.

Using the data in Fig. 3-1, the first method is pictured in Fig. 3-4(a) and the second method is pictured in Fig. 3-4(b).



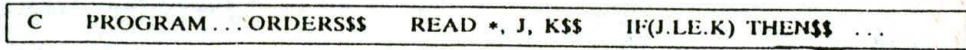
(a) Records with sentinels.



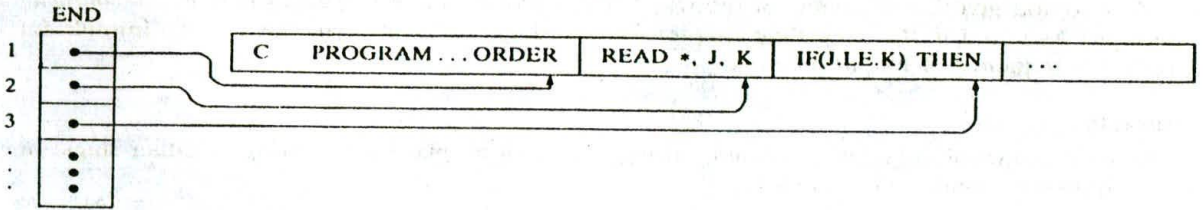
(b) Records whose lengths are listed.

Fig. 3-4

Remark: One might be tempted to store strings one after another by using some separation marker, such as the two dollar signs (\$\$) in Fig. 3-5(a), or by using a pointer array giving the location of the strings, as in Fig. 3-5(b). These ways of storing strings will obviously save space and are sometimes used in secondary memory when records are relatively permanent and require little change. However, such methods of storage are usually inefficient when the strings and their lengths are frequently being changed.



(a)



(b)

Fig. 3-5 Records stored one after another.

Linked Storage

Computers are being used very frequently today for word processing, i.e., for inputting, processing and outputting printed matter. Therefore, the computer must be able to correct and modify the printed matter, which usually means deleting, changing and inserting words, phrases, sentences and even paragraphs in the text. However, the fixed-length memory cells discussed above do not easily lend themselves to these operations. Accordingly, for most extensive word processing applications, strings are stored by means of linked lists. Such linked lists, and the way data are inserted and deleted in them, are discussed in detail in Chap. 5. Here we simply look at the way strings appear in these data structures.

By a (one-way) linked list, we mean a linearly ordered sequence of memory cells, called *nodes*, where each node contains an item, called a *link*, which points to the next node in the list (i.e., which contains the address of the next node). Figure 3-6 is a schematic diagram of such a linked list.

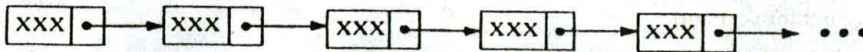
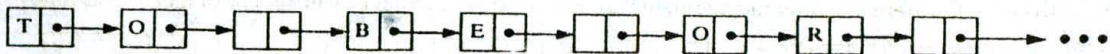
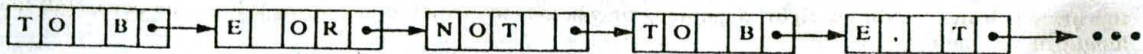


Fig. 3-6 Linked list.



(a) One character per node.



(b) Four characters per node.

Fig. 3-7

Strings may be stored in linked lists as follows. Each memory cell is assigned one character or a fixed number of characters, and a link contained in the cell gives the address of the cell containing the next character or group of characters in the string. For example, consider this famous quotation:

To be or not to be, that is the question.

Figure 3-7(a) shows how the string would appear in memory with one character per node, and Fig. 3-7(b) shows how it would appear with four characters per node.

3.4 CHARACTER DATA TYPE

This section gives an overview of the way various programming languages handle the *character* data type. As noted in the preceding chapter (in Sec. 2.7), each data type has its own formula for decoding a sequence of bits in memory.

Constants

Many programming languages denote string constants by placing the string in either single or double quotation marks. For example,

'THE END' and 'TO BE OR NOT TO BE'

are string constants of lengths 7 and 18 characters respectively. Our algorithms will also define character constants in this way.

Variables

Each programming language has its own rules for forming character variables. However, such variables fall into one of three categories: static, semistatic and dynamic. By a *static* character variable, we mean a variable whose length is defined before the program is executed and cannot change throughout the program. By a *semistatic* character variable, we mean a variable whose length may vary during the execution of the program as long as the length does not exceed a maximum value determined by the program before the program is executed. By a *dynamic* character variable, we mean a variable whose length can change during the execution of the program. These three categories correspond, respectively, to the ways the strings are stored in the memory of the computer as discussed in the preceding section.

EXAMPLE 3.2

- (a) Many versions of FORTRAN use static CHARACTER variables. For example, consider the following FORTRAN program segment:

```
CHARACTER ST1*10, ST2*14
ST1 = 'THE END'
ST2 = 'TO BE OR NOT TO BE'
```

The first statement declares ST1 and ST2 to be CHARACTER variables with lengths 10 and 14, respectively. After both assignment statements are executed, ST1 and ST2 will appear in memory as follows:

ST1	T	H	E		E	N	D									ST2	T	O		B	E		O	R		N	O	T		T
-----	---	---	---	--	---	---	---	--	--	--	--	--	--	--	--	-----	---	---	--	---	---	--	---	---	--	---	---	---	--	---

That is, a string is stored left-justified in memory. Either blank spaces are added on the right of the string, or the string is truncated on the right, depending on whether the length of the string is less than or exceeds the length of the memory location.

- (b) BASIC defines character variables as those variables whose name ends with a dollar sign. Generally speaking, the variables are semistatic ones whose lengths cannot exceed a fixed bound. For example, the BASIC program segment

```
A$ = 'THE BEGINNING'
B$ = 'THE END'
```

defines A\$ and B\$ to be character variables. When the segment is executed, the lengths of A\$ and B\$ will be 13 and 7, respectively.

Also, BASIC uses double quotation marks to denote string constants.

- (c) SNOBOL uses dynamic character variables. For example, the SNOBOL program segment

```
WORD = 'COMPUTER'
TEXT = 'IN THE BEGINNING'
```

defines WORD and TEXT as character variables. When the segment is executed, the lengths of WORD and TEXT will be 8 and 16, respectively. However, the lengths may change later in the program.

- (d) PL/1 uses both static and semistatic CHARACTER variables. For example, the PL/1 statement

```
DECLARE NAME CHARACTER(20),
        WORD CHARACTER(15) VARYING;
```

designates NAME as a static CHARACTER variable of length 20 and designates WORD as a semistatic CHARACTER variable whose length may vary but may not exceed 15.

- (e) In Pascal, a character variable (abbreviated CHAR) can represent only a single character, and hence a string is represented by a linear array of characters. For example,

```
VAR WORD: ARRAY[1..20] OF CHAR
```

declares WORD to be a string of 20 characters. Furthermore, WORD[1] is the first character of the string, WORD[2] the second character and so on. In particular, CHAR arrays have fixed lengths and hence are static variables.

3.5 STRING OPERATIONS

Although a string may be viewed simply as a sequence or linear array of characters, there is a fundamental difference in use between strings and other types of arrays. Specifically, groups of consecutive elements in a string (such as words, phrases and sentences), called *substrings*, may be units unto themselves. Furthermore, the basic units of access in a string are usually these substrings, not individual characters.

Consider, for example, the string

'TO BE OR NOT TO BE'

We may view the string as the 18-character sequence T, O, \square , B, . . . , E. However, the substrings TO, BE, OR, . . . have their own meaning.

On the other hand, consider an 18-element linear array of 18 integers,

4, 8, 6, 15, 9, 5, 4, 13, 8, 5, 11, 9, 9, 13, 7, 10, 6, 11


The basic unit of access in such an array is usually an individual element. Groups of consecutive elements normally do not have any special meaning.

For the above reason, various string operations have been developed which are not normally used with other kinds of arrays. This section discusses these string-oriented operations. The next section shows how these operations are used in word processing. (Unless otherwise stated or implied, we assume our character-type variables are dynamic and have a variable length determined by the context in which the variable is used.)

Substring

Accessing a substring from a given string requires three pieces of information: (1) the name of the string or the string itself, (2) the position of the first character of the substring in the given string and

(3) the length of the substring or the position of the last character of the substring. We call this operation SUBSTRING. Specifically, we write

 SUBSTRING(string, initial, length)

to denote the substring of a string S beginning in a position K and having a length L.

EXAMPLE 3.3

(a) Using the above function we have:


SUBSTRING('TO BE OR NOT TO BE', 4, 7) = 'BE OR N'
SUBSTRING('THE END', 4, 4) = '□END'

(b) Our function SUBSTRING(S, 4, 7) is denoted in some programming languages as follows:

PL/1:	SUBSTR(S, 4, 7)
FORTRAN 77:	S(4:10)
UCSD Pascal:	COPY(S, 4, 7)
BASIC:	MID\$(S, 4, 7)

Indexing

Indexing, also called *pattern matching*, refers to finding the position where a string pattern P first appears in a given string text T. We call this operation INDEX and write

 INDEX(text, pattern)

If the pattern P does not appear in the text T, then INDEX is assigned the value 0. The arguments "text" and "pattern" can be either string constants or string variables.

EXAMPLE 3.4

(a) Suppose T contains the text

'HIS FATHER IS THE PROFESSOR'

Then,

INDEX(T, 'THE'), INDEX(T, 'THEN') and INDEX(T, '□THE□')

have the values 7, 0 and 14, respectively.

(b) The function INDEX(text, pattern) is denoted in some of the programming languages as follows:

PL/1:	INDEX(text, pattern)
UCSD Pascal:	POS(pattern, text)

Observe the reverse order of the arguments in UCSD Pascal.

Concatenation

Let S_1 and S_2 be strings. Recall (Sec. 3.2) that the *concatenation* of S_1 and S_2 , which we denote by $S_1 // S_2$, is the string consisting of the characters of S_1 followed by the characters of S_2 .

EXAMPLE 3.5

(a) Suppose $S_1 = \text{'MARK'}$ and $S_2 = \text{'TWIN'}$. Then:

$S_1 // S_2 = \text{'MARKTWIN'}$ but $S_1 // \text{'□'} // S_2 = \text{'MARK TWIN'}$

(b) Concatenation is denoted in some programming languages as follows:

PL/1:	$S_1 S_2$
FORTRAN 77:	$S_1 // S_2$
BASIC:	$S_1 + S_2$
NOBOL:	$S_1 S_2$ (juxtaposition with a blank space between S_1 and S_2)

Length

The number of characters in a string is called its length. We will write

$LENGTH(string)$

for the length of a given string. Thus

$LENGTH('COMPUTER') = 8$ $LENGTH('□') = 0$ $LENGTH('') = 0$

Some of the programming languages denote this function as follows:

PL/1:	$LENGTH(string)$
BASIC:	$LEN(string)$
UCSD Pascal:	$LENGTH(string)$
NOBOL:	$SIZE(string)$

FORTRAN and standard Pascal, which use fixed-length string variables, do not have any built-in $LENGTH$ functions for strings. However, such variables may be viewed as having variable length if one ignores all trailing blanks. Accordingly, one could write a subprogram $LENGTH$ in these languages so that

$LENGTH('MARC ') = 4$

In fact, SNOBOL has a built-in string function $TRIM$ which omits trailing blanks:

$TRIM('ERIK ') = 'ERIK'$

This $TRIM$ function is occasionally used in our algorithms.

3.6 WORD PROCESSING

In earlier times, character data processed by the computer consisted mainly of data items, such as names and addresses. Today the computer also processes printed matter, such as letters, articles and reports. It is in this latter context that we use the term "word processing."

Given some printed text, the operations usually associated with word processing are the following:

- Replacement.* Replacing one string in the text by another.
- Insertion.* Inserting a string in the middle of the text.
- Deletion.* Deleting a string from the text.

The above operations can be executed by using the string operations discussed in the preceding section. This we show below when we discuss each operation separately. Many of these operations are built into or can easily be defined in each of the programming languages that we have cited.

Insertion

Suppose in a given text T we want to insert a string S so that S begins in position K . We denote this operation by

$INSERT(text, position, string)$

For example,

$$\begin{aligned}\text{INSERT}('ABCDEFGF', 3, 'XYZ') &= 'ABXYZCDEFG' \\ \text{INSERT}('ABCDEFGF', 6, 'XYZ') &= 'ABCDEXYZFG'\end{aligned}$$

This INSERT function can be implemented by using the string operations defined in the previous section as follows:

$$\text{INSERT}(T, K, S) = \text{SUBSTRING}(T, 1, K - 1) // S // \text{SUBSTRING}(T, K, \text{LENGTH}(T) - K + 1)$$

That is, the initial substring of T before the position K , which has length $K - 1$, is concatenated with the string S , and the result is concatenated with the remaining part of T , which begins in position K and has length $\text{LENGTH}(T) - (K - 1) = \text{LENGTH}(T) - K + 1$. (We are assuming implicitly that T is a dynamic variable and that the size of T will not become too large.)

Deletion

Suppose in a given text T we want to delete the substring which begins in position K and has length L . We denote this operation by

$$\text{DELETE}(\text{text}, \text{position}, \text{length})$$

For example,

$$\begin{aligned}\text{DELETE}('ABCDEFGF', 4, 2) &= 'ABCFGF' \\ \text{DELETE}('ABCDEFGF', 2, 4) &= 'AFGF'\end{aligned}$$

We assume that nothing is deleted if position $K = 0$. Thus

$$\text{DELETE}('ABCDEFGF', 0, 2) = 'ABCDEFGF'$$

The importance of this "zero case" is seen later.

The DELETE function can be implemented using the string operations given in the preceding section as follows:

$$\text{DELETE}(T, K, L) = \text{SUBSTRING}(T, 1, K - 1) // \text{SUBSTRING}(T, K + L, \text{LENGTH}(T) - K - L + 1)$$

That is, the initial substring of T before position K is concatenated with the terminal substring of T beginning in position $K + L$. The length of the initial substring is $K - 1$, and the length of the terminal substring is:

$$\text{LENGTH}(T) - (K + L - 1) = \text{LENGTH}(T) - K - L + 1$$

We also assume that $\text{DELETE}(T, K, L) = T$ when $K = 0$.

Now suppose text T and pattern P are given and we want to delete from T the first occurrence of the pattern P . This can be accomplished by using the above DELETE function as follows:

$$\text{DELETE}(T, \text{INDEX}(T, P), \text{LENGTH}(P))$$

That is, in the text T , we first compute $\text{INDEX}(T, P)$, the position where P first occurs in T , and then we compute $\text{LENGTH}(P)$, the number of characters in P . Recall that when $\text{INDEX}(T, P) = 0$ (i.e., when P does not occur in T) the text T is not changed.

EXAMPLE 3.6

(a) Suppose $T = 'ABCDEFGF'$ and $P = 'CD'$. Then $\text{INDEX}(T, P) = 3$ and $\text{LENGTH}(P) = 2$. Hence

$$\text{DELETE}('ABCDEFGF', 3, 2) = 'ABEFGF'$$

- (b) Suppose $T = \text{'ABCDEFGG'}$ and $P = \text{'DC'}$. Then $\text{INDEX}(T, P) = 0$ and $\text{LENGTH}(P) = 2$. Hence, by the "zero case,"

$$\text{DELETE}(\text{'ABCDEFGG'}, 0, 2) = \text{'ABCDEFGG'}$$

as expected.

Suppose after reading into the computer a text T and a pattern P , we want to delete every occurrence of the pattern P in the text T . This can be accomplished by repeatedly applying

$$\text{DELETE}(T, \text{INDEX}(T, P), \text{LENGTH}(P))$$

until $\text{INDEX}(T, P) = 0$ (i.e., until P does not appear in T). An algorithm which accomplishes this follows.

Algorithm 3.1: A text T and a pattern P are in memory. This algorithm deletes every occurrence of P in T .

1. [Find index of P .] Set $K := \text{INDEX}(T, P)$.
2. Repeat while $K \neq 0$:
 - (a) [Delete P from T .]
Set $T := \text{DELETE}(T, \text{INDEX}(T, P), \text{LENGTH}(P))$
 - (b) [Update index.] Set $K := \text{INDEX}(T, P)$.
- [End of loop.]
3. Write: T .
4. Exit.

We emphasize that after each deletion, the length of T decreases and hence the algorithm must stop. However, the number of times the loop is executed may exceed the number of times P appears in the original text T , as illustrated in the following example.

EXAMPLE 3.7

- (a) Suppose Algorithm 3.1 is run with the data

$$T = \text{XABYABZ}, \quad P = \text{AB}$$

Then the loop in the algorithm will be executed twice. During the first execution, the first occurrence of AB in T is deleted, with the result that $T = \text{XYABZ}$. During the second execution, the remaining occurrence of AB in T is deleted, so that $T = \text{XYZ}$. Accordingly, XYZ is the output.

- (b) Suppose Algorithm 3.1 is run with the data

$$T = \text{XAAABBBY}, \quad P = \text{AB}$$

Observe that the pattern AB occurs only once in T but the loop in the algorithm will be executed three times. Specifically, after AB is deleted the first time from T we have $T = \text{XAABBY}$, and hence AB appears again in T . After AB is deleted a second time from T , we see that $T = \text{XABY}$ and AB still occurs in T . Finally, after AB is deleted a third time from T , we have $T = \text{XY}$ and AB does not appear in T , and thus $\text{INDEX}(T, P) = 0$. Hence XY is the output.

The above example shows that when a text T is changed by a deletion, patterns may occur that did not appear originally.

Replacement

Suppose in a given text T we want to replace the first occurrence of a pattern P_1 by a pattern P_2 . We will denote this operation by

$$\text{REPLACE}(\text{text}, \text{pattern}_1, \text{pattern}_2)$$

For example

```
REPLACE('XABYABZ', 'AB', 'C') = 'XCYABZ'
REPLACE('XABYABZ', 'BA', 'C') = 'XABYABZ'
```

In the second case, the pattern BA does not occur, and hence there is no change.

We note that this REPLACE function can be expressed as a deletion followed by an insertion if we use the preceding DELETE and INSERT functions. Specifically, the REPLACE function can be executed by using the following three steps:

```
K := INDEX(T, P1)
T := DELETE(T, K, LENGTH(P1))
INSERT(T, K, P2)
```

The first two steps delete P₁ from T, and the third step inserts P₂ in the position K from which P₁ was deleted.

Suppose a text T and patterns P and Q are in the memory of a computer. Suppose we want to replace every occurrence of the pattern P in T by the pattern Q. This might be accomplished by repeatedly applying

REPLACE(T, P, Q)

until INDEX(T, P) = 0 (i.e., until P does not appear in T). An algorithm which does this follows.

Algorithm 3.2: A text T and patterns P and Q are in memory. This algorithm replaces every occurrence of P in T by Q.

1. [Find index of P.] Set K := INDEX(T, P).
2. Repeat while K ≠ 0:
 - (a) [Replace P by Q.] Set T := REPLACE(T, P, Q).
 - (b) [Update index.] Set K := INDEX(T, P).
- [End of loop.]
3. Write: T.
4. Exit.

Warning: Although this algorithm looks very much like Algorithm 3.1, there is no guarantee that this algorithm will terminate. This fact is illustrated in Example 3.8(b). On the other hand, suppose the length of Q is smaller than the length of P. Then the length of T after each replacement decreases. This guarantees that in this special case where Q is smaller than P the algorithm must terminate.

EXAMPLE 3.8

(a) Suppose Algorithm 3.2 is run with the data

T = XABYABZ, P = AB, Q = C

Then the loop in the algorithm will be executed twice. During the first execution, the first occurrence of AB in T is replaced by C to yield T = XCYABZ. During the second execution, the remaining AB in T is replaced by C to yield T = XCYCZ. Hence XCYCZ is the output.

(b) Suppose Algorithm 3.2 is run with the data

T = XAY, P = A, Q = AB

Then the algorithm will never terminate. The reason for this is that P will always occur in the text T, no matter how many times the loop is executed. Specifically,

T = XABY at the end of the first execution of the loop
 T = XAB²Y at the end of the second execution of the loop

 T = XABⁿY at the end of the *n*th execution of the loop

(The infinite loop arises here since P is a substring of Q.)

3.7 PATTERN MATCHING ALGORITHMS

Pattern matching is the problem of deciding whether or not a given string pattern P appears in a string text T. We assume that the length of P does not exceed the length of T. This section discusses two pattern matching algorithms. We also discuss the complexity of the algorithms so we can compare their efficiencies.

Remark: During the discussion of pattern matching algorithms, characters are sometimes denoted by lowercase letters (*a, b, c, . . .*) and exponents may be used to denote repetition; e.g.,

$$a^2b^3ab^2 \text{ for } aabbbabb \quad \text{and} \quad (cd)^3 \text{ for } cdcdcd$$

In addition, the empty string may be denoted by Λ , the Greek letter lambda, and the concatenation of strings X and Y may be denoted by $X \cdot Y$ or, simply, XY.

First Pattern Matching Algorithm

The first pattern matching algorithm is the obvious one in which we compare a given pattern P with each of the substrings of T, moving from left to right, until we get a match. In detail, let

$$W_k = \text{SUBSTRING}(T, K, \text{LENGTH}(P))$$

That is, let W_k denote the substring of T having the same length as P and beginning with the *K*th character of T. First we compare P, character by character, with the first substring, W_1 . If all the characters are the same, then $P = W_1$ and so P appears in T and $\text{INDEX}(T, P) = 1$. On the other hand, suppose we find that some character of P is not the same as the corresponding character of W_1 . Then $P \neq W_1$ and we can immediately move on to the next substring, W_2 . That is, we next compare P with W_2 . If $P \neq W_2$, then we compare P with W_3 , and so on. The process stops (a) when we find a match of P with some substring W_k and so P appears in T and $\text{INDEX}(T, P) = K$, or (b) when we exhaust all the W_k 's with no match and hence P does not appear in T. The maximum value MAX of the subscript K is equal to $\text{LENGTH}(T) - \text{LENGTH}(P) + 1$.

Let us assume, as an illustration, that P is a 4-character string and that T is a 20-character string, and that P and T appear in memory as linear arrays with one character per element. That is,

$$P = P[1]P[2]P[3]P[4] \quad \text{and} \quad T = T[1]T[2]T[3] \cdots T[19]T[20]$$

Then P is compared with each of the following 4-character substrings of T:

$$W_1 = T[1]T[2]T[3]T[4], \quad W_2 = T[2]T[3]T[4]T[5], \quad \dots, \quad W_{17} = T[17]T[18]T[19]T[20]$$

Note that there are $\text{MAX} = 20 - 4 + 1 = 17$ such substrings of T.

A formal presentation of our algorithm, where P is an *r*-character string and T is an *s*-character string, is shown in Algorithm 3.3.

Observe that Algorithm 3.3 contains two loops, one inside the other. The outer loop runs through each successive R-character substring

$$W_k = T[K]T[K + 1] \cdots T[K + R - 1]$$

of T. The inner loop compares P with W_k , character by character. If any character does not match, then control transfers to Step 5, which increases K and then leads to the next substring of T. If all the R

Algorithm 3.3: (Pattern Matching) P and T are strings with lengths R and S, respectively, and are stored as arrays with one character per element. This algorithm finds the INDEX of P in T.

1. [Initialize.] Set $K := 1$ and $MAX := S - R + 1$.
2. Repeat Steps 3 to 5 while $K \leq MAX$:
3. Repeat for $L = 1$ to R : [Tests each character of P.]
 If $P[L] \neq T[K + L - 1]$, then: Go to Step 5.
 [End of inner loop.]
4. [Success.] Set $INDEX = K$, and Exit.
5. Set $K := K + 1$.
 [End of Step 2 outer loop.]
6. [Failure.] Set $INDEX = 0$.
7. Exit.

characters of P do match those of some W_K , then P appears in T and K is the INDEX of P in T. On the other hand, if the outer loop completes all of its cycles, then P does not appear in T and so $INDEX = 0$.

The complexity of this pattern matching algorithm is measured by the number C of comparisons between characters in the pattern P and characters of the text T. In order to find C, we let N_k denote the number of comparisons that take place in the inner loop when P is compared with W_k . Then

$$C = N_1 + N_2 + \cdots + N_L$$

where L is the position L in T where P first appears or $L = MAX$ if P does not appear in T. The next example computes C for some specific P and T where $LENGTH(P) = 4$ and $LENGTH(T) = 20$ and so $MAX = 20 - 4 + 1 = 17$.

EXAMPLE 3.9

(a) Suppose $P = aaba$ and $T = cdcd \cdots cd = (cd)^{10}$. Clearly P does not occur in T. Also, for each of the 17 cycles, $N_k = 1$, since the first character of P does not match W_k . Hence

$$C = 1 + 1 + 1 + \cdots + 1 = 17$$

(b) Suppose $P = aaba$ and $T = ababaaba \dots$. Observe that P is a substring of T. In fact, $P = W_5$ and so $N_5 = 4$. Also, comparing P with $W_1 = abab$, we see that $N_1 = 2$, since the first letters do match; but comparing P with $W_2 = baba$, we see that $N_2 = 1$, since the first letters do not match. Similarly, $N_3 = 2$ and $N_4 = 1$. Accordingly,

$$C = 2 + 1 + 2 + 1 + 4 = 10$$

(c) Suppose $P = aaab$ and $T = aa \cdots a = a^{20}$. Here P does not appear in T. Also, every $W_k = aaaa$; hence every $N_k = 4$, since the first three letters of P do match. Accordingly,

$$C = 4 + 4 + \cdots + 4 = 17 \cdot 4 = 68$$

In general, when P is an r -character string and T is an s -character string, the data size for the algorithm is

$$n = r + s$$

The worst case occurs when every character of P except the last matches every substring W_k , as in Example 3.9(c). In this case, $C(n) = r(s - r + 1)$. For fixed n , we have $s = n - r$, so that

$$C(n) = r(n - 2r + 1)$$

The maximum value of $C(n)$ occurs when $r = (n + 1)/4$. (See Prob. 3.19.) Accordingly, substituting this value for r in the formula for $C(n)$ yields

$$C(n) = \frac{(n + 1)^2}{8} = O(n^2)$$

The complexity of the average case in any actual situation depends on certain probabilities which are usually unknown. When the characters of P and T are randomly selected from some finite alphabet, the complexity of the average case is still not easy to analyze, but the complexity of the average case is still a factor of the worst case. Accordingly, we shall state the following: *The complexity of this pattern matching algorithm is equal to $O(n^2)$.* In other words, the time required to execute this algorithm is proportional to n^2 . (Compare this result with the one on page 57.)

Second Pattern Matching Algorithm

The second pattern matching algorithm uses a table which is derived from a particular pattern P but is independent of the text T. For definiteness, suppose

$$P = aaba$$

First we give the reason for the table entries and how they are used. Suppose $T = T_1T_2T_3 \dots$, where T_i denotes the i th character of T; and suppose the first two characters of T match those of P; i.e., suppose $T = aa \dots$. Then T has one of the following three forms:

- (i) $T = aab \dots$, (ii) $T = aaa \dots$, (iii) $T = aax$

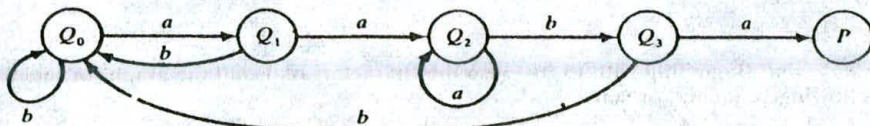
where x is any character different from a or b . Suppose we read T_3 and find that $T_3 = b$. Then we next read T_4 to see if $T_4 = a$, which will give a match of P with W_1 . On the other hand, suppose $T_3 = a$. Then we know that $P \neq W_1$; but we also know that $W_2 = aa \dots$, i.e., that the first two characters of the substring W_2 match those of P. Hence we next read T_4 to see if $T_4 = b$. Last, suppose $T_3 = x$. Then we know that $P \neq W_1$, but we also know that $P \neq W_2$ and $P \neq W_3$, since x does not appear in P. Hence we next read T_4 to see if $T_4 = a$, i.e., to see if the first character of W_4 matches the first character of P.

There are two important points to the above procedure. First, when we read T_3 we need only compare T_3 with those characters which appear in P. If none of these match, then we are in the last case, of a character x which does not appear in P. Second, after reading and checking T_3 , we next read T_4 ; we do not have to go back again in the text T.

Figure 3-8(a) contains the table that is used in our second pattern matching algorithm for the pattern $P = aaba$. (In both the table and the accompanying graph, the pattern P and its substrings Q

	a	b	x
Q ₀	Q ₁	Q ₀	Q ₀
Q ₁	Q ₂	Q ₀	Q ₀
Q ₂	Q ₂	Q ₃	Q ₀
Q ₃	P	Q ₀	Q ₀

(a) Pattern matching table.



(b) Pattern matching graph.

Fig. 3-8

will be represented by italic capital letters.) The table is obtained as follows. First of all, we let Q_i denote the initial substring of P of length i ; hence

$$Q_0 = \Lambda, \quad Q_1 = a, \quad Q_2 = a^2, \quad Q_3 = a^2b, \quad Q_4 = a^2ba = P$$

(Here $Q_0 = \Lambda$ is the empty string.) The rows of the table are labeled by these initial substrings of P , excluding P itself. The columns of the table are labeled a , b and x , where x represents any character that doesn't appear in the pattern P . Let f be the function determined by the table; i.e., let

$$f(Q_i, t)$$

denote the entry in the table in row Q_i and column t (where t is any character). This entry $f(Q_i, t)$ is defined to be the largest Q that appears as a terminal substring in the string $Q_i t$, the concatenation of Q_i and t . For example,

a^2 is the largest Q that is a terminal substring of $Q_2 a = a^3$, so $f(Q_2, a) = Q_2$

Λ is the largest Q that is a terminal substring of $Q_1 b = ab$, so $f(Q_1, b) = Q_0$

a is the largest Q that is a terminal substring of $Q_0 a = a$, so $f(Q_0, a) = Q_1$

Λ is the largest Q that is a terminal substring of $Q_3 x = a^2bx$, so $f(Q_3, x) = Q_0$

and so on. Although $Q_1 = a$ is a terminal substring of $Q_2 a = a^3$, we have $f(Q_2, a) = Q_2$ because Q_2 is also a terminal substring of $Q_2 a = a^3$ and Q_2 is larger than Q_1 . We note that $f(Q_i, x) = Q_0$ for any Q_i , since x does not appear in the pattern P . Accordingly, the column corresponding to x is usually omitted from the table.

Our table can also be pictured by the labeled directed graph in Fig. 3-8(b). The graph is obtained as follows. First, there is a node in the graph corresponding to each initial substring Q_i of P . The Q_i 's are called the *states* of the system, and Q_0 is called the *initial state*. Second, there is an arrow (a directed edge) in the graph corresponding to each entry in the table. Specifically, if

$$f(Q_i, t) = Q_j$$

then there is an arrow labeled by the character t from Q_i to Q_j . For example, $f(Q_2, b) = Q_3$, so there is an arrow labeled b from Q_2 to Q_3 . For notational convenience, we have omitted all arrows labeled x , which must lead to the initial state Q_0 .

We are now ready to give the second pattern matching algorithm for the pattern $P = aaba$. (Note that in the following discussion capital letters will be used for all single-letter variable names that appear in the algorithm.) Let $T = T_1 T_2 T_3 \cdots T_N$ denote the n -character-string text which is searched for the pattern P . Beginning with the initial state Q_0 and using the text T , we will obtain a sequence of states S_1, S_2, S_3, \dots as follows. We let $S_1 = Q_0$ and we read the first character T_1 . From either the table or the graph in Fig. 3-8, the pair (S_1, T_1) yields a second state S_2 ; that is, $F(S_1, T_1) = S_2$. We read the next character T_2 . The pair (S_2, T_2) yields a state S_3 , and so on. There are two possibilities:

- (1) Some state $S_k = P$, the desired pattern. In this case, P does appear in T and its index is $k - \text{LENGTH}(P)$.
- (2) No state S_1, S_2, \dots, S_{N+1} is equal to P . In this case, P does not appear in T .

We illustrate the algorithm with two different texts using the pattern $P = aaba$.

EXAMPLE 3.10

(a) Suppose $T = aabcaba$. Beginning with Q_0 , we use the characters of T and the graph (or table) in Fig. 3-8 to obtain the following sequence of states:

$$Q_0 \xrightarrow{ca} Q_1 \xrightarrow{ca} Q_2 \xrightarrow{cb} Q_3 \xrightarrow{cc} Q_0 \xrightarrow{ca} Q_1 \xrightarrow{cb} Q_0 \xrightarrow{ca} Q_1$$

We do not obtain the state P , so P does not appear in T .

(b) Suppose $T = abcaabaca$. Then we obtain the following sequence of states:

$$Q_0 \xrightarrow{ca} Q_1 \xrightarrow{cb} Q_0 \xrightarrow{cc} Q_0 \xrightarrow{ca} Q_1 \xrightarrow{ca} Q_2 \xrightarrow{cb} Q_3 \xrightarrow{ca} P$$

Here we obtain the pattern P as the state S_8 . Hence P does appear in T and its index is $8 - \text{LENGTH}(P) = 4$.
 The formal statement of our second pattern matching algorithm follows:

Algorithm 3.4: (Pattern Matching). The pattern matching table $F(Q_i, T)$ of a pattern P is in memory, and the input is an N-character string $T = T_1 T_2 \dots T_N$. This algorithm finds the INDEX of P in T.

1. [Initialize.] Set $K := 1$ and $S_1 = Q_0$.
2. Repeat Steps 3 to 5 while $S_K \neq P$ and $K \leq N$.
3. Read T_K .
4. Set $S_{K+1} := F(S_K, T_K)$. [Finds next state.]
5. Set $K := K + 1$. [Updates counter.]
- [End of Step 2 loop.]
6. [Successful?]
 - If $S_K = P$, then:
 - INDEX = $K - \text{LENGTH}(P)$.
 - Else:
 - INDEX = 0.
 - [End of If structure.]
7. Exit.

The running time of the above algorithm is proportional to the number of times the Step 2 loop is executed. The worst case occurs when all of the text T is read, i.e., when the loop is executed $n = \text{LENGTH}(T)$ times. Accordingly, we can state the following: *The complexity of this pattern matching algorithm is equal to $O(n)$.*

Remark: A combinatorial problem is said to be *solvable in polynomial time* if there is an algorithmic solution with complexity equal to $O(n^m)$ for some m , and it is said to be *solvable in linear time* if there is an algorithmic solution with complexity equal to $O(n)$, where n is the size of the data. Thus the second of the two pattern matching algorithms described in this section is solvable in linear time. (The first pattern matching algorithm was solvable in polynomial time.)

Solved Problems

TERMINOLOGY; STORAGE OF STRINGS

3.1 Let W be the string ABCD. (a) Find the length of W. (b) List all substrings of W. (c) List all the initial substrings of W.

(a) The number of characters in W is its length, so 4 is the length of W.

(b) Any subsequence of characters of W is a substring of W. There are 11 such substrings:

Substrings:	<u>ABCD</u> ,	<u>ABC, BCD</u> ,	<u>AB, BC, CD</u> ,	<u>A, B, C, D</u> ,	<u>Λ</u>
Lengths:	4	3	2	1	0

(Here Λ denotes the empty string.)

(c) The initial substrings are ABCD, ABC, AB, A, Λ ; that is, both the empty string and those substrings that begin with A.

- 3.2** Assuming a programming language uses at least 48 characters—26 letters, 10 digits and a minimum of 12 special characters—give the minimum number and the usual number of bits to represent a character in the memory of the computer.

Since $2^5 < 48 < 2^6$, one requires at least a 6-bit code to represent 48 characters. Usually a computer uses a 7-bit code, such as ASCII, or an 8-bit code, such as EBCDIC, to represent characters. This allows many more special characters to be represented and processed by the computer.

- 3.3** Describe briefly the three types of structures used for storing strings.

- Fixed-length-storage structures. Here strings are stored in memory cells that are all of the same length, usually space for 80 characters.
- Variable-length storage with fixed maximums. Here strings are also stored in memory cells all of the same length; however, one also knows the actual length of the string in the cell.
- Linked-list storage. Here each cell is divided into two parts; the first part stores a single character (or a fixed small number of characters), and the second part contains the address of the cell containing the next character.

- 3.4** Find the string stored in Fig. 3-9, assuming the link value 0 signals the end of the list.

	START		CHAR	LINK
	4	1	OY F	10
		2	ING	7
		3		
		4	A TH	2
		5		
		6	ER.	0
		7	OF B	11
		8	A J	1
		9		
		10	OREV	6
		11	EAUT	12
		12	Y IS	8

Fig. 3-9

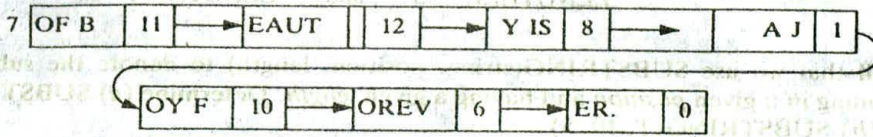
Here the string is stored in a linked-list structure with 4 characters per node. The value of START gives the location of the first node in the list:

4	A TH	2
---	------	---

The link value in this node gives the location of the next node in the list:

2	ING	7
---	-----	---

Continuing in this manner, we obtain the following sequence of nodes:



Thus the string is:

A THING OF BEAUTY IS A JOY FOREVER.

3.5 Give some (a) advantages and (b) disadvantages of using linked storage for storing strings.

- (a) One can easily insert, delete, concatenate and rearrange substrings when using linked storage.
- (b) Additional space is used for storing the links. Also, one cannot directly access a character in the middle of the list.

3.6 Describe briefly the meaning of (a) static, (b) semistatic and (c) dynamic character variables.

- (a) The length of the variable is defined before the program is executed and cannot change during the execution of the program.
- (b) The length of the variable may vary during the execution of the program, but the length cannot exceed a maximum value defined before the program is executed.
- (c) The length of the variable may vary during the execution of the program.

3.7 Suppose MEMBER is a character variable with fixed length 20. Assume a string is stored left-justified in a memory cell with blank spaces padded on the right or with the right-most characters truncated. Describe MEMBER (a) if 'JOHN PAUL JONES' is assigned to MEMBER and (b) if 'ROBERT ANDREW WASHINGTON' is assigned to MEMBER.

The data will appear in MEMBER as follows:

(a) MEMBER J O H N P A U L J O N E S [] [] [] [] [] [] [] [] [] []

(b) MEMBER R O B E R T A N D R E W W A S H I N [] [] [] [] [] [] [] [] [] []

STRING OPERATIONS

In Probs. 3.8 to 3.11 and 3.13, let S and T be character variables such that

S = 'JOHN PAUL JONES'

T = 'A THING OF BEAUTY IS A JOY FOREVER.'

3.8 Recall that we use LENGTH(string) for the length of a string.

- (a) How is this function denoted in (i) PL/I, (ii) BASIC, (iii) UCSD Pascal, (iv) SNOBOL and (v) FORTRAN?
- (b) Find LENGTH(S) and LENGTH(T).
- (a) (i) LENGTH(string). (ii) LEN(string). (iii) LENGTH(string). (iv) SIZE(string). (v) FORTRAN has no length function for strings, since the language uses only fixed-length variables.

(b) Assuming there is only one blank space character between words,

$$\text{LENGTH}(S) = 15 \quad \text{and} \quad \text{LENGTH}(T) = 35$$

3.9 Recall that we use $\text{SUBSTRING}(\text{string}, \text{position}, \text{length})$ to denote the substring of *string* beginning in a given *position* and having a given *length*. Determine (a) $\text{SUBSTRING}(S, 4, 8)$ and (b) $\text{SUBSTRING}(T, 10, 5)$.

(a) Beginning with the fourth character and recording 8 characters, we obtain

$$\text{SUBSTRING}(S, 4, 8) = \text{'NPAUL'}$$

(b) Similarly,

$$\text{SUBSTRING}(T, 10, 5) = \text{'FBEAU'}$$

3.10 Recall that we use $\text{INDEX}(\text{text}, \text{pattern})$ to denote the position where a pattern first appears in a text. This function is assigned the value 0 if the pattern does not appear in the text. Determine (a) $\text{INDEX}(S, \text{'JO'})$, (b) $\text{INDEX}(S, \text{'JOY'})$, (c) $\text{INDEX}(S, \text{'JO'})$, (d) $\text{INDEX}(T, \text{'A'})$, (e) $\text{INDEX}(T, \text{'A'})$ and (f) $\text{INDEX}(T, \text{'THE'})$.

(a) $\text{INDEX}(S, \text{'JO'}) = 1$, (b) $\text{INDEX}(S, \text{'JOY'}) = 0$, (c) $\text{INDEX}(S, \text{'JO'}) = 10$, (d) $\text{INDEX}(T, \text{'A'}) = 1$, (e) $\text{INDEX}(T, \text{'A'}) = 21$ and (f) $\text{INDEX}(T, \text{'THE'}) = 0$. (Recall that \square is used to denote a blank space.)

3.11 Recall that we use $S_1 // S_2$ to denote the concatenation of strings S_1 and S_2 .

(a) How is this function denoted in (i) PL/1, (ii) FORTRAN, (iii) BASIC, (iv) SNOBOL and (v) UCSD Pascal?

(b) Find (i) $\text{'THE'} // \text{'END'}$ and (ii) $\text{'THE'} // \text{' } // \text{'END'}$.

(c) Find (i) $\text{SUBSTRING}(S, 11, 5) // \text{' } // \text{SUBSTRING}(S, 1, 9)$ and (ii) $\text{SUBSTRING}(T, 28, 3) // \text{'GIVEN'}$.

(a) (i) $S_1 // S_2$, (ii) $S_1 // S_2$, (iii) $S_1 + S_2$, (iv) $S_1 S_2$ (juxtaposition with a blank space between S_1 and S_2) and (v) $\text{CONCAT}(S_1, S_2)$.

(b) $S_1 // S_2$ refers to the string consisting of the characters of S_1 followed by the characters of S_2 . Hence, (i) THEEND and (ii) THE END .

(c) (i) JONES, JOHN PAUL and (ii) FORGIVEN .

3.12 Recall that we use $\text{INSERT}(\text{text}, \text{position}, \text{string})$ to denote inserting a *string* S in a given *text* T beginning in *position* K .

(a) Find (i) $\text{INSERT}(\text{'AAAAA'}, 1, \text{'BBB'})$, (ii) $\text{INSERT}(\text{'AAAAA'}, 3, \text{'BBB'})$ and (iii) $\text{INSERT}(\text{'AAAAA'}, 6, \text{'BBB'})$.

(b) Suppose T is the text $\text{'THE STUDENT IS ILL.}'$ Use INSERT to change T so that it reads: (i) The student is very ill. (ii) The student is ill today. (iii) The student is very ill today.

(a) (i) BBBAAAAA , (ii) AABBBAAA and (iii) AAAAABBB .

(b) Be careful to include blank spaces when necessary. (i) $\text{INSERT}(T, 15, \text{'VERY'})$, (ii) $\text{INSERT}(T, 19, \text{'TODAY'})$, (iii) $\text{INSERT}(\text{INSERT}(T, 19, \text{'TODAY'}), 15, \text{'VERY'})$ or $\text{INSERT}(\text{INSERT}(T, 15, \text{'VERY'}), 24, \text{'TODAY'})$.

3.13 Find

(a) $\text{DELETE}(\text{'AAABBB'}, 2, 2)$ and $\text{DELETE}(\text{'JOHN PAUL JONES'}, 6, 5)$

(b) $\text{REPLACE}(\text{'AAABBB'}, \text{'AA'}, \text{'BB'})$ and $\text{REPLACE}(\text{'JOHN PAUL JONES'}, \text{'PAUL'}, \text{'DAVID'})$

(a) DELETE(T, K, L) deletes from a text T the substring which begins in position K and has length L. Hence the answers are

ABBB and JOHN JONES

(b) REPLACE(T, P₁, P₂) replaces in a text T the first occurrence of the pattern P₁ by the pattern P₂. Hence the answers are

BBABBB and JOHN DAVID JONES

WORD PROCESSING

In Probs. 3.14 to 3.17, S is a short story stored in a linear array LINE with n elements such that each LINE[K] is a static character variable storing 80 characters and representing a line of the story. Also, LINE[1], the first line, contains only the title of the story, and LINE[N], the last line, contains only the name of the author. Furthermore, each paragraph begins with 5 blank spaces, and there is no other indention except possibly the title in LINE[1] or the name of the author in LINE[N].

3.14 Write a procedure which counts the number NUM of paragraphs in the short story S.

Beginning with LINE[2] and ending with LINE[N - 1], count the number of lines beginning with 5 blank spaces. The procedure follows.

Procedure P3.14: PAR(LINE, N, NUM)

1. Set NUM := 0 and BLANK := '□□□□□'.
2. [Initialize counter.] Set K := 2.
3. Repeat Steps 4 and 5 while K ≤ N - 1.
4. [Compare first 5 characters of each line with BLANK.]
If SUBSTRING(LINE[K], 1, 5) = BLANK, then:
Set NUM := NUM + 1.
[End of If structure.]
5. Set K := K + 1. [Increments counter.]
[End of Step 3 loop.]
6. Return.

3.15 Write a procedure which counts the number NUM of times the word "the" appears in the short story S. (We do not count "the" in "mother," and we assume no sentence ends with the word "the.")

Note that the word "the" can appear as THE□ at the beginning of a line, as □THE at the end of a line, or as □THE□ elsewhere in a line. Hence we must check these three cases for each line. The procedure follows.

Procedure P3.15: COUNT(LINE, N, NUM)

1. Set WORD := 'THE' and NUM := 0.
2. [Prepare for the three cases.]
Set BEG := WORD // '□', END := '□' // WORD and
MID := '□' // WORD // '□'.
3. Repeat Steps 4 through 6 for K = 1 to N:
4. [First case.] If SUBSTRING(LINE[K], 1, 4) = BEG, then:
Set NUM := NUM + 1.
5. [Second case.] If SUBSTRING(LINE[K], 77, 4) = END, then:
Set NUM := NUM + 1.
6. [General case.] Repeat for J = 2 to 76.
If SUBSTRING(LINE[K], J, 5) = MID, then:
Set NUM := NUM + 1.
[End of If structure.]
[End of Step 6 loop.]
- [End of Step 3 loop.]
7. Return.

- 3.16** Discuss the changes that must be made in Procedure P3.15 if one wants to count the number of occurrences of an arbitrary word W with length R .

There are three basic types of changes.

- (a) Clearly, 'THE' must be changed to W in Step 1.
 (b) Since the length of W is r and not 3, appropriate changes must be made in Steps 3 to 6.
 (c) One must also consider the possibility that W will be followed by some punctuation, e.g.,

$W, \quad W; \quad W. \quad W?$

Hence more than the three cases must be treated.

- 3.17** Outline an algorithm which will interchange the k th and l th paragraphs in the short story S .

The algorithm reduces to two procedures:

Procedure A. Find the values of arrays BEG and END where

$LINE[BEG[K]]$ and $LINE[END[K]]$

contain, respectively, the first and last lines of paragraph K of the story S .

Procedure B. Using the values of $BEG[K]$ and $END[K]$ and the values of $BEG[L]$ and $END[L]$, interchange the block of lines of paragraph K with the block of lines of paragraph L .

PATTERN MATCHING

- 3.18** For each of the following patterns P and texts T , find the number C of comparisons to find the INDEX of P in T using the "slow" algorithm, Algorithm 3.3:

- (a) $P = abc, T = (ab)^5 = ababababab$ (c) $P = aaa, T = (aabb)^3 = aabbaabbaabb$
 (b) $P = abc, T = (ab)^{2n}$ (d) $P = aaa, T = abaabbaaabbbaaaabbbb$

Recall that $C = N_1 + N_2 + \dots + N_n$ where N_k denotes the number of comparisons that take place in the inner loop when P is compared with W_k .

- (a) Note first that there are

$$LENGTH(T) - LENGTH(P) + 1 = 10 - 3 + 1 = 8$$

substrings W_k . We have

$$C = 2 + 1 + 2 + 1 + 2 + 1 + 2 + 1 = 4(3) = 12$$

and $INDEX(T, P) = 0$, since P does not appear in T .

- (b) There are $2n - 3 + 1 = 2(n - 1)$ subwords W_k . We have

$$C = 2 + 1 + 2 + 1 + \dots + 2 + 1 = (n + 1)(3) = 3n + 3$$

and $INDEX(T, P) = 0$.

- (c) There are $12 - 3 + 1 = 10$ subwords W_k . We have

$$C = 3 + 2 + 1 + 1 + 3 + 2 + 1 + 1 + 3 + 2 = 19$$

and $INDEX(T, P) = 0$.

- (d) We have

$$C = 2 + 1 + 3 + 2 + 1 + 1 + 3 = 13$$

and $INDEX(T, P) = 7$.

3.19 Suppose P is an r -character string and T is an s -character string, and suppose $C(n)$ denotes the number of comparisons when Algorithm 3.3 is applied to P and T . (Here $n = r + s$.)

- (a) Find the complexity $C(n)$ for the best case.
- (b) Prove that the maximum value of $C(n)$ occurs when $r = (n + 1)/4$.
- (a) The best case occurs when P is an initial substring of T , or, in other words, when $\text{INDEX}(T, P) = 1$. In this case $C(n) = r$. (We assume $r \leq s$.)
- (b) By the discussion in Sec. 3.7,

$$C = C(n) = r(n - 2r + 1) = nr - 2r^2 + r$$

Here n is fixed, so $C = C(n)$ may be viewed as a function of r . Calculus tells us that the maximum value of C occurs when $C' = dC/dr = 0$ (here C' is the derivative of C with respect to r). Using calculus, we obtain:

$$C' = n - 4r + 1$$

Setting $C' = 0$ and solving for r gives us the required result.

3.20 Consider the pattern $P = aaabb$. Construct the table and the corresponding labeled directed graph used in the "fast," or second pattern matching, algorithm.

First list the initial segments of P :

$$Q_0 = \Lambda, \quad Q_1 = a, \quad Q_2 = a^2, \quad Q_3 = a^3, \quad Q_4 = a^3b, \quad Q_5 = a^3b^2$$

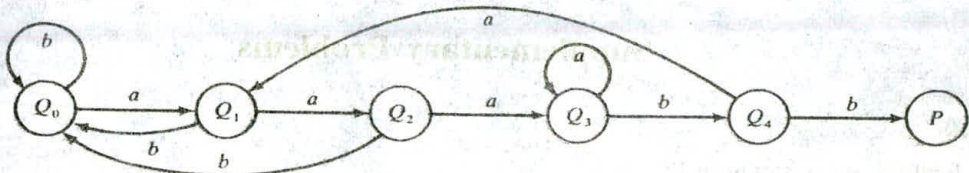
For each character t , the entry $f(Q_i, t)$ in the table is the largest Q which appears as a terminal substring in the string $Q_i t$. We compute:

$$\begin{aligned} f(\Lambda, a) &= a, & f(a, a) &= a^2, & f(a^2, a) &= a^3, & f(a^3, a) &= a^3, & f(a^3b, a) &= a \\ f(\Lambda, b) &= \Lambda, & f(a, b) &= \Lambda, & f(a^2, b) &= \Lambda, & f(a^3, b) &= a^3b, & f(a^3b, b) &= P \end{aligned}$$

Hence the required table appears in Fig. 3-10(a). The corresponding graph appears in Fig. 3-10(b), where there is a node corresponding to each Q and an arrow from Q_i to Q_j labeled by the character t for each entry $f(Q_i, t) = Q_j$ in the table.

	a	b
Q_0	Q_1	Q_0
Q_1	Q_2	Q_0
Q_2	Q_3	Q_0
Q_3	Q_3	Q_4
Q_4	Q_1	P

(a)



(b)

Fig. 3-10

- 3.21 Find the table and corresponding graph for the second pattern matching algorithm where the pattern is $P = ababab$.

The initial substrings of P are:

$$Q_0 = \Lambda, \quad Q_1 = a, \quad Q_2 = ab, \quad Q_3 = aba, \quad Q_4 = abab, \quad Q_5 = ababa, \quad Q_6 = ababab = P$$

The function f giving the entries in the table follows:

$$f(\Lambda, a) = a$$

$$f(a, a) = a$$

$$f(ab, a) = aba$$

$$f(aba, a) = a$$

$$f(abab, a) = ababa$$

$$f(ababa, a) = a$$

$$f(\Lambda, b) = \Lambda$$

$$f(a, b) = ab$$

$$f(ab, b) = \Lambda$$

$$f(aba, b) = abab$$

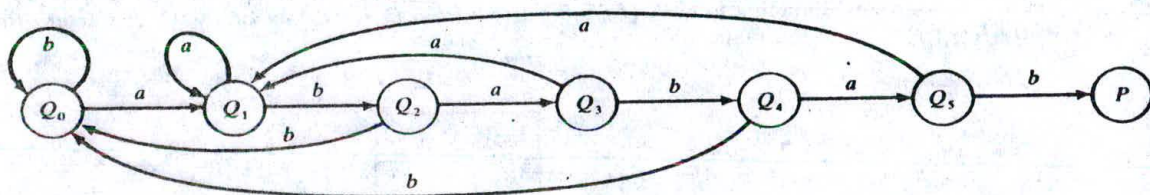
$$f(abab, b) = \Lambda$$

$$f(ababa, b) = P$$

The table appears in Fig. 3-11(a) and the corresponding graph appears in Fig. 3-11(b).

	a	b
Q_0	Q_1	Q_0
Q_1	Q_1	Q_2
Q_2	Q_3	Q_0
Q_3	Q_1	Q_4
Q_4	Q_5	Q_0
Q_5	Q_1	P

(a)



(b)

Fig. 3-11

Supplementary Problems

STRINGS

3.22 Find the string stored in Fig. 3-12.

3.23 Consider the string $W = 'XYZST'$. List (a) all substrings of W and (b) all initial substrings of W .

3.24 Suppose W is a string of length n . Find the number of (a) substrings of W and (b) initial substrings of W .

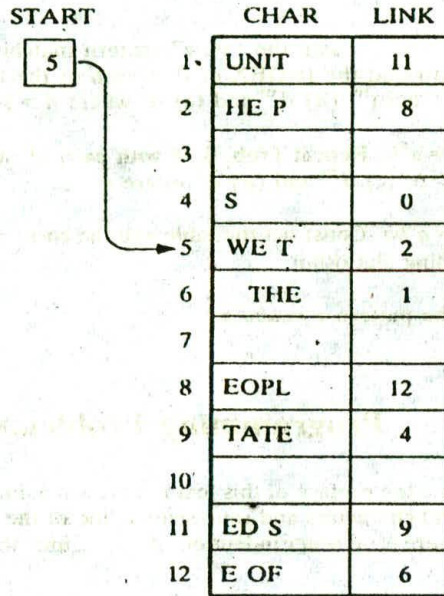


Fig. 3-12

- 3.25 Suppose STATE is a character variable with fixed length 12. Describe the contents of STATE after the assignment (a) STATE := 'NEW YORK', (b) STATE := 'SOUTH CAROLINA' and (c) STATE := 'PENNSYLVANIA'

STRING OPERATIONS

In Probs. 3.26 to 3.31, let S and T be character variables such that

$$S = \text{'WE THE PEOPLE'} \quad \text{and} \quad T = \text{'OF THE UNITED STATES'}$$

- 3.26 Find the length of S and T.
- 3.27 Find (a) SUBSTRING(S, 4, 8) and (b) SUBSTRING(T, 10, 5).
- 3.28 Find (a) INDEX(S, 'P'), (b) INDEX(S, 'E'), (c) INDEX(S, 'THE'), (d) INDEX(T, 'THE'), (e) INDEX(T, 'THEN') and (f) INDEX(T, 'TE').
- 3.29 Using $S_1 // S_2$ to stand for the concatenation of S_1 and S_2 , find (a) 'NO' // 'EXIT', (b) 'NO' // ' ' // 'EXIT' and (c) SUBSTRING(S, 4, 10) // ' ARE ' // SUBSTRING(T, 8, 6).
- 3.30 Find (a) DELETE('AAABBB', 3, 3), (b) DELETE('AAABBB', 1, 4), (c) DELETE(S, 1, 3) and (d) DELETE(T, 1, 7).
- 3.31 Find (a) REPLACE('ABABAB', 'B', 'BAB'), (b) REPLACE(S, 'WE', 'ALL') and (c) REPLACE(T, 'THE', 'THESE').
- 3.32 Find (a) INSERT('AAA', 2, 'BBB'), (b) INSERT('ABCDE', 3, 'XYZ') and (c) INSERT('THE BOY', 5, 'BIG').
- 3.33 Suppose U is the text 'MARC STUDIES MATHEMATICS.' Use INSERT to change U so that it reads: (a) MARC STUDIES ONLY MATHEMATICS. (b) MARC STUDIES MATHEMATICS AND PHYSICS. (c) MARC STUDIES APPLIED MATHEMATICS.

PATTERN MATCHING

- 3.34** Consider the pattern $P = abc$. Using the "slow" pattern matching algorithm, Algorithm 3.3, find the number C comparisons to find the INDEX of P in each of the following texts T :
 (a) a^{10} , (b) $(aba)^{10}$, (c) $(cbab)^{10}$, (d) d^{10} and (e) d^n where $n > 3$.
- 3.35** Consider the pattern $P = a^5b$. Repeat Prob. 3.34 with each of the following texts T :
 (a) a^{20} , (b) a^n where $n > 6$, (c) d^{20} and (d) d^n where $n > 6$.
- 3.36** Consider the pattern $P = a^3ba$. Construct the table and the corresponding labeled directed graph used in the "fast" pattern matching algorithm.
- 3.37** Repeat Prob. 3.36 for the pattern $P = aba^2b$.

Programming Problems

In Probs. 3.38 to 3.40, assume the preface of this text is stored in a linear array $LINE$ such that $LINE[K]$ is a static character variable storing 80 characters and represents a line of the preface. Assume that each paragraph begins with 5 blank spaces and there is no other indentation. Also, assume there is a variable NUM which gives the number of lines in the preface.

- 3.38** Write a program which defines a linear array PAR such that $PAR[K]$ contains the location of the K th paragraph, and which also defines a variable $NPAR$ which contains the number of paragraphs.
- 3.39** Write a program which reads a given $WORD$ and then counts the number C of times $WORD$ occurs in $LINE$. Test the program using (a) $WORD = 'THE'$ and (b) $WORD = 'HENCE'$.
- 3.40** Write a program which interchanges the J th and K th paragraphs. Test the program using $J = 2$ and $K = 4$.

In Probs. 3.41 to 3.46, assume the preface of this text is stored in a single character variable $TEXT$. Assume 5 blank spaces indicates a new paragraph.

- 3.41** Write a program which constructs a linear array PAR such that $PAR[K]$ contains the location of the K th paragraph in $TEXT$, and which finds the value of a variable $NPAR$ which contains the number of paragraphs. (Compare with Prob. 3.38.)
- 3.42** Write a program which reads a given $WORD$ and then counts the number C of times $WORD$ occurs in $TEXT$. Test the program using (a) $WORD = 'THE'$ and (b) $WORD = 'HENCE'$. (Compare with Prob. 3.39.)
- 3.43** Write a program which interchanges the J th and K th paragraphs in $TEXT$. Test the program using $J = 2$ and $K = 4$. (Compare with Prob. 3.40.)
- 3.44** Write a program which reads words $WORD1$ and $WORD2$ and then replaces each occurrence of $WORD1$ in $TEXT$ by $WORD2$. Test the program using $WORD1 = 'HENCE'$ and $WORD2 = 'THUS'$.
- 3.45** Write a subprogram $INST(TEXT, NEW, K)$ which inserts a string NEW into $TEXT$ beginning at $TEXT[K]$.
- 3.46** Write a subprogram $PRINT(TEXT, K)$ which prints the character string $TEXT$ in lines with at most K characters. No word should be divided in the middle and appear on two lines, so some lines may contain trailing blank spaces. Each paragraph should begin with its own line and be indented using 5 blank spaces. Test the program using (a) $K = 80$, (b) $K = 70$ and (c) $K = 60$.