



**M.SC(IT) PART-2 (SEM-3)**

**PAPER: MITM2102T**

**DATA AND FILE STRUCTURE**

UNIT No.1

Center for Distance and Online  
Education,  
PunjabiUniversity, Patiala

Lesson No:

- 1.1 : **INTRODUCTION TO DATA STRUCTURES**
- 1.2 : **ALGORITHM ANALYSIS**
- 1.3 : **Arrays**
- 1.4 : **Stacks**
- 1.5 : **APPLICATIONS OF STACKS**
- 1.6 : **QUEUES**
- 1.7 : **TYPES OF QUEUES**
- 1.8 : **LINKED LIST**
- 1.9 : **TYPES OF LINKED LIST**
- 1.10 : **LINKED REPRESENTATION OF STACKS AND QUEUES**
- 1.11 : **TREES**
- 1.12 : **BINARY SEARCH TREE**
- 1.13 : **HEIGHT BALANCED TREE**
- 1.14 : **HEAPS**

**(Syllabus)**

**MITM2102T: Data and File Structures**

**Maximum Marks: 70**

**Maximum Time: 3 Hrs.**

**Minimum Pass Marks: 35%**

**Course Objective:** This course is designed to explore computing and to show students the art of practical implementation and usage of Algorithms and Data Structures. On completion of this course, the students will be able to

- Be familiar with basic data structure of algorithms.
- Design and analyze programming problem statements
- Choose appropriate data structures and algorithms and use it to design algorithms for a specific problem.
- Handle operations like searching, insertion, deletion and traversing mechanism
- Come up with analysis of efficiency and proofs of correctness

**Course Content**

**SECTION A**

Data Structure: Introduction to data structure and algorithm, Algorithm analysis: Time space trade off algorithms and Big O notation. Arrays: Introduction, one dimensional and multidimensional arrays, memory representation of arrays, operations on arrays, sparse arrays and sparse matrices and their implementation, Advantages and limitation of arrays.

Stacks: Introduction; Operation on stacks; Implementation of stacks, Application of stacks: matching parenthesis, evaluation of arithmetic expressions, conversion from infix to postfix, recursion.

Queues: Introduction, operation on queues, circular queue, memory representation of queues, dequeues, priority queues, application of queues.

Linked List: Introduction; operation on linked list, circular linked list, doubly linked list, header linked list, implementation of linked list, application of linked lists.

Trees: Introduction; Binary Tree; Threaded Binary Trees; Binary Search Tree; Balanced Trees; B-Trees; Heap

## SECTION B

Graphs: Introduction Graph: Graph terminology, Memory Representation of Graphs: adjacency matrix representation of graphs, adjacency list or linked representation of graphs, Operations performed on graphs, Application of graphs

Sorting: Selection Sort, Insertion Sort, Merge Sort, Bucket Sort, Radix Sort, Quick Sort and Heap Sort

Hashing: Hashing techniques; Collision resolution; Deleting items from a hash table; Application of hashing

File Organization: Introduction, External Storage Device: Sequential Access Storage Device (SASD), Direct Access Storage Device (DASD) Sequential File Organization: processing sequential files, operations on sequential files, advantages and disadvantages of sequential file organization Direct File Organization: introduction, processing of direct files, advantages and disadvantages of direct organization Indexed Sequential Organization: introduction, processing of indexed sequential files, advantages and disadvantages of indexed sequential organization

### **Pedagogy:**

The Instructor is expected to use leading pedagogical approaches in the class room situation, research-based methodology, innovative instructional methods, extensive use of technology in the class room, online modules of MOOCS, and comprehensive assessment practices to strengthen teaching efforts and improve student learning outcomes.

The Instructor of class will engage in a combination of academic reading, analyzing case studies, preparing the weekly assigned readings and exercises, encouraging in class discussions, and live project based learning.

**Text and Readings:** Students should focus on material presented in lectures. The text should be used to provide further explanation and examples of concepts and techniques discussed in the course:

- A. Tanenbaum, Y. Lanhgsam and A.J. Augenstein, "Data Structures Using C", PHI.
- Loomis, Marry, "Data Management and File Structures", PHI
- Seymour Lipschultz, "Theory and Practice of Data Structures", McGraw-Hill.
- E. Horowitz and S. Sahni, "Data Structures with Pascal", Galgotia.
- M. J. Folk, B. Zoellick, G Riccardi, "File Structures", Pearson Education.

### **Scheme of Examination**

- English will be the medium of instruction and examination.

- Written Examinations will be conducted at the end of each Semester as per the Academic Calendar notified in advance
- Each course will carry 100 marks of which 30 marks shall be reserved for internal assessment and the remaining 70 marks for written examination to be held at the end of each semester.
- The duration of written examination for each paper shall be three hours.
- The minimum marks for passing the examination for each semester shall be 35% in aggregate as well as a minimum of 35% marks in the semester-end examination in each paper.
- A minimum of 75% of classroom attendance is required in each subject.

#### **Instructions to the External Paper Setter**

The question paper will consist of three Sections: A, B and C. Sections A and B will have four questions each from the respective section of the syllabus and will carry 10.5 marks for each question. Section C will consist of 7-15 short answer type questions covering the entire syllabus uniformly and will carry a total of 28 marks.

#### **Instructions for candidates**

- Candidates are required to attempt five questions in all, selecting two questions each from section A and B and compulsory question of section C.
- Use of non-programmable scientific calculator is allowed.

## INTRODUCTION TO DATA STRUCTURES

- 1.1 Objectives**
- 1.2 Introduction**
- 1.3 Definition of Data Structures & Algorithms**
- 1.4 Some other Definitions**
  - 1.4.1 Algorithm
  - 1.4.2 Data Type
  - 1.4.3 Four Fundamental data structures
  - 1.4.4 Recursion
- 1.5 Types of Data Structures**
- 1.6 Characteristics of Data Structures**
- 1.7 Operations on Data Structures**
- 1.8 Abstract Data Types**
- 1.9 Summary**
- 1.10 Questions**
- 1.11 Suggested readings**

### **1.1 Objectives**

In this lesson we will discuss the meaning, types, characteristics and operations on data structures.

### **1.2 Introduction**

Software engineering is the study of ways in which to create large and complex computer applications and that generally involve many programmers and designers. At the

heart of software engineering is with the overall design of the applications and on the creation of a design that is based on the needs and requirements of end users. While software engineering involves the full life cycle of a software project, it includes many different components - specification, requirements gathering, design, verification, coding, testing, quality assurance, user acceptance testing, production and ongoing maintenance.

Having an in-depth understanding on every component of software engineering is not mandatory; however, it is important to understand that the subject of data structures and algorithms is concerned with the coding phase. The use of data structures and algorithms is the nuts-and-bolts used by programmers to store and manipulate data.

When thinking about a particular application or programming problem, many developers find themselves most interested in writing the algorithm to tackle the problem at hand or adding cool features to the application to enhance the user's experience. Rarely, if ever, will you hear someone excited about what type of data structure they are using. However, the data structures used for a particular algorithm can greatly impact its performance. A common example is finding an element in a data structure. With an array, this process takes time proportional to the number of elements in the array. With binary search trees or SkipLists, the time required is sub-linear. When searching large amounts of data, the data structure chosen can make a difference in the application's performance that can be visibly measured in seconds or even minutes.

Since the data structure used by an algorithm can greatly affect the algorithm's performance, it is important that there exists a rigorous method by which to compare the efficiency of various data structures. What we, as developers utilizing a data structure, are primarily interested in is how the data structure's performance changes as the amount of data stored increases. That is, for each new element stored by the data structure, how are the running times of the data structure's operations affected?

### 1.3 Definition of Data Structures & Algorithms

**A data structure is an arrangement of data in a computer's memory or even disk storage.** An example of several common data structures are arrays, linked lists, queues, stacks, binary trees and hash tables. Algorithms, on the other hand are used to manipulate the data contained in these data structures as in searching and sorting.

The choice of the data structure often begins from the choice of an abstract data type. A well-designed data structure allows a variety of critical operations to be performed, using as few resources, both execution time and memory space, as possible. Data structures are implemented by a programming language as data types and the references and operations they provide.

Different kinds of data structures are suited to different kinds of applications, and some are highly specialized to certain tasks. For example, B-trees are particularly well-suited for implementation of databases, while networks of machines rely on routing tables to function.

In the design of many types of computer program, the choice of data structures is a primary design consideration. Experience in building large systems has shown that the difficulty of implementation and the quality and performance of the final result depends heavily on choosing the best data structure. After the data structures are chosen, the algorithms to be used often become relatively obvious. Sometimes things work in the opposite direction — data structures are chosen because certain key tasks have algorithms that work best with particular data structures. In either case, the choice of appropriate data structures is crucial.

Many algorithms apply directly to a specific data structures. When working with certain data structures you need to know how to insert new data, search for a specified item and deleting a specific item.

## 1.4 Some other Definitions

**1.4.1 Algorithm:** A finite sequence of instructions, each of which has a clear meaning and can be executed with a finite amount of effort in finite time. Whatever the input values, an algorithm will definitely terminate after executing a finite number of instructions.

**1.4.2 Data Type:** Data type of a variable is the set of values that the variable may assume.

Basic data types in C:

- int, char, float and double

**1.4.3 Four Fundamental Data Structures:** The following four data structures are used ubiquitously in the description of algorithms and serve as basic building blocks for realizing more complex data structures.

- Sequences (also called as lists)
- Dictionaries
- Priority Queues
- Graphs

Dictionaries and priority queues can be classified under a broader category called *dynamic sets*. Also, binary and general trees are very popular building blocks for implementing dictionaries and priority queues.

**1.4.4 Recursion:** It is a method call in which the method being called is the same as the one making the call. The recursive algorithm is implemented by using a method that makes recursive calls to itself. It can be:

- **Direct recursion:** Recursion in which a method directly calls itself
- **Indirect recursion:** Recursion in which a chain of two or more method calls returns to the method that originated the chain

## 1.5 Types of Data Structures

Data structures can be broadly divided into three categories

**1.5.1 Base data structures:** Basic data structures are the building blocks for composite and complex data structures.

<b>General type</b>	<b>Specific types</b>
Primitive types	<ul style="list-style-type: none"> <li>• Character</li> <li>• Integer</li> <li>• String</li> <li>• Double</li> <li>• Float</li> </ul>
Composite types	<ul style="list-style-type: none"> <li>• Struct</li> <li>• Bit field</li> <li>• Union</li> </ul>

**1.5.2 Linear Data Structures:** Where data is stored in sequential or linear order. The most commonly used linear data structures are given below:

General type	<ul style="list-style-type: none"> <li>• Specific types</li> <li>• Array               <ul style="list-style-type: none"> <li>○ Bitmaps</li> <li>○ Dynamic array</li> <li>○ Sparse array</li> <li>○ Matrix                   <ul style="list-style-type: none"> <li>▪ Sparse matrix</li> </ul> </li> </ul> </li> </ul>
List (or vector or sequence)	<ul style="list-style-type: none"> <li>• Linked list               <ul style="list-style-type: none"> <li>○ Doubly linked list</li> </ul> </li> <li>• Buffer               <ul style="list-style-type: none"> <li>○ Stack</li> <li>○ Queue                   <ul style="list-style-type: none"> <li>▪ Priority queue                       <ul style="list-style-type: none"> <li>▪ Double-ended priority queue</li> </ul> </li> </ul> </li> <li>○ Deque</li> <li>○ Circular buffer</li> </ul> </li> </ul>
Associative array (a.k.a. dictionary or map)	<ul style="list-style-type: none"> <li>• Hash table</li> <li>• Self-balancing binary search tree</li> <li>• Skip list</li> </ul>

**1.5.3 Non-linear Data Structures:** Where data is not stored in sequential or linear order like trees and graphs etc.



## General type

## Specific types

### Graph data structures

- Adjacency list
- Adjacency matrix
- Disjoint-set data structure
- Graph-structured stack
- Scene graph
- M-Way Tree
- B-tree
  - B+ tree
  - Generalized search tree
- B\* tree
- K-ary tree

### Tree data structures

- Binary tree
  - Binary heap
  - Binary search trees (each tree node compares entire key values)
    - Self-balancing binary search trees
      - AVL tree
      - Red-black tree
      - Scapegoat tree
      - Splay tree
      - Interval tree
      - Treap
  - Exponential tree
- Heap
  - Binary heap
  - Binomial heap
  - Fibonacci heap
- Syntax tree
  - Abstract syntax tree
  - Parse tree
- Decision theory
  - Binary decision diagram
  - Decision tree
    - Alternating decision tree
    - Minimax tree
    - Expectiminimax tree

## 1.6 Characteristics of Data Structures

<b>Data Structure</b>	<b>Advantages</b>	<b>Disadvantages</b>
<b>Array</b>	Quick inserts Fast access if index known	Slow search Slow deletes Fixed size
<b>Ordered Array</b>	Faster search than unsorted array	Slow inserts Slow deletes Fixed size
<b>Stack</b>	Last-in, first-out access	Slow access to other items
<b>Queue</b>	First-in, first-out access	Slow access to other items
<b>Linked List</b>	Quick inserts Quick deletes	Slow search
<b>Binary Tree</b>	Quick search Quick inserts Quick deletes (If the tree remains balanced)	Deletion algorithm is complex
<b>Red-Black Tree</b>	Quick search Quick inserts Quick deletes (Tree always remains balanced)	Complex to implement
<b>2-3-4 Tree</b>	Quick search Quick inserts Quick deletes	Complex to implement

	(Tree always remains balanced) (Similar trees good for disk storage)	
<b>Hash Table</b>	Very fast access if key is known Quick inserts	Slow deletes  Access slow if key is not known Inefficient memory usage
<b>Heap</b>	Quick inserts  Quick deletes  Access to largest item	Slow access to other items
<b>Graph</b>	Best models real-world situations	Some algorithms are slow and very complex
NOTE: The data structures shown above (with the exception of the array) can be thought of as Abstract Data Types (ADTs).		

Commonly used algorithms include are useful for:

Searching for a particular data item (or record).

Sorting the data. There are many ways to sort data. Simple sorting, Advanced sorting

Iterating through all the items in a data structure. (Visiting each item in turn so as to display it or perform some other action on these items)

### 1.7 Operation on Data Structures

The data present in data structure is processed by means of certain operations. The choice of a particular data structures depends specifically on the frequency of specific operations. The following are some common operations performed on the data structures:

- **Insertion:** Adding new elements to the data structure
- **Deletion:** Removing an element from the data structure
- **Traversal:** Accessing each record only once so that certain items in the record may be processed. (This is also termed as visit the record)
- **Searching:** Finding the location of a record with a given key or finding the locations of all records satisfying one or more given conditions.

Some times two or more operations are used in a given situation, e.g. for deleting a record with a specific key requires first finding the location of the record and then deleting it.

The following operations are used in some specific situations only.

- **Sorting:** Arranging the elements in some logical order, e.g. alphabetically according to some name or in numerical order according to some number key or chronologically according to dates of events etc.)
- **Merging:** Combining two sets of records into one.

## 1.8 Abstract Data Types

**An Abstract Data Type (ADT)** is more a way of looking at a data structure: focusing on what it does and ignoring how it does its job. A stack or a queue is an example of an ADT. It is important to understand that both stacks and queues can be implemented using an array. It is also possible to implement stacks and queues using a linked list. This demonstrates the "abstract" nature of stacks and queues: how they can be considered separately from their implementation.

To best describe the term Abstract Data Type, it is best to break the term down into "data type" and then "abstract".

### 1.8.1 data type

When we consider a primitive type we are actually referring to two things: a data item with certain characteristics and the permissible operations on that data. An int in Java, for example, can contain any whole-number value from -2,147,483,648 to +2,147,483,647. It can also be used with the operators +, -, \*, and /. The data type's permissible operations are an inseparable part of its identity; understanding the type means understanding what operations can be performed on it.

In C++, any class represents a data type, in the sense that a class is made up of data (fields) and permissible operations on that data (methods). By extension, when a data storage structure like a stack or queue is represented by a class, it too can be referred to as a data type. A stack is different in many ways from an int but they are both defined as a certain arrangement of data and a set of operations on that data.

### 1.8.2 abstract

Now lets look at the "abstract" portion of the phrase. The word abstract in our context stands for "considered apart from the detailed specifications or implementation". In C++, an Abstract Data Type is a class considered without regard to its implementation. It can be thought of as a "description" of the data in the class and a list of operations that can be carried out on that data and instructions on how to use these operations. What is excluded though is the details of how the methods carry out their tasks. An end user (or class user), you should be told what methods to call, how to call them and the results that should be expected, but not HOW they work.

We can further extend the meaning of the ADT when applying it to data structures such as a stack and queue. In C++, as with any class, it means the data and the operations that can be performed on it. In this context, although, even the fundamentals of how the data is stored should be invisible to the user. Users not only should not know

how the methods work, they should also not know what structures are being used to store the data.

Consider for example the stack class. The end user knows that push() and pop() (among other similar methods) exist and how they work. The user doesn't and shouldn't have to know how push() and pop() work, or whether data is stored in an array, a linked list, or some other data structure like a tree.

### **1.8.3 The Interface**

The ADT specification is often called an interface. It's what the user of the class actually sees. In C++, this would often be the public methods. Consider for example, the stack class - the public methods push() and pop() and similar methods from the interface would be published to the end user.

## **1.9 Summary**

A data structure is an arrangement of data in a computer's memory or even disk storage. Algorithms, on the other hand are used to manipulate the data contained in these data structures as in searching and sorting. The data structure used by an algorithm can greatly affect the algorithm's performance, therefore, it is important that there exists a rigorous method by which to compare the efficiency of various data structures.

Data structures can be basic (character, integer, float etc.), linear (array, stacks, queues, linked lists etc.) or non-linear (trees, graphs etc). The operations that are performed on data structures include insertion, deletion, searching, traversing, sorting and merging.

### **1.10 Questions**

1. Define data structures.
2. What are the characteristics of data structures?
3. What are the operations performed on data structures?
4. Define Abstract Data Structure.

### **1.11 Suggested readings**

1. A. Tanenbaum, Y. Lanhgsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley

## ALGORITHM ANALYSIS

- 2.1 Objectives
- 2.2 Introduction
- 2.3 Asymptotic notation
- 2.4 Big 'O' notation
- 2.5 Algorithm analysis
- 2.6 Time Space Tradeoffs
- 2.7 Algorithm efficiency
- 2.8 Summary
- 2.9 Questions
- 2.10 Suggested readings

### 2.1 Objectives

In this lesson we will discuss the various operations applied on the data structures, time and space tradeoffs using different types of data structures and the representation of algorithm complexity in terms of big 'O' notation.

### 2.2 Introduction

In computer science, often the question is not how to solve a problem, but how to solve a problem well. For instance, take the problem of sorting. Many sorting algorithms are well-known; the problem is not to find a way to sort words but to find a way to efficiently sort words. This article is about understanding how to compare the relative efficiency of algorithms and why it's important to do so.

If it's possible to solve a problem by using a brute force technique, such as trying out all possible combinations of solutions (for instance, sorting a group of words by trying all possible orderings until you find one that is in order), then why is it necessary to find a better approach? The simplest answer is, if you had a fast enough computer, maybe it

wouldn't be. But as it stands, we do not have access to computers fast enough. For instance, if you were to try out all possible orderings of 100 words, that would require  $100!$  (100 factorial) orders of words. (Explanation) That's a number with a 158 digits; were you to compute 1,000,000,000 possibilities were second, you would still be left with the need for over  $1 \times 10^{149}$  seconds, which is longer than the expected life of the universe. Clearly, having a more efficient algorithm to sort words would be handy and of course there are many that can sort 100 words within the life of the universe.

Before going further, it's important to understand some of the terminology used for measuring algorithmic efficiency. Usually, the efficiency of an algorithm is expressed as how long it runs in relation to its input. For instance, in the above example, we showed how long it would take our naive sorting algorithm to sort a certain number of words. Usually we refer to the length of input as  $n$ ; so, for the above example, the efficiency is roughly  $n!$ . You might have noticed that it's possible to come up with the correct order early on in the attempt -- for instance, if the words are already partially ordered, it's unlikely that the algorithm would have to try all  $n!$  combinations. Often we refer to the average efficiency, would be in this case  $n!/2$ . But because the division by two is nearly insignificant as  $n$  grows larger (half of 2 billion is, for instance, still a very large number), we usually ignore constant terms (unless the constant term is zero).

Now that we can describe any algorithm's efficiency in terms of its input length (assuming we know how to compute the efficiency), we can compare algorithms based on their "order". Here, "order" refers to the mathematical method used to compare the efficiency -- for instance,  $n^2$  is "order of  $n$  squared," and  $n!$  is "order of  $n$  factorial." It should be obvious that an order of  $n^2$  algorithm is much less efficient than an algorithm of order  $n$ . But not all orders are polynomial -- we've already seen  $n!$ , and some are order of  $\log n$ , or order  $2^n$ .

Often order is abbreviated with a capital O: for instance,  $O(n^2)$ . This notation, known as big-O notation (explain later in this lesson) is a typical way of describing algorithmic efficiency; note that big-O notation typically does not call for inclusion of constants. Also, if you are determining the order of an algorithm and the order turns out to be the sum of several terms, you will typically express the efficiency as only the term with the highest order. For instance, if you have an algorithm with efficiency  $n^2 + n$ , then it is an algorithm of order  $O(n^2)$ .

### **2.3 Asymptotic Notation**

A problem may have numerous algorithmic solutions. In order to choose the best algorithm for a particular task, you need to be able to judge how long a particular solution will take to run. Or, more accurately, you need to be able to judge how long two solutions will take to run and choose the better of the two. You don't need to know how many minutes and seconds they will take but you do need some way to compare algorithms against one another.

**Asymptotic complexity is a way of expressing the main component of the cost of an algorithm, using idealized units of computational work.** Consider, for example, the algorithm for sorting a deck of cards, which proceeds by repeatedly searching through the deck for the lowest card. The asymptotic complexity of this algorithm is the square of the number of cards in the deck. This quadratic behavior is the main term in the complexity formula, it says, e.g., if you double the size of the deck, then the work is roughly quadrupled.

The exact formula for the cost is more complex and contains more details than are needed to understand the essential complexity of the algorithm. With our deck of cards, in the worst case, the deck would start out reverse-sorted, so our scans would have to go all the way to the end. The first scan would involve scanning 52 cards, the next would take 51, etc. So the cost formula is  $52 + 51 + \dots + 1$ . generally, letting  $N$  be the number of cards, the formula is  $1 + 2 + \dots + N$ , which equals  $(N+1) * (N/2) = (N^2 + N)/2 = (1/2)N^2 + N/2$ . But the  $N^2$  term dominates the expression and this is what is key for comparing algorithm costs. (This is in fact an expensive algorithm; the best sorting algorithms run in sub-quadratic time.)

Asymptotically speaking, in the limit as  $N$  tends towards infinity,  $1 + 2 + \dots + N$  gets closer and closer to the pure quadratic function  $(1/2) N^2$ . And what difference does the constant factor of  $1/2$  make at this level of abstraction. So the behavior is said to be  $O(n^2)$ .

Now let us consider how we would go about comparing the complexity of two algorithms. Let  $f(n)$  be the cost, in the worst case, of one algorithm, expressed as a function of the input size  $n$  and  $g(n)$  be the cost function for the other algorithm. E.g., for sorting algorithms,  $f(10)$  and  $g(10)$  would be the maximum number of steps that the algorithms would take on a list of 10 items. If, for all values of  $n \geq 0$ ,  $f(n)$  is less than or equal to  $g(n)$ , then the algorithm with complexity function  $f$  is strictly faster. But, generally speaking, our concern for computational cost is for the cases with large inputs; so the comparison of  $f(n)$  and  $g(n)$  for small values of  $n$  is less significant than the "long term" comparison of  $f(n)$  and  $g(n)$ , for  $n$  larger than some threshold.

Note that we have been speaking about bounds on the performance of algorithms, rather than giving exact speeds. The actual number of steps required to sort our deck of cards (with our naive quadratic algorithm) will depend upon the order in which the cards begin. The actual time to perform each of our steps will depend upon our processor speed, the condition of our processor cache, etc., etc. It's all very complicated in the concrete details and moreover not relevant to the essence of the algorithm.

## 2.4 Big-O Notation

**Big-O notation is used to express an ordering property among functions.** In the context of our study of algorithms, the functions are the amount of resources consumed by an algorithm when the algorithm is executed. This function is usually denoted  $T(N)$ .



$T(N)$  is the amount of the resource (usually time or the count of some specific operation) consumed when the input to the algorithm is of size  $N$ .

Sometimes it is not obvious what the "size" of the input is, so here are few examples:

- **Multiplication:** If we multiply two numbers together using the algorithm we learned in grade school, the number of single-digit multiplications and additions that we do depends on the number of digits in the two numbers being multiplied. (If we multiply two numbers together using a calculator, then the number of keystrokes we need to type also depends on the number of digits in the two numbers.) Therefore, for an algorithm that performs multiplication (or any other arithmetic operation) the "size" of the input is usually the number of digits in the input numbers.
- **Searching:** If we have a set of  $N$  unknown items, arranged in an unknown order and we want to do something involving all of them, then the size of the problem is  $N$ . For example, if we want to find the smallest element in an array of  $N$  numbers, then the size of the problem is  $N$ .

Big-O is the formal method of expressing the upper bound of an algorithm's running time. It's a measure of the longest amount of time it could possibly take for the algorithm to complete.

More formally, for non-negative functions,  $f(n)$  and  $g(n)$ , if there exists an integer  $n_0$  and a constant  $c > 0$  such that for all integers  $n > n_0$ ,  $f(n) \leq cg(n)$ , then  $f(n)$  is Big O of  $g(n)$ . This is denoted as " $f(n) = O(g(n))$ ". If graphed,  $g(n)$  serves as an upper bound to the curve you are analyzing,  $f(n)$ .

### Theory Examples

So, let's take an example of Big-O. Say that  $f(n) = 2n + 8$ , and  $g(n) = n^2$ . Can we find a constant  $c$ , so that  $2n + 8 \leq cn^2$ ? The number 4 works here, giving us  $16 \leq 16$ . For any number  $c$  greater than 4, this will still work. Since we're trying to generalize this for large values of  $n$  and small values (1, 2, 3) aren't that important, we can say that  $f(n)$  is generally faster than  $g(n)$ , i.e.,  $f(n)$  is bound by  $g(n)$  and will always be less than it. It could then be said that  $f(n)$  runs in  $O(n^2)$  time: "f-of-n runs in Big-O of n-squared time".

To find the upper bound - the Big-O time - assuming we know that  $f(n)$  is equal to (exactly)  $2n + 8$ , we can take a few shortcuts. For example, we can remove all constants from the runtime; eventually, at some value of  $c$ , they become irrelevant. This makes  $f(n) = 2n$ . Also, for convenience of comparison, we remove constant multipliers; in this case, the 2. This makes  $f(n) = n$ . It could also be said that  $f(n)$  runs in  $O(n)$  time; that lets us put a tighter (closer) upper bound onto the estimate.

### Practical Examples

$O(n)$ : printing a list of  $n$  items to the screen, looking at each item once.  $O(\ln n)$ : also "log  $n$ ", taking a list of items, cutting it in half repeatedly until there's only one item left.  $O(n^2)$ : taking a list of  $n$  items and comparing every item to every other item.

### **Big-Omega Notation**

For non-negative functions,  $f(n)$  and  $g(n)$ , if there exists an integer  $n_0$  and a constant  $c > 0$  such that for all integers  $n > n_0$ ,  $f(n) \geq cg(n)$ , then  $f(n)$  is omega of  $g(n)$ . This is denoted as " $f(n) = \Omega(g(n))$ ".

This is almost the same definition as Big Oh, except that " $f(n) \geq cg(n)$ ", this makes  $g(n)$  a lower bound function, instead of an upper bound function. It describes the best that can happen for a given data size.

### **Theta Notation**

For non-negative functions,  $f(n)$  and  $g(n)$ ,  $f(n)$  is theta of  $g(n)$  if and only if  $f(n) = O(g(n))$  and  $f(n) = \Omega(g(n))$ . This is denoted as " $f(n) = \Theta(g(n))$ ".

This is basically saying that the function,  $f(n)$  is bounded both from the top and bottom by the same function,  $g(n)$ .

### **Little-O Notation**

For non-negative functions,  $f(n)$  and  $g(n)$ ,  $f(n)$  is little o of  $g(n)$  if and only if  $f(n) = O(g(n))$ , but  $f(n) \neq \Theta(g(n))$ . This is denoted as " $f(n) = o(g(n))$ ". This represents a loose bounding version of Big O.  $g(n)$  bounds from the top but it does not bound the bottom.

### **Little Omega Notation**

For non-negative functions,  $f(n)$  and  $g(n)$ ,  $f(n)$  is little omega of  $g(n)$  if and only if  $f(n) = \Omega(g(n))$ , but  $f(n) \neq \Theta(g(n))$ . This is denoted as " $f(n) = \omega(g(n))$ ". Much like Little Oh, this is the equivalent for Big Omega.  $g(n)$  is a loose lower boundary of the function  $f(n)$ ; it bounds from the bottom, but not from the top.

### **How asymptotic notation relates to analyzing complexity**

Temporal comparison is not the only issue in algorithms. There are space issues as well. Generally, a trade off between time and space is noticed in algorithms. Asymptotic notation empowers you to make that trade off. If you think of the amount of time and space your algorithm uses as a function of your data over time or space (time and space are usually analyzed separately), you can analyze how the time and space is handled when you introduce more data to your program.

This is important in data structures because you want a structure that behaves efficiently as you increase the amount of data it handles. Keep in mind though that algorithms that are efficient with large amounts of data are not always simple and efficient for small amounts of data. So if you know you are working with only a small amount of

data and you have concerns for speed and code space, a trade off can be made for a function that does not behave well for large amounts of data.

### **A few examples of asymptotic notation**

Generally, we use asymptotic notation as a convenient way to examine what can happen in a function in the worst case or in the best case. For example, if you want to write a function that searches through an array of numbers and returns the smallest one:

```
function find-min(array a[1..n])  
  
  let j :=  $\infty$   
  
  for i := 1 to n:  
  
    j := min(j, a[i])  
  
  repeat  
  
  return j  
  
end
```

Regardless of how big or small the array is, every time we run find-min, we have to initialize the i and j integer variables and return j at the end. Therefore, we can just think of those parts of the function as constant and ignore them.

So, how can we use asymptotic notation to discuss the find-min function? If we search through an array with 87 elements, then the for loop iterates 87 times, even if the very first element we hit turns out to be the minimum. Likewise, for n elements, the for loop iterates n times. Therefore, we say the function runs in time  $O(n)$ .

What about this function:

```
function find-min-plus-max(array a[1..n])  
  
  // First, find the smallest element in the array  
  
  let j :=  $\infty$ ;  
  
  for i := 1 to n:  
  
    j := min(j, a[i])  
  
  repeat  
  
  let minim := j
```

```

// Now, find the biggest element, add it to the smallest and
j :=  $-\infty$  ;
for i := 1 to n:
    j := max(j, a[i])
repeat
let maxim := j
// return the sum of the two
return minim + maxim;
end

```

What's the running time for find-min-plus-max? There are two for loops, that each iterate  $n$  times, so the running time is clearly  $O(2n)$ . Because 2 is a constant, we throw it away and write the running time as  $O(n)$ . Why can you do this? If you recall the definition of Big-O notation, the function whose bound you're testing can be multiplied by some constant. If  $f(x) = 2x$ , we can see that if  $g(x) = x$ , then the Big-O condition holds. Thus  $O(2n) = O(n)$ . This rule is general for the various asymptotic notations.

### Finding the Big-O of a Function

Although the definitions given above completely define what the big-O of a function is, it is not always immediately obvious how to use them to actually discover the big-O of a function. Finding the order of many useful functions can be a challenging task, and there are even functions that have defied all attempts at a complete analysis! Luckily, since there has been so much work done in the area already, finding the big-O of many functions is simply a matter of finding a known result for a similar function and using it.

The following rules provide a jumping off point:

1. If  $T_1(N) = O(f_1(N))$  and  $T_2(N) = O(f_2(N))$  then:
2.  $T_1(N) + T_2(N) = \max(O(f_1(N)), O(f_2(N)))$
3.  $T_1(N) \times T_2(N) = O(f_1(N) \times f_2(N))$
4. If  $T(N)$  is a polynomial of degree  $k$ , then  $T(N) = \Theta(N^k)$

The fact that  $T(N) = O(N^k)$  follows from the previous rule. Showing that  $T(N) = \Theta(N^k)$  is slightly more complicated, but can be proven by showing that there is always a choice of  $c$  and  $n_0$  that satisfies the definition of  $\Theta$ .

5.  $\log^k N = O(N)$ , for any constant  $k$ .

This is true, but of relatively limited use. Since  $O(N)$  is an upper bound, it is OK to say that anything with a lower order is  $O(N)$ . This is handy as a last resort, when trying to find the big-O of a nasty function that includes logarithms, but whenever possible it is useful to draw a distinction between  $O(N)$  and  $O(\log^k N)$ . For example, an algorithm that has a big-O of  $O(N \log N)$  is (usually) much better than an algorithm that is  $O(N^2)$ .

Another method of finding the big-O of a function is to find the dominant term of the function and find its order. The order of the dominant term will also be the order of the function.

The dominant term is the term that grows most quickly as  $n$  becomes large. Some rules of dominance include the following:

- Any exponential function of  $n$  dominates any polynomial function of  $n$ .
- Any polynomial function of  $n$  dominates any logarithmic function of  $n$ .
- Any logarithmic function of  $n$  dominates a constant term.
- A polynomial of degree  $k$  dominates a polynomial of degree  $l$  if and only if  $k > l$ .

There are additional rules that you can discover for yourself. The key is that term  $x(n)$  dominates function  $y(n)$  if and only if  $x(n)/y(n)$  grows as  $n$  grows large.

Using these rules can sometimes make finding the order of a complicated function easy. For example, the function  $f(n) = n^4 + 2n + \log n$  is clumsy to manipulate algebraically. However, thanks to the rules of dominance, we know that the  $2n$  term will dominate the order of this function, so we can simply find the order of  $2n$ , which is itself. Therefore, we can immediately say that  $f(n) = O(2n)$ .

There are other methods that can be used find the order of functions. but we will not introduce them yet. Later in the semester, when we analyze several recursive algorithms, methods for analyzing the order of recursive functions will be introduced.

## Properties of Orders

Since the orders of functions look like ordinary functions, it is tempting to treat them as such and manipulate them as though they were ordinary functions. This is usually a huge mistake.

For example, let  $f(n) = n^2$  and  $g(n) = 2n^2$ . It should be clear that the both  $f(n)$  and  $g(n)$  are  $O(n^2)$ . What is the big-O of  $g(n) - f(n)$ ? If we compute the order of the difference as the difference of the orders, then we would mistakenly conclude that  $(g(n) - f(n)) =$  the order of  $f(n)$  less the order of  $g(n)$ , or  $(n^2 - n^2) = 0$ . This answer is obviously wrong. If we actually compute  $g(n) - f(n)$ , however, we arrive at the correct conclusion:

$$(g(n) - f(n)) = (2n^2 - n^2) = n^2 = O(n^2).$$

## 2.5 Algorithm analysis

Now that we have seen the basics of big-O notation, it is time to relate this to the analysis of algorithms.

In our study of algorithms, nearly every function whose order we are interested in finding is a function that defines the quantity of some resource consumed by a particular algorithm in relationship to the parameters of the algorithm. (This function is often referred to as a **complexity** of the algorithm or less frequently as the *cost function* of the algorithm.) This function is usually not the same as the algorithm itself but is a property of the algorithm. For example, when we are analyzing an algorithm that multiplies two numbers, the functions we might be interested in are the relationships between the number of digits in each number and the length of time or amount of memory required by the algorithm.

Although big-O notation is a way of describing the order of a function. It is also often meant to represent the time complexity of an algorithm. This is sloppy use of the mathematics, but unfortunately not uncommon. (For example, at the bottom of page 22 in Weiss, a factorial function is described as being  $O(N)$ . Clearly, the factorial function itself is not  $O(N)$ , it is  $O(N!)$  according to its definition. What is meant is that the particular algorithm given for computing the factorial of  $N$  requires  $O(N)$  time to execute.)

Note that this notation, like the orders themselves, doesn't tell us how quickly or slowly the algorithms actually execute for a given input. This information can be extremely important in practice- during the semester. We will study several algorithms that address the same problems and have the same order running time, but take substantially different amounts of time to execute.

Similarly, the order doesn't tell us how fast an algorithm will run for a small  $N$ . This may also be quite important in practice- during the semester, we will study at least one group of algorithms where the choice of the best algorithm depends on  $N$ - when  $N$  is small, one algorithm is best, but when  $N$  is large, a different algorithm is much better.

## 2.6 Time Space Tradeoffs

**In computer science, a space-time or time-memory tradeoff is a situation where the memory use can be reduced at the cost of slower program execution or vice versa**, the computation time can be reduced at the cost of increased memory use. As the relative costs of CPU cycles, RAM space and hard drive space change - hard drive space has for some time been getting cheaper at a much faster rate than other components of computers - the appropriate choices for space-time tradeoffs have changed radically. Often, by exploiting a space-time tradeoff, a program can be made to run much faster.

The most common situation is an algorithm involving a lookup table. An implementation can include the entire table, which reduces computing time but increases the amount of memory needed or it can compute table entries as needed, increasing computing time but reducing memory requirements.

A space-time tradeoff can be applied to the problem of data storage. If data is stored uncompressed, it takes more space but less time than if the data were stored compressed (since compressing the data reduces the amount of space it takes but it takes time to run the compression algorithm). Depending on the particular instance of the problem, either way is practical. Another example is displaying mathematical formulae on primarily text-based websites. Storing only the LaTeX source and rendering it as an image every time the page is requested would be trading time for space - more time used but less space. Rendering the image when the page is changed and storing the rendered images would be trading space for time - more space used but less time. Note that there are also rare instances where it is possible to directly work with compressed data, such as in the case of compressed bitmap indices, where it is faster to work with compression than without compression.

Larger code size can be traded for higher program speed when applying loop unwinding. This technique makes the code longer for each iteration of a loop but saves the computation time required for jumping back to the beginning of the loop at the end of each iteration.

In the field of cryptography, where the adversary is trying to do better than the exponential time required for a brute force attack, the advantages of a space-time tradeoff can readily be seen. Rainbow tables use partially precomputed values in the hash space of a cryptographic hash function to crack passwords in minutes instead of weeks. Decreasing the size of the rainbow table increases the time required to iterate over the hash space. The meet-in-the-middle attack uses a space-time tradeoff to find the cryptographic key in only  $2n + 1$  encryptions (and  $O(2n)$  space) versus the expected  $2^{2n}$  encryptions (but only  $O(1)$  space) of the naive attack.

Dynamic programming is another example where the time complexity of a problem can be reduced significantly by using more memory.

## 2.7 Algorithm efficiency

The ability to analyze a piece of code or an algorithm and understand its efficiency is vital for understanding computer science.

One approach to determining an algorithm's order is to start out assuming an order of  $O(1)$ , an algorithm that doesn't do anything and immediately terminates no matter what the input. Then, find the section of code that you expect to have the highest order. From there, work out the algorithmic efficiency from the outside in -- figure out the efficiency of the outer loop or recursive portion of the code, then find the efficiency of the inner code; the total efficiency is the efficiency of each layer of code multiplied together. For instance, to compute the efficiency of a simple selection sort

```
for(int x=0; x<n; x++)
```

```
{
```

```

int min = x;

for(int y=x; y<n; y++)
{
    if(array[y]<array[min])
        min=y;
}

int temp = array[x];

array[x] = array[min];

array[min] = temp;
}

```

We compute that the order of the outer loop (for(int x = 0; ..)) is  $O(n)$ ; then, we compute that the order of the inner loop is roughly  $O(n)$ . Note that even though its efficiency varies based on the value of  $x$ , the average efficiency is  $n/2$  and we ignore the constant, so it's  $O(n)$ . After multiplying together the order of the outer and the inner loop, we have  $O(n^2)$ .

In order to use this approach effectively, you have to be able to deduce the order of the various steps of the algorithm. And you won't always have a piece of code to look at; sometimes you may want to just discuss a concept and determine its order. Some efficiencies are more important than others in computer science and on the next page, you'll see a list of the most important and useful orders of efficiency, along with examples of algorithms having that efficiency.

### **Algorithmic Efficiency -- Various Orders and Examples**

Below is a list of some of the common orders and with example algorithms.

#### **$O(1)$**

An algorithm that runs the same no matter what the input. For instance, an algorithm that always returns the same value regardless of input could be considered an algorithm of efficiency  $O(1)$ .

#### **$O(\log n)$**

Algorithms based on binary trees are often  $O(\log n)$ . This is because a perfectly balanced binary search tree has  $\log n$  layers and to search for any element in a binary search tree requires traversing a single node on each layer.

The binary search algorithm is another example of a  $O(\log n)$  algorithm. In a binary search, one is searching through an ordered array and beginning in the middle of the



remaining space to be searched, whether to search the top or the bottom half. You can divide the array in half only  $\log n$  times before you reach a single element, which is the element being searched for, assuming it is in the array.

### **$O(n)$**

Algorithms of efficiency  $n$  require only one pass over an entire input. For instance, a linear search algorithm, which searches an array by checking each element in turn is  $O(n)$ . Often, accessing an element in a linked list is  $O(n)$  because linked lists do not support random access.

### **$O(n \log n)$**

Often, good sorting algorithms are roughly  $O(n \log n)$ . An example of an algorithm with this efficiency is merge sort, which breaks up an array into two halves, sorts those two halves by recursively calling itself on them and then merging the result back into a single array. Because it splits the array in half each time, the outer loop has an efficiency of  $\log n$  and for each "level" of the array that has been split up (when the array is in two halves, then in quarters, and so forth), it will have to merge together all of the elements, an operations that has order of  $n$ .

### **$O(n^2)$**

A fairly reasonable efficiency, still in the polynomial time range, the typical examples for this order come from sorting algorithms, such as the selection sort example on the previous page.

### **$O(2^n)$**

The most important non-polynomial efficiency is this exponential time increase. Many important problems can only be solved by algorithms with this (or worse) efficiency. One example is factoring large numbers expressed in binary; the only known way is by trial and error, and a naive approach would involve dividing every number less than the number being factored into that number until one divided in evenly. For every increase of a single digit, it would require twice as many tests.

## **2.8 Summary**

Efficiency of algorithm is of prime concern in computer science. The complexity of an algorithm can be measured in terms of number of operations that are required to be performed for solving the problem called time complexity or the space required by the algorithm called space complexity. The notation used for representing complexity is Big-O. The time and space trade off of an algorithm maintains a balance in the time and space complexity of the algorithm.

## **2.9 Questions**

1. What is the need of writing efficient algorithms?
2. Define asymptotic notation.

3. What do you mean by Big-O notation? How is it computed?
4. What are the various notations used for representing algorithm complexity?
5. What do you mean by algorithm analysis?
6. Define time and space trade off.
7. Define algorithm efficiency. How efficiency of an algorithm is measured?

### **2.10 Suggested readings**

1. A. Tanenbaum, Y. Lanhgsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley

## ARRAYS

### 3.1 Objectives

### 3.2 Introduction

### 3.3 One Dimensional Arrays

### 3.4 Multi-dimensional Arrays

### 3.5 Operations on Arrays

### 3.6 Sparse Arrays and Matrices

### 3.7 Summary

### 3.8 Questions

### 3.9 Suggested readings

### 3.1 Objectives

We will explore arrays, multidimensional arrays, sparse arrays and matrices in this lesson.

### 3.2 Introduction

**An Array is a collection of variables having the same data type. It is a way to reference a series of memory locations using the same name.** Each memory location is represented by an array element. An array element is similar to one variable except it is identified by an index value instead of a name. An index value is a number used to identify an array element.

Now we'll show you what an array looks like, with the three array elements shown next. The array is called grades. The first array element is called grades[0]. The zero is the index value. The square bracket tells the computer that the value inside the square bracket is an index.

grades[0]

```
grades[1]
```

```
grades[2]
```

Each array element is like a variable name. For example, the following variables are equivalent to array elements. There is no difference between array elements and variables—well, almost no difference but we'll get to the differences in a moment. For now, let's explore how they are the same. Here are three integer variables:

```
int maryGrade;
```

```
int bobGrade;
```

```
int amberGrade;
```

You probably recall from your programming class that you store a value into a memory location by using an assignment statement. Here are two assignment statements. The first assigns a value to a variable and the other assigns a value to an array element. Notice that these statements are practically the same except reference is made to the index of the array element in the second statement:

```
int grades[1];
```

```
maryGrade = 90;
```

```
grades[0] = 90;
```

Suppose you want to use the value stored in a memory location. There are a number of ways to do this in a program but a common way is to use another assignment statement like the ones shown in the next example. The first assignment statement uses two variables, the next assignment statement uses two array elements and the last assignment statement assigns the value referenced by a variable name and assigns that value to an array element:

```
bobGrade = maryGrade;
```

```
grades[0] = grades[1];
```

```
grades[0] = bobGrade;
```

You've probably noticed a pattern developing. You use an array element the same way you use a variable.

There are two important differences between an array element and a variable, and those differences make working with large amounts of data a breeze. Suppose you had to work with 100 grades to calculate the average grade. How would you do this?

The challenge isn't applying the formula for calculating an average. You know how that's done. The challenge is to come up with 100 variable names and then reference all those variable names in a program. Ouch!

First, you'd need to sum all the grades by writing a statement similar to the following. (We'll stop at three variables because it's difficult to identify 100 variables—and we'd run out of space on this page.)

```
sum = maryGrade + bobGrade + amberGrade;
```

Now, here's how a smart programmer meets this challenge using an array:

```
sum = 0;

for (int i = 0; i < 100; i++)

    sum = sum + grades[i];
```

Big difference. The control variable of the for loop is the index for the array element, enabling the program to quickly walk through all array elements in two lines of code. (The first statement has nothing to do with walking through all the array elements. It only initializes the sum variable with the total grades.)

The other difference between an array and a variable is that all the array elements are next to each other in memory. Variables can be anywhere in memory. For example, `grades[0]` is next to `grades[1]` in memory, `grades[1]` is next to `grades[2]` in memory, and so on. In contrast, `maryGrade` and `bobGrade` variables can be anywhere in memory, even if they are declared in the same declaration statement.

The location of array elements is important when pointers are used to manipulate data stored in memory. It is more efficient to point to array elements than variables because the computer moves to the next memory location when you point to the next array element.

Some programmers might say that arrays are the backbone of data structures because an array enables a programmer to easily reorganize hundreds of values stored in memory by using an array of pointers to pointers.

So we drew a picture to show you the importance of arrays in data structures. Figure 3.1 shows memory. Each block is a byte. We'll say that two bytes are needed to store a memory address in memory. You need to store a memory address in memory because you'll use it to refer to other memory addresses in the "An Array of Pointers" section of this lesson.

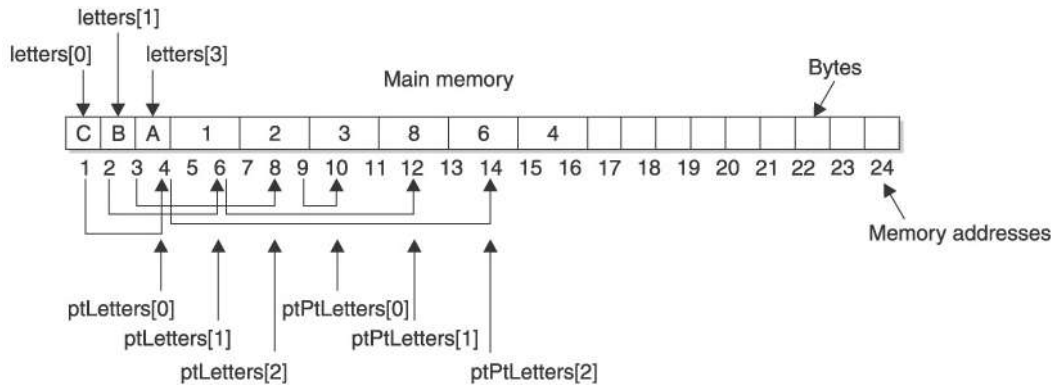


Figure 3.1: Elements of an array are stored sequentially in memory.

First, create an array called letters and assign characters to it, as shown here:

```
char letters[3];

letters[0] = 'C';

letters[1] = 'B';

letters[2] = 'A';
```

You'll notice in Figure 3.1 that each letter appears one after the other in memory. This is because these values are assigned to elements of an array and each array element is placed sequentially in memory.

Next, create an array of pointers. A pointer is a variable that contains a memory address of another variable. In this example, you'll use an array of pointers instead of a pointer variable.

An array of pointers is nearly identical to a pointer variable except each array element contains a memory address. Assign the memory address of each element of the letters array to elements of the ptLetters array, which is an array of pointers. Here's how this is done in C and C++:

```
char * ptLetters[3];

for (int i = 0; i < 3; i++)

    ptLetters[i] = &letters[i];
```

The ampersand (&), which is called the address operator, tells the computer to assign the memory address of the element of the letters array and not the contents of the element.

The final step is to create an array of pointers to pointers and then use it to change the order of the letters array when printing the letters array on the screen. A pointer to a pointer is a variable that contains the address of another pointer. In the example, you use an array of pointers to a pointer where each element of the array is like a pointer to a pointer variable. That is, each element is assigned an address of a pointer.

Use the following code to assign the memory address of each element of the `ptLetters` pointer array to the `ptPtLetters` pointer to pointer array. Notice that you don't use a for loop. This is because you need to change the order of the letters array without changing the letters array itself. Figure 3.1 shows how memory looks after the following code executes. If you printed elements of the `ptPtLetters` array, what would be displayed on the screen?

```
char ** ptPtLetters[3];  
  
ptPtLetters[0] = &ptLetters[2];  
  
ptPtLetters[1] = &ptLetters[1];  
  
ptPtLetters[2] = &ptLetters[0];
```

Here is the code that prints the `ptPtLetters` array:

```
for ( i = 0; i <3; i++)  
  
    cout << **ptPtLetters[i] << endl;
```

The answer to the question is A B C. Follow Figure 3.1 as we explain how this works. The first element of the `ptPtLetters` array is located at memory address 10. The content of memory address 10 is 8, which is memory address 8 because memory address 10 is the last element of the array `ptLetters`—a pointer. The value of memory address 8 is 3, which is the memory address of the third element of the array `letters`.

When the computer sees the `ptPtLetters[i]` statement for the first time, it goes to the array element `ptPtLetters[0]` and reads its value, which is 8. The computer then goes to memory address 8 and reads its content because memory address 8 is a pointer. The content of memory address 8 is 3, which is the memory address of the third element of the `letters` array. The computer reads the content of memory address 3 and displays the content on the screen.

This can be a bit tricky to follow unless you use Figure 3.1 as a guide; you can also use Figure 3.1 to explain how the computer displays the other letters.

The importance of using arrays for data structures is that you can easily change the order of data by using pointers and pointers to pointers without having to touch the original data. Some smart programmer might tell you that you're not saving any time or memory by using pointers and pointers to pointers to rearrange an array of characters.

The programmer is correct. However, we're juggling characters to illustrate how arrays and pointers to pointers work. In the real world, pointers typically point to a whole group of information such as a client's name, address, phone number and other pertinent data. Instead of juggling all that information, you need only to juggle memory addresses.

### 3.3 One Dimensional Arrays

The way to declare an array depends on the programming language used to write your program. In Java, there are two techniques for declaring an array. You can declare and initialize an array either where memory is allocated at compile time or where memory is dynamically allocated at runtime. Allocation is another way of saying reserving memory.

Let's begin by declaring an array where memory is reserved when you compile your program. This technique is similar in Java, C and C++, except in Java you must initialize the array when the array is declared. There are four components of a statement that declares an array. These components are a data type, an array name, the total number of array element to create and a semicolon (;). The semicolon tells the computer that the preceding is a statement. Here's the declaration statement in C and C++:

```
int grades[10];
```

In Java, you must initialize the array when the array is declared as shown here. The size of the array is automatically determined by counting the number of values within the braces. Therefore, there is n't any need to place the size of the array within the square brackets:

```
int[] grades = {0,0,0,0,0,0,0,0,0,0};
```

The data type is a keyword that tells the computer the amount of memory to reserve for each element of the array. In this example, the computer is told to reserve enough memory to store an integer for each array element.

The array name is the name you use within a program to reference an array element. The array name in this example is grades. The number within the square brackets is the total number of elements that will be in the array. The previous statements tell the computer to create an array of 10 elements. Previously in this chapter you learned that the index for the first array element is zero, not one. Therefore, the tenth array element has the index value 9, not 10.

Some programs confuse an index with the total number of array elements. That is, they use the value 9 within the square brackets when declaring an array because they assume they are declaring 10 elements. In reality, they are declaring an array of 9 elements. The confusion stems from the fact that 9 is the index to reference the tenth array element.

With a little practice you can avoid making this mistake. Remember that the value within the square brackets in the statement that creates an array is not an index,



although it resembles an index. This value is the number of array elements you need. That is, you insert the number 10 within the square brackets if you need 10 array elements. You use the index value of 9 if you want to access the tenth element.

In order to allocate memory at compile time, you must know the number of array elements that you need. Sometimes you don't know this, especially if your program loads the array with data stored in a database. The amount of data stored in a database typically fluctuates.

### 3.4 Multi-dimensional Arrays

The array described in this lesson is referred to as a one-dimensional array because the array consists of one series of elements. However, **an array can have more than one series of elements. This is called a multidimensional array.**

A multidimensional array consists of two or more arrays defined by sets of array elements, as shown in Figure 3.2. Each set of array elements is an array. The first set of array elements is considered the primary array and the second and subsequent sets of array elements are considered subarrays.

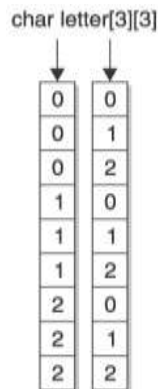


Figure 3.2: A two-dimensional array is a multidimensional array consisting of two arrays.

There are two arrays in the multidimensional array shown in Figure 3.2. Each element of the first array points to a corresponding array. For example, letters[1] in Figure 3.2 points to the array beginning with array element letters[1][0] where the zero is the first element of the second array.

Although you can create an array with any size multidimension, many programmers limit an array to two dimensions. Any size greater than two dimensions become unwieldy to manage.

An analogy we find helpful is visualizing a table (rows and columns) for a two-dimensional array and a cube (or similar figure) for a three-dimensional array.

#### 3.4.1 Need of a Multidimensional Array

A multidimensional array can be useful to organize subgroups of data within an array. Let's say that a student has three grades, a mid-term grade, a final exam grade, and a final grade. You can store all three grades for an endless number of students in a two-dimensional array, as shown in Figure 3.3.

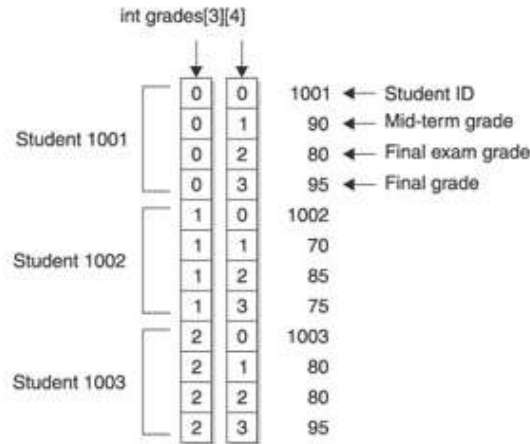


Figure 3.3: All three grades can be stored in a multidimensional array.

Figure 3.3 declares a multidimensional array of integers. The first set of array elements contains three array elements, one for each student. The second set of array elements has four array elements. The first of the four elements contains the student ID and the other three contain the three grades for that student ID.

Data stored in a multidimensional array is stored sequentially by sets of elements, as shown in Figure 3.4. The first set of four array elements is placed in memory, followed by the second set of four array elements, and so on.

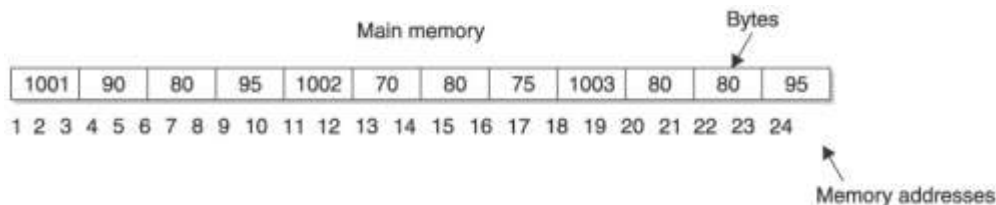


Figure 3.4: Elements of a multidimensional array are stored sequentially in memory.

The name of a multidimensional array references the memory address of the first element of the first set of four elements. That is, grades is the equivalent of using memory address 1 in Figure 3.4. You can use the name of a multidimensional array as a pointer to the entire array.

The index of the first element of the first set of array elements points to the memory address where values assigned to array elements are stored.

Referencing the index of the first dimension points to the memory address of the first element of that dimension. For example, referencing `grades[1]` points to memory address 9 in Figure 3.4. Memory address 9 is the first memory address of contiguous memory where values of the second set of array elements that are associated with `grades[1]` are stored.

A multidimensional array is declared similar to the way you declare a one-dimensional array except you specify the number of elements in both dimensions. For example, the multidimensional array shown in Figure 3.3 is declared as follows in C or C++:

```
int grades[3][4];
```

The first bracket (`[3]`) tells the compiler that you're declaring 3 pointers, each pointing to an array. This concept might be confusing because the term "pointer" may make some programmers think of pointer variable or pointer array, which you'll learn about later in this lesson. However, we are not talking about a pointer variable or pointer array. Instead, we are saying that each element of the first dimension of a multidimensional array reference a corresponding second dimension, which is an array.

In this example, all the arrays pointed to by the first index are of the same size. The second index can be of variable size. For example, the previous statement declares a two-dimensional array where there are 3 elements in the first dimension and 4 elements in the second dimension.

The element `grades[0]` is said to "point" (just as you use your finger to point) to the second dimension of the array, which is referenced as `grades[0][0]`. The second dimension is considered an array. Therefore, programmers say that the first element of a multidimensional array points to another array (that is, the second dimension).

The data type tells the computer that each element of the array will contain an integer data type. The data type is followed by the array name and two values that indicate the size of each dimension used for the array. In this case, there are three sets of four array elements.

You declare a multidimensional array and initialize its elements by using French braces, as shown in Figure 3.5. There are three sets of inner French braces. Each of these sets represents the first dimension of the array. There are four values within each set of inner French braces. These values are assigned to each element of the second dimension of the array.

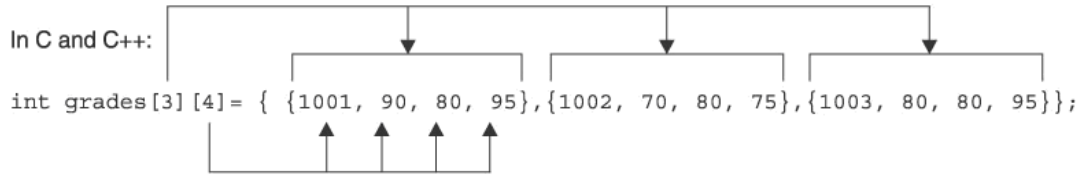


Figure 3.5: Braces define sets of values to be assigned to array elements.

You assign a value to an element of a multidimensional array with an assignment statement similar to the assignment statement that assigns a value to a single-dimensional array, as shown here:

```
grades[0][0] = 1001;
```

You must specify the index for both dimensions. In this example, the integer 1001, which is a student ID, is assigned to the first element of the first set of elements in the grades array.

### 3.4.2 Referencing Contents of a Multidimensional Array

The contents of elements of a multidimensional array can be used in a program by referencing the index of both dimensions of the array element. Figure 3.6 shows you how to display the final exam grade for the second student.

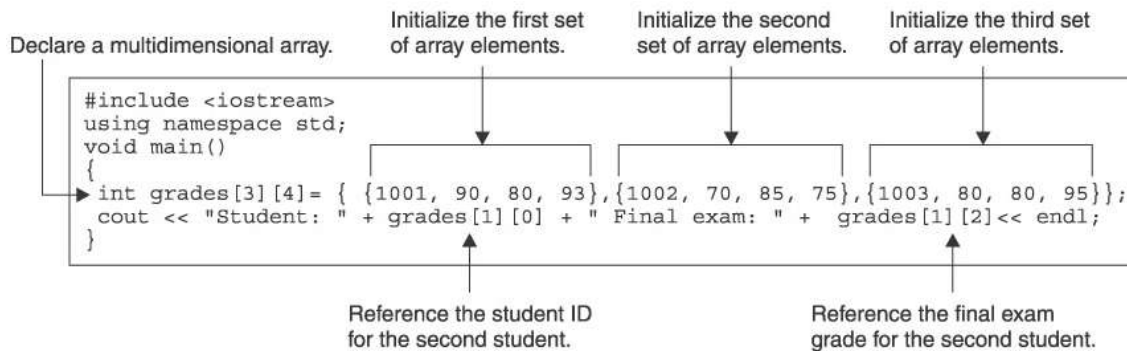


Figure 3.6: Display the contents of array elements by referencing the index of both sets of array elements.

In this example, the student ID is displayed by referencing the first element of the second set, and the grade for the final exam is displayed by referencing the third element of the second set.

#### Index of the first element

The index of the first element (sometimes called the "origin") varies by language. There are three main implementations: zero-based, one-based and n-based arrays, for which the first element has an index of zero, one or a programmer-specified value. The zero-based array is more natural in the root machine language and was popularized by the C programming language, where the abstraction of array is very weak and an index  $n$  of a one-dimensional array is simply the offset of the element accessed from the address of the first (or "zeroth") element (scaled by the size of the element). One-based arrays are based on traditional mathematics notation for matrices and most but not all, mathematical sequences. n-based is made available so the programmer is free to choose the lower bound, which may even be negative, which is most naturally suited for the problem at hand.

### **Index of the last element**

The relation between numbers appearing in an array declaration and the index of that array's last element also varies by language. In some languages (e.g. C) the number of elements contained in the arrays must be specified, whereas in others (e.g. Visual Basic .NET) the numeric value of the index of the last element must be specified.

### **3.5 Operations on Arrays**

An array is the simplest of all the data structures to implement and use. The various operation that can be applied on array are:

- i. Insertion of an element:** An element can be inserted in the array by specifying the location where it is to be stored. The location should not be greater than the size of the array. e.g.  $A[i] = 10$  this will insert value 10 at the  $i^{\text{th}}$  location in the array and value of  $i$  in this case should not be greater than size of the array.
- ii. Deletion of an element:** Since all element of an array are store in sequential order in the memory. Deletion of an element can be performed by sifting the element appearing after the element to be deleted from the array.
- iii. Searching:** Searching on array can be performed in linear order if the array elements are not arranged in ascending or descending order. If elements are arranged then binary search can be applied which is far superior to linear search. Complexity of linear search in worst case is  $O(n)$  as compared to that of binary search which is  $\log_2 n$ . Consider search for a number, which is not in the array, in an array of 10000 numbers. Linear search needs 10000 comparisons to declare that the number is not present. Compare this to the comparisons needed by binary search, which will need about 15 comparisons to make the same declaration.
- iv. Sorting:** It means arranging the elements of the array in ascending or descending order. There are various sorting algorithms for performing this task. These algorithms have been discussed in detail in lessons 17 and 18.

### 3.6 Sparse Arrays and Matrices

**A sparse array is an array in which most of the elements have the same value (known as the default value -- usually 0 or null).**

A naive implementation of an array may allocate space for the entire array but in the case where there are few non-default values, this implementation is inefficient. Typically the algorithm used instead of an ordinary array is determined by other known features (or statistical features) of the array, for instance if the sparsity is known in advance, or if the elements are arranged according to some function (e.g. occur in blocks).

As an example, a spreadsheet containing 100x100 mostly empty cells might be more efficiently stored as a linked list rather than an array containing ten thousand array elements.

#### 3.6.1 Representation of Sparse matrix

The naive data structure for a matrix is a two-dimensional array. Each entry in the array represents an element  $a_{i,j}$  of the matrix and can be accessed by the two indices  $i$  and  $j$ . For an  $m \times n$  matrix we need at least enough memory to store  $(m \times n)$  entries to represent the matrix.

Many if not most entries of a sparse matrix are zeros. The basic idea when storing sparse matrices is to store only the non-zero entries as opposed to storing all entries. Depending on the number and distribution of the non-zero entries, different data structures can be used and yield huge savings in memory when compared to a naïve approach.

A simple method of representing sparse matrix is a two dimensional array of  $NNZ+1 \times 3$ , where  $NNZ$  is the number of non-zero elements and number of columns is fixed to 3.

The first row of the matrix contains number of rows, number of columns and number of non-zero elements in the source matrix. The remaining  $NNZ$  rows store the row, column and value of the non-zero element.

For example, consider the source matrix

$$\begin{bmatrix} 1 & 2 & 0 & 0 \\ 0 & 3 & 9 & 0 \\ 0 & 1 & 4 & 0 \end{bmatrix}$$

The simple sparse representation of this matrix is

$$\begin{bmatrix} 3 & 4 & 6 \\ 1 & 1 & 1 \\ 1 & 2 & 2 \end{bmatrix}$$

[ 2 2 3 ]

[ 2 3 9 ]

[ 3 2 1 ]

[ 3 3 4 ]

Another example of such a sparse matrix format is the Yale Sparse Matrix Format. It stores an initial sparse  $m \times n$  matrix,  $M$ , in row form using three one-dimensional arrays. Let  $NNZ$  denote the number of nonzero entries of  $M$ . The first array is  $A$ , which is of length  $NNZ$  and holds all nonzero entries of  $M$  in left-to-right top-to-bottom order. The second array is  $IA$ , which is of length  $m + 1$  (i.e., one entry per row, plus one).  $IA(i)$  contains the index in  $A$  of the first nonzero element of row  $i$ . Row  $i$  of the original matrix extends from  $A(IA(i))$  to  $A(IA(i+1)-1)$ . The third array,  $JA$ , contains the column index of each element of  $A$ , so it also is of length  $NNZ$ .

Consider the source matrix given above, the Yale sparse matrix format is

$A = [ 1 2 3 9 1 4 ]$

$IA = [ 1 3 5 7 ]$

$JA = [ 1 2 2 3 2 3 ]$

### 3.7 Summary

Array is a linear data structure. Elements of an array are stored sequentially in memory. Each element of an array can be addressed by an index. Therefore, accessing an element of array does not require any traversal through the other element of it. All languages provide arrays as inbuilt data structures.

Arrays can be one dimensional or multi-dimensional. Two dimensional array is called matrix, which has wide application in mathematics as well as computer science.

A matrix with very few non-zero elements is called a sparse matrix and there are different methods of storing the sparse matrix other than the normal two dimensional matrix.

### 3.8 Questions

1. What is the difference between an array element and a variable?
2. How do you declare an array?
3. How are elements of an array stored in memory?
4. Define a matrix.
5. Define sparse matrix. What are the various methods of representing a sparse matrix?
6. Write the program for converting a two dimensional matrix to Yale's form of sparse matrix representation.
7. WAP to traverse an Array.
8. WAP to merge two sorted Array.
9. WAP to transpose an Array.

### **3.9 Suggested readings**

1. A. Tanenbaum, Y. Lanhgsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley.



## STACKS

### 4.1 Objectives

### 4.2 Introduction

### 4.3 Operations on Stacks

4.3.1 Push

4.3.2 Pop

4.3.3 Checking overflow and underflow

### 4.4 Creating a Stack in C++

4.4.1 Implementing IsFull

4.4.2 Implementing IsEmpty

4.4.3 Creating a Push Member Function in C++

4.4.4 Creating a Pop Member Function in C++

### 4.5 An Example of Stack Implementation

### 4.6 Summary

### 4.7 Questions

### 4.8 Suggested readings

### 4.1 Objectives

A stack is a linear data structure which follows a particular order in which operations are performed. A stack is the way you group things together by placing one thing on top of another and then removing things one at a time from the top of the stack. It is amazing that something this simple is a critical component of nearly every program that is written. In this lesson, you'll learn how to create and use a stack in your programs.

### 4.2 Introduction

A stack is a type of data structure – a means of storing information in a computer. When a new object is entered in a stack, it is placed on the top of all the previously entered objects. When you hear the term “stack” used outside the context of computer programming, you might envision a stack of dishes on your kitchen counter. This organization is structured in a particular way: the newest dish is on top and the oldest is on the bottom of the stack.

Each dish in a stack is accessed using fifo: first in, first out. The only way to access each dish is from the top of the stack. If you want the third dish (the third oldest on the stack), then you must remove the first two dishes from the top of the stack. This places the third dish at the top of the stack making it available to be removed.

There’s no way to access a dish unless the dish is at the top of the stack. You might be thinking stacks are inefficient and you’d be correct if the objective was to randomly access things on the stack. There are other data structures that are ideal for random access, which you’ll learn about throughout these lessons.

However, if the object is to access things in the order in which they were placed on the stack, such as computer instructions, stacks are efficient. In these situations, using a stack makes a lot of sense.

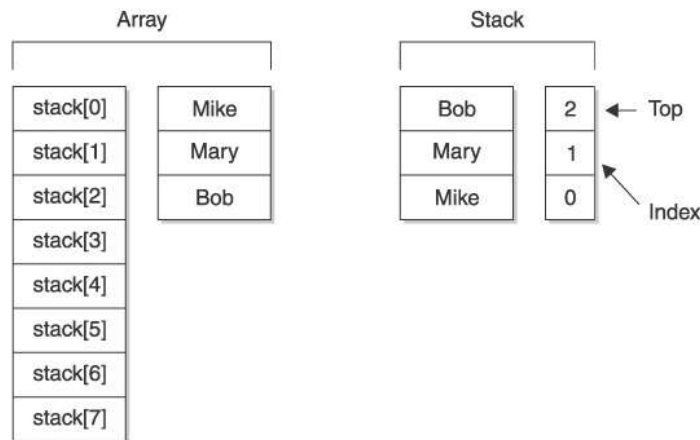


Figure 4.1: A stack and an array are two different things: an array stores values in memory; a stack tracks which of the array elements is at the top of the stack.

### 4.3 Operations on Stacks

Programmers use arrays to store values that are referenced by a stack. An array consists of a series of array elements, each of which is similar in concept to a variable. The stack contains the index of the array element that is at the top of the stack.

Figure 4.1 is the way some programmers envision an array used with a stack. This example shows an array called stack with 8 array elements. The entire array contains values that are referenced by the stack. Three array elements are assigned values, while

the other array elements are empty and can be used when new items are placed on the stack (see the upcoming section “Push”).

Mike is the first value placed on the stack. You know this because Mike is at the bottom of the stack. Bob is the last item placed on the stack because Bob is the top item on the stack.

### 4.3.1 Push

**Programmers use the term “push” to mean placing an item on a stack.** Push is the direction that data is being added to the stack. Think of this as pushing items down on the stack to move the items already on the stack down to make room for the next item.

Here’s what actually happens. The new value is assigned to the next available array element and the index of that array element becomes the top of the stack, as shown in Figure 4.2. The program increments the current index of the stack by 1. In this example, the index is incremented by 1, resulting in index 3 being at the top of the stack, which is the index of the new values assigned to the array.

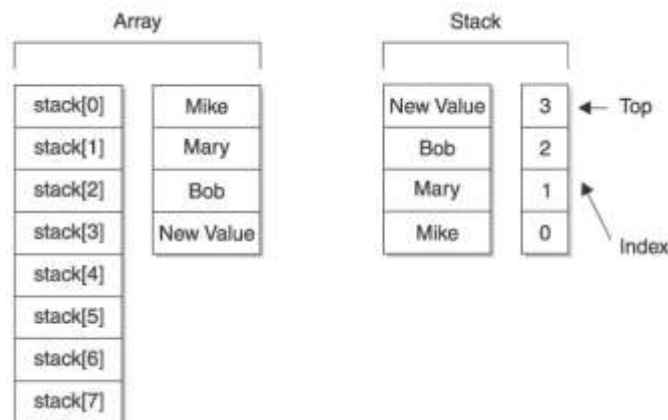


Figure 4.2: The new value is assigned to the next array element and its index becomes the top of the stack.

### 4.3.2 Pop

**Popping is the reverse process of pushing: it removes an item from the stack.** It is important to understand that popping an item off the stack doesn’t copy the item. Once an item is popped from the stack, the item is no longer available on the stack, although the value remains in the array.

Here’s what really happens. Remember that the top of the stack contains an index of the array element whose value is at the top of the stack. In Figure 4.2, index 3 is at the top of the stack, which means New Value in array element 3 is at the top of the stack.

When you pop New Value from the stack, you decrement the index at the top of the stack. That is, you make its index 2 instead of 3. This makes Bob the new value at the top

of the stack (see Figure 4.3). Notice that New Value and array element 3 remain untouched in the array because popping a value from the stack only alters the stack, not the underlying array.

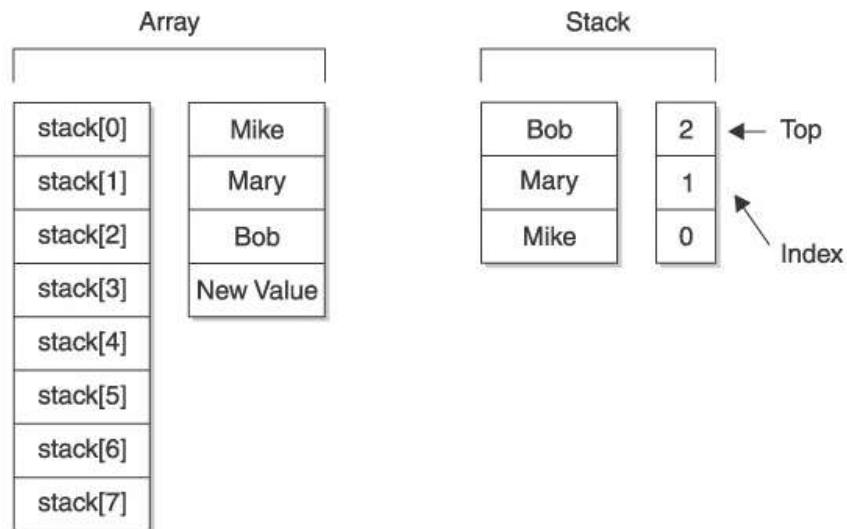


Figure 4.3: All values move toward the top of the stack when the top item is popped off the stack.

### 4.3.3 Checking overflow and underflow

Deletion of stack element is not possible when the stack is empty. This is called underflow checking. Similarly an element can not be inserted when the stack is full. This is called overflow checking.

While implementing Push operation, before actual insertion of the element in the stack, its overflow condition must be checked. And while implementing Pop operation, before actual deletion of stack element, its underflow condition must be checked.

## 4.4 Creating a Stack in C++

You can create a stack in C++ by defining a stack class and declaring an instance of the class. The Stack class requires three attributes and several member functions, which are defined as you learn about them. You'll begin by defining a basic stack class that has only the components needed to create the stack.

The class is called Stack but you can call it any name you wish. This class definition is divided into a private access specifier section and a public access specifier section. The private access specifier section has attributes and member functions (although not in this example) that are accessible only by a member function defined in this class. The public access specifier section has attributes (although not in this example) and member functions that are accessible by using an instance of the class.

The private access specifier section of the Stack class defines three attributes: size, top, and values, all of which are integers. The size attribute stores the number of elements in the stack, the top attribute stores the index of the top element of the stack, and the values attribute is a pointer to the stack, which is an array. The stack in this example is a stack of integer values, but you can use an array of any data type, depending on the nature of your program.

Only one member function is defined in the Stack class, although we'll define other member functions for the class in upcoming sections of this lesson. For now, let's keep the example simple and easy to understand. This member function is called Stack, which is the constructor of the class. A constructor is a member function that has the same name as the class and is called when an instance of the class is created. The code for this is on the next page.

Several things are happening in the constructor. First, the constructor receives an integer as an argument that is passed when the instance of the Stack class is declared. The integer determines the number of elements in the stack and is assigned to the size variable.

The first statement might look a bit confusing. It appears that the value of the size variable from the argument list is being assigned to itself but that's not the case. Actually, the size variable from the argument list is local to the Stack member function. The `this->size` combination refers to the size attribute of the Stack class, as shown here:

```
    this->size = size;
```

Programmers use the `this` pointer within a member function to refer to the current instance of the class. In this example, the 'this' pointer uses the pointer reference (`->`) to tell the computer to use the size attribute of the class. As you'll remember from your C++ programming class, the pointer reference is used when indirectly working with a class member and the dot operator is used when you are directly working with a class member.

This allows the compiler to distinguish between a class variable and local variable that have the same name. This means that the value of the size variable that is passed as an argument to the Stack member function is assigned to the size attribute, making the value available to other members of the Stack class.

You can see how the size attribute is used in the next statement. This statement does two things. First, it allocates memory for the stack by using the new operator (`new int[size]`). The new operator returns a pointer to the reserved memory location. The size is the size attribute of the class and determines the size of the array. The array is an array of integers.

Next, the pointer to the array of integers is assigned to the values attribute of the class. The values attribute is a pointer variable that is defined in the private attribute section of the Stack class.

The last statement in the Stack member function assigns a -1 to the top attribute. The value of the top attribute is the index of the top element of the stack. A -1 means that that stack doesn't have any elements. Remember from your programming class that index values are memory offsets from the start of the array. Index 0 means "move 0 bytes from the start of the array." So index -1 is just a convenient way to say that the stack is empty.

We'll expand on the definition of the Stack class in the next section but for now let's create an instance of the Stack class. The instance is declared within the main() function of this example. Three things are happening here. First, the new operator is creating an instance of the stack in memory. The new operator returns a pointer to that memory location.

Next, the statement declares a reference to the stack, which is called myStack. The reference is a pointer. The final step is to assign the pointer returned by the new operator to the reference. You then use the reference (myStack) as the name of the instance of the Stack class throughout the program.

```
public class Stack{
    private:
        int size;
        int top;
        int* values;
    public:
        Stack(int size){
            this->size = size;
            values = new int[size];
            top = -1;
        }
};

void main(){
    Stack *myStack = new Stack(10);
}
```

#### 4.4.1 Implementing IsFull

We'll create different member functions for each step, beginning by defining a member function that determines if there is room on the stack i.e. it check the overflow condition. We'll call it `IsFull()` and define it in the following code. The `IsFull()` member function is simple. It compares the value of the `top` attribute with the one less than the value of the `size` attribute.

```
bool IsFull(){
    if(top < size-1){
        return false;
    }
    else{
        return true;
    }
}
```

With the `IsFull()` member function defined, we'll move on to defining the `push()` member function, as shown in section 4.4.3.

#### **4.4.2 Implementing IsEmpty**

The `IsEmpty()` member function determines if there are any values on the stack i.e. the underflow condition. Let's define the `IsEmpty()` member function in this next example. The `IsEmpty()` member function contains an `if` statement. The condition expression of the `if` statement compares the value of the `top` attribute to `-1`. Remember that `-1` is the initial value of the `top` attribute when the instance of a stack is declared. If the `top` attribute is equal to `-1`, then a `true` is returned because the stack is empty; otherwise, a `false` is returned.

```
bool IsEmpty(){
    if(top == -1){
        return true;
    }
    else{
        return false;
    }
}
```

#### **4.4.3 Creating a Push Member Function in C++**

Now that you've seen how to define a class that creates a stack, we'll show you how to define additional member functions that enable the class to push values onto the stack.

Pushing a value onto the stack is a two-step process. First, you must determine if there is room on the stack for another value. If there is, you push the value onto the stack; otherwise, you don't.

The value of the `top` attribute is `-1` when the instance of the stack is declared. Suppose the value of `size` is `10`. The condition expression in the `if` statement of the `IsFull()` member function determines if the value of `top`, which is `-1`, is 1 less than the value of `size`. Since the value of `size` is `10`, the condition expression compares `-1 < 9`. If `top` is greater than or equal to `9`, then a `true` is returned; otherwise, a `false` is returned.

Why do you subtract 1 from the size of the stack? The value of the `top` attribute is an index of an array element. Remember that the index begins with zero. In contrast, the `size` is actually the number of array elements in the stack. Therefore, the tenth array element on the stack has an index of `9`.

The `push()` member function pushes a value onto the stack. The value being pushed onto the stack is passed as an argument to the `push()` member function and is assigned to the variable `x` in this example.

Before doing anything else, the `push()` member function determines if there is room on the stack by calling the `IsFull()` member function in the condition expression of the `if` statement. The condition expression might look a little strange because the call is preceded by an exclamation point (!) so we'll take apart the condition expression to explain what is really happening here.

Remember from your programming classes that statements within an `if` statement execute only if the condition expression is `true`. This means the condition expression must be `true` for the value passed to the `push()` member function to be placed on the stack.

Here's a slight problem. We're calling the `IsFull()` member function to determine if there is room on the stack for another value. However, the `IsFull()` member function returns `false` if there is room and `true` if there isn't room. A `false` causes the `push()` member function to skip statements that place the value on the stack. We really need the `IsFull()` member function to return a `true` if there is room available, not a `false`. Rather than rewrite the `IsFull()` member function, we use the exclamation point to reverse the logic. As you remember from your programming class, the exclamation point is the not operator—that is, a `false` is treated as a `true`, which causes the value to be placed on the stack.

There are two statements within the `if` statement. The first statement increments the value of the `top` attribute, which is the index of the last value placed on the stack. If the stack is empty, then the current value of the `top` attribute is `-1`. Incrementing `-1` changes the value of the `top` attribute to `0`, which is the index of the first array element of the stack. The last statement in the `if` statement assigns the value passed to the `push()` member function to the next available array element.



```

void push(int x){
    if(!IsFull()){
        top++;
        values[top] = x;
    }
}

```

#### 4.4.4 Creating a Pop Member Function in C++

The pop() member function of the Stack class has the job of changing the index that is at the top of the stack and returning the value of the corresponding array to the statement that calls the pop() member function. The next example defines the pop() member function.

The first statement in the definition declares an integer variable called retVal that stores the value returned by the pop() member function. The retVal is initialized to zero. Next, the IsEmpty() member function is called in the condition expression of the if statement to determine if there is a value at the top of the stack. Notice the exclamation point reverses the logic as it did in the pop() member function.

Statements within the if statement should only execute if the IsEmpty() member function returns a false, meaning the stack is not empty. Therefore, we need to use the exclamation point to reverse the logic of the condition expression to make the condition expression true if the IsEmpty() member function returns a false.

Two steps occur within the if statement. First, the value at the top of the stack is assigned to the retVal variable by referencing the values array using the index contained in the top attribute. Next, the value of the top attribute is decremented. The return retVal is then returned by the pop() member function.

```

int pop(){
    int retVal = 0;
    if(!IsEmpty()){
        retVal = values[top];
        top--;
    }
    return retVal;
}

```

#### 4.5 An Example of Stack Implementation

Now that you understand how to create and use a stack, we'll pick up the pace and explore an industrial-strength stack. You've may have heard the term industrial strength used in relation to programming and may be curious what this really means.

Industrial strength is a term used in industry that implies a product is designed to withstand stress. Industrial strength can be used to describe any kind of product but in this case the product is the program that creates and uses a stack.

Programs used to illustrate the concepts of a stack in this lesson are bare bones and lack the robust features that are found in industrial-strength programs. A bare-bones program is what you need when you're learning the concepts of stacks and other data structures because the program contains only statements that pertain to what you are learning. However, once you learn the concept, you need to see how it's applied in a real-world program. In this section, you'll take a look at how a stack is created and used in an industrial-strength C++ program.

We'll use as an example an industrial-strength C++ program that creates and uses a stack. The program is contained within three files, `stack.h`, `stack.cpp` and `stackDemo.cpp`. The `stack.h` file is a header file that contains the definition of the `Stack` class, which is the "blueprint" of the `Stack` class. The `stack.cpp` file is a source code file that contains the implementation of the member functions of the `Stack` class. The `stackDemo.cpp` file contains the source code for the C++ program that declares the instance of the `Stack` class and calls its member functions. Let's begin by taking a look at the `stack.h` header file, which is shown in the next code example. As you'll recall from your C++ classes, a header file typically contains definitions and preprocessor instructions. A preprocessor is a program that applies preprocessor instructions to source code before the code is compiled.

The `stack.h` header file contains one preprocessor instruction, `#define`, which defines a symbol. Here we've defined the symbol `DEFAULT_SIZE` and given it a value of 10. The preprocessor then replaces all occurrences of `DEFAULT_SIZE` with 10 before the code is compiled. The `DEFAULT_SIZE` is the default size of the stack if the no argument is passed to the constructor. Function parameters in C and C++ can be assigned default values in the function prototype as long as those arguments are at the end of the argument list. If the size value is not passed in, it gets defaulted to the value of `DEFAULT_SIZE`, which is 10 in our example.

The `stack.h` file also contains the definition of the `Stack` class. The `Stack` class definition has the same size, top and values attributes you saw in the previous C++ example. However, the definition of member functions is different from what you saw because member functions are implemented outside the class definition in the `stack.cpp` source code file. The header file contains only the prototype of the functions, which make up the blueprint for the class.

From your C++ class, you'll remember that only the prototype or signature of a member function needs to be included in a class definition. The implementation of the member function can be outside the class definition. There are two important reasons for keeping the definition (header file) and implementation (source) in separate files:

It allows you to provide a commercial software application programmer interface to a programmer without handing over your source code. You provide the programmer with

your header files, which they will use to compile their code (they only need header files to compile the code). You provide your source code in the form of precompiled libraries that are referenced by the programmer's program during linking.

The class definition contains signatures of six member functions. The first member function is called Stack, which is the constructor that you learned about previously in this lesson. Previously, you learned that the constructor is passed an integer representing the size of the stack. In the real-world version, the program sets a default size that can be overridden when an instance of the class is created in the program. The default size is specified by using the DEFAULT\_SIZE, which is 10 (see #define). The next member function is ~Stack() and is the destructor of the class. A destructor is the last member function that is called when the instance of the class goes out of scope and dies. A constructor must always be the same name as the class and begin with a tilde (~). By definition, destructors cannot accept any arguments. The purpose of the destructor is to free memory that is used by the stack or do any other sort of cleanup that's required.

The remaining member functions are the same functions that you learned about previously in this lesson.

```
//stack.h

#define DEFAULT_SIZE 10

class Stack {

private:

    int size;

    int top;

    int* values;

public:

    Stack(int size = DEFAULT_SIZE);

    virtual ~Stack();

    bool IsFull();

    bool IsEmpty();

    void push(int);

    int pop();

};
```

The `stack.cpp` file is a source code file that contains the implementation of the `Stack` class's member functions. We placed these in a different file from the class definition because it is easier to read and maintain as well as for other reasons explained previously. The file begins with the preprocessor instruction `#include` that tells the computer to evaluate the contents of the `stack.h` file before compiling the `stack.cpp` file so it "knows" about the `Stack` class definition before compiling the program. Member functions in the `stack.cpp` file will be familiar to you because all except one are the same member functions that you learned about previously in the lesson. However, the names of the member functions might be confusing at first glance because each name begins with the name of the class followed by two colons (`::`). The two colons are called the scope resolution operator.

You must precede the name of a member function with the class name and scope resolution operator if the member function is defined outside the class definition. Think of this as telling the computer that the member function belongs to the `Stack` class. The `~Stack()` member function frees memory used by the stack. It does this by using the `delete` operator and referencing the name of the array used for the stack. In this example, `values` is the name of the array.

To avoid memory leaks, freeing memory is important whenever memory is dynamically allocated. The square brackets (`[]`) are used with `delete` because the object being removed from memory was dynamically created.

The `stack.cpp` is compiled as you would compile any source code. The result is an object file that is joined together with the compiled `stackDemo.cpp` source code file by the linker to create an executable program called a load module.

```
//stack.cpp

#include "stack.h"

Stack::Stack(int size){

    this->size = size;

    values = new int[size];

    top = -1;

}

Stack::~~Stack(){

    delete[] values;

}

bool Stack::IsFull(){
```

```

    if(top < size-1){
        return false;
    }
    else{
        return true;
    }
}

bool Stack::IsEmpty(){
    if(top == -1){
        return true;
    }
    else{
        return false;
    }
}

void Stack::push(int x){
    if(!IsFull()){
        top++;
        values[top] = x;
    }
}

int Stack::pop(){
    int retVal = 0;
    if(!IsEmpty()){
        retVal = values[top];
        top--;
    }
}

```

```

    }

    return retVal;

}

```

Finally, we come to the stackDemo.cpp program, which is the C++ program that creates the instance of the Stack class. The first statement creates the stack in a three-step process. The first step is to use the new operator to allocate space in memory for the Stack class by calling the constructor of the class. The new operator returns the memory location of the stack. The second step is to declare a pointer called stack. The last step is to assign the memory location returned by the new operator to the stack pointer. In this example, we used the default size for the stack, which is 10 elements. We can pass the Stack() constructor an integer to change the size of the stack. The push() member function is called three times. Each time a different value is placed on the stack. Notice that the -> pointer is used instead of the dot operator. You must do this because stack is a pointer to an instance of the class and not the instance itself.

The last portion of the stackDemo.cpp program calls the pop() member three times. Each time a value is removed from the top of the stack and displayed on the screen.

```

//stackDemo.cpp

void main() {

    Stack *stack = new Stack();

    stack->push(10);

    stack->push(20);

    stack->push(30);

    for(int i=0; i<3; i++) {

        cout << stack->pop() << endl;

    }

}

```

## 4.6 Summary

Stack is a linear data structures extensively used in varied application of computer programming. It is based in Last In First Out (LIFO) principle. A top pointer maintains the count of elements present in the stack. Push operation is used to insert elements in stack and Pop operation is used to remove an element from top of the stack. Underflow and overflow conditions are check using IsEmpty and IsFull functions, respectively.

#### **4.7 Questions**

1. What is a stack?
2. What is the purpose of the push() member method?
3. What is the purpose of the pop() member method?
4. What is the purpose of the IsFull() member method?
5. What is the purpose of the IsEmpty() member method?
6. What kind of value is assigned to the top attribute?
7. Why is the top attribute initialized to -1?

#### **4.8 Suggested readings**

1. A. Tanenbaum, Y. Lanhsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley.

## APPLICATIONS OF STACKS

### 5.1 Objectives

### 5.2 Introduction

### 5.3 Parenthesis Matching

### 5.4 Conversion from Infix to Postfix expression

### 5.5 Evaluation of Postfix Expression

### 5.6 Recursion

### 5.7 Summary

### 5.8 Questions

### 5.9 Suggested readings

### 5.1 Objectives

Stacks have wide applications in computers programming. In this lesson we shall discuss some most common applications of stacks.

### 5.2 Introduction

Stacks are used in various applications of computers. These may be used for reversal of a string, matching or parenthesis in expression or programs, conversion of infix expression to postfix, evaluation of postfix expression, simulating recursion etc.

### 5.3 Parenthesis Matching

**The simplest application of stacks deals with removing the parentheses, this task ensures that all paired symbols match correctly, including parentheses ( ), brackets [ ] and braces { }.** This program will read strings of input from the user, indicating whether the parentheses, brackets and braces match up correctly. Use a character stack to store the necessary characters and make comparisons. All characters other than parentheses, brackets and braces are ignored for purposes of matching. Many



applications of stacks deal with determining whether an input string (or file or sequence of symbols) is a member of a context-free language. Many programming languages are context free.

As an example, consider the language of balanced parentheses. A sequence of symbols involving  $()^{+*}$  and integers, is said to have balanced-parentheses if each right parenthesis has a corresponding left parenthesis that occurs before it. For example: These expressions have balanced parentheses:

`2*7 // no parens - still balanced`

`(1+3)`

`((2*16)+1)*(44+(17+9))`

These expressions do not:

`(44+38`

`) // a right paren with no left`

`(55+(12*11) // missing a )`

How do we tell if a sequence of characters represents a balanced-parentheses expression? Use stacks. Idea:

- Start with an empty stack.
- For each left parenthesis, push.
- For each right parenthesis, pop.
- For each non-parenthesis character, do nothing.
- The expression is balanced if the stack is empty and there are no more characters to process.
- It is not balanced if either after the last character the stack is not empty (too many left parens) or if the stack is empty and a right paren is encountered (a right without a left).

What items do we push onto the stack? If we are just interested in knowing if the expression is balanced, it does not matter. For clarity, we might choose a stack of characters, pushing "(" onto the stack for each left parenthesis.

We can do more. We can match multiple types of delimiters. For example, we might want to match `()`, `[]` and `{}`. In that case we can push the left delimiter onto the stack and when we pop, do a check, as in the pseudocode:

```
if (next character is a right delimiter) then {  
    if (stack is empty) then  
        return false
```

```

else {

    pop the stack

    if (right delimiter is corresponding version
        of what was popped off the stack) then

        continue processing

    else

        return false

}

}

```

Other variations on parentheses matching include:

- valid prefix: return true if the expression is a valid prefix for a balanced parentheses expression, for example " $((3+4)+$ " is a valid prefix. Same idea, but we can return true as long as there is no attempt to pop an empty stack.
- identifying matches: it is often useful to not only know that an expression is balanced, but to see which pair correspond. For example, in some editors (like Emacs), when you type a right paren the corresponding left paren is momentarily highlighted. Think about how you might adapt a paren matching function to do this.

```

int main(){
    char c;
    int k;
    ifstream f;
    f.open( "input.txt" );
    if( ! f ){
        cout << "Error opening file. Quitting.\n";
        exit(-1);
    }
    //initialize
    stk_init();
    c = getchar();
    putchar(c);
    // loop through file one char at a time
    while (c!=EOF)
    {
        // process one character push (, pop if ),
        // check stack empty if \n, ignore others
        // if error, move to next line
    }
}

```

```

if (c == '{'){
    stk_push(c);
}
else if (c == ')') {
    stk_top (c);
    if (c == '{')
        stk_pop (c);
    else// pop stack, if empty have missing (
}
else if (c == '\n'){
    // check on missing ) set up for next line
    stk_init();
}
if (stk_error()) //skip rest of line
else{ // get next char
}
} // endwhile
return 0;
}

```

#### 5.4 Conversion from Infix to Postfix expression

Any expression in the standard form like "2\*3-4/5" is an Infix (In order) expression. The Postfix (Post order) form of the above expression is "23\*45/-". In a postfix expression the operands precede the operator.

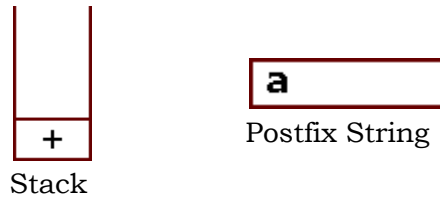
In normal algebra we use the infix notation like  $a+b*c$ . The corresponding postfix notation is  $abc*+$ . The algorithm for the conversion is as follows:

- Scan the Infix string from left to right.
- Initialise an empty stack.
- If the scanned character is an operand, add it to the Postfix string. If the scanned character is an operator and if the stack is empty Push the character to stack.
- If the scanned character is an Operand and the stack is not empty, compare the precedence of the character with the element on top of the stack (topStack). If topStack has higher precedence over the scanned character Pop the stack else Push the scanned character to stack. Repeat this step as long as stack is not empty and topStack has precedence over the character.  
Repeat this step till all the characters are scanned.
- (After all characters are scanned, we have to add any character that the stack may have to the Postfix string.) If stack is not empty add topStack to Postfix string and Pop the stack. Repeat this step as long as stack is not empty.
- Return the Postfix string.

Example :

Let us see how the above algorithm will be implemented using an example.  
Infix String :  $a+b*c-d$

Initially the Stack is empty and our Postfix string has no characters. Now, the first character scanned is 'a'. 'a' is added to the Postfix string. The next character scanned is '+'. It being an operator, it is pushed to the stack.



Next character scanned is 'b' which will be placed in the Postfix string. Next character is '\*' which is an operator. Now, the top element of the stack is '+' which has lower precedence than '\*', so '\*' will be pushed to the stack.



The next character is 'c' which is placed in the Postfix string. Next character scanned is '-'. The topmost character in the stack is '\*' which has a higher precedence than '-'. Thus '\*' will be popped out from the stack and added to the Postfix string. Even now the stack is not empty. Now the topmost element of the stack is '+' which has equal priority to '-'. So pop the '+' from the stack and add it to the Postfix string. The '-' will be pushed to the stack.



Next character is 'd' which is added to Postfix string. Now all characters have been scanned so we must pop the remaining elements from the stack and add it to the Postfix string. At this stage we have only a '-' in the stack. It is popped out and added to the

Postfix string. So, after all characters are scanned, this is how the stack and Postfix string will be:



End result :

Infix String : a+b\*c-d

Postfix String : abc\*+d-

**Example program:**

```
#include <iostream.h>
#include <conio.h>
#include <string.h>
#include <ctype.h>
#include "CHRStack.h"
#define MAXSIZE 100
int IsOperator(char c){
    if (ch == '^' || ch == '*' || ch == '/' || ch == '+' || ch == '-')
        return 1;
    return 0;
}
void main(){
    Stack s;
    char c,ch;
    char str[MAXSIZE];
    int l,i;
```

```

int HP;

clrscr();

cout << "Enter the input expression in infix notation -> ";

cin.getline(str, MAXSIZE);

l = strlen(str);

str[l++] = ')';

str[l] = '\0';

s.Push('(');

cout << "The corresponding expression in postfix notation is -> ";

for (i = 0; i < l; i++){

    if (isalpha(str[i])){

        cout << str[i];

        continue;

    }

    if (str[i] == '{'){

        s.Push(str[i]);

        continue;

    }

    if (IsOperator(str[i])){

        HP = 1;

        while (HP){

            ch = s.GetTopElement();

            if (IsOperator(ch)){

                switch(str[i]){

                    case '^':

                        if (ch == '^'){

```

```

        cout << s.Pop();

        HP = 1;
    }
    else HP = 0;

    break;
case '*':
case '/':
    if (ch == '^' || ch == '*' || ch == '/') {
        cout << s.Pop();

        HP = 1;
    }
    else HP = 0;

    break;
case '+':
case '-':
    if (IsOperator(ch)) {
        cout << s.Pop();

        HP = 1;
    }
    else HP = 0;

    break;
}
}
else HP = 0;
}
s.Push(str[i]);

```

```

    }
    if (str[i] == '){
        while (1){
            if (IsOperator(s.GetTopElement())) cout << s.Pop();
            else break;
        }
        s.Pop();
    }
}
getch();
}

```

## 5.5 Evaluation of Postfix Expression

In normal algebra we use the infix notation like  $a+b*c$ . The corresponding postfix notation is  $abc*+$ . The algorithm for the conversion is as follows:

- Scan the Postfix string from left to right.
- Initialise an empty stack.
- If the scanned character is an operand, add it to the stack. If the scanned character is an operator, there will be atleast two operands in the stack.
- If the scanned character is an Operator, then we store the top most element of the stack(`topStack`) in a variable `temp`. Pop the stack. Now evaluate `topStack(Operator)temp`. Let the result of this operation be `retVal`. Pop the stack and Push `retVal` into the stack.

Repeat this step till all the characters are scanned.

- After all characters are scanned, we will have only one element in the stack. Return `topStack`.

### Example:

Let us see how the above algorithm will be implemented using an example.  
Postfix String:  $123*+4-$

Initially the Stack is empty. Now, the first three characters scanned are 1, 2 and 3, which are operands. Thus they will be pushed into the stack in that order.





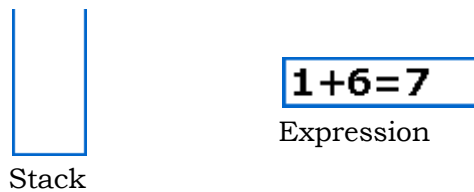
Next character scanned is "\*", which is an operator. Thus, we pop the top two elements from the stack and perform the "\*" operation with the two operands. The second operand will be the first element that is popped.



The value of the expression(2\*3) that has been evaluated(6) is pushed into the stack.



Next character scanned is "+", which is an operator. Thus, we pop the top two elements from the stack and perform the "+" operation with the two operands. The second operand will be the first element that is popped.



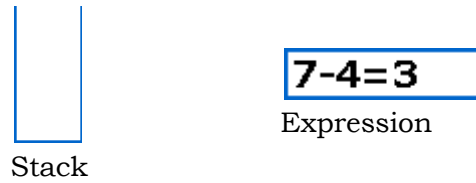
The value of the expression(1+6) that has been evaluated(7) is pushed into the stack.



Next character scanned is "4", which is added to the stack.



Next character scanned is "-", which is an operator. Thus, we pop the top two elements from the stack and perform the "-" operation with the two operands. The second operand will be the first element that is popped.



The value of the expression(7-4) that has been evaluated(3) is pushed into the stack.



Now, since all the characters are scanned, the remaining element in the stack (there will be only one element in the stack) will be returned.

**End result:**

Postfix String: 123\*+4-

Result: 3

**Example Program**

```
#include <iostream.h>
#include <conio.h>
#include <string.h>
#include <ctype.h>
#include "INTStack.h"
#define MAXSIZE 100
int IsOperator(char c){
```

```

if (ch == '^' || ch == '*' || ch == '/' || ch == '+' || ch == '-')
    return 1;
return 0;
}

void main(){
    Stack s;

    int i,l,t,a,b;

    char str[MAXSIZE];

    clrscr();

    cout << "Enter the expression in postfix (use comma as separator) -> ";
    cin.getline(str,MAXSIZE);

    l = strlen(str);

    for (i = 0;i < l;i++){
        if (isdigit(str[i])){
            s.Push(str[i++]-'0');

            while (isdigit(str[i])){ //For handling multi-digit numbers
                int t = s.Pop() * 10 + str[i] - '0';
                s.Push(t);
                i++;
            }
            continue;
        }
        if (IsOperator(str[i])) {
            b = s.Pop();
            a = s.Pop();
            switch(str[i]){

```

```

        case '^':    t = a^b;
                    break;
        case '*':    t = a*b;
                    break;
        case '/':    t = a/b;
                    break;
        case '+':    t = a+b;
                    break;
        case '-':    t = a-b;
                    break;
    }
    s.Push(t);
}
}
cout << "Result of expression is -> " << s.Pop();
getch();
}

```

## 5.6 Recursion

**Recursion is calling a procedure from itself, directly (simple recursion) or via calls to other procedures (mutual recursion).** The power of recursion evidently lies in the possibility of defining an infinite set of objects by a finite statement. In the same manner, an infinite number of computations can be described by a finite recursive program, even if this program contains no explicit repetitions.

Recursion is usually less computationally efficient, compiler generated procedure calls introduce overhead when storing all variables before a procedure call and loading them back upon entering the procedure.

Automatic optimizations reduce the overhead associated with procedure calls by eliminating redundant store/load operations to some degree but of course hand optimizing may improve upon that. Efficiency is usually not the most important consideration when

choosing an algorithm, in the name of good programming we already agreed to gladly suffer some measure of inefficiency.

## Factorial

A classic example of a recursive procedure is the function used to calculate the factorial of an integer.

Function definition:

$$fact(n) = \begin{cases} 1 & \text{if } n = 0 \\ n \times fact(n - 1) & \text{if } n > 0 \end{cases}$$

```
Pseudocode (recursive):  
  
function      factorial      is:  
input: integer n such that n >= 1  
output: [n × (n-1) × (n-2) × ... × 1]  
  
    1. if n is 0, return 1  
  
    2. otherwise, return [ n × factorial(n-1) ]  
  
end factorial
```

## Simulating Recursion

Knowing the process by which recursion passes data upward and downward through the called modules, you can isolate and preserve the variables unique to each recursive step and simply loop a given piece of code to achieve simulated recursion. Since looping would start the code process at the same place each time, you must keep a place marker, indicating just where to begin the processing of the current iteration.

The data used by each iterative step can be in any form and is a function of the programming requirements. A binary switch is probably the best form for a place marker. A combination of the required variable data and a place marker form a snapshot of the conditions during any step of the recursion. You should construct a table where snapshots can be stored. A pointer to this stack of snapshots allows the program to select the appropriate set of values for a particular iteration of the recursive process. The depth of the stack must be such that it can contain all the steps necessary to achieve the objective of the program being executed.

## Example Program

```
#include <iostream.h>
```

```

#include <conio.h>

#include <string.h>

#include "intstack.h"

void main(){

    Stack s;

    int n,t;

    clrscr();

    cout << "Enter the number for computing its factorial (less than 8) -> ";

    cin >> n;

    s.Push(1);

    while (n != 0){

        t = s.Pop() * n--;

        s.Push(t);

    }

    cout << "Factorial of n " << n << " is n is -> " << s.Pop();

    getch();

}

```

## 5.7 Summary

Stack has many applications, including parenthesis matching, infix to postfix conversion, postfix evaluation and simulation of recursion among others. For matching program delimiters stacks are extensively used.

## 5.8 Questions

1. What do you mean by infix and postfix expressions?
2. Convert the following infix expressions to postfix
  - a.  $a+b*c+d$
  - b.  $a*b+c*d$
3. Write the pseudo code for converting an infix expression to postfix.

4. Write the pseudo code for evaluating a postfix expression.
5. Write the pseudo code for simulating recursion using stacks for computing factorial.

### **5.9 Suggested readings**

1. A. Tanenbaum, Y. Lanhgsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley

## QUEUES

### 6.1 Objectives

### 6.2 Introduction

### 6.3 Operations on Queues

6.3.1 Inserting Element in Queue

6.3.2 Deleting Element from Queue

6.3.3 Checking Underflow and Overflow

### 6.4 Queues Using an Array in C++

6.4.1 Implementing IsFull

6.4.2 Implementing IsEmpty

6.4.3 Implementing Enqueue

6.4.4 Implementing Dequeue

### 6.5 Summary

### 6.6 Questions

### 6.7 Suggested readings

### 6.1 Objectives

You probably never thought that waiting in line in the supermarket would help you become a whiz at data structures but it's a big help. The checkout line at a supermarket is similar to the way data structures are organized. We're the "things" organized by the supermarket line and the same kind of organization is used for data within your program. The checkout line in your program is called a queue. In this lesson, you'll learn the ins and outs of implementing a queue within your program.



## 6.2 Introduction

**A queue is a container of objects (a linear collection) that are inserted and removed according to the FIFO principle.** A queue is like the checkout line at the supermarket where the first customer is at the front of the line, the second customer is next in line and so on until you reach the last customer who is at the back of the line. Customers check out of the supermarket in the order they arrive in the line. That is, the first customer is the first one to check out. This is referred to as first in, first out (FIFO).

The same concept applies to a queue in your program. A queue is a sequential organization of data. Data is accessible using FIFO. That is, the first data in the queue is the first data that is accessible by your program. In this lesson, you will explore the simplest type of queue, a fixed size, first in, first out queue using an array.

Programmers use one of two kinds of queues depending in the objective of the program, a simple queue or a priority queue. A simple queue organizes items in a line where the first item is at the beginning of the line and the last item is at the back of the line. Each item is processed in the order in which it appears in the queue. The first item in line is processed first, followed by the second item and then the third until the last item on the line is processed. There isn't any way for an item to cut the line and be processed out of order.

A priority queue is similar to a simple queue in that items are organized in a line and processed sequentially. However, items on a priority queue can jump to the front of the line if they have priority. Priority is a value that is associated with each item placed in the queue. The program processes the queue by scanning the queue for items with high priority. These are processed first regardless of their position in the line. All the other items are then processed sequentially after high priority items are processed.

Queues are very important in business applications that require items to be processed in the order they are received. The supermarket checkout line is a queue that most of us have experienced but you won't be creating a supermarket checkout line in a program unless the program is designed to simulate a checkout line.

In the real world, queues are used in programs that process transactions. A transaction is a set of information such as an order form. Transaction information is received by a program and then placed in a simple queue waiting to be processed by another part of the program.

Many other applications use a simple queue to maintain the order in which to process items. These include programs that process stock and bond trades and those that process students registering for a course. Queues are also used within a computer to manage printing.

## 6.3 Operations on Queues

Data organized by a queue may be stored in an array. The queue determines the array element that is at the front and back of the queue. The array is not the queue. Likewise, the queue is not the array. Both are two separate things. This is an important concept to grasp and one that may be difficult to understand at first.

Take a look at Figure 6.1 and you'll see how an array and a queue are different and yet are linked together to organize data. The array is pictured as a block of elements. The queue is pictured as a circle. The empty boxes are where values are stored in the queue and the numbers correspond to the index of the array that is associated with the queue. To the right of the circle are three values. The front and back values store the index of the front and back of the queue. The size value is the number of elements in the queue, which is 8 in this example.

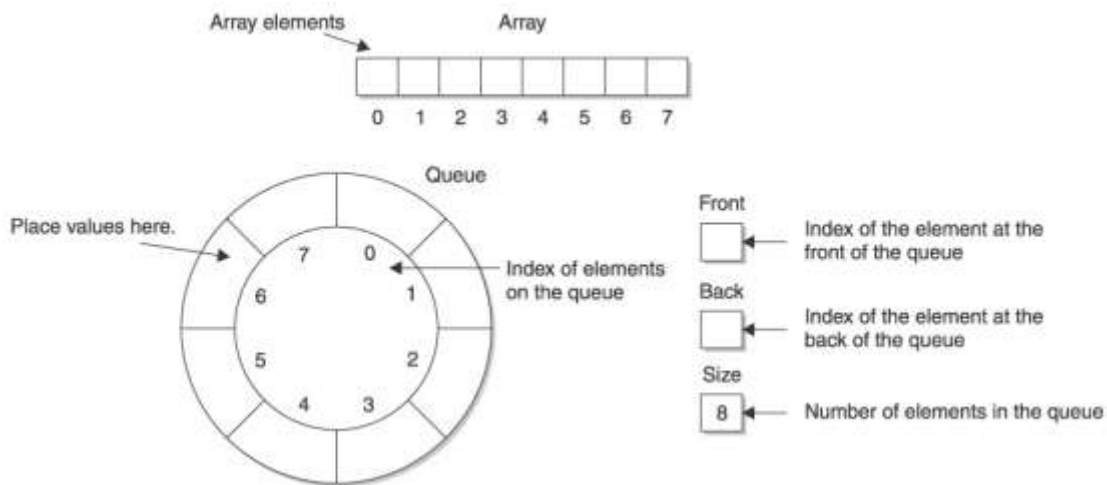


Figure 6.1: The queue is different from the array to store data that appears in the queue.

### 6.3.1 Inserting Element in Queue

A value is placed in the queue by performing the enqueue process, which consists of two steps. The first step is to identify the array element that is at the back of the queue. However, this is not necessarily the last element of the array. Remember that the queue is not the array. The back of the queue is calculated by using the following formula:

$$\text{back} = (\text{back} + 1) \% \text{size}$$

Figure 6.2 shows how to use the formula and gives the values for the front, back and size of the queue. The front and back variables are set to zero because the queue is empty and size is set to 8 because the array has 8 elements.

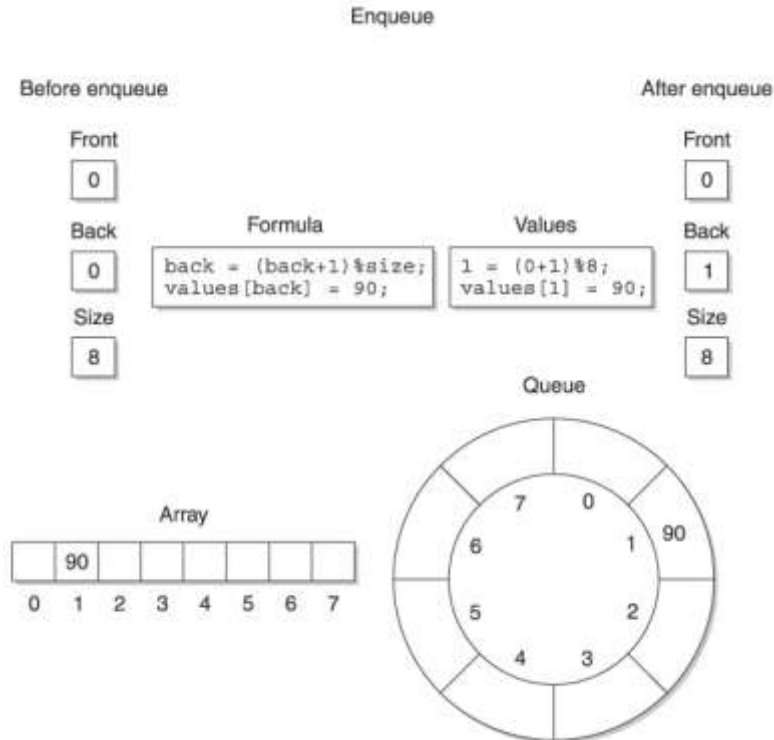


Figure 6.2: The enqueue process places a new value at the back of the queue.

The next box shows the formula that identifies the back of the queue and assigns it the value 90. To the right of this box is the same formula with variable names replaced by actual values. Let's take a closer look at this and see how the back of the queue is calculated.

The first operation occurs within the parentheses, where 1 is added to the value of the back variable. The modulus operator determines where the next element should be placed in the queue by performing an integer division and returning the remainder of the division. Although we've described a queue as a checkout line in the supermarket, a queue is actually circular. This is illustrated in the calculation used to determine the back of the queue, as shown here:

$$(7 + 1) \% 8$$

When you get to the last element in the array at index 7, the calculation returns 0 (8 divided by 8 is 1 and the remainder is 0). So after the last element in the array, you come around to the beginning of the array as the back of the queue. You check to see if you're at the front of the queue before placing an item at the back of the queue so you don't overwrite the item at the front and corrupt the queue.

The second step is to assign the value 90 to array element 1. That is, place the value 90 at the back of the queue. Remember that values are added to the queue from the

back just as you go to the back of the checkout line to wait your turn at the supermarket. Notice that the value 90 is assigned to the array in Figure 6.2.

### 6.3.2 Deleting Element from Queue

Dequeue is the process that removes a value from the front of the queue. It is important to understand that the value is removed from the queue, not the array. The value always remains assigned to the array until the value is either overwritten or the queue is abandoned.

There are two steps in the dequeue process, as illustrated in Figure 6.3.

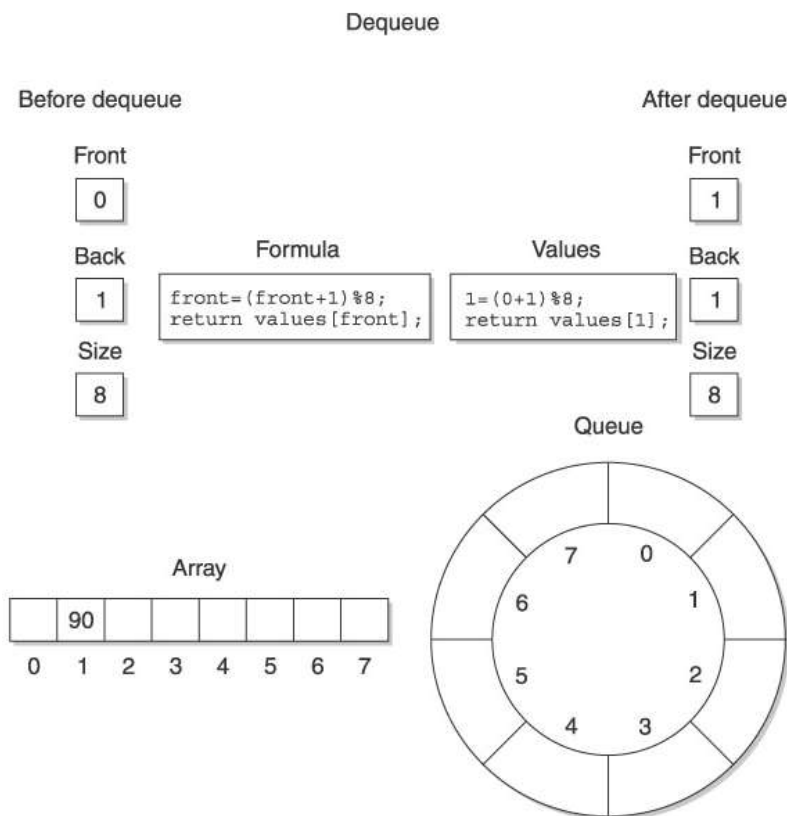


Figure 6.3: The dequeue process removes an item from the front of a queue.

The initial step is to calculate the index of the array element at the front of the queue using the following expression:

$$\text{front} = (\text{front} + 1) \% 8$$

Notice that this expression is very similar to the expression used in the enqueue process to calculate the index of the array element at the back of the queue. The first operation in this expression increments the value of the front variable. As you can see in Figure 6.3, the front variable is assigned the initial value zero. Therefore, the result of the

first operation is 1. The next operation is to apply the modulus operator, which is identical to the modulus operation performed in the enqueue process. The result of this operation is 1, meaning that the front of the queue is the array element whose index is 1. This value is then assigned to the front variable. Previously in this lesson, you learned that if you were at index 7 in the array, the result of this calculation would be 0  $((7+1)\%8 = 0)$ , so you would chase the queue around in a circle.

The final step in the dequeue process is to use the value located at the front of the queue. Typically, the dequeue process is a method and the front of the queue is returned to the statement that called the method.

In Figure 6.3, the array element `values[1]` is at the front of the queue. The value assigned to this element is 90, which was placed at the back of the queue by the previous enqueue process.

Notice that the value 90 remains assigned to the `values[1]` array element in Figure 6.3 because values assigned to the array associated with a queue are not affected when a value is removed from the front of the queue. The queue keeps track of array elements that are at the front and back of the queue, not the front or back of the array. In this case, we're using a simple integer array to illustrate the principles behind implementing the queue data structure. You may come across more complex implementations, where each element in the array is a pointer to a class object or structure. In these cases, you should be concerned about memory management when you perform enqueue and dequeue operations.

### **6.3.3 Checking Underflow and Overflow**

Deletion of queue element is not possible when the queue is empty. This is called underflow checking. Similarly an element can not be inserted when the queue is full. This is called overflow checking.

While implementing Enqueue operation, before actual insertion of the element in the queue, its overflow condition must be checked. And while implementing Dequeue operation, before actual deletion of queue element, its underflow condition must be checked.

## **6.4 Queues Using an Array in C++**

Now that you understand how queues work with an array, it is time to create a real queue. In this section, you'll create a queue using C++.

The C++ queue program is organized into three files: the `queue.h` file, the `queue.cpp` file and the `queueProgram.cpp` file. The `queue.h` file, shown next, sets the default size of the array and defines the `Queue` class. The `Queue` class declares `size`, `front` and `back` attributes that store the array size and the index of the front and back of the queue. The `Queue` class also declares a pointer that will point to the array. In addition to these, the

Queue class defines a set of member functions that manipulate the queues. These are explained later in this section.

```
//queue.h

#define DEFAULT_SIZE 8

class Queue{

private:

    const int size;

    int front;

    int back;

    int* values;

public:

    Queue(int size = DEFAULT_SIZE);

    virtual ~Queue();

    bool IsFull();

    bool IsEmpty();

    void enqueue(int);

    int dequeue();

};
```

The queue.cpp file contains the implementation of the member functions for the Queue class. There are six member functions defined in this file: Queue(), ~Queue(), IsFull(), IsEmpty(), enqueue() and dequeue().

The Queue() member function is a constructor, which is passed the size of the array when an instance of the Queue class is declared. If the constructor is called with no parameters, then the default size is used, otherwise, the value passed to the constructor is used. The value of the array size is assigned to the attribute size by the first statement within the constructor.

The second statement uses the new operator to declare an array of integers whose size is determined by the size passed to the constructor. The new operator returns a pointer to the array, which is assigned to the values pointer. The last two statements in the constructor initialize the front and back attributes to zero.

The `~Queue()` member function is the destructor and uses the delete operator to remove the array from memory when the instance of the Queue class goes out of scope.

### 6.4.1 Implementing IsFull

The `IsFull()` member function (see Figure 6.4) determines if there is room available in the queue by comparing the calculated value of the back of the queue with the value of the front of the queue, as in shown Figure 6.4. Notice that the expression that calculates the back of the queue is very similar to the expression in the enqueue process (see the “Enqueue” section of this lesson) and both produce the same result. The queue is full when the back index is 1 behind the front. Placing another element in the queue would overwrite the front element and corrupt the queue. The modulus operator is used again to make this a circular queue, so when you’re at element 7 on the back, the next element to look at is element 0.

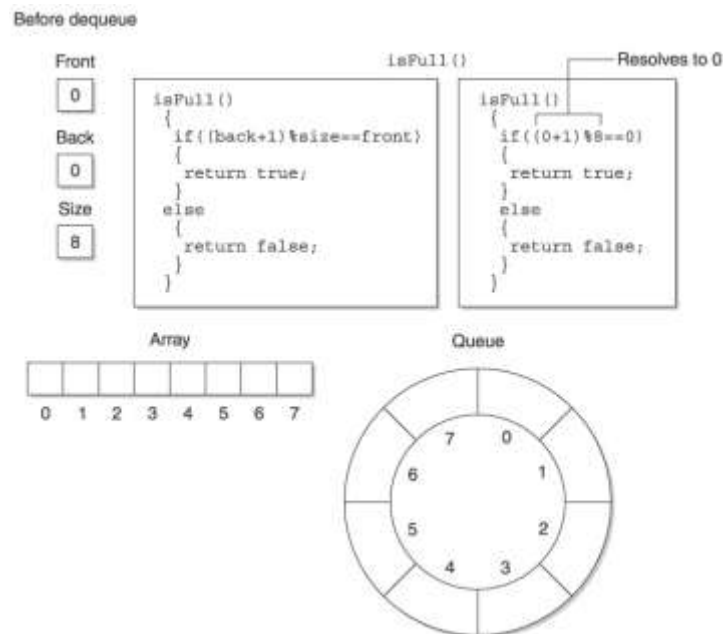


Figure 6.4: The `IsFull()` member function determines if there is room to place another item on the back of the queue.

The `IsFull()` member function is called by the `enqueue()` member function before an attempt is made to place a value on the back of the queue. The `IsFull()` member function returns a `true` if no more room is available in the queue or a `false` if there is room available.

### 6.4.2 Implementing IsEmpty

The `IsEmpty()` member function determines (see Figure 6.5) if the queue is empty by comparing the back and front variables. If they have the same values, a `true` is returned,

otherwise, a false is returned. The `IsEmpty()` member function is called within the `dequeue()` member function before it attempts to remove the front item from the queue.

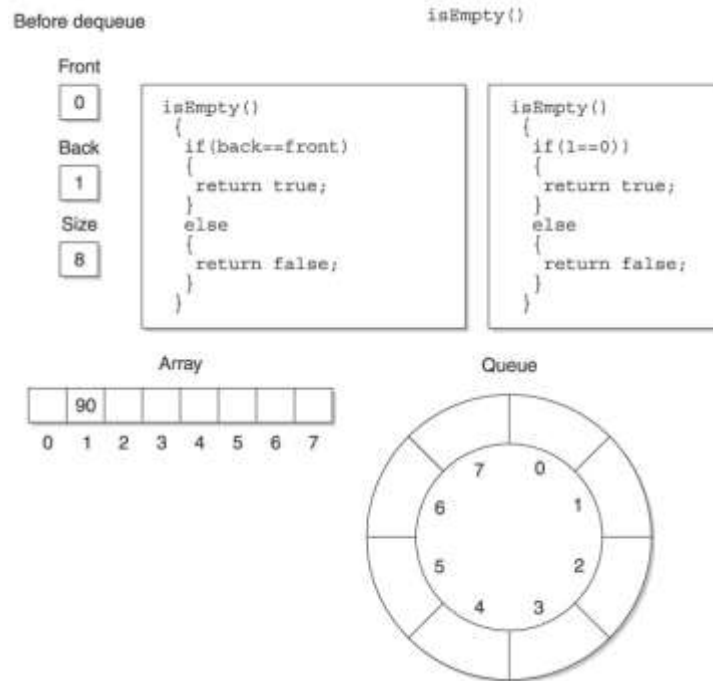


Figure 6.5: The `IsEmpty()` member function determines if the queue contains any values.

### 6.4.3 Implementing Enqueue

The `enqueue()` member function places an item at the back of the queue, as described in the “Enqueue” section of this lesson. The `enqueue()` member function is passed the value that is to be placed in the queue. However, before doing so, the `IsFull()` member function is called to determine if there is room available in the queue. Notice in the following example that the `IsFull()` member function is called as the condition expression of the `if` statement. Also notice that the `not` operator reverses the bool value returned by the `IsFull()` method. That is, a `false` is returned by the `IsFull()` member function if room is available in the queue. The condition expression in the `if` statement reverses this logic to `true` so that statements within the `if` statement execute to place the new item on the back of the queue.

### 6.4.4 Implementing Dequeue

The `dequeue()` member function removes an item from the queue and returns that item to the statement within the program that calls the `dequeue()` member function. However, the `IsEmpty()` member function is called in the condition expression of the `if` statement within the `dequeue()` member function, as shown in the next code listing.



The not operator in this expression reverses the logic returned by the `IsEmpty()` member function. The `IsEmpty()` member function returns a false if the queue is not empty. The not operator changes this to true, enabling statements within the if statement to remove the front item from the queue and return it to the statement that calls the `dequeue()` member function.

```
//queue.cpp
#include "queue.h"
Queue::Queue(int size){
    this->size = size;
    values = new int[size];
    front = 0;
    back = 0;
}
Queue::~Queue(){
    delete[] values;
}
bool Queue::IsFull(){
    if( (back+1) % size == front)
        return true;
    else
        return false;
}
bool Queue::IsEmpty(){
    if(back == front)
        return true;
    else
        return false;
}
```

```

void Queue::enqueue(int x){
    if(!IsFull()){
        back = (back+1) % size;
        values[back] = x;
    }
}

int Queue::dequeue(){
    if(!IsEmpty()){
        front = (front+1) % size;
        return queue[front];
    }
    return 0;
}

```

The `queueProgram.cpp` is where all the action takes place. It is here that an instance of the `Queue` class is declared and manipulated. As you can see in the next example, the first statement in the program uses the new operator to declare an instance of the `Queue` class and set the size to 8 elements. The new operator returns a pointer that is assigned to a pointer to an instance of the `Queue` class.

The next three statements call the `enqueue()` member function three times to place the values 10, 20 and 30 in the queue, respectively. The program concludes by calling the `dequeue()` member function three times to display the contents of the queue. Figure 6.6 shows the queue and the array after the last call to the `enqueue()` member function is made.

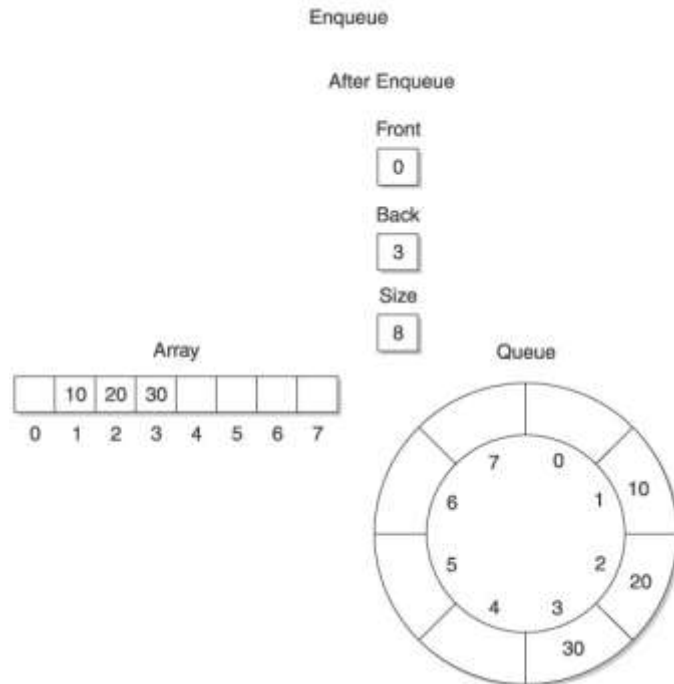


Figure 6.6: Here's the queue and the array after the last call to the enqueue() member function is made.

```
//queueProgram.cpp
#include <iostream>
using namespace std;
void main(){
    Queue *queue = new Queue(8);
    queue->enqueue(10);
    queue->enqueue(20);
    queue->enqueue(30);
    for(int i=0; i<3; i++){
        cout << queue->dequeue() << endl;
    }
}
```

## 6.5 Summary

Queue is a linear data structures extensively used in varied application of computer programming. It is based in First In First Out (FIFO) principle. A front pointer is used for deleting element from the beginning of the stack and a rear pointer is maintained for deleting elements from end of the stack. Enqueue operation is used to insert elements in the stack and dequeue operation is used for deleting an element from the stack. Underflow and overflow conditions are check using IsEmpty and IsFull functions, respectively.

## 6.6 Questions

1. What is a queue?
2. What is the relationship between a queue and its underlying array?
3. Explain how the index of the front and back of the queue is calculated.
4. What is the purpose of the enqueue process?
5. What is the purpose of the dequeue process?
6. Why is the IsFull() member method called?
7. Why is the IsEmpty() member method called?
8. What happens to the data stored on the array when the data is removed from the queue?
9. What is the purpose of setting the default size of the queue?
10. WAP to implement queue using Array.
11. WAP to implement queue using linked list.

## 6.7 Suggested readings

1. A. Tanenbaum, Y. Lanhsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley

## TYPES OF QUEUES

- 7.1 Objectives
- 7.2 Introduction
- 7.3 Circular Queue
- 7.4 Double Ended Queue (de-queues)
- 7.5 Priority Queue
- 7.6 Application of Queues
- 7.7 Summary
- 7.8 Questions
- 7.9 Suggested Readings

### 7.1 Objectives

In this lesson, we will discuss the various variants of queues including circular, double ended and priority queue. We shall also discuss the various applications of queues.

### 7.2 Introduction

There are different variants of queues for diverse areas of applications. These include circular queue, double ended queue also called dequeue or deque and priority queue as used in operating systems.

### 7.3 Circular Queue

A circular queue is a Queue but a particular implementation of a queue. It is very efficient. It is also quite useful in low level code, because insertion and deletion are totally independent, which means that you don't have to worry about an interrupt handler trying to do an insertion at the same time as your main code is doing a deletion.

Algorithm for Insertion:-

Step-1: If "rear" of the queue is pointing to the last position then go to step-2 or else step-3

Step-2: make the "rear" value as 0

Step-3: increment the "rear" value by one

Step-4:

1. if the "front" points where "rear" is pointing and the queue holds a not NULL value for it, then its a "queue overflow" state, so quit; else go to step-4.2
2. insert the new value for the queue position pointed by the "rear"

Algorithm for deletion:-

Step-1: If the queue is empty then say "empty queue" and quit; else continue

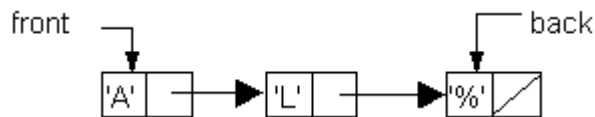
Step-2: Delete the "front" element

Step-3: If the "front" is pointing to the last position of the queue then step-4 else step-5

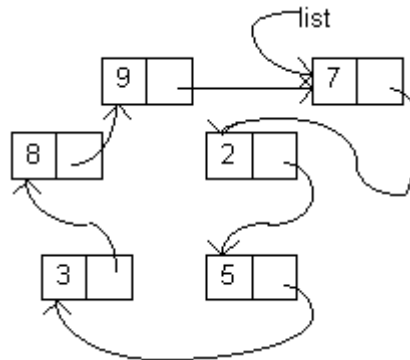
Step-4: Make the "front" point to the first position in the queue and quit

Step-5: Increment the "front" position by one

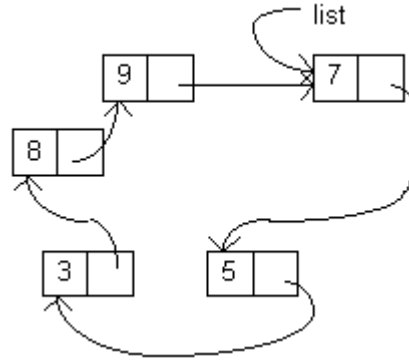
One of the nice applications of pointers is in implementing a circular queue. First let us consider the fact, that we only need one pointer and that is to the last element in the queue. In the example below, the node containing 2 is the head of the list and the node containing 9 is the rear of the list.



To enqueue a node containing 7 the queue would look like this:

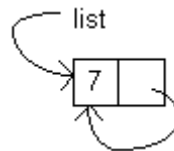


And then if we were to dequeue a node, we would dequeue the node containing the 2 and the result would look like this.



In this implementation, an empty list would be list = NULL.

A list containing one element would look like this. In coding enqueue, care must be taken with the first element.



To get rid of a queue or makeEmpty, the nodes must be disposed so that the memory is free again. Use the myqueue Class, a pointer implementation of a queue, to modify the code to be a circular queue. In your client program, to show that each function works, write a non-destructive function show that will display the contents of the queue.

A key advantage of a circular queue is that it has a static size and elements need not be shuffled around when a portion of the queue is used. This means that only new data is written to the queue and the computational cost is independent of the length of the queue.

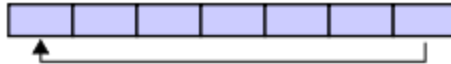
For example, in implementing a Transmission Control Protocol stack the windows used to hold the data can be circular queues. When the receiver acknowledges the reception of a packet then the amount of data acknowledged is used to advance the start of the queue. This allows the stack to restrict the amount of unacknowledged data by the transmitter. (This is an example of when it is not permitted for the start of the queue to be overwritten.)

An example that could possibly use an overwriting circular queue is with multimedia. If the queue is used as the bounded queue in the producer-consumer problem then it is probably desired for the producer (e.g., an audio generator) to overwrite old data if the consumer (e.g., the sound card) is unable to momentarily keep up. Another example is the digital waveguide synthesis method which uses circular queues to efficiently simulate the sound of vibrating strings or wind instruments.

The "prized" attribute of a circular queue is that it does not need to have its elements shuffled around when one is consumed. (If a non-circular queue were used then

it would be necessary to shift all elements when one is consumed.) In other words, the circular queue is well suited as a FIFO queue while a standard.

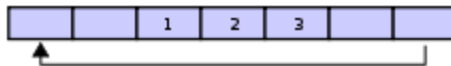
A circular queue first starts empty and of some predefined length. For example, this is a 7-element queue:



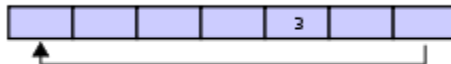
Assume that a 1 is written into the middle of the queue (exact starting location does not matter in a circular queue):



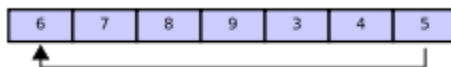
Then assume that two more elements are added — 2 & 3 — which get appended after the 1:



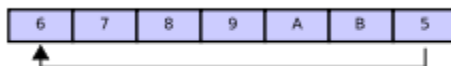
If two elements are then removed from the queue then they come from the end. The two elements removed, in this case, are 1 & 2 leaving the queue with just a 3:



If the queue has 7 elements then it is completely full:



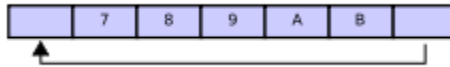
A consequence of the circular queue is that when it is full then a subsequent write is performed then it starts overwriting the oldest data. In this case, two more elements — A & B — are added and they overwrite then 3 & 4:



Alternatively, the routines that manage the queue could easily not allow data to be overwritten and return an error or raise an exception. Whether or not data is overwritten is up to the semantics of the queue routines or the application using the circular queue.



Finally, if after overwriting elements two elements are removed then what would be returned is not 3 & 4 but 5 & 6 because A & B overwrote the 3 & the 4 yielding the queue with:



### Circular queue mechanics

What is not shown in the example above is the mechanics of how the circular queue is managed.

### Start / End Pointers

Generally, a circular queue requires three pointers:

- one to the actual queue in memory
- one to point to the start of valid data
- one to point to the end of valid data

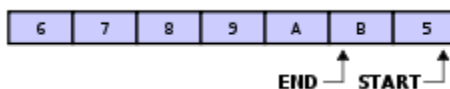
Alternatively, a fixed-length queue with two integers to keep track of indices can be used in languages that do not have pointers.

Taking a couple of examples from above. (While there are numerous ways to label the pointers and exact semantics can vary, this is one way to do it.)

This image shows a partially-full queue:



This image shows a full queue with two elements having been overwritten:

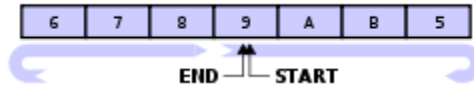


What to note about the second one is that after each element is overwritten then the start pointer is incremented as well.

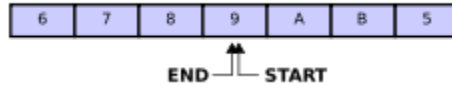
### Difficulties

### Full / Empty Queue Distinction

Some small disadvantage of relying on pointers or relative indices of the start and end of data is, that in the case the queue is entirely full, both pointers pointing at the same element:



This is exactly the same situation as when the queue is empty:



To solve this problem there are a number of solutions:

- Always keep one byte open.
- Use a fill count to distinguish the two cases.
- Use read and write counts to get the fill count from.
- Use absolute indices.

### **Always Keep One Byte Open**

This simple solution always keeps one byte unallocated. A full queue has at most (size - 1) bytes. If both pointers are pointing at the same location, the queue is empty.

The advantages are:

- Very simple and robust.
- You need only the two pointers.

The disadvantages are:

- You can never use the entire queue.
- If you cannot read over the queue border, you get a lot of situations where you can only read one element at once.

### **Use a Fill Count**

The second simplest solution is to use a fill count. The fill count is implemented as an additional variable which keeps the number of readable bytes in the queue. This variable has to be increased if the write (end) pointer is moved and to be decreased if the read (start) pointer is moved. In the situation if both pointers pointing at the same location, you consider the fill count to distinguish if the queue is empty or full.

The advantages are:

- Simple.

- Needs only one additional variable.

The disadvantage is:

- You need to keep track of a third variable. This can require complex logic, especially if you are working with different threads.

Alternately, you can replace the second pointer with the fill count and generate the second pointer as required by incrementing the first pointer by the fill count.

The advantages are:

- Simple.
- No additional variables.

The disadvantage is:

- Additional overhead when generating the write pointer.

### **Read / Write Counts**

Another solution is to keep counts of the number of items written to and read from the circular queue. Both counts are stored in unsigned integer variables with numerical limits larger than the number of items that can be stored and are allowed to wrap freely from their limit back to zero.

The unsigned difference ( $\text{write\_count} - \text{read\_count}$ ) always yields the number of items placed in the queue and not yet retrieved. This can indicate that the queue is empty, partially full, completely full (without waste of a storage location) or in a state of overrun.

The advantage is:

- The source and sink of data can implement independent policies for dealing with a full queue and overrun while adhering to the rule that only the source of data modifies the write count and only the sink of data modifies the read count. This can result in elegant and robust circular queue implementations even in multi-threaded environments.

The disadvantage is:

- You need two additional variables.

### **7.4 Double Ended Queue (de-queues)**

This differs from the queue abstract data type or First-In-First-Out List (FIFO), where elements can only be added to one end and removed from the other. This general data class has some possible sub-types:

- An input-restricted deque is one where deletion can be made from both ends, but input can only be made at one end.
- An output-restricted deque is one where input can be made at both ends, but output can be made from one end only.

Both the basic and most common list types in computing, the queues and stacks can be considered specializations of deques, and can be implemented using deques.

A deque (short for double-ended queue—usually pronounced deck) is an abstract list type data structure also called a head-tail linked list, for which elements can be added to or removed from the front (head) or back (tail).

Deque is sometimes written dequeue, but this use is generally deprecated in technical literature or technical writing because dequeue is also a verb meaning "to remove from a queue". Nevertheless, several libraries and some writers, such as Aho, Hopcroft, and Ullman in their textbook *Data Structures and Algorithms*, spell it dequeue. DEQ and DQ are also used.

This differs from the queue abstract data type or First-In-First-Out List (FIFO), where elements can only be added to one end and removed from the other. This general data class has some possible sub-types:

- An input-restricted deque is one where deletion can be made from both ends, but input can only be made at one end.
- An output-restricted deque is one where input can be made at both ends, but output can be made from one end only.

Both the basic and most common list types in computing, the queues and stacks can be considered specializations of deques, and can be implemented using deques.

## Operations

The following operations are possible on a deque:

operation	C++
insert element at back	push_back
insert element at front	push_front
remove last element	pop_back
remove first element	pop_front
examine last element	back
examine first element	front

## Implementations

There are at least two common ways to efficiently implement a deque: with a modified dynamic array or with a doubly-linked list. For information about doubly-linked lists, see the linked list article.

### **Dynamic array implementation**

Deques are often implemented using a variant of a dynamic array that can grow from both ends, sometimes called array deques. These array deques have all the properties of a dynamic array, such as constant time random access, good locality of reference, and inefficient insertion/removal in the middle, with the addition of amortized constant time insertion/removal at both ends, instead of just one end. Two common implementations include:

- Storing deque contents in a circular queue, and only resizing when the queue becomes completely full. This decreases the frequency of resizings, but requires an expensive branch instruction for indexing.
- Allocating deque contents from the center of the underlying array, and resizing the underlying array when either end is reached. This approach may require more frequent resizings and waste more space, particularly when elements are only inserted at one end.

### **Complexity**

- In a doubly-linked list implementation, the time complexity of all operations is  $O(1)$ , except for accessing an element in the middle using only its index, which is  $O(n)$ .
- In a growing array, the amortized time complexity of all operations is  $O(1)$ , except for removals and inserts into the middle of the array, which are  $O(n)$ .

## **7.5 Priority Queue**

**A priority queue is similar to a simple queue in that items are organized in a line and processed sequentially. However, items on a priority queue can jump to the front of the line if they have priority.** Priority is a value that is associated with each item placed in the queue. The program processes the queue by scanning the queue for items with high priority. These are processed first regardless of their position in the line. All the other items are then processed sequentially after high priority items are processed.

Unlike a standard queue where items are ordered in terms of who arrived first, a priority queue determines the order of its items by using a form of custom comparer to see which item has the highest priority. Other than the items in a priority queue being ordered by priority it remains the same as a normal queue, you can only access the item at the front of the queue.

A sensible implementation of a priority queue is to use a heap data structure. Using a heap we can look at the first item in the queue by simply returning the item at index 0 within the heap array. A heap provides us with the ability to construct a priority queue

where the items with the highest priority are either those with the smallest value or those with the largest.

## 7.6 Application of Queues

An operating system often maintains a FIFO queue of processes that are ready to execute or that are waiting for a particular event to occur. The programmer who creates the operating system can use a Queue ADT to implement this.

Computer systems must often provide a “holding area” for messages between two processes, two programs or even two systems. This holding area is usually called a “buffer” and is often implemented as a FIFO queue. For example, if a large number of mail messages arrive at a mail server at about the same time, the messages are held in a buffer until the mail server can get around to processing them. It processes them in the order they arrived—in “first in, first out” order. (Some mail servers may be designed to handle the messages based on a priority system. In that case, the priority queue would be a more appropriate data structure.)

To demonstrate the use of queues, we look at a simpler problem i.e. identifying palindromes. A palindrome is a string that reads the same forwards as backwards. While we are not sure of their general usefulness, identifying them provides us with a good example for the use of both queues and stacks. Besides, palindromes can be entertaining.

Some famous palindromes are:

- A tribute to Teddy Roosevelt, who orchestrated the creation of the Panama Canal: “A man, a plan, a canal—Panama!”
- Allegedly muttered by Napoleon Bonaparte upon his exile to the island of Elba (although this is hard to believe since Napoleon mostly spoke French!): “Able was I ere, I saw Elba.”
- Overheard in a Chinese restaurant: “Won ton? Not now!”
- And possibly the world’s first palindrome: “Madam, I’m Adam.”
- Followed immediately by one of the world’s shorted palindromes: “Eve.”

As you can see, the rules for what is a palindrome are somewhat lenient. Typically, we do not worry about punctuation, spaces or matching the case of letters. We again follow the input/output model we have established for our test drivers—the same model used for the balanced parentheses example in the section on stacks.

An input file holds a separate string on each line. A corresponding output file is created, repeating each of the input lines and stating whether or not it is a palindrome.

This program assumes that no input line is more than 180 characters in length. If an input line is longer than that it is skipped. The names of the input and output files are

passed to the program on the command line. Summary statistics are written to an output frame.

The program reads a line of input and checks to see how long it is. If it is too long, it moves on to the next line of input. Otherwise, it creates a new stack and a new queue and it repeatedly pushes each letter from the input line onto the stack and also enqueues it onto the queue. To simplify comparison later, the actual characters pushed and enqueued are the lowercase versions of the characters in the string. When all of the characters of the line have been processed, the program repeatedly pops a letter from the stack and dequeues a letter from the queue. As long as these letters match each other for the entire way through this process, we have a palindrome. Can you see why? Since the queue is a “first in, first out” structure, the letters are returned from the queue in the same order they appear in the string. But the letters taken from the stack, a “last in, first out” structure, are returned in the opposite order from the way they appear in the string. So, we are comparing the letters from the forward view of the string to the letters from the backward view of the string.

## **7.7 Summary**

There are different variants of queues for diverse areas of applications. These include circular queue, double ended queue also called dequeue or deque and priority queue as used in operating systems. A circular queue is a Queue but a particular implementation of a queue. It is very efficient.

A double ended queue or dequeue may be:

- An input restricted where deletion can be made from both ends, but input can only be made at one end.
- An output restricted where input can be made at both ends but output can be made from one end only.

A priority queue is similar to a simple queue in that items are organized in a line and processed sequentially.

Queues are used by operating system for process management, by computer for buffering etc.

## **7.8 Questions**

1. What are the various types of queues?
2. Define circular queue. What are its properties? Explain.
3. Write the procedure of implementing circular queue.
4. Discuss in detail the concept of double ended queue.
5. What are priority queues? How these are implemented? Explain.

6. WAP to implement circular queue using Array.

### **7.9 Suggested Readings**

1. A. Tanenbaum, Y. Lanhsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley.



---

**LINKED LIST****8.1 Objectives****8.2 Introduction****8.3 Structure of a Linked List****8.4 Operations on Linked List**

8.4.1 The Linked List Class

8.4.2 Creation of a Linked List

8.4.3 Insertion in a Linked List

8.4.4 Displaying elements of a Linked List

8.4.5 Destroying the Linked List

**8.5 Implementing Linked List in C++****8.6 Summary****8.7 Questions****8.8 Suggested readings****8.1 Objectives**

In this lesson, we shall discuss meaning of linked list, creation of linked list and various operations like insertion, deletion, traversal and displaying of linked list.

**8.2 Introduction**

**A linked list is a data structure that makes it easy to rearrange data without having to move data in memory.** Sound a little confusing? If so, picture a classroom of students who are seated in no particular order. A unique number identifies each seat, as shown in Figure 8.1. We've also included the relative height of each student, which we'll use in the next exercise.

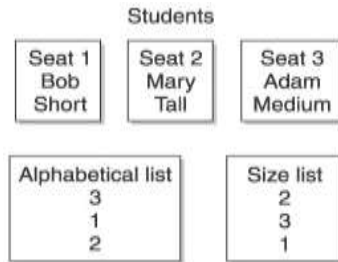


Figure 8.1: Students are seated in a classroom in random order.

Let's say that a teacher needs to place students names in alphabetical order so she can easily find a name on the list. One option is to have students change their seats so that Adam sits in seat 1, Bob sits in seat 2 and Mary in seat 3. However, this can be chaotic if there are a lot of students in the class.

Another option is to leave students seated and make a list of seat numbers that corresponds to the alphabetical order of students. The list would look something like this: 3, 1 and 2 as shown in Figure 8.1. The student in seat 3 is the first student who appears in alphabetical order followed by the student seated in seat 1 and so on. Notice how this option doesn't disrupt the class.

Suppose you want to rearrange students in size order. There's a pretty good chance that you won't move students about the classroom. Instead, you'd probably create another list of seat numbers that reflect each student's height. Here's the list: 1, 3 and 2 which is illustrated in Figure 8.1. The list can be read from bottom to top for the shortest to tallest or vice versa for tallest to shortest.

Once the list is created, the teacher can simply go down the list to see which seat contains the next student. To quiz students in alphabetical order, the teacher would use the alphabetical list to see that the student sitting in seat 3 is alphabetically first, followed by seat 1. The teacher can be tricky and call on the previous student by looking at the list to determine the student's seat.

Programmers call this sort of list a linked list because each item on the list is linked to the previous item and the next item. That is, the seat of the current student is linked to the seat of the previous student and to the seat of the next student by the list.

It is very important to keep the real world in mind as you learn how to use a linked list, otherwise, you'll fall into the trap of thinking that a linked list is an abstract concept that has little use in the real world. Actually, linked lists play a critical role in applications that help companies and governments manage data dynamically.

There are two versions of a linked list, a single link and a double link. A single link list enables a program to move through the list in one direction, which is usually from the front of the list moving to the end of the list. A doubly linked list enables the program to move through the list in both directions. We'll focus on the doubly linked list for most of

the examples in this lesson and then discuss the single link list toward the end of the lesson.

Although we've mentioned that an entry in a linked list contains data and pointers to the previous and next entries in the list, this is an over simplification. The data we're talking about is typically a set of data such as customer information. Customer information could be a customer ID, customer first name, customer last name, customer street address, customer city, customer state, customer ZIP and so on. Programmers call this a record. This means that an entry in a linked list may contain several data elements. In our example, however, we'll store only a single value of an integer so that we can focus on the principle of how a linked list works. In reality, you can add as many additional attributes to each node as you need.

Programmers choose linked lists over an array because linked lists can grow or shrink in size during runtime. Another entry can be placed at the end of the last entry on the linked list simply by assigning reference to the new entry in the last entry on the linked list. Likewise, the last entry can be removed from the linked list by simply removing reference to it in the next element of the second-to-last entry on the linked list. This is more efficient than using an array and resizing at runtime.

If you change the size of the array, the operating system tries to increase the array by using memory alongside the array. If this location is unavailable, then the operating system finds another location large enough to hold elements of the array and new array elements. Elements of the array are then copied to the new location.

If you change the size of a linked list, the operating system changes references to the next item on the list, which is fewer steps than changing the size of an array.

### **8.3 Structure of a Linked List**

Each entry in a linked list is called a node. Think of a node as an entry that has two subentries. One subentry contains the data, which may be one attribute or many attributes and the other points to the next node. When you enter a new item on a linked list, you allocate the new node and then set the pointer to next nodes.

Programmers create a node in C++ by using either a structure or a class object; our example uses a structure. As you'll recall from your C++ programming course, a structure is a user-defined data type. The following example is a structure used to define a node. Figure 8.2 shows a node.

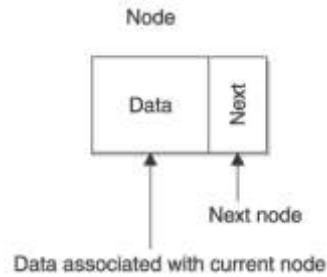


Figure 8.2: A node contains reference to the next node and the previous node in the linked list and contains data that is associated with the current node.

```

struct Node{
    int data;
    Node *next;
}*List;

```

The structure might look a bit strange even if you are familiar with structures because this example uses a pointer to the structure itself as one of its attribute. We'll clear up any confusion by taking apart this example. The structure is called Node. The name of the structure creates an instance of the structure similar to how you use a constructor to create an instance of a class and data type.

The first statement of the structure declares an integer that stores the current data of the node. The next statement declares pointers to the next node in the linked list.

#### 8.4 Operations on Linked List

The various operations that can be performed on linked list are:

- i. Creation (CreateList())
- ii. Insertion of a node
  - a. Insertion at beginning (InsertAtBeg())
  - b. Insertion at end (InsertAtEnd())
  - c. Insertion at a position in the list (InsertAtPos())
  - d. Insertion after a specific value (InsertAfterVal())
- iii. Deletion of a node
  - a. Deletion from beginning (DeleteFromBeg())
  - b. Deletion from end (DeleteFromEnd())

- iv. Traversal of linked list or Displaying elements of a linked list (Disp())
- v. Reverse nodes of the list (Reverse())
- vi. Destroy list (Destroy())

#### 8.4.1 The Linked List Class

Programmers use a `LinkedList` class to create and manage a linked list. C++ programmers define their own `LinkedList` class. For now, we'll focus on defining a `LinkedList` class in C++.

The `LinkedList` class definition consists of two data members and eleven function members, as shown in the example in this section. The two data members are: **data** which stores the actual data of the node and the pointer, **next** which references the next node on the linked list if it exists or `NULL`.

The twelve member functions manipulate the linked list. The first member function is the constructor of the `LinkedList` class and is called when an instance of the class is declared. Following the constructor is the destructor. When you return memory to the operating system by using the delete operator, the destructor is called. If you don't call the delete operator, then the destructor never gets called and the application causes a memory leak.

- The `CreateList(int val)` member function create a new list with the value of the first node supplied as argument to the function.
- The `InsertAtBeg(int val)` function inserts a new node in the beginning of the list making it the first node of the list.
- The `InsertAtEnd(int val)` function inserts a new node at the end of the list.
- The `InsertAtPos(int pos, int val)` function inserts a new node at the specified position in the list.
- The `InsertAfterVal(int aval, int val)` function inserts a new node after a specific value already existing in the list.
- The `int DeleteFromBeg()` function deletes the first node of the list.
- The `int DeleteFromEnd()` function deletes the last node of the list.
- The `void Disp()` functions displays the values of all nodes of the list.
- The `void Destroy()` function destroys the list.

The `LinkedList` class specification is defined in the header file, and the implementation is defined in the source file. We'll take a closer look at the implementation of these member functions in the next few sections of this lesson.

```
class LinkedList{  
    struct Node{  
        int data;
```

```

        Node *next;

}*List;

public:

    LinkedList(){NODE = NULL;}

    ~LinkedList(){Destroy();}

    void CreateList(int val);

    void InsertAtBeg(int val);

    void InsertAtEnd(int val);

    void InsertAtPos(int pos, int val);

    void InsertAfterVal(int aval, int val);

    int DeleteFromBeg();

    int DeleteFromEnd();

    void Reverse();

    void Disp();

    void Destroy();

};

```

#### 8.4.2 Creation of a Linked List

The LinkedList constructor is a member function that is called when an instance of the LinkedList is declared. The purpose of the constructor in the linked list example is to initialize the List pointer also known as header node, as shown in the following definition. The List pointer is assigned a NULL value, which is used by member functions to determine if the linked list is empty. You'll see how this is done in the next section.

```

LinkedList(){
    List = NULL;
}

```

The destructor is a member function called when the instance of the LinkedList class is deleted using the delete operator. In the example shown next, the destructor contains one statement that calls the Destroy() member function.

The Destroy() member function deletes the contents of the linked list but does not delete the linked list itself. That is, it removes all the nodes from the linked list. The Destroy() also resets the List pointer to NULL, signifying the linked list is empty of nodes.

The destructor is responsible for deallocating all the memory that was allocated for the linked list. In this case, it would be all the nodes.

You might be wondering why we defined two member functions to perform basically the same task. We do so to enable the programmer to empty the linked list. This way you can reset the contents of the linked list without destroying the instance of the LinkedList class.

```
~LinkedList(){
    Destroy();
}
```

### 8.4.3 Insertion in a Linked List

Insertion in a linked list can be made by different ways. We shall consider four methods of inserting a node in the linked list.

- i. Insert at beginning:** The InsertAtBeg() member function places a new node at beginning of the list. If the list is empty the next pointer the newly inserted node points to NULL else the next node points to node of the list which was first node before insertion.

```
void LinkedList::InsertAtBeg(int val)
{
    Node *p;
    p = new Node;    //Create a new node
    p->data = val;
    p->next = List; //Make it first node of the list
    List = p;      //Point List to the new node
}
```

An argument val is passed to the function which is the data for the new node to be inserted. The function, first, creates a new node and points the next pointer of the new node to the beginning of the List and then makes it the first node of the list.

- ii. Insert at end:** The InsertAtEnd() member function insert a node at the end of the linked list. There are several steps that must be performed in order to add the node to the end of the list. These are shown in the following definition of the member function:

```
void LinkedList::InsertAtEnd(int val)
{
```

```

Node *p, *Start;

Start = List;

p = new Node;

p->data = val;

p->next = NULL;

if (Start == NULL)

    List = p;

else

{

    while (Start->next != NULL)

        Start = Start->next;

    Start->next = p;

}

}

```

The `InsertAtEnd()` member function requires one argument called `val`, which is the current data for the node. The first statement in the `InsertAtEnd()` member function declares two instances of the `List` structure, one (`p`) for creating the new node and the other (`Start`) for moving to the end of the `List`.

Once the new node is created, the `InsertAtEnd()` member function positions the new node in the linked list. First, it determines if the linked list is empty by comparing the `List` pointer to `NULL`. As you'll recall, the `List` pointer is assigned a `NULL` when an instance of the `LinkedList` class is declared and when the `Destroy()` member function removes all the nodes from the list. If the linked list is empty (checked by `if` condition), then the new node is assigned to the `List` pointer. This means that the linked list contains one node after the `InsertAtEnd()` member function is called, which is the new node. However, if there is at least one node on the linked list, then a little shifting of pointers must be performed. The `while` statement performs the required shifting and the new node is then assigned to the `next` pointer, making the new node the last node on the linked list. This can be a little confusing, so take a look at Figure 8.3. Figure 8.3 shows nodes of the linked list. Assume that the linked list has two nodes before the new node is appended to the list. This is represented in the top block.



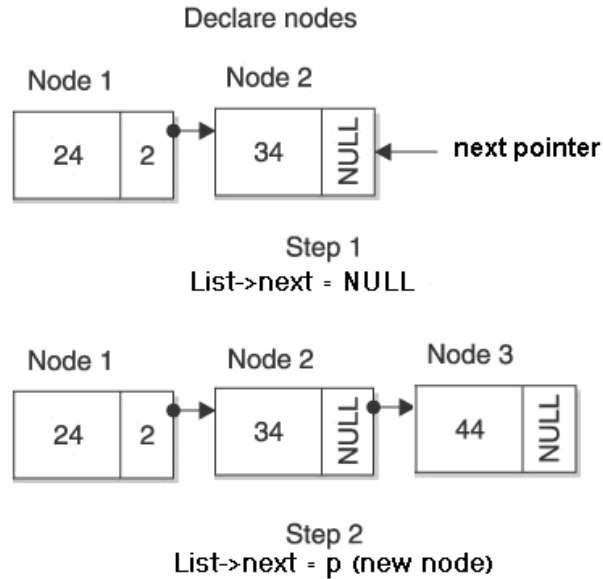


Figure 8.3: The appendNode() member function changes what nodes are pointed to in the linked list.

The first step moves the pointer to the end of the list, which is shown in the second block of memory in Figure 8.3. The second step assigns the memory address of the new node to the next member of the last node. This links both nodes.

**iii. Insert at position:** The InsertAtPos() function insert an element at the specified position in the list. If the specified position is greater than the number of nodes in the list then the node is inserted at the end.

```
void List::InsertAtPos(int pos, int val)
{
    int count = 1;

    Node *p, *q = List;

    p = new Node;
    p->data = val;

    while (q->next != NULL && count < pos-1)
    {
        count++;
        q = q->next;
    }
}
```

```

    }

    p->next = q->next;

    q->next = p;

}

```

The pointer is first moved to the specified position using the while loop. If the number of nodes in the list is less than the specified position then the node is inserted at the end of the List.

**iv. Insert after value:** The InsertAfterVal() function inserts a new node after a node having the specified value.

```

void List::InsertAfterVal(int aval, int val)
{
    Node *p, *q = List;

    p = new Node;

    p->data = val;

    while (q != NULL && q->data != aval)
        q = q->next;

    p->next = q->next;

    q->next = p;

}

```

The function is just like the previous function, the only difference is that the node is inserted after a specified value instead of position. If the specified value does not exist in the list then the new node is inserted at the end.

### 8.4.3 Deletion from beginning of the List

We have specified two different function for deleting a node from the List.

**i. Delete node from beginning:** The function DeleteFromBeg() deletes the first node of the List and the second node becomes the first node if it exists otherwise the list becomes empty.

```

int List::DeleteFromBeg()
{

    int val;

```

```

Node *p = List;

List = List->next;

val = p->data;

delete p;

return val;

}

```

- ii. Deletion from end of the list:** This is an interesting situation. You will never cut the stem of the tree while sitting on it or break the stair while standing on it. Similarly for deleting a node from the end of the list we need to reach the last but one node of the list. The function DeleteFromEnd() does the same.

```

int List::DeleteFromEnd()

{

    int val;

    Node *p,*q = List;

    while (q->next->next != NULL)

        q = q->next;

    p = q->next;

    q->next = NULL;

    val = p->data;

    delete p;

    return val;

}

```

The above function reaches the last but one node of the list and from there deletes the last node.

#### 8.4.5 Displaying elements of a Linked List

The Disp() member function displays each node of the linked list, beginning with the node at the beginning of the list and ending with the node at the end of the list. This is shown in the next example:

```
void List::Disp()
{
    Node *p = List;
    while (p != NULL)
    {
        cout << p->data << "\t";
        p = p->next;
    }
    cout << endl;
}
```

The Disp() member function begins by declaring a pointer p to the List. Before attempting to display data assigned to the node, Disp() determines if there is a node at the next position of the linked list. It does so by determining if the node pointed to by the p pointer is NULL. If so, the linked list is empty and there is nothing to display. If not, the member function proceeds and displays the data assigned to the node of the linked list.

The Disp() member function uses the node's next member to assign the pointer to the next node to the p pointer. The process continues by first determining if the node isn't NULL before displaying the data assigned to the node.

This process ends after the node at the end of the linked list is displayed because the next member of the node at the end of the list is NULL.

#### **8.4.5 Destroying the Linked List**

The Destroy() member function removes nodes from the linked list without removing the linked list itself, as shown in the following example. Each node is declared dynamically using the new operator, as you learned previously in this lesson. This enables you to remove the node by using the delete operator.

```
void List::Destroy()
{
    Node *p;
```

```

        while (List != NULL)
        {
            p = List;
            List = List->next;
            delete p;
        }
        List = NULL;
    }

```

The Destroy() member function begins by declaring a temporary pointer that is assigned the pointer to the node that is to be deleted and then the List points to the next node. However, before the node is removed, the member function determines if there is a node at the current position of the linked list by testing whether the List pointer is NULL. If so, then the Destroy() member function assumes there are no nodes on the linked list. If the List pointer isn't NULL, then the member function proceeds to delete the node.

This process continues until all the nodes are removed from the linked list. The final step in the destroyList() member function is to assign NULL values to the front and back of the linked list, which indicates that the linked list is empty of any nodes.

## 8.5 Implementing Linked Lists in C++

Now that you know the parts of a linked list and how to create and manipulate the linked list using a class, we'll put those parts together and create a real-world C++ application that uses a linked list.

Professional programmers organized a linked list C++ application into three files. The first file is the header file that contains the definition of the NODE structure and the LinkedList class definition. The second file is a source code file containing the implementation of member functions of the LinkedList class. The last file is the application file that contains code that creates and uses the LinkedList class.

Let's begin with the header file. LinkedList.h, shown in the code in this section, is the header file for the C++ linked list example. This file contains the definition of the LinkedList class, which programmers called a class specification.

You'll notice that the LinkedList class definition does not contain the implementation of member functions. Instead, it contains prototypes of member functions that are implemented in the source file. Keeping the specifications and implementation in separate header and source files is common practice. Parts of the program that use the class only care about the interface functions defined in the header file; they don't care

about the implementation. This also allows you to precompile your source code into library modules so the users of this class need only the headers and modules.

```
//LinkedList.h

#include <iostream.h>

class LinkedList{

    struct Node{

        int data;

        Node *next;

    }*List;

public:

    LinkedList(){List = NULL;}

    ~LinkedList(){Destroy();}

    void CreateList(int val);

    void InsertAtBeg(int val);

    void InsertAtEnd(int val);

    void InsertAtPos(int pos, int val);

    void InsertAfterVal(int aval, int val);

    int DeleteFromBeg();

    int DeleteFromEnd();

    void Disp();

    void Destroy();

};
```

Definitions of member functions for the LinkedList class are contained in the LinkedList.cpp file as shown in the next code in this section. The file begins with a preprocessor statement that tells the preprocessor to reference the contents of the LinkedList.h file during preprocessing. The LinkedList.h file contains the LinkedList class definition, which is required to resolve statements in the LinkedList.cpp that refer to the class.

Each member function definition in this example is practically the same definition as those discussed in the last several sections of this lesson. The only exception is that reference is made to the LinkedList class in the name of each member function definition. This associates each definition with the LinkedList class for the compiler.

```
//LinkedList.cpp

#include "LinkedList.h"

void LinkedList::CreateList(int val)
{
    Node *p;
    p = new Node;
    p->data = val;
    p->next = NULL;
    List = p;
}

void LinkedList::InsertAtBeg(int val)
{
    Node *p;
    p = new Node;
    p->data = val;
    p->next = List;
    List = p;
}

void LinkedList::InsertAtEnd(int val)
{
    Node *p, *Start;
    Start = List;
    p = new Node;
```

```

    p->data = val;
    p->next = NULL;
    if (Start == NULL)
        List = p;
    else
    {
        while (Start->next != NULL)
            Start = Start->next;
        Start->next = p;
    }
}

void List::InsertAtPos(int pos, int val)
{
    int count = 1;
    Node *p, *q = List;
    p = new Node;
    p->data = val;

    while (q->next != NULL && count < pos-1)
    {
        count++;
        q = q->next;
    }
    p->next = q->next;
    q->next = p;
}

```



```

void List::InsertAfterVal(int aval, int val)
{
    Node *p, *q = List;
    p = new Node;
    p->data = val;
    while (q != NULL && q->data != aval)
        q = q->next;
    p->next = q->next;
    q->next = p;
}

```

```

int List::DeleteFromBeg()

```

```

{
    int val;
    Node *p = List;
    List = List->next;
    val = p->data;
    delete p;
    return val;
}

```

```

int List::DeleteFromEnd()

```

```

{
    int val;
    Node *p,*q = List;

    while (q->next->next != NULL)
        q = q->next;
}

```

```

        p = q->next;
        q->next = NULL;
        val = p->data;
        delete p;
        return val;
    }
void List::Disp()
{
    Node *p = List;
    while (p != NULL)
    {
        cout << p->data << "\t";
        p = p->next;
    }
    cout << endl;
}
void List::Destroy()
{
    Node *p;
    while (List != NULL)
    {
        p = List;
        List = List->next;
        delete p;
    }
    List = NULL;
}

```

```
}
```

The last file is the C++ application that uses the linked list. We call the file `LinkedListDemo.cpp`, which is shown next. It is amazing that the application itself is so small when compared to all the code used to define the `NODE` structure and the `LinkedList` class.

```
//LinkedListDemo.cpp

#include <iostream>

using namespace std;

void main()
{
    LinkedList list;

    clrscr();

    list.CreateList(10);
    list.InsertAtBeg(5);
    list.InsertAtEnd(20);
    list.InsertAtPos(3,15);
    list.InsertAfterVal(15,18);

    list.Disp();

    cout << "Element deleted is -> " << list.DeleteFromBeg() << endl;
    cout << "Element deleted is -> " << list.DeleteFromEnd() << endl;

    cout << "Final list is -> \n";

    list.Disp();

    list.Destroy();
}
```

The application begins by declaring an instance of the `LinkedList` class. As you recall from earlier in this lesson, the constructor initializes the front and back pointers. Next, the `CreateList()`, `InsertAtBeg()`, `InsertAtEnd()`, `InsertAfterPos()` and `InsertAfterVal()` member functions are called. Next the `Disp()` function displays the values of nodes of the list, which is:

5 10 15 18 20

The next two statements delete nodes from the list. First DeleteFromBeg() function is called which deleted the first node of the list (node with value 5), then DeleteFromEnd() function is called which deletes the last node of the list (node with value 20). Then the Disp() function again displays the nodes of the list, which are:

10 15 18

Finally the list is destroyed by calling Destroy() function, which successively deleted all nodes of the list.

## 8.6 Summary

Linked list is a linear data structure, in which each node is a combination of two value, one for storing the data of the node and the other for pointing to the next node if there is any or else NULL. A linked list is a dynamic data structure in which nodes can be inserted or deleted during execution of the program. Size of the linked list is constrained by the amount of memory available. Insertion in the list can be done at any point. If a node is to be inserted within a linked list then no shifting of other nodes is required.

## 8.7 Questions

1. What is a linked list?
2. What is the benefit of using a linked list?
3. What is a node?
4. What are the elements of a node?
5. What advantage does a linked list have over an array?
6. Can a node reference more than one data element?
7. How can a node be inserted in the middle of a linked list?
8. Create a function which inserts nodes in ascending order in the list.
9. WAP to traverse a linked list.
10. WAP to insert item after a specific node.

## 8.8 Suggested readings

1. A. Tanenbaum, Y. Lanhsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.

## TYPES OF LINKED LIST

- 9.1 Objectives**
- 9.2 Introduction**
- 9.3 Circular Linked List**
- 9.4 Doubly Linked List**
- 9.5 Linked List with Headers and Trailers**
- 9.6 Linked List as an arrays of nodes**
- 9.7 Applications of Linked Lists.**
- 9.8 Summary**
- 9.9 Questions**
- 9.10 Suggested readings**

### **9.1 Objectives**

This lesson begins with three new implementations of reference based lists i.e. circular linked lists, doubly linked lists and lists with headers and trailers. We then introduce an array-based approach to implementing a linked list. This implementation is widely used in operating systems software.

### **9.2 Introduction**

There are many variation of linked list discussed in lesson 8. The most common of these are circular linked list and doubly linked list. Linked list can be implemented using arrays as well.

### **9.3 Circular Linked List**

The linked lists that we implemented in lesson 8 are characterized by a linear (linelike) relationship between the elements. Each element (except the first one) has a unique predecessor and each element (except the last one) has a unique successor. Let's consider a small change to our linked list approach and see how much it would affect our

implementation and use of the Sorted List ADT. Suppose we change the linear list slightly, making the next reference of the last node point back to the first node instead of containing null (Figure 9.1). Now our list is a circular linked list rather than a linear linked list. We can start at any node in the list and traverse the whole list. Of course, we must now ensure that all of our list operations maintain this new property of the list that after the execution of any list operation, the last node continues to point to the front node. A quick consideration of each of the operations should convince us that we could continue to efficiently support all of them except when an operation changes the first element on the list. Consider, for example, if we try to delete the first element. Our previous delete approach would simply change the list reference to point to the second element on the list, effectively removing the first element. Now, however, we must also update the reference in the last element on the list, so that it points to the new first element. The only way to do that is to traverse the entire list to obtain access to the last element and then make the change. A similar problem arises if we insert an item into the front of the list. Circular linked list is a list in which every node has a successor the “last” element is succeeded by the “first” element.

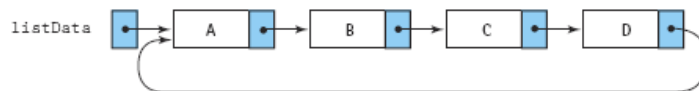


Figure 9.1 A circular linked list

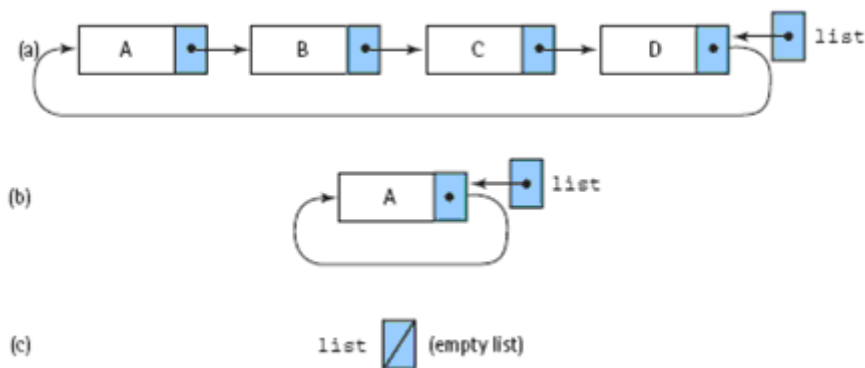


Figure 9.2 Circular linked lists with the external pointer pointing to the rear element

Inserting and deleting elements at the front of a list might be a common operation for some applications. Our linear linked list approach supported these operations very efficiently, but our circular linked list approach does not. We can fix this problem by letting our list reference point to the last element in the list rather than the first; now we have direct access to both the first and the last elements in the list (See Figure 9.2 where

list.info references the information in the last node and list.next.info references the information in the first node.)

### **Deleting from a Circular List**

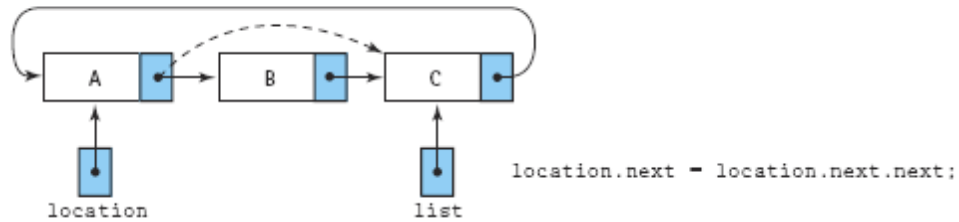
We can use the same basic approach to deleting an element from a circular list as we used for a linear list. First, find the element that matches the targeted item and then delete it. To delete it we unlink it from the chain of elements by setting the next reference of the element previous to the identified element to reference the element after the identified element. Thus, when we delete an element we need a reference to the element that precedes it. Recall the “trick” we used for the linear list delete method, where we always looked at the element in the position after our current location. That way, when we found the element to delete, location held a reference to the previous node and we could “jump over” the node to be deleted with the statement

```
location.next = location.next.next;
```

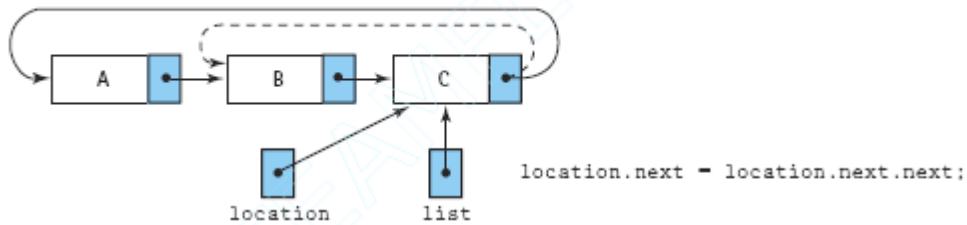
In fact, using our trick with the circular list works very nicely. Since the original value of location is the node at the end of the list, the first information that we check is actually associated with the first node on the list. As with the linear approach, we are guaranteed to find our targeted element. When we do, location is referencing its predecessor on the list. To remove the targeted element from the list, we simply reset location next as described above. We set it to jump over the node we are deleting. That works for the general case, at least (see Figure 9.3a).

However, the primary reason that there was a special case was that the overall list reference pointed to the first list element and had to be updated if that element was deleted. In the circular version the overall list reference points to the last list element, so it is very possible that deleting the first element is not a special case. Figure 9.3(b) shows that guess to be correct. However, deleting the only node in a circular list is a special case, as we see in Figure 9.3(c). The reference to the list must be set to null to indicate that the list is now empty. We can detect this situation by checking, at the start of the method, whether location is equal to location.next. If it is, since we know from the method preconditions that the item we are deleting is on the list, in the case of the single element list we can simply delete it immediately. It must be the element that we wish to delete. We delete it by setting the list reference to null.

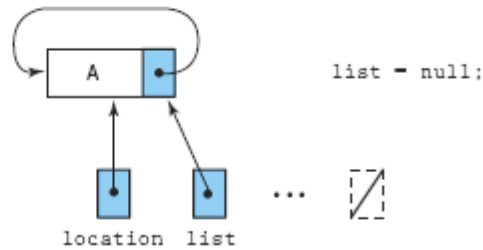
(a) The general case (delete B)



(b) Special case (?): deleting the smallest item (delete A)



(c) Special case: deleting the only item (delete A)



(d) Special case: deleting the largest item (delete C)

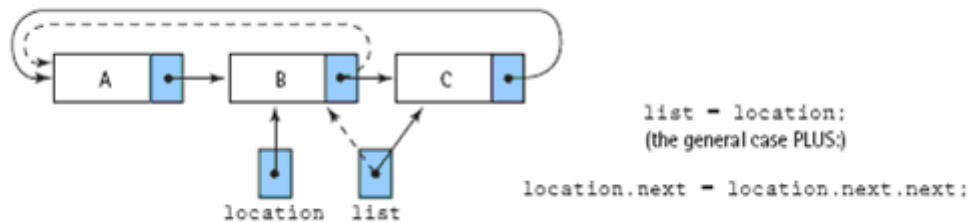


Figure 9.3 Deleting from a circular linked list

We might also guess that deleting the largest list element (the last node) from a circular list is a special case. After all, our reference to the list points to the last element, so if we delete it we must change our reference. As Figure 9.3(d) illustrates, when we delete the last node, we first update the overall list reference to point to the preceding element.



We can detect this situation by checking whether `location.next` equals `list` after the search phase.

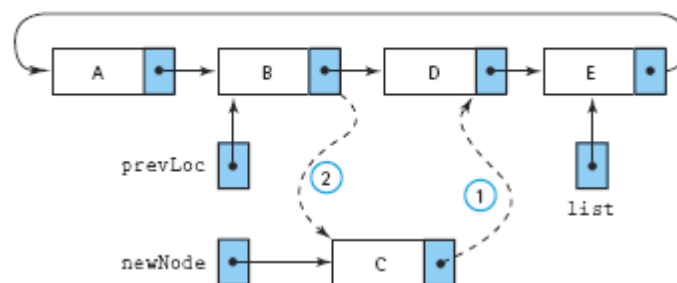
### The insert Method

The algorithm to insert an element into a circular linked list is also similar to its linear list counterpart. Essentially, we find the insertion location by performing a search and insert the new item by rearranging some references. To do the insertion we need to have access to both the node preceding the insertion point and the node following the insertion point. And we need to handle special cases carefully. The task of creating a new node is the same as for the linear list. We allocate space for the node using the `new` operator and then store a copy of `item` into `newNode.info`.

The next task is one that we are used to by now. We search through the list maintaining two references, `location` and `prevLoc`, until we find an element larger than `item` or reach the end of the list. The new node is linked into the list immediately after `prevLoc`. To put the new element into the list we store `location` into `newNode.next` and `newNode` into `prevLoc.next`.

The general case is illustrated in Figure 9.4(a). What are the special cases? First, we have the case of inserting the first element into an empty list. In this case, we want to make `list` point to the new node and to make the new node point to itself (Figure 9.4b). We handle this special case first, before doing any other processing. In the insertion algorithm for the linear linked list we also had a special case when the new element key was smaller than any other key in the list. Because the new node became the first node in the list, we had to change the reference to point to the new node. The reference to a circular list, however, doesn't point to the first node in the list—it points to the last node. Therefore, inserting the smallest list element is not a special case for a circular linked list (Figure 9.4c). However, inserting the largest list element at the end of the list is a special case. In addition to linking the node to its predecessor (previously the last list node) and its successor (the first list node), we must modify the list reference to point to `newNode`—the new last node in the circular list (Figure 9.4d).

(a) The general case (insert C)



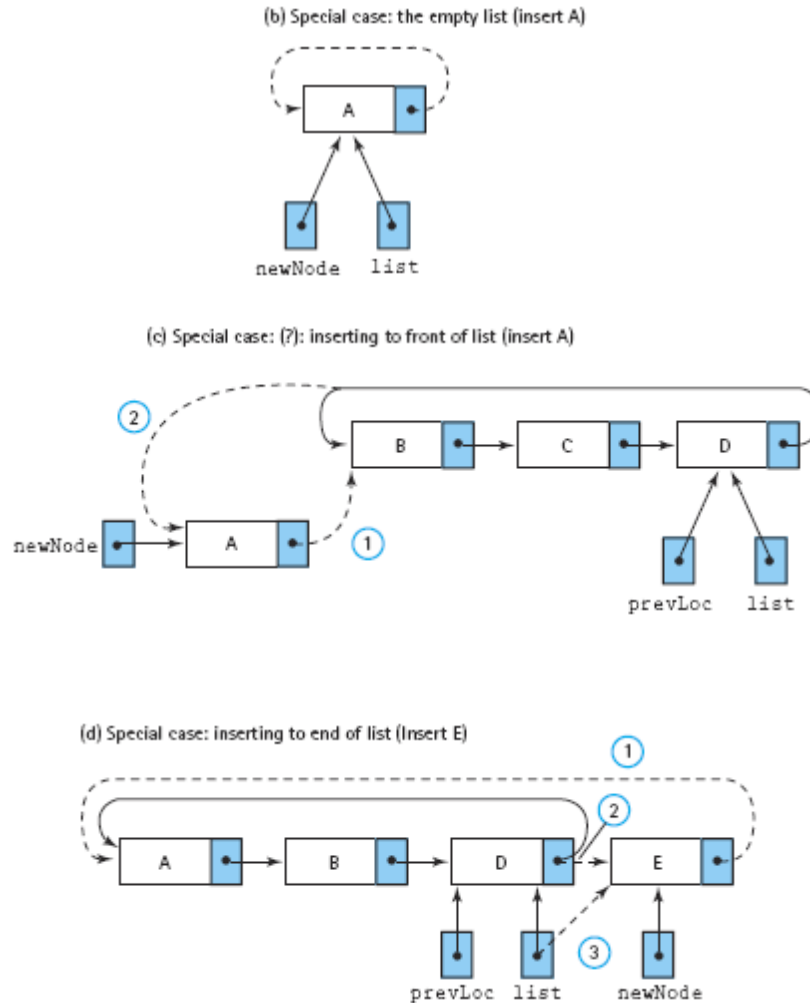


Figure 9.4 Inserting into a circular linked list

The statements to link the new node to the end of the list are the same as the general case, plus the assignment of the reference list. Rather than checking for this special case before the search, we can treat it together with the general case. We search for the insertion place and link in the new node. Then, if we detect that we have added the new node to the end of the list, we reassign list to point to the new node. To detect this condition, we compare item to list.info.

### insert (item)

Create a node for the new list element. Find the place where the new element belongs. Put the new element into the list Increment the number of items

### Circular Versus Linear

Studying circular linked lists provided good practice with using references and self-referential structures. You may have noticed that the only operation that is simpler for the circular approach, as compared to the linear approach, is `getNextItem`; that minimal advantage is counterbalanced by a more complicated `reset` operation. Why then might we want to use a circular, rather than linear, linked list? Circular lists are good for applications that require access to both ends of the list. Our Circular Sorted Linked List class could be used as the basis for other classes that include operations that can take advantage of the new implementation. Perhaps we need a “maximum” operation that returns the largest list element; with the circular approach we have easy access to the largest element (through `list`). Or suppose we need an operation `inBetween` that returns a boolean value indicating whether a parameter item is “in between” the largest and smallest element of the list; as just mentioned, with the circular approach we have easy access to the largest element (through `list`) and we also have easy access to the smallest element (through `list.next`). Therefore, with the circular list, we could implement `inBetween` in  $O(1)$ , whereas with our linear approach it would take  $O(N)$ .

In addition, it is common for the data we want to add to a sorted list to already be in order. Some times people manually sort raw data before turning it over to a data entry clerk. Data produced by other programs are often in sorted order. Given a Sorted List ADT and sorted input data, we always insert at the end of the list, the most expensive place to insert in terms of machine time. It is ironic that the work done manually to order the data now results in maximum insertion times. A circular list with the list reference to the end of the list, as developed in this section, can be designed to avoid this execution overhead. You may have realized that many of the benefits described here for circular lists could also be obtained by using the linear linked list defined in Lesson 8 augmented with a reference to the last element of the list. This is yet another list variation as with the circular list, this variation requires changes to some of the linear list methods. We ask you to explore this variation in the exercises. The existence of many list variations is another reason for studying circular lists. It helps us understand how small changes in the structure underlying an ADT can require many subtle changes in the implementations of the ADT operations.

#### **9.4 Doubly Linked List**

We have discussed using circular linked lists to enable us to reach any node in the list from any starting point. Although this structure has advantages over a simple linear linked list for some applications, it is still too limited for others. Suppose we want to be able to delete a particular node in a list, given only a reference to that node. This task involves changing the next reference of the node preceding the targeted node. However, given only a reference to a node, it is not easy to access its predecessor in the list.

Another task that is difficult to perform on a linear linked list (or even a circular linked list) is traversing the list in reverse. For instance, suppose we have a list of student records, sorted by grade point average (GPA) from lowest to highest. The Dean of Students might want a printout of the students’ records, sorted from highest to lowest, to use in

preparing the Dean’s List. Consider the Real Estate application where the user can step through a list of house information, viewing the information house by house on the screen, by pressing a “next” button. Suppose the user requests an enhancement to the interface—the idea is to include a “previous” button so that the user can browse through the houses in either direction.

In cases like these, where we need to be able to access the node that precedes a given node, a doubly linked list is useful. In a doubly linked list, the nodes are linked in both directions. Each node of a doubly linked list contains three parts:

- info: the data stored in the node
- next: the reference to the following node
- back: the reference to the preceding node

A linear doubly linked list is pictured in Figure 9.5. Note that the back reference of the first node, as well as the next reference of the last node, contains a null. The following definition might be used to declare the nodes in such a list:

**Doubly linked list** A linked list in which each node is linked to both its successor and its predecessor



Figure 9.5 A linear doubly linked list

### The Insert and Delete Operations

Using the definition of DLLListNode, let’s discuss the corresponding insert and delete methods. The first step for both is to find the location to do the insertion or deletion. This step was complicated in the singly linked list situation by the need to hold onto a reference to the previous location during the search. That is why we created our inchworm search approach. That approach is no longer needed; instead, we can get the predecessor to any node through its back reference. This means we can revert to the simpler search approaches used with our array-based lists.

Although our search phase is simpler, the algorithms for the insertion and deletion operations on a doubly linked list are somewhat more complicated than for a singly linked list. The reason is clear: There are more references to keep track of in a doubly linked list.

For example, consider insert. To link a new node newNode, after a given node referenced by prevLoc, in a singly linked list, we need to change two references: newNode.next and prevLoc.next (see Figure 9.6a). The same operation on a doubly linked list requires four reference changes (see Figure 9.6b).

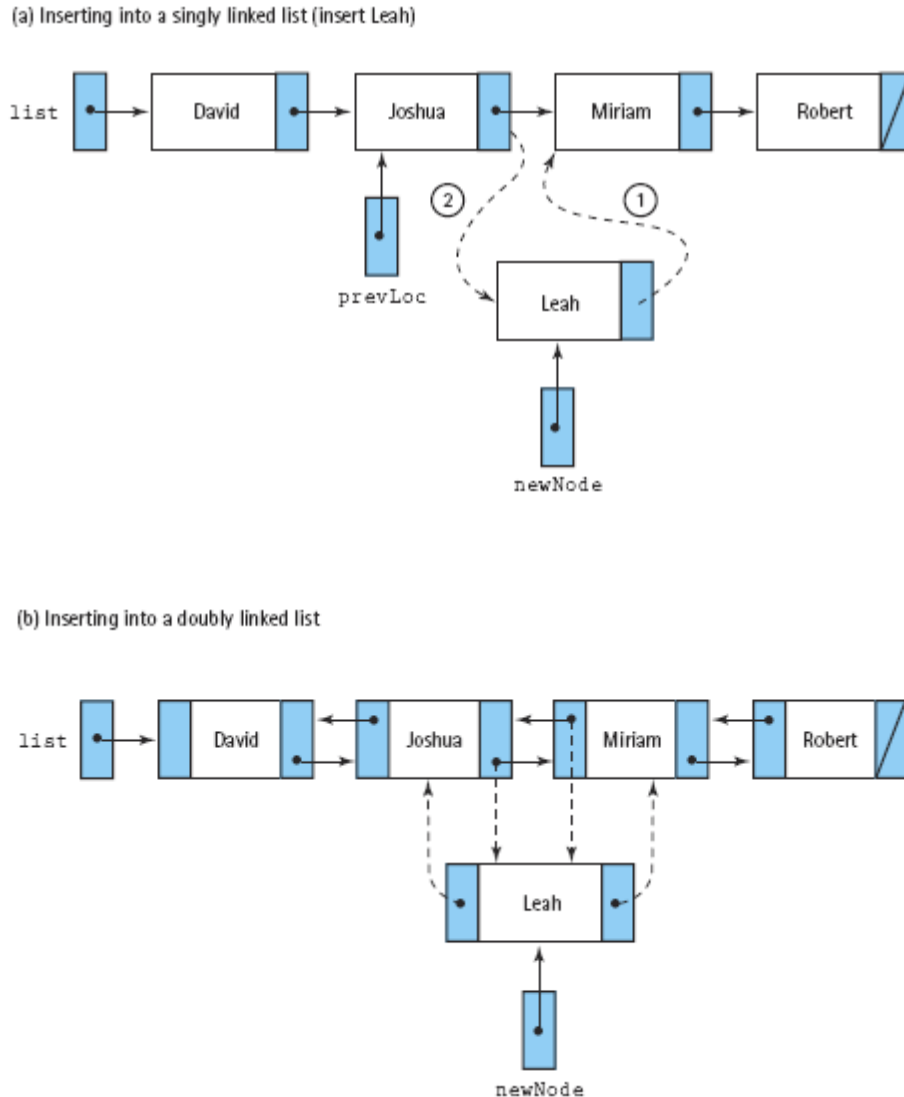


Figure 9.6 Insertions into single and doubly linked lists

To insert a new node we allocate space for the new node and search the list to find the insertion point. The result of our search is that location references the node that should follow the new node. Now we are ready to link the new node into the list. Because of the complexity of the operation, it is important to be careful about the order in which you change the references. For instance, when inserting the new node before location, if we change the reference in location.back first, we lose our reference to the node that is to precede the new node. The correct order for the reference changes is illustrated in Figure 9.7. The corresponding code would be

```
newNode.back = location.back;
newNode.next = location;
location.back.next = newNode;
```

```
location.back = newNode;
```

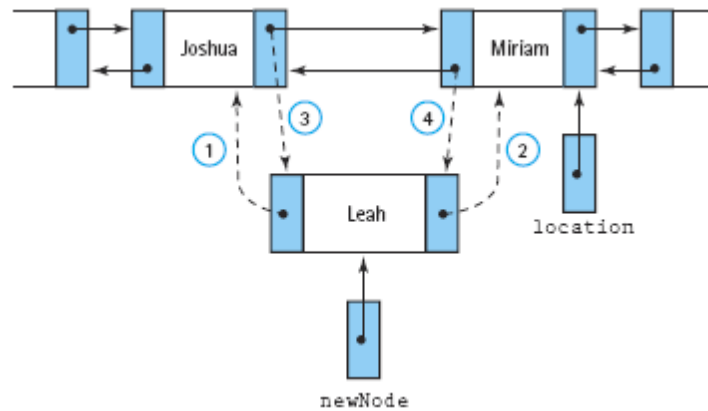


Figure 9.7 Inserting into a doubly linked list

We do have to be careful about inserting into an empty list, as it is a special case. Now let's consider the delete method. One of the useful features of a doubly linked list is that we don't need a reference to a node's predecessor in order to delete the node. Through the back reference, we can alter the next variable of the preceding node to make it jump over the unwanted node. Then we make the back reference of the succeeding node point to the preceding node. This operation is pictured in Figure 9.8. We do, however, have to be careful about the end cases.

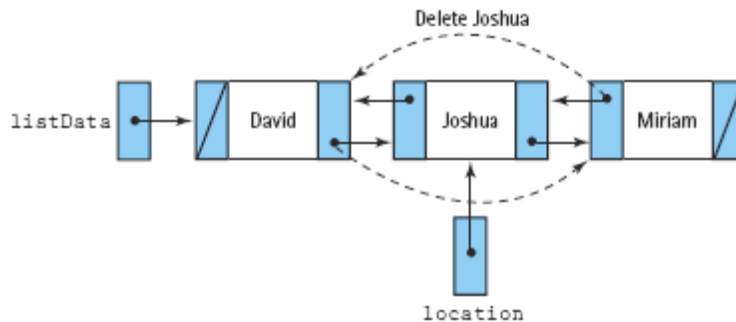


Figure 9.8 Deleting from a doubly linked list

If `location.back` is null, we are deleting the first node; if `location.next` is null, we are deleting the last node. If both `location.back` and `location.next` are null, we are deleting the only node. We leave the complete coding of the insert and delete methods for the doubly linked list as an exercise.

### The List Framework

Before leaving the topic of doubly linked lists we should address the question of how they fit into our list framework. Although a doubly linked list is a form of linked list, it is based on a different underlying logical structure than our other linked lists. The relationship among its list elements is different from the case of the singly linked list.

Therefore, it does not make sense to have it extend the abstract `LinkedList` class as we did for the other implementations. Instead, there are several other viable options. We could create a new class, perhaps called `DoublyLinkedList`, and require it to implement our current `ListInterface` interface. In this case, we would probably want to add some additional methods to `DoublyLinkedList` that are not required by the interface. For example, we could add a `getPreviousItem` method; otherwise, what is the benefit of having the double links?

Alternately, we could create a new interface, perhaps called `TwoWayListInterface` that defines what we expect of lists that can be traversed in two directions. Then we could create a class based on doubly linked lists that implements the new interface. The new interface would probably require all of the operations that are part of our current `ListInterface` interface, plus a few more. Certainly, we would add the `getPreviousItem` operation. Of course, a doubly linked list is not the only possible way to implement such an interface. An array-based approach would also work well.

The fact that the proposed new interface would require all of the operations of our current `ListInterface` interface raises another possible approach. Java supports the inheritance of interfaces. (In fact, the language supports multiple inheritance of interfaces, so that a single interface can extend any number of other interfaces.) A good approach would be to define a new interface that extends `ListInterface` and adds a `getPreviousItem` method.

### **Single Linked List vs. Doubly Linked List**

Programmers call this a doubly linked list or bidirectional (Figure 8.3) because each node contains reference to the previous and next node on the linked list. This enables the programmer to traverse the linked list in both directions by referencing the previous and next nodes. The node can be transformed into a single linked list (Figure 9.9) by only having one pointer in the structure that contains the address of the next node. Typically, a node in a single linked list references the next node and not the previous node, although nothing stops you from creating a backward reference by using only the previous node reference.

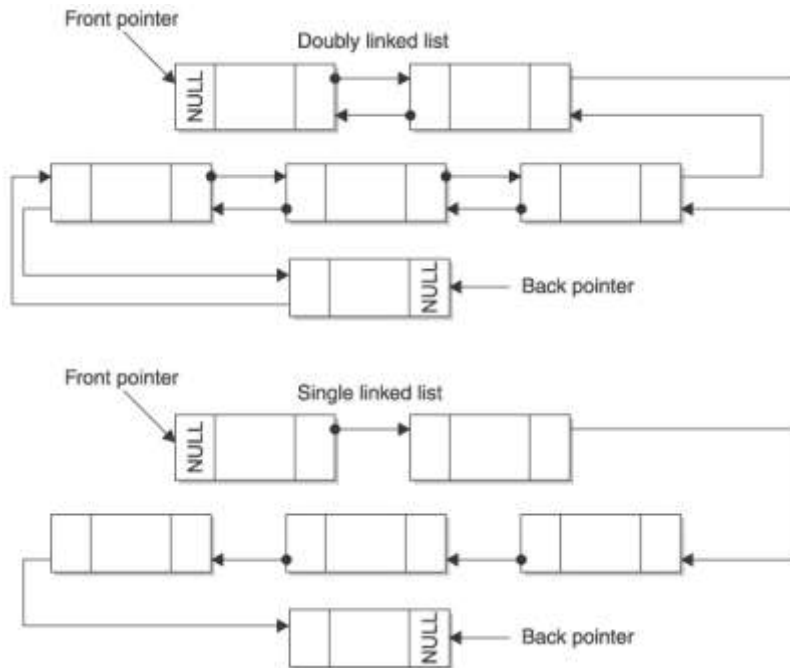


Figure 9.9 A doubly linked list contains next and previous members and a single linked list contains only a next member.

The following example is nearly the same as the previous example except this is a single direction node. You'll notice that reference to the previous node is missing. This means a programmer can only move down the linked list and not in both directions.

### 9.5 Linked List with Headers and Trailers

In writing the insert and delete algorithms for all implementations of linked lists, we see that special cases arise when we are dealing with the first node or the last node. One way to simplify these algorithms is to make sure that we never insert or delete at the ends of the list.

How can this be accomplished? Recall that the elements in the sorted linked list are arranged according to the value in some key—for example, numerically by identification number or alphabetically by name. If the range of possible values for the key can be determined, it is often a simple matter to set up dummy nodes with values outside of this range. A header node, containing a value smaller than any possible list element key, can be placed at the beginning of the list. A trailer node, containing a value larger than any legitimate element key, can be placed at the end of the list.

The header and the trailer are regular nodes of the same type as the real data nodes in the list. They have a different purpose, however, instead of storing list data, they act as placeholders.

**Header node** A placeholder node at the beginning of a list used to simplify list processing



**Trailer node** A placeholder node at the end of a list, used to simplify list processing

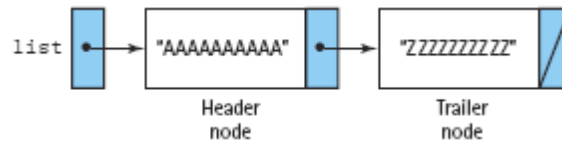


Figure 9.10 An “empty” list with a header and a trailer

If a list of students is sorted by last name, for example, we might assume that there are no students named “AAAAAAAAAA” or “ZZZZZZZZZZ”. We could, therefore, initialize our linked list to contain header and trailer nodes with these values as the keys, see Figure 9.10. How can we implement a general list ADT if we must know the minimum and maximum key values? We can use a parameterized class constructor and let the user pass as arguments elements containing the dummy keys.

## 9.6 Linked List as an arrays of nodes

We tend to think of linked structures as being dynamically allocated as needed, using self-referential nodes as illustrated in Figure 9.11(a) but this is not a requirement. A linked list could be implemented in an array; the elements might be stored in the array in any order and “linked” by their indexes (see Figure 9.11b). In this section, we develop an array-based linked-list implementation.

For the array-based linked representation developed in this section, we can no longer rely on dynamic memory management support. Instead, we predetermine the maximum list size and instantiate an array of list nodes of that size. We then directly manage the nodes in the array. We keep a separate list of the available nodes and write routines to allocate and deallocate nodes, from and to this free list.

We have seen that dynamic allocation of list nodes has many advantages, so why would we even discuss using an array-of-nodes implementation instead? Remember that dynamic allocation is only one advantage of choosing a linked implementation another advantage is the efficiency of the insert and delete algorithms. Most of the algorithms that we have discussed for operations on a linked structure can be used for either an array-based or a reference-based implementation. The main difference is the requirement that we manage our own free space in an array-based implementation. Sometimes, managing the free space ourselves gives us greater flexibility.

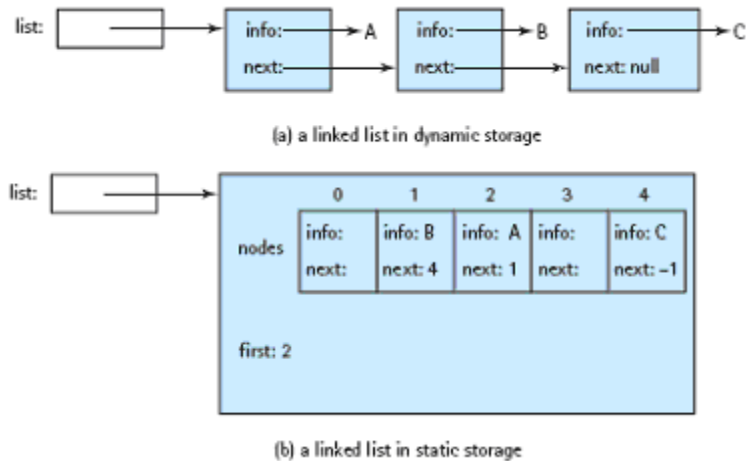


Figure 9.11 Linked lists in dynamic and static storage

Another reason to use an array of nodes is that there are programming languages that do not support dynamic allocation or reference types. You can still use linked structures if you are programming in one of these languages, using the techniques presented in this section.

With some languages, using references presents a problem when we need to save the information in a data structure between runs of a program. Suppose we want to write all the nodes in a list to a file and then use this file as input the next time we run the program to recreate the list. If the links are reference values—containing memory addresses—they are meaningless on the next run of the program because the program may be placed somewhere else in memory the next time. We must save the user data part of each node in the file and then rebuild the linked structure the next time we run the program. An array index, however, is still valid on the next run of the program. We can store the array information, including the next data index and then read it back in the next time we run the program.

Most importantly, there are times when dynamic allocation isn't possible or feasible or when dynamic allocation of each node, one at a time, is too costly in terms of time—especially in real-time system software such as operating systems, air traffic controllers and automotive systems. In such situations, an array-based linked approach provides the benefits of linked structures without the runtime costs.

Let's get back to our discussion of how a linked list can be implemented in an array. We can associate a next variable with each array node to indicate the array index of the succeeding node. The beginning of the list is accessed through a "reference" that contains the array index of the first element in the list. Figure 9.12 shows how a sorted list containing the elements "David," "Joshua," "Leah," "Miriam" and "Robert" might be stored in an array of nodes called nodes. Do you see how the order of the elements in the list is explicitly indicated by the chain of next indexes?

nodes	.info	.next
[0]	David	4
[1]		
[2]	Miriam	6
[3]		
[4]	Joshua	7
[5]		
[6]	Robert	-1
[7]	Leah	2
[8]		
[9]		

list	0
------	---

Figure 9.12 A sorted list stored in an array of nodes

What goes in the next index of the last list element? Its “null” value must be an invalid address for a real list element. Because the nodes array indexes begin at 0, the value `_1` is not a valid index into the array, that is, there is no `nodes[_1]`. Therefore, `_1` makes an ideal value to use as a “null” address. We could use the literal value `_1` in our programs:

```
while (location != -1)
```

but it is better programming style to declare a named constant. We use the identifier

```
NUL and define it to be _1:
```

```
private static final int NUL = -1;
```

When an array-of-nodes implementation is used to represent a linked list, the programmer must write routines to manage the free space available for new list elements. Where is this free space? Look again at Figure 9.12. All of the array elements that do not contain values in the list constitute free space. Instead of the built-in allocator `new`, which allocates memory dynamically, we must write our own method to allocate nodes from the

free space. We call this method `getNode`. We use `getNode` when we insert new items onto the list.

When elements are deleted from the list, we need to free the node space, that is, to return the deleted node to the free space, so it can be used again later. We can't depend on a garbage collector the node we delete remains in the allocated array so it is not reclaimed by the run-time engine. We write our own method, `freeNode`, to put a node back into the pool of free space.

We need a way to track the collection of nodes that are not being used to hold list elements. We can link this collection of unused array elements together into a second list, a linked list of free nodes. Figure 9.13 shows the array nodes with both the list of elements and the list of free space linked through their next values. The list variable is a reference to a list that begins at index 0 (containing the value David). Following the links in next, we see that the list continues with the array slots at index 4 (Joshua), 7 (Leah), 2 (Miriam) and 6 (Robert), in that order. The free list begins at free, at index 1.

nodes	.info	.next
[0]	David	4
[1]		5
[2]	Miriam	6
[3]		8
[4]	Joshua	7
[5]		3
[6]	Robert	NUL
[7]	Leah	2
[8]		9
[9]		NUL

list	0
free	1

Figure 9.13 An array with linked list of values and free space

Following the next links, we see that the free list also includes the array slots at index 5, 3, 8 and 9. You see two NUL values in the next column because there are two

linked lists contained in the nodes array so there are two end-of-list values in the array. There are two approaches to using an array-of-nodes implementation for linked structures. The first is to simulate dynamic memory with a single array. One array is used to store many different linked lists, just as the computer's free space can be dynamically allocated for different lists. In this approach, the references to the lists are not part of the storage structure but the reference to the list of free nodes is part of the structure. Figure 9.14 shows an array that contains two different lists. The list indicated by list1 contains the values "John," "Nell," "Susan" and "Susanne" and the list indicated by list2 contains the values "Mark," "Naomi" and "Robert." The remaining three array slots in Figure 9.14 are linked together in the free list.

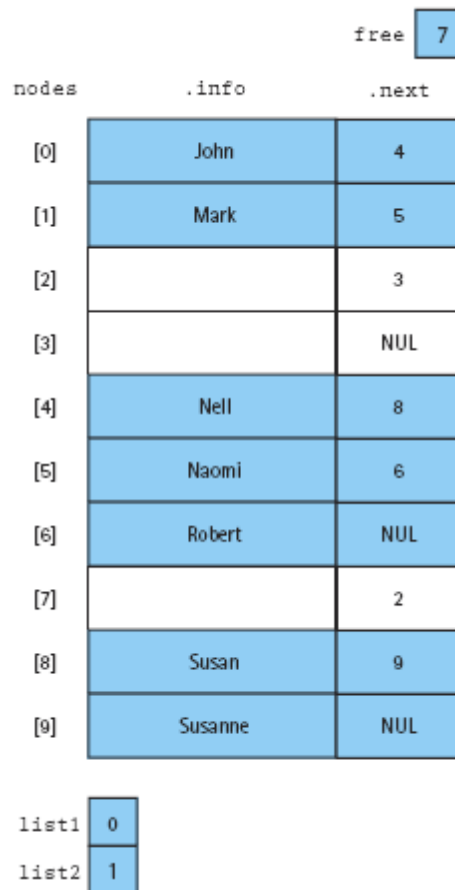


Figure 9.14 An array with three lists (including the free list)

Another approach is to have one array of nodes for each list. In this approach, the reference to the list is part of the storage structure itself (see Figure 9.15). This works since there is only one list. The list constructor has a parameter that specifies the maximum number of items to be on the list. This parameter is used to dynamically allocate an array of the appropriate size.

	free	1
	list	0
nodes	.info	.next
[0]	David	4
[1]		5
[2]	Miriam	6
[3]		8
[4]	Joshua	7
[5]		4
[6]	Robert	NUL
[7]	Leah	2
[8]		9
[9]		NUL

Figure 9.15 List and link structure are together

In this section, we implement this second approach. We call our new class `ArrayLinkedList`. In implementing our class methods, we need to keep in mind that there are two distinct processes going on within the array of nodes: bookkeeping relating to the space (such as initializing the array of nodes, getting a new node and a node) and the operations on the list that contains the user's data. The bookkeeping operations are transparent to the user. Our list interface does not change. In fact, our new class implements the `ListInterface` interface, just as all of our other list implementations have done. The private data, however, change. We need to include the array of nodes. Let's call this array `nodes` and have it hold elements of the class `AListNode`.

Objects of the `AListNode` class, then, contain two attributes: `info` of type `Listable` that holds a reference to a copy of the user's data and `next`, of the primitive type `int`, that holds the index of the next element on the list. In addition to the array of nodes, we need an integer "reference" to the first node of the list and another to the first free node. We call these `list` and `free`. And of course, we still need our `numItems` and `currentPos` variables. Next is the beginning of our class file.

The class constructors for class `ArrayLinkedList` must allocate the storage for the array of nodes and initialize all of the private instance variables. They must also set up the initial free list of nodes. At the time of instantiation, all of the nodes are on the free list. So, the variable `free` is set to 0 to “reference” the first array node and the next value of that node is set to 1 and so on until all of the nodes are chained together. This initialization can be handled by a for loop, followed by a single assignment statement to set the last next value to `NUL`. To be consistent with our past array based implementation we should provide two constructors, one that accepts a size parameter and one that uses a default maximum size.

The methods that do the bookkeeping, `getNode` and `freeNode` are auxiliary (“helper”) methods and therefore are private methods. The `getNode` method must return the index of the next free node. The easiest node to use is the one at the beginning of the free list, so `getNode` returns the value of `free`. Therefore, `getNode` must also update the value of `free` to indicate the next node on the free list. Other than the fact that we must be careful of our order of operations and use a temporary variable to hold the index we need to return.

The `freeNode` method must take the node index received as an argument and insert the corresponding node into the list of free nodes. As the first element in the list is the one that is directly accessible, we have `freeNode` insert the node being returned at the beginning of the free list, in the variable `free`.

The public methods are very similar to their reference-based linked list counterparts. From the point of view of the algorithm used, they are identical. Now we have a third implementation approach, the array-based linked approach. The following table shows equivalent expressions in each of our views of a list. It also shows the expressions for creating and deleting nodes, where appropriate.

Design Terminology	Array-Based	Reference-Based	Array Index Links
<code>location.node()</code>	<code>list[location]</code>	<code>location</code>	<code>nodes[location]</code>
<code>location.info()</code>	<code>list[location]</code>	<code>location.info</code>	<code>nodes[location].info</code>
<code>location.next()</code>	<code>list[location+1]</code>	<code>location.next</code>	<code>nodes[location].next</code>
allocate a node N	done by constructor	<code>ListNode N = new ListNode</code>	<code>N = getNode()</code>
free a node N	not applicable	remove links, garbage collector	<code>freeNode(N)</code>

We look here at some of the methods and leave the rest as an exercise. First an easy one is full. We have two ways of determining whether or not the list is full. As for our array-based list implementation, we can compare the number of items on the list to the size of the underlying array. If they are equal then the list is full. But there is an even easier way. If the entire array is being used to hold our list, then the list of free space must be empty.

The remaining methods can be implemented following the same scheme devised for the reference-based approach. You must be careful, however, to correctly transform the implementation. And don't forget that you have to handle the memory management yourself. Let's look at the delete method.

Think for a minute about how you would represent this statement using the approach of this section. It is not as simple as it might first appear. You need to compare the info value of the next element to item. Using the table above, you see that you access the info attribute of a location with `nodes[location].info`. However, you don't want the info value of the location, you want the info value of the next location. So, you must replace `location` with the expression that stands for the next location. In this case, that is `nodes[location].next`. Putting this altogether (see Figure 6.15), the corresponding code is:

```
while (item.compareTo(nodes[nodes[location].next].info) != 0)
    location = nodes[location].next;
```

## 9.7 Applications of Linked Lists.

The main Applications of Linked Lists are

- For representing Polynomials  
It means in addition/subtraction /multiplication of two polynomials.

Eg: $p_1=2x^2+3x+7$  and  $p_2=3x^3+5x+2$

$p_1+p_2=3x^3+2x^2+8x+9$

- In Dynamic Memory Management
  - In Symbol Tables
- In Dynamic Memory Management
  - Allocation and releasing memory at runtime.
  - In Symbol Tables
  - Balancing parenthesis
- Representing Sparse Matrix
- Representation and manipulation of long numbers.

## 9.8 Summary

The idea of linking the elements in a list has been extended to include lists with header and trailer nodes, circular lists and doubly linked lists. In addition to using dynamically allocated nodes to implement a linked structure, we looked at a technique for implementing linked structures in an array of nodes. In this technique the links are not references into the free store but indexes into the array of nodes. This type of linking is used extensively in systems software.

Although a linked list can be used to implement almost any list application, its real strength is in applications that largely process the list elements in order. This is not to say



that we cannot do “random access” operations on a linked list. Our specifications include operations that logically access elements in random order. For instance, public methods retrieve and delete manipulate an arbitrary element in the list. However, at the implementation level the only way to find an element is to search the list, beginning at the first element and continuing sequentially to examine element after element. This search is  $O(N)$ , because the amount of work required is directly proportional to the number of elements in the list. A particular element in a sequentially sorted list in an array, in contrast, can be found with a binary search, decreasing the search algorithm to  $O(\log_2 N)$ . For a large list, the  $O(N)$  sequential search can be quite time-consuming. There is a linked structure that supports  $O(\log_2 N)$  searches the binary search tree. We discuss the binary search tree data structure in detail in Lesson 12.

### **9.9 Questions**

1. Define circular linked list. What are its properties?
2. How insertion and deletion operations are performed on circular linked list?
3. Define doubly linked list. What are its properties?
4. How insertion and deletion operations are performed on doubly linked list?
5. What are the properties of linked list with header and trailers?
6. How linked list is implemented using arrays? Explain.
7. What are the various applications of linked list?

### **9.10 Suggested readings**

1. A. Tanenbaum, Y. Lanhsam and A. J. Augenstein, “Data Structures Using C++”, PHI.
2. M. A. Weiss, “Data Structures and Algorithm Analysis in C++”, Pearson Education.
3. R. Sedgewick, “Algorithms in C++”, Pearson Education.
4. S. Lipschutz, “Data Structures”, Tata McGraw Hill.
5. Donald E. Knuth, “The Art Of Computer Programming”, Vol 1-4. 3rd ed., Addison Wesley.

## LINKED REPRESENTATION OF STACKS AND QUEUES

### 10.1 Objectives

### 10.2 Introduction

### 10.3 Linked representation of Stack

- 10.3.1 The StackLinkedList Class
- 10.3.2 StackLinkedList Constructor and Destructor
- 10.3.3 Pushing a Node onto a Stack-Linked List
- 10.3.4 Popping a Node from a Stack-Linked List
- 10.3.5 Determine If the Stack Is Empty

### 10.4 StackLinked List Using C++

- 10.4.1 LinkedList Header File and LinkedList Functions
- 10.4.2 StackLinkedList Header File and StackLinkedList Source File
- 10.4.3 StackLinkedList Application

### 10.5 Linked representation of Queue

- 10.5.1 The QueueLinkedList class
- 10.5.2 Enqueue
- 10.5.3 Dequeue

### 10.6 Linked List Queue Using C++

### 10.7 Summary

### 10.8 Questions

### 10.9 Suggested readings

## 10.1 Objectives

You learned about stacks back in lesson 4 and about queues in lesson 6 when you discovered how to use an array to create your own stack or queue. However, using arrays presents a problem: you cannot adjust the size of the stack or queue when the program runs. The solution is to use a linked list to create a stack or queue. You learned about linked lists in general in lesson 8. In this lesson, you'll learn how to use a linked list to create a stack and a queue.

## 10.2 Introduction

As you'll recall from Lesson 4, a stack is a data structure that organizes data similar to how you organize dishes in a stack on your kitchen counter. The newest dish is on top and the oldest is on the bottom of the stack.

When accessing dishes, the last dish on the stack is the first dish removed from the stack. If you want the third dish, you must remove the first two dishes from the top of the stack first so that the third dish becomes the top of the stack and you can remove it. There is no way to remove a dish from anywhere other than the top of the stack. You'd need to use a different kind of data structure (or stacking system) if you wanted to randomly access dishes.

A stack is useful whenever you need to store and retrieve data in last in, first out order. For example, your computer processes instructions using a stack in which the next instruction to execute is at the top of the stack.

In Lesson 6, you learned that a queue is a sequential organization of data where data is accessible on a first in, first out (fifo) basis, which is similar to the line that you stand in to buy concert tickets.

The queue in Lesson 6 was created using an array to store data. As you'll recall, the array is separate from the queue. Data is assigned to elements of the array. The queue itself consists of two variables called front and back. Each points to the array element that is at the front of the queue or at the back of the queue. When data is removed from the front of the queue, the program changes the value of the front variable to point to the next array element. However, the data removed from the queue remains assigned to the array. That is, data isn't removed from memory.

There is a serious problem with using arrays to store data for queues: you must know the size of the array when you write the program. An array can store only a specific maximum number of elements at any point in time, similar to an architect designing a specific space for a box office that can accommodate a maximum number of fans at any point in time.

However, there is a difference between exceeding the number of array elements and overflowing the space around the box office: unlike the stadium, there is no parking lot for fans to gather in while waiting to get in the queue to purchase tickets inside a computer.

Programmers work around the size issue by using a linked list instead of an array when creating a queue. As you learned in previous lessons, a linked list can grow and shrink at runtime based on the needs of the application.

### 10.3 Linked representation of Stack

Although we discuss data as being stacked like a stack of dishes, it isn't physically stacked at all. Instead, data is linked together sequentially in a list, where the last data always appears at the front of the list. Data is removed only from the front of the list.

You create this sequential list by using a linked list. In lesson 8, you learned that a linked list contains entries called nodes. A node has three subentries, data and two pointers. The data subentry is the data stored on the stack. Pointers point to the previous node and the next node (Figure 10.1). When you enter a new item on a linked list, you allocate the new node and then set the pointers to the previous and next nodes.

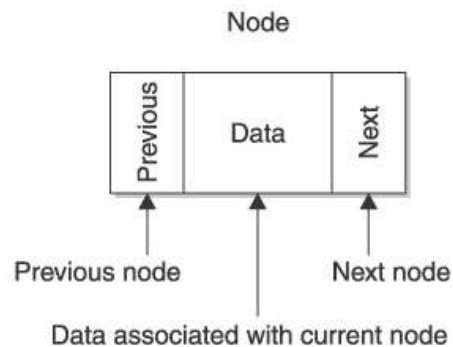


Figure 10.1: A node contains references to the previous node and the next node in the linked list and contains data that is associated with the current node.

A node is defined in C++ by using a structure, which is a user-defined data type. The following structure defines a node:

```
typedef struct Node
{
    struct Node(int data)
    {
        this->data = data;
        previous = NULL;
    }
};
```

```

        next = NULL;
    }

    int data;

    struct Node* previous;

    struct Node* next;

} NODE;

```

The structure is called Node. The name of the structure creates an instance of the structure similar to the way a constructor creates an instance of a class and data type. Let's skip the second definition of the structure and look at the last three statements within the structure because statements at the beginning of the structure actually create an instance of the structure and don't define the structure. The first statement declares an integer that stores the current data of the node. The next two statements declare pointers to the previous and next nodes in the linked list.

The constructor initializes elements of the node when the node is created, which is similar to the way constructors work in a class definition. You provide the current data to the structure when you create a new node. This data is assigned to data in the argument list, which is then assigned to the data element of the instance of the structure. The previous and next nodes are initialized to NULL, which indicates there are no other elements of the linked list. The NULL is replaced with a reference to a node when a new node is added to the linked list.

A LinkedList class is defined to create and manage a linked list. There are two data members and six function members defined in the LinkedList class. Data members are pointers to instances of the Node structure. The first pointer is called front, and it references the first node on the linked list. The second pointer is called back, and it references the last node on the linked list.

Both the front and back pointers are declared in the protected access specifier area of the class definition because the LinkedList class is inherited by the StackLinkedList class, which you'll learn about in the "The StackLinkedList Class" section of this lesson. The StackLinkedList class uses the front and back pointers.

The six member functions manipulate the linked list. These function members are the constructor, destructor, appendNode(), displayNodes(), displayReverseNodes(), and destroyNodes().

Here is the LinkedList class definition. The front and back pointers are defined in the protected access specifier area because the StackLinkedList class will use them:

```

class LinkedList{

```

```

protected:
    NODE* front;

    NODE* back;

public:
    LinkedList();
    ~LinkedList();

    void appendNode(int);
    void displayNodes();
    void displayNodesReverse();
    void destroyList();

};

```

### 10.3.1 The StackLinkedList Class

An efficient programmer does not repeat code if possible and instead inherits attributes and behaviors of another class, defining a LinkedList class to create and manipulate a linked list. An efficient programmer might also define a StackLinkedList class to create and manipulate a stack-linked list. The StackLinkedList class inherits attributes and behaviors of the LinkedList class and then defines other behaviors that are necessary to work with a stack-linked list.

In addition to the attributes and behaviors defined in the LinkedList class, the StackLinkedList class requires five behaviors defined as member functions: a constructor and destructor, push(), pop(), and isEmpty(). The StackLinkedList class definition is shown here:

```

class StackLinkedList : public LinkedList{
public:
    StackLinkedList();
    virtual ~StackLinkedList();
    void push(int);
    int pop();
    bool isEmpty();
};

```

### 10.3.2 StackLinkedList Constructor and Destructor

The constructor and destructor of the StackLinkedList class may be confusing the first time you look at them because both are empty and there aren't any instructions specified in the body of the constructor and destructor, as shown here:

```
StackLinkedList() {  
    }  
~StackLinkedList() {  
    }  
}
```

The constructor is empty because the constructor of the `LinkedList` class is called before the constructor of the `StackLinkedList` class. You'll recall that the `StackLinkedList` class inherits the `LinkedList` class. The `LinkedList` class constructor initializes the front and back pointers of the linked list to `NULL`. Therefore, there is nothing else for the `StackLinkedList` class constructor to do.

Likewise, the destructor of the `LinkedList` class is called before the destructor of the `StackLinkedList` class. The `LinkedList` class constructor deletes all memory that is associated with the nodes of the linked list. Therefore, the destructor of the `StackLinkedList` class also has nothing to do.

### 10.3.3 Pushing a Node onto a Stack-Linked List

In lesson 4, you learned that data is placed at the top of the stack and removed from the top of the stack. Programmers call this pushing data onto the stack and popping data off the stack. The same steps occur when using a linked list for the stack, but instead of placing data at the next available index in an array, it is placed at the back of the linked list.

You'll need to define a `push()` member function for the `StackLinkedList` class that is called whenever data is added to the stack. Remember that you are really adding a node to the linked list and not simply data. Data is contained in the node.

To add a node to the stack, you use the same steps you use to add a node to a linked list. This means that the `appendNode()` member function of the `LinkedList` class can be used to place a new node on the stack. Therefore, all you need is to call the `appendNode()` member function from the `push()` member function. Because `appendNode()` is public, you could just call `appendNode` directly to push a node onto the stack, but putting a `push()` function in the stack class makes this more intuitive to somebody using this class. This also helps hide the underlying implementation so using the class is a little more straightforward.

The `appendNode()` member function requires one argument, which is the data that is assigned to the new node. You must define the `push()` member function to accept the same data as its argument in order to pass this data to the `appendNode()` member function. This is illustrated in the following example. The `push()` member function requires an integer passed as an argument. The integer is then passed to the `appendNode()` member function within the body of the `push()` function definition.

```
void push(int x) {  
    appendNode(x);  
}
```

### 10.3.4 Popping a Node from a Stack-Linked List

You'll also need to define a member function to pop a node from the stack. In this example, we'll call it `pop()`. Because you're using the linked list as the stack, the `pop()` member function must remove the node from the back of the linked list.

Unfortunately, you cannot simply call a member function of the `LinkedList` class to pop the node off the stack because the `LinkedList` class doesn't define a member function that removes a node from the linked list. If you had a member function in the base class for `removeBack()`, you could call that to pop a node off the list. In this case, you'll need to define a `pop()` function in the `StackLinkedList` class to do this. This will give you a last in, first out access to the stack.

Here is the definition of the `pop()` member function. Refer to the picture of the stacked linked list in Figure 10.2 as you read to help you understand how the `pop()` member function works.

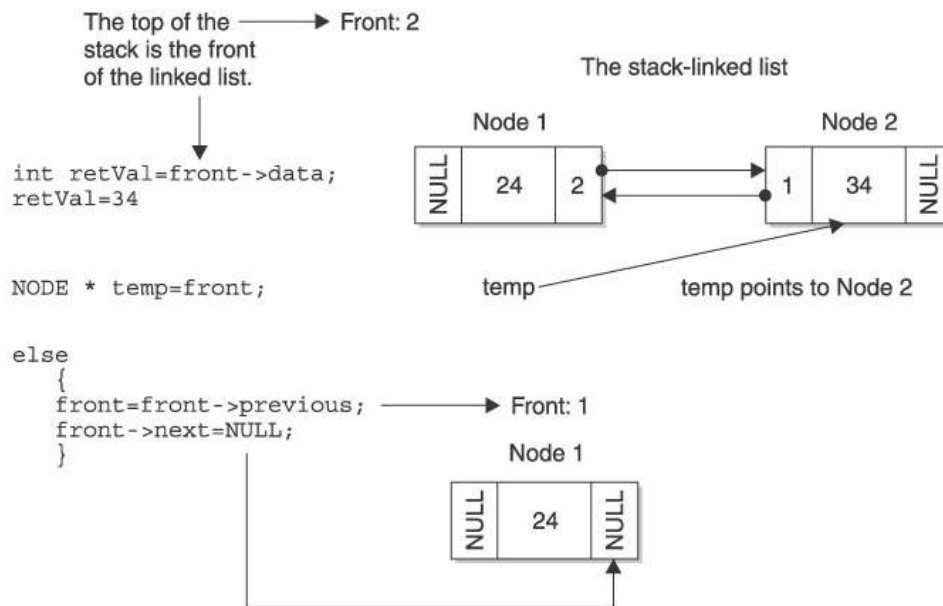


Figure 10.2: The `pop()` member function removes the node at the top of the stack, which is the node at the front of the linked list.

```
int pop()
{
  if (isEmpty())
  {
    return -1;
  }
  int retVal = back->data;
  NODE * temp = back;
  if (back->previous == NULL)
```



```

    {
        back = NULL;
        front = NULL;
    }
    else
    {
        back = back->previous;
        back->next = NULL;
    }
    delete temp;
    return retVal;
}

```

First, you must determine if there is anything on the stack by calling the `isEmpty()` member function. We'll show you how this member function works later in this section. For now, understand that the `isEmpty()` member function returns a Boolean `true` if the stack is empty or a Boolean `false` if it is not. You can see in Figure 10.2 that the stack has two nodes on the stack, so it is not empty. Therefore, the return statement in the `if` statement is not executed.

The `pop()` member function refers to the `back` attribute of the `LinkedList` class. It is important to remember that the `back` attribute refers to the top of the stack. Nodes will be removed from the back of the linked list to do a `pop` operation. Therefore, the value of `front` is Node 2.

The value of Node 2 is assigned to the `retVal` variable, which is the value returned by the `pop()` member function if there is a node on the stack. This pops the value from the stack.

Next, the address of the `back` node, which is Node 2, is assigned to a temporary pointer. The node that the temporary pointer points to is removed from memory with the `delete` operator at the end of the `pop()` member function.

Next, you determine if the node at the back of the stack was the only node on the linked list. You make this determination by seeing if the `previous` attribute of the node is `NULL`. If the `previous` pointer on the back of the list is `NULL`, this indicates that there's only one node in the linked list.

Be careful when analyzing the `pop()` member function. Remember that the back of the linked list is the top of the stack and that the bottom of the stack is the top of the linked list. If the `pop()` member function is removing the only node on the stack, then the `front` and `back` attributes of the `LinkedList` class are set to `NULL`, indicating there are no nodes left on the linked list after the `pop()` is executed.

However, if there is at least one node on the stack, statements within the `else` statement are executed, as in Figure 10.2. The first statement within the `else` statement assigns the `previous` attribute of the `back` attribute as the new `back`. In Figure 10.2, the

previous attribute is 1. This tells you that Node 1 comes before Node 2. You then assign Node 1 as the new back of the stack. Remember that there isn't a next node on the stack because you are always working with the back of the linked list. Therefore, you need to assign NULL to the next attribute of the back node, which is Node 1. This makes Node 1 at the top of the stack.

The next to last statement removes the node from memory using the delete operator. The temporary pointer points to the memory address of the node that was removed from the stack. The last statement returns the value of the node that was removed from the stack.

### 10.3.5 Determine If the Stack Is Empty

The pop() member function must determine if the stack is empty, or it will attempt to remove a node that isn't on the stack. The pop() member function determines if the stack is empty by calling the isEmpty() member function, which you must define as part of the StackLinkedList class.

The isEmpty() member function is a simple function, as shown next. It determines if the stack is empty by seeing if the value of the front attribute of the LinkedList class is NULL. If so, then a Boolean true is returned; otherwise, a Boolean false is returned. If the stack is empty, both front and back are equal to NULL but you only need to check one of them.

```
bool isEmpty()
{
    if(front == NULL)
    {
        return true;
    }
    else
    {
        return false;
    }
}
```

## 10.4 StackLinked List Using C++

Now that you have an understanding of components, you need to create a stack-linked list. In this section, we'll focus on assembling them into a working C++ application. Some programmers organize components of a stack-linked list into five files: LinkedList.h, LinkedList.cpp, StackLinkedList.h, StackLinkedList.cpp, and StackLinkedListDemo.cpp. All these files are joined together at compile time to create the executable.

The LinkedList.h file is the header file that contains the definition of the Node structure and the definition of the LinkedList class. The LinkedList.cpp is a source code file that contains the implementation of member functions of the LinkedList class. The

StackLinkedList.h file is the header that contains the definition of the StackLinkedList class. The StackLinkedList.cpp is the source code file that contains the implementation of member functions of the StackLinkedList class.

The StackLinkedListDemo.cpp contains the application. It is here where an instance of the StackLinkedList class is declared and member functions are called.

#### 10.4.1 LinkedList Header File and LinkedList Functions

The LinkedList.h file and the LinkedList.cpp file are shown in the following code. The front and back attributes defined in the LinkedList class in the LinkedList.h file are defined within the protected access specifier section of the class definition. The StackLinkedList class needs access to these variables, so you protect them so they're visible to the subclass.

```
//LinkedList.h
typedef struct Node
{
    struct Node(int data)
    {
        this->data = data;
        previous = NULL;
        next = NULL;
    }
    int data;
    struct Node* previous;
    struct Node* next;
} NODE;
class LinkedList
{
protected:
    NODE* front;
    NODE* back;
public:
    LinkedList();
    ~LinkedList();
    void appendNode(int);
    void displayNodes();
    void displayNodesReverse();
    void destroyList();
};
//LinkedList.cpp
#include "LinkedList.h"
LinkedList::LinkedList()
{
    front = NULL;
```

```

        back = NULL;
    }
LinkedList::~LinkedList()
{
    destroyList();
}
void LinkedList::appendNode(int data)
{
    NODE* n = new NODE(data);
    if(back == NULL)
    {
        back = n;
        front = n;
    }
    else
    {
        back->next = n;
        n->previous = back;
        back = n;
    }
}
void LinkedList::displayNodes()
{
    cout << "Nodes:";
    NODE* temp = front;
    while(temp != NULL)
    {
        cout << " " << temp->data;
        temp = temp->next;
    }
}
void LinkedList::displayNodesReverse()
{
    cout << "Nodes in reverse order:";
    NODE* temp = back;
    while(temp != NULL)
    {
        cout << " " << temp->data;
        temp = temp->previous;
    }
}
void LinkedList::destroyList()
{
    NODE* temp = back;
    while(temp != NULL)

```

```

    {
        NODE* temp2 = temp;
        temp = temp->previous;
        delete temp2;
    }
    back = NULL;
    front = NULL;
}

```

#### 10.4.2 StackLinkedList Header File and StackLinkedList Source File

The StackLinkedList.h file contains the definition of the StackLinkedList class, as shown next. Below the StackLinkedList.h file is the StackLinkedList.cpp file that contains the definitions of member functions.

The class definition and each member function were explained in the “The StackLinkedList Class” section of this chapter.

```

//StackLinkedList.h
class StackLinkedList : public LinkedList
{
public:
    StackLinkedList();
    virtual ~StackLinkedList();
    void push(int);
    int pop();
    bool isEmpty();
};
//StackLinkedList.cpp
StackLinkedList.h
StackLinkedList::StackLinkedList()
{
}
StackLinkedList::~~StackLinkedList()
{
}
void StackLinkedList::push(int x)
{
    appendNode(x);
}
int StackLinkedList::pop()
{
    if(isEmpty())
    {
        return -1;
    }
    int retVal = back->data;

```

```

    NODE* temp = back;
    if(back->previous == NULL)
    {
        back = NULL;
        front = NULL;
    }
    else
    {
        back = back->previous;
        back->next = NULL;
    }
    delete temp;
    return retVal;
}

bool StackLinkedList::isEmpty()
{
    if(front == NULL)
    {
        return true;
    }
    else
    {
        return false;
    }
}

```

### 10.4.3 StackLinkedList Application

The StackLinkedListDemo.cpp file contains the actual stack application, as shown in the following code listing. The application begins by declaring an instance of the StackLinkedList class. Remember that this statement also indirectly calls the constructor of the LinkedList class, which is inherited by the StackLinkedList class.

The application then calls the push() member function to push the values 10, 20, and 30 onto the stack. The displayNodes() member function is then called to display the values on the stack.

The pop() member function is then called to remove the last node on the stack, which is then displayed on the screen (see Figure 10.3). The program then calls the delete operator to remove the stack from memory.

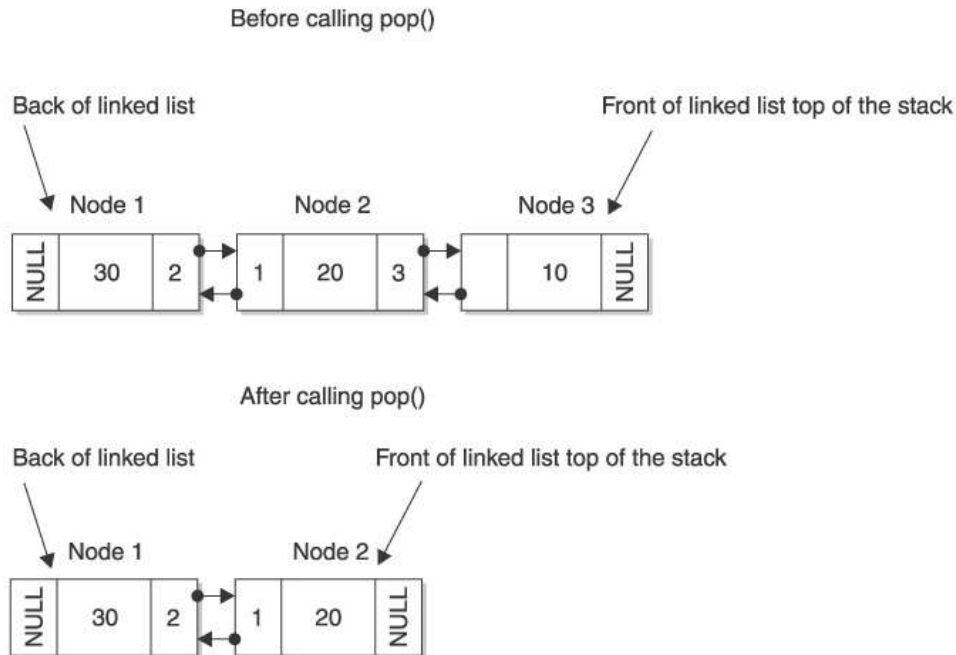


Figure 10.3: Before the pop() member function is called, there are three nodes on the stack. Two nodes remain after pop() is called.

Here's the output of this program:

Nodes: 10 20 30 10

```
//StackLinkedListDemo.cpp
#include <iostream>
using namespace std;
void main(){
    StackLinkedList* stack = new StackLinkedList();
    stack->push(10);
    stack->push(20);
    stack->push(30);
    stack->displayNodes();
    cout << stack->pop() << endl;
    delete stack;
}
```

## 10.5 Linked representation of Queue

Conceptually, a linked list queue is the same as a queue built using an array. Both store data. Both place data at the front of the queue and remove data from the front of the queue. However, in an array queue, data is stored in an array element. In a linked list queue, data is stored in a node of a linked list. The linked list queue consists of three major components: the node, the LinkedList class definition, and the QueueLinkedList class definition. Collectively, they are assembled to organized data into a queue.

Node contains the data and pointers to the previous node and the next node on the linked list (Figure 10.4). The next code snippet is the user-defined data type structure

node. You'll be using the following user-defined data type structure in this chapter to create the linked list queue.

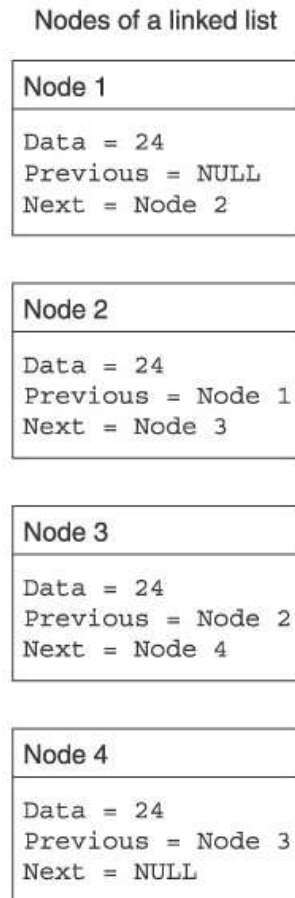


Figure 10.4: Each node points to the previous node and the next node.

The name of the user-defined data structure is called `Node` in this example and is used within the `LinkedList` class definition to declare instances of the node. The last three statements in the structure declare an integer that stores the current data and declares two pointers to reference the previous node and the next node on the linked list.

Each time a node is created, the user-defined structure is passed data for the node. Pointers to the previous node and to the next node are assigned `NULL`, which indicates there isn't a previous node or next node. `NULL` is replaced with reference to a node once the new node is added to the linked list.

```
typedef struct Node
{
    struct Node(int data)
    {
        this->data = data;
        previous = NULL;
        next = NULL;
    }
};
```



```

    }
    int data;
    struct Node* previous;
    struct Node* next;
} NODE;

```

The `LinkedList` class creates and manages the linked list. In addition, the `LinkedList` class defines member functions that manage the linked list. These are the same member functions described earlier in this lesson, a constructor and destructor, `appendNode()`, `displayNodes()`, `displayNodesReverse()`, and `destroyList()`. Here is the `LinkedList` class definition that you'll use to create the linked list queue:

```

class LinkedList
{
protected:
    NODE* front;
    NODE* back;
public:
    LinkedList();
    ~LinkedList();
    void appendNode(int);
    void displayNodes();
    void displayNodesReverse();
    void destroyList();
};

```

Programmers usually place the node structure and the `LinkedList` class definition in the same header file, `LinkedList.h`. Placing the code needed to create a linked list in one file like this helps keep it organized. Programmers then use the preprocessor directive `#include` to include `LinkedList.h` in any program that uses a linked list. The last component of the linked list queue is the `QueueLinkedList` class definition. The `QueueLinkedList` class inherits the `LinkedList` class and then defines member functions that are specifically designed to manage a queue.

You might wonder why you don't simply define one class that combines the `LinkedList` class and the `QueueLinkedList` class. Intuitively, this seems to be a good idea because everything needed to create a linked list queue is contained in one file. However, doing so repeats code, which is something programmers avoid if possible.

For example, definitions of a node and the `LinkedList` class would be located in two places. If you needed to upgrade either definition, you'd need to remember all the places where they are defined in your code. A better approach is to place each definition in its own file (for example, `LinkedList.h`, `QueueLinkedList.h`) so code won't be repeated.

### 10.5.1 The `QueueLinkedList` class

Here is the definition of the `QueueLinkedList` class that you'll use to create a queue. Programmers save this definition in a file called `QueueLinkedList.h`. The `QueueLinkedList`

class has five member functions: a constructor and destructor, enqueue(), dequeue(), and isEmpty().

```
//QueueLinkedList.h
#include "LinkedList.h"
class QueueLinkedList : public LinkedList
{
    public:
        QueueLinkedList();
        virtual ~QueueLinkedList();
        void enqueue(int);
        int dequeue();
        bool isEmpty();
};
```

The constructor and destructor of the QueueLinkedList class are empty, as shown in the next code snippet. The constructor typically initializes data members of an instance of the class. In the case of the linked list queue, initialization is performed by the constructor of the LinkedList class which is called before the constructor of the QueueLinkedList class. This means there isn't anything for the constructor of the QueueLinkedList class to do.

The destructor typically frees memory used by an instance of a class. The linked list used for the queue is removed by the destructor of the LinkedList class, which is also called before the destructor of the QueueLinkedList class. Therefore, there isn't anything for the destructor of the QueueLinkedList to do either.

```
QueueLinkedList::QueueLinkedList() {
}
QueueLinkedList::~~QueueLinkedList() {
}
```

### 10.5.2 Enqueue

The enqueue() member function of the QueueLinkedList class is called whenever a new node is placed on the queue. As you see from the function definition in the next code snippet, the enqueue() member function is sparse because it contains only one statement, which calls the appendNode() member function of the LinkedList class.

You don't have to include additional statements in the enqueue() member function because placing a node on the queue is the same process as appending a node to the linked list. Each new node is placed at the back of the linked list. Therefore, the appendNode() member function is all you need.

You may wonder why the new node is being placed on the back of the queue, but it's just because you're reusing the same code in the LinkedList class. The new node will be placed on the back of the queue like a line at the grocery store. Nodes will be pulled off the front.

The enqueue() member function has one argument, which is the data that is being assigned to the new node. In this example, the node is used to store an integer. However, you can store any type of data in a node. In fact, the data can be a pointer to a set of data such as student information. To change this example from integer data to another type of data, you'd need to change the data element in the Node structure to reflect the type of data you want to store in the node.

Data received by the enqueue() member function is passed to the appendNode() member function. Figure 10.5 illustrates how the appendNode() member function places a new node at the back of the linked list. At the top of the illustration is a linked list that contains two nodes. The appendNode() is then called to add a new node to the back of this linked list.

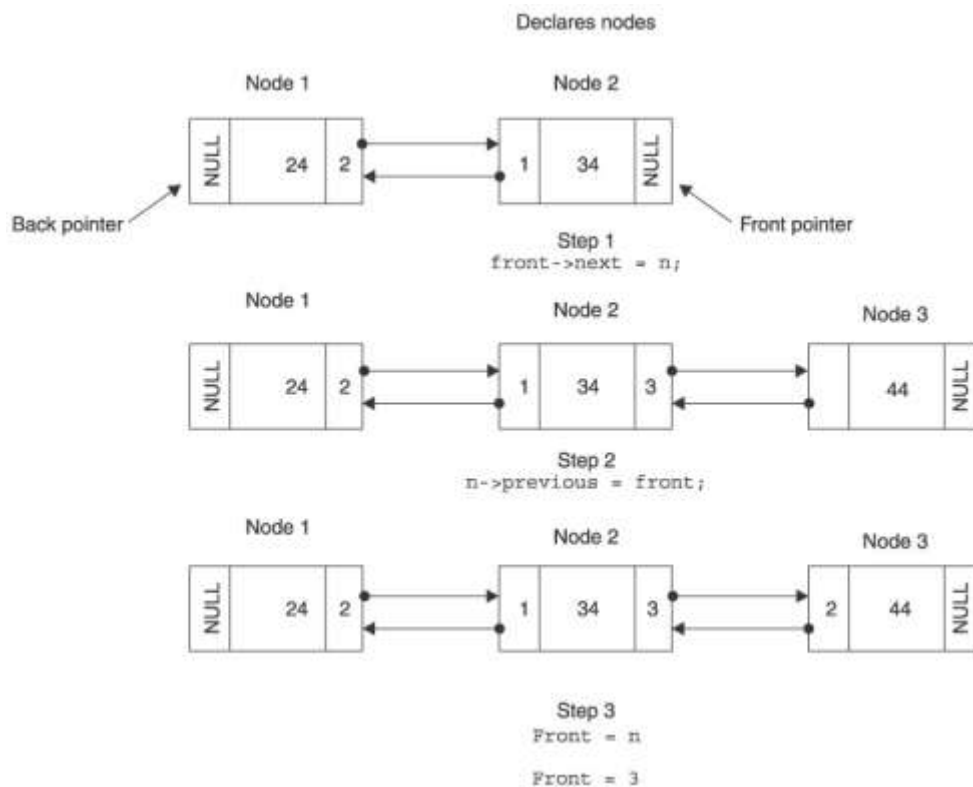


Figure 10.5: A new node is added to the queue at the back of the linked list.

The first step in this process assigns a reference to the new node to the next member of the front node. The front node is Node 2 and is assigned the reference Node 3 as the value of the next node in the linked list. This makes Node 3 the back of the linked list.

The second step assigns reference to Node 2 as the value of the previous node in Node 3. This means the program looks at the value of the previous node of Node 3 to know

which node comes before Node 3 in the linked list. The last step is to assign Node 3 as the new value of the back data member of the LinkedList class.

```
void enqueue(int x){
    appendNode(x);
}
```

### 10.5.3 Dequeue

The dequeue() member function of the QueueLinkedList class removes a node from the front of the queue. Unfortunately, there aren't any member functions in the LinkedList class that remove a node from the back of the linked list. Therefore, the dequeue() member function must do the job.

The dequeue() function begins by determining if there are any nodes on the queue by calling the isEmpty(). The isEmpty() member function returns a Boolean true if the queue is empty, in which case the dequeue() returns a -1. A Boolean false is returned if there is at least one node on the queue.

Figure 10.6 shows how the dequeue() member function works. You'll notice there are three nodes on the queue, so the isEmpty() member function returns a Boolean false, causing the program to remove the front node from the queue.

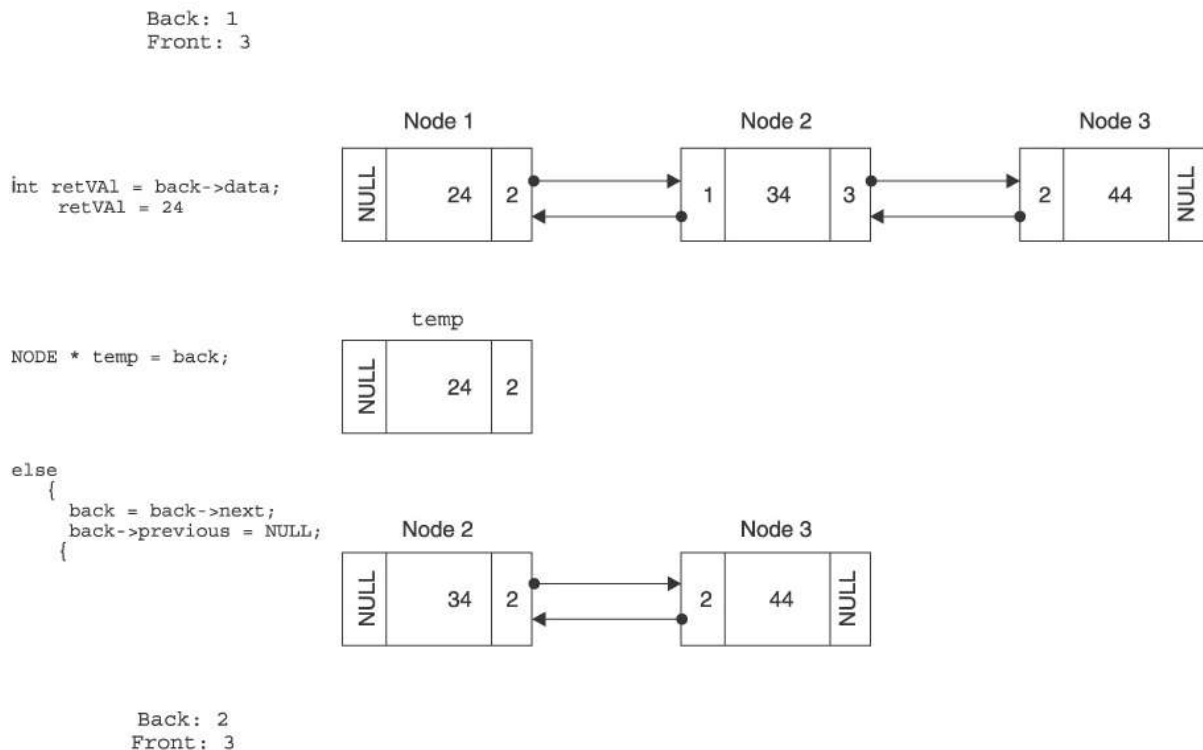


Figure 10.6: Node 1 is removed from the back of the queue by the dequeue() member function.

The removal process starts by assigning the data of the node at the front of the queue to a variable called `retVal`. The value of the `retVal` is returned by the `dequeue()` member function in the last statement of the function. Next, reference to the front node is assigned to the `temp` variable pointer. The `delete` operator later in the function uses the `temp` variable to remove the back node from memory.

Next, the function determines if there is another node on the queue by examining the value of the `next` member of the front node. If the value of the `next` member is `NULL`, there aren't any other nodes on the queue. In this case, the `front` and `back` members of the `LinkedList` class are set to `NULL`, indicating that the queue is empty. However, if the `next` member of the front node is not `NULL`, the value of the `next` member of the front node is assigned to the `front` member of the `LinkedList` class. In this example, Node 2 is the next node following Node 1. Node 2 becomes the new front of the queue.

Notice that the previous member of Node 2 is set to Node 1. However, Node 1 no longer exists. Therefore, the previous member must be set to `NULL` because there isn't a previous node. Node 2 is the front of the queue.

The `temp` node is then deleted from memory. Remember that the `temp` node is a pointer that points to Node 1 and Node 1 no longer exists in memory. The final statement returns the value of the `retVal` variable, which is the data that was stored in Node 1.

```
int dequeue() {
    if(isEmpty()){
        return -1;
    }
    int retVal = front->data;
    NODE* temp = front;
    if(front->next == NULL){
        back = NULL;
        front = NULL;
    }
    else{
        front = front->next;
        front->previous = NULL;
    }
    delete temp;
    return retVal;
}
```

The `isEmpty()` member function determines if there are any nodes on the queue, which is called by the `dequeue()` member function. The `isEmpty()` member function examines the value of the `front` data member of the `LinkedList` class. If the value of `front` is `NULL`, then the queue is empty; otherwise, the queue has at least one node.

The `isEmpty()` member function returns a Boolean `true` if the value of `front` is `NULL`, otherwise a Boolean `false` is returned as shown in the definition of the `isEmpty()` here:

```

bool isEmpty(){
    if(front == NULL){
        return true;
    }
    else{
        return false;
    }
}

```

## 10.6 Linked List Queue Using C++

Now that you understand how to create a queue using a linked list, let's assemble all the pieces and build a working queue in C++. Programmers organize an application into several files, each containing a distinct component of the application.

In the case of the demo queue application illustrated next, there are five distinct components: the driver file (QueueLinkedListDemo.cpp), the header file that contains the definition of the node and the LinkedList class (LinkedList.h), the file that contains the implementation of member functions of the LinkedList class (LinkedList.cpp), the header file that contains the definition of the QueueLinkedList class (QueueLinkedList.h) and the file that contains the implementation of member functions of the QueueLinkedList class (QueueLinkedList.cpp).

The application is called QueueLinkedListDemo, and it uses a linked list to create a queue, as shown in the next code. The application begins by declaring an instance of the QueueLinkedList class using the new operator. It then declares a pointer to an instance of the QueueLinkedList. The pointer is called queue, which is assigned a reference to the instance created by the new operator.

The enqueue() member function is then called three times, each time another node is placed on the queue. The queue shown in Figure 10.7 depicts the queue after the last time the enqueue() method is called.

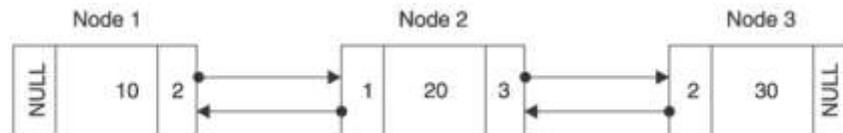


Figure 10.7: The queue after all three values are placed on the queue.

The dequeue() member function is then called to remove the first node from the queue and display its data member on the screen. Figure 10.8 shows the queue after the dequeue() member function is called.

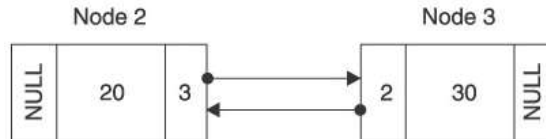


Figure 10.8: The queue after the dequeue() member function is called

The last statement in the program removes the queue from memory.

Each of the remaining components of the application was discussed in the previous section.

```
//QueueLinkedListDemo.cpp
#include <iostream>
using namespace std;
void main(){
    QueueLinkedList* queue = new QueueLinkedList();
    queue->enqueue(10);
    queue->enqueue(20);
    queue->enqueue(30);
    cout << queue->dequeue() << endl;
    delete queue;
}
//LinkedList.h
typedef struct Node
{
    struct Node(int data)
    {
        this->data = data;
        previous = NULL;
        next = NULL;
    }
    int data;
    struct Node* previous;
    struct Node* next;
} NODE;
class LinkedList
{
protected:
    NODE* front;
    NODE* back;
public:
    LinkedList();
    ~LinkedList();
    void appendNode(int);
    void displayNodes();
};
```

```

        void displayNodesReverse();
        void destroyList();
};
//LinkedList.cpp
#include "LinkedList.h"
LinkedList::LinkedList()
{
    front = NULL;
    back = NULL;
}
LinkedList::~LinkedList()
{
    destroyList();
}
void LinkedList::appendNode(int data)
{
    NODE* n = new NODE(data);
    if(front == NULL)
    {
        back = n;
        front = n;
    }
    else
    {
        back->next = n;
        n->previous = back;
        back = n;
    }
}
void LinkedList::displayNodes()
{
    cout << "Nodes:";
    NODE* temp = front;
    while(temp != NULL)
    {
        cout << " " << temp->data;
        temp = temp->next;
    }
}
void LinkedList::displayNodesReverse()
{
    cout << "Nodes in reverse order:";
    NODE* temp = back;
    while(temp != NULL)
    {

```



```

        cout << " " << temp->data;
        temp = temp->previous;
    }
}
void LinkedList::destroyList()
{
    NODE* temp = back;
    while(temp != NULL)
    {
        NODE* temp2 = temp;
        temp = temp->previous;
        delete temp2;
    }
    back = NULL;
    front = NULL;
}
//QueueLinkedList.h
#include "LinkedList.h"
class QueueLinkedList : public LinkedList
{
public:
    QueueLinkedList();
    virtual ~QueueLinkedList();
    void enqueue(int);
    int dequeue();
    bool isEmpty();
};
//QueueLinkedList.cpp
#include "QueueLinkedList.h"
QueueLinkedList::CQueueLinkedList()
{
}
QueueLinkedList::~~CQueueLinkedList()
{
}
void QueueLinkedList::enqueue(int x)
{
    appendNode(x);
}
int QueueLinkedList::dequeue()
{
    if(isEmpty())
    {
        return -1;
    }
}

```

```

    int retVal = front->data;
    NODE* temp = front;
    if(front->next == NULL)
    {
        back = NULL;
        front = NULL;
    }
    else
    {
        front = front->next;
        front->previous = NULL;
    }
    delete temp;
    return retVal;
}
bool QueueLinkedList::isEmpty()
{
    if(front == NULL)
    {
        return true;
    }
    else
    {
        return false;
    }
}

```

## 10.7 Summary

There is a serious problem with using arrays to store data for queues: you must know the size of the array when you write the program. An array can store only a specific maximum number of elements at any point in time. Programmers work around the size issue by using a linked list instead of an array when creating a stack or a queue. A linked list can grow and shrink at runtime based on the needs of the application.

Although we discuss data as being stacked like a stack of dishes, it isn't physically stacked at all. Instead, data is linked together sequentially in a list, where the last data always appears at the front of the list. Data is removed only from the front of the list. You create this sequential list by using a linked list.

Conceptually, a linked list queue is the same as a queue built using an array. Both store data. Both place data at the front of the queue and remove data from the front of the queue. However, in an array queue, data is stored in an array element. In a linked list queue, data is stored in a node of a linked list.

## 10.8 Questions

1. What is a stack-linked list?
2. How does a stack-linked list differ from a linked list?
3. What is the benefit of using a stack-linked list?
4. Where is the front of the stack in a stack-linked list?
5. What is the maximum number of nodes that you can have on a stack-linked list?
6. Can a node on a stack-linked list have more than one data element?
7. Why does the StackLinkedList class inherit the LinkedList class?
8. What happens when you push a new node onto a stack?
9. What is a queue linked list?
10. How does a queue linked list differ from an array queue?
11. What is the benefit of using a queue linked list?
12. Where are new nodes added to the queue?
13. Which node is removed from the queue when the dequeue() member method is called?
14. Can a node on a queue linked list have more than one data element?
15. What form of access is used to add and remove nodes from a queue?
16. Why does the QueueLinkedList class inherit the LinkedList class?
17. What happens when dequeue() is called?

### **10.9 Suggested readings**

1. A. Tanenbaum, Y. Lanhgsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley.

## TREES

### 11.1 Objectives

### 11.2 Introduction

### 11.3 Structure of a binary tree

### 11.4 Terminology

### 11.5 Need of a binary tree

### 11.6 Types of binary trees

11.6.1 Binary search tree

11.6.2 Height balanced tree (AVL)

11.6.3 Heaps

### 11.7 Traversing a binary tree

### 11.8 Summary

### 11.9 Questions

### 11.10 Suggested readings

### 11.1 Objectives

In this lesson, you'll learn about the concept of trees, different terms related with trees and types of trees.

### 11.2 Introduction

When you read the word "tree" in the title of this lesson, you probably imagined your favorite tree covered with a full coat of green leaves. However, this doesn't truly represent the tree that we'll be talking about. Instead, envision a tree barren of leaves, where all you can see are branches stretched in all directions. Each stem terminates with only two branches, similar to a fork in the road. Those branches lead to other stems and other forks.

This type of tree is a binary tree. Binary means two, as you learned when you studied the binary numbering systems in your first computer course. The binary numbering system consists of two digits, zero and one.

A **binary tree is a tree where each stem has not more than two branches** (see Figure 11.1). Typically, the stem has two branches, but there can be situations when the stem has one branch or simply terminates, resulting in no additional branches.

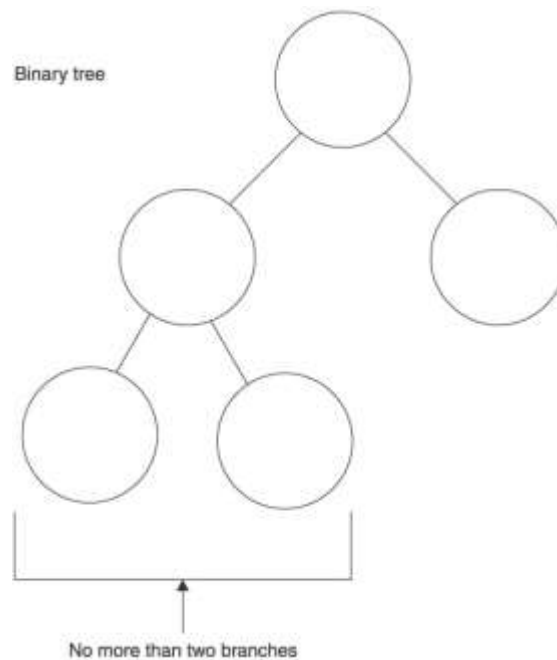


Figure 11.1: A binary tree is a tree where each stem has no more than two branches.

Programmers use a binary tree as a model to create a data structure to encode logic used to make complex decisions. Here’s how this works. Let’s say that a stem consists of a set of program instructions. At the end of the stem, the program evaluates a binary expression. You’ll recall that a binary expression evaluates to either a Boolean true or false. Based on the evaluation, the program proceeds down one of two branches. Each branch has its own set of program instructions.

The basic concept of a binary tree isn’t new to you because it uses Boolean logic that you learned to implement using an if statement in your program. An if statement evaluates an expression that results in a Boolean value. Depending on the Boolean value, the if statement executes one of two sets of instructions. However, a binary tree is much more than an if statement, as you’ll learn in this lesson.

### 11.3 Structure of a binary tree

Although we introduced the concept of a binary tree with terms commonly used when referring to a tree, programmers established different terms to refer to parts of a binary tree. Let’s take a moment to become familiar with those terms.

“Node” is the term used to describe a termination point. There are three kinds of termination points in a binary tree (see Figure 11.2): the starting node, the ending node,

and the branch node. The starting node is called the root node, which is the top-level node in the tree. The stem leading from the root node leads to the branch node. The branch node is the fork in the road that links the root node to two branches. Each branch terminates with a child node. Programmers call these the left branch and the right branch.

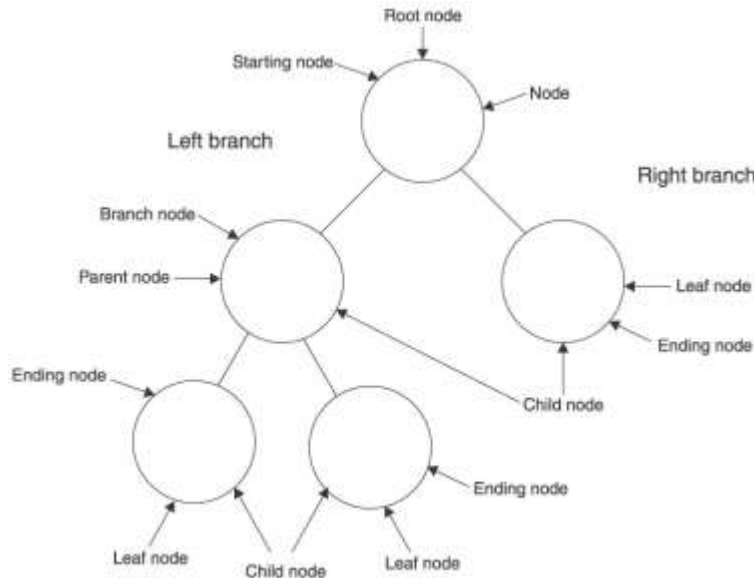


Figure 11.2: A binary tree is comprised of several nodes, each of which are related to other nodes on the tree.

As you can see, a binary tree defines a strong parent-child relationship among nodes. A parent-child relationship is relative to a node. All nodes except the root node have a parent node. However, some nodes have no children, while other nodes have one or two child nodes. Programmers determine the parent-child relationship by selecting a node, which is called the current node. The node that spawns the current node is called the current node's parent node. The node or nodes spawned by the current node is called the child node.

The child node is also referred to as the left node or the right node, depending on which direction the node branches from the current node. If the current node doesn't have any child nodes, then the current node is referred to as the leaf node. A leaf node is located at the bottom of the tree. If you look at your favorite tree, you'll notice that the end of nearly every branch is either another branch (child node) or a leaf, and that's why programmers call a node with no child nodes a leaf node.

### Depth and Size

A binary tree is described by using two measurements: tree depth and tree size (see Figure 11.3). Tree depth is the number of levels in the tree. A new level is created each time a current node branches to a child node. For example, one level is created when the root node branches into child nodes.

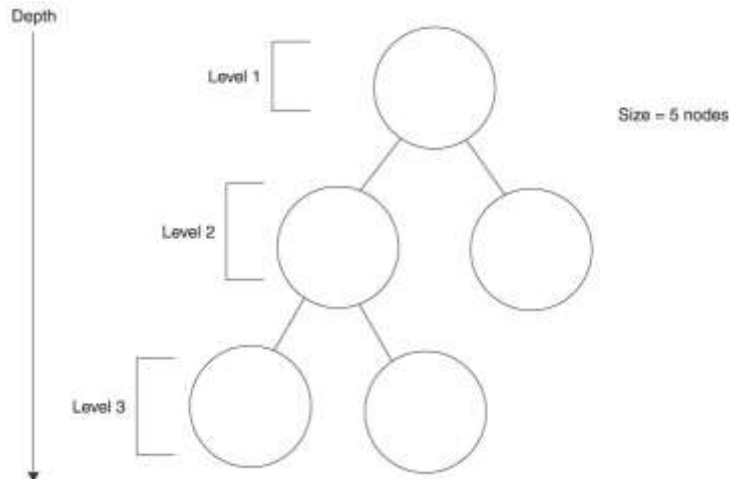


Figure 11.3: The number of levels in a tree defines a tree's depth, and the number of nodes defines the size of the tree

The size of a tree is the number of nodes in the tree. For example, the first level in Figure 11.3 has one node, which is the root node. The second level has up to two nodes, which are the child nodes of the root. The third level may have up to four nodes. Programmers estimate the size of a tree by using the following formula.

$$\text{size} = 2^{\text{depth}}$$

Let's say the binary tree has five levels, which is a depth of 5. Here's how you estimate the size of the tree:  $32 = 2^5$ .

The size is an approximation because the tree may or may not be balanced. A balanced tree is a binary tree where each current node has two child nodes. An unbalanced tree is a binary tree where one or more current nodes have fewer than two child nodes. This formula gives you a rough idea of how well balanced the tree is. Binary trees are usually used for very large data sets. The formula is not terribly accurate for the small tree shown in Figure 11.3.

#### 11.4 Terminology

Terminology describing family relationships is frequently used to describe relationships between the nodes of a tree  $T$ . Suppose  $N$  is a node of the tree with left successor  $S_L$  and right successor  $S_R$ . Then  $N$  is called parent of  $S_L$  and  $S_R$ .  $S_L$  is also called left child of  $N$  and  $S_R$  is also called right child of  $N$ . Furthermore  $S_L$  and  $S_R$  are called siblings. Each node  $N$  in a binary tree has a unique parent, except the root node, called the predecessor of  $N$ . The terms descendant and ancestor have their usual meaning. That is, a node  $L$  is called a descendant of node  $N$  (and  $N$  is called an ancestor of  $L$ ) if there is a succession of children from  $N$  to  $L$ . In particular,  $L$  is called left or right descendant of  $N$  according to whether  $L$  belongs to the left subtree or right subtree of  $N$ .

Terminology from graph theory and horticulture is also used with a binary tree  $T$ . The line drawn from node  $N$  of  $T$  to a successor is called an edge and a sequence of consecutive edges is called a path. A terminal node is called leaf node and a path ending in a leaf node is called a branch.

As discussed above each node of tree  $T$  can have at most two children and that level  $r$  of  $T$  can have at most  $2^r$  nodes. The tree  $T$  is said to be complete if all its levels, except possibly the last, have the maximum number of possible nodes and if all the nodes at the last level appear as far left as possible. Tree in figure 11.3 is a complete binary tree but the one in figure 11.4 is not.

### **11.5 Need of a binary tree**

Programmers use a binary tree to quickly look up data stored at each node on the binary tree. Let's say that you need to find a student ID in a list of a million student IDs. What is the maximum number of comparisons that you'll need to make before finding the student ID?

You could make a maximum of a million comparisons if you sequentially searched the list of a million student IDs. More than a million comparisons are necessary if you randomly selected student IDs from the list and then replaced those that didn't match back into the list.

However, you'd need to make a maximum of only 20 comparisons if student IDs were stored in a binary tree. This is because of the way data is organized in a binary tree. Data stored on the left node is less than data stored on all the right nodes at any current node.

This might sound a little confusing, but an example will make this concept clear. Suppose you had a list of five student IDs: 101, 102, 103, 104 and 105. These student IDs are stored in a binary tree so that the center student ID is the root node, the student ID that is less than the current node is stored on the left child node, and the student ID that is more than the current node is stored on the right child node, as shown in Figure 11.4.



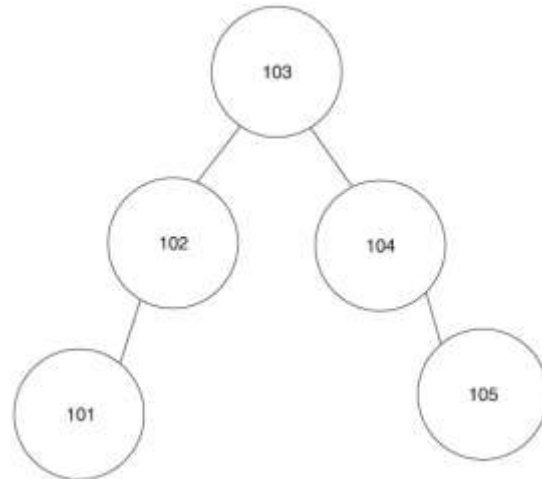


Figure 11.4: The left child node is always less than the parent node, and the right child node is always greater than the parent node.

The same pattern is applied to each child node. Thus, student ID 101 is the left child of the node that contains student ID 102 because student ID 101 is less than student ID 102. Likewise, student ID 105 is the right node of student ID 104 because student ID is greater than student ID 105.

Let's say you want to locate student ID 101 on the binary tree. First, you compare the value of the root node to student ID 101. There's no match. Because student ID 101 is less than student ID 103, your next comparison uses the left child node. This eliminates the need to compare all the nodes to the right of the node that contains student ID 103. You can ignore half the student IDs because you know that student ID 101 isn't the right node or a child of the right node.

After comparing student ID 101 to student ID 102, you notice two things. First, they don't match. Second, student ID 102 is greater than student ID 101. This means you compare the left child node to student ID 101. You ignore the right child node and subsequent child nodes because they are greater than student ID 101. There aren't any right child nodes of student ID 102 in this example. The next comparison results in a match. So, in a large binary tree, each time you do a comparison, you eliminate another half of the remaining nodes from the search. If you had 1 million nodes in the tree, you would divide 1 million by 2 about 20 times to reduce it down to one node ( $2^{20}$  is about 1 million). This way, you can find the node you're looking for by doing about 20 comparisons.

Programmers think of every node as the root node and all subsequent nodes as its own subtree. They approach nodes in this way because functions that deal with trees are recursive. A function works on a child node and performs the same functionality as if the child node is the root node of the entire tree. That is, the value of the child node is

compared to the value of its left node and right node to determine which branch of the tree to pursue.

### The Key

Each node of a tree contains a key that is associated with a value similar to the relationship between a primary key of a database and a row within a table of the database. A key is a value that is compared to search criteria. If the index and the search criteria match, then the application retrieves data stored in the row that corresponds to the key. The data is referred to as the value of the node, as shown in Figure 11.5.

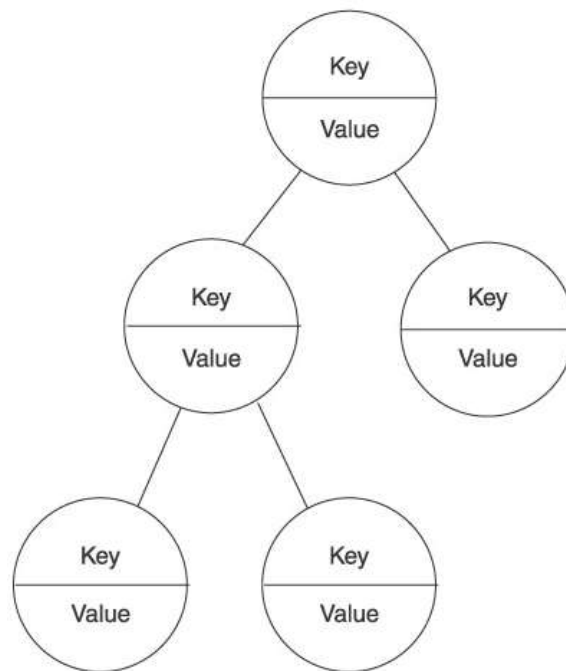


Figure 11.5: Each node has an index and a value; the index uniquely identifies the node and retrieves the value of a node.

Any data type can be used as the key. We use a string as the key in examples in this lesson, but we could have chosen any data type. Unlike the primary key of a database, the key to the tree doesn't have to be in natural order. That is, the key doesn't have to be in alphabetical order or numerical order. In a typical tree implementation, you would define a comparator to tell the tree how to order the nodes. In our case, we'll use the natural ordering sequence of a string so we can keep our focus on the internal workings of the tree.

### 11.6 Types of binary trees

Tree is just an abstract data type. For practical implementation and use we need to define some variant of the binary tree. Following are some of the types of binary trees:

### 11.6.1 Binary search tree

A binary tree which conforms to the following properties is called a **binary search tree**.

Properties:

- Each value (key) in the tree exists at most once (i.e. no duplicates).
- The "greater-than" and "less-than" relations are well defined for the data value.
- Sorting constraints:- for every node  $n$  :
  - o All data in the left subtree of  $n$  is **less than** the data in the root of that subtree.
  - o All data in the right subtree of  $n$  is **greater than** the data in the root of that subtree.

The examples of binary search trees are given in figure 11.6

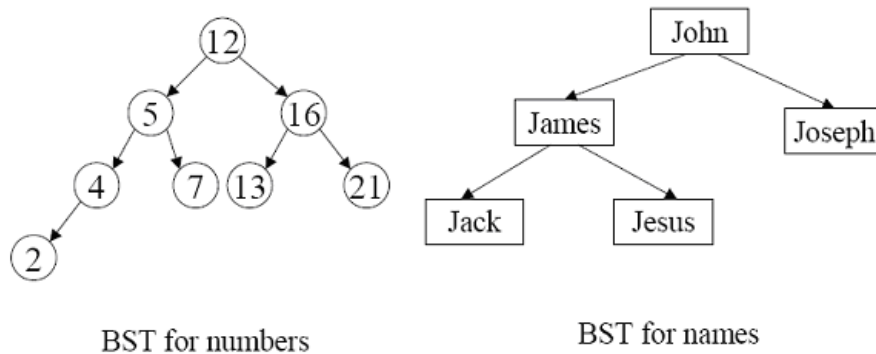


Figure 11.6: Examples of binary search tree

### 11.6.2 Height balanced tree (AVL)

A height balanced tree also known as an AVL tree is a special type of binary tree that is always "partially" balanced. The criteria that is used to determine the "level" of "balanced-ness" is the difference between the heights of subtrees of a root in the tree. The "height" of tree is the "number of levels" in the tree. Or to be more formal, the height of a tree is defined as follows:

1. The height of a tree with no elements is 0
2. The height of a tree with 1 element is 1
3. The height of a tree with > 1 element is equal to 1 + the height of its tallest subtree.

An AVL tree is a binary tree in which the difference between the height of the right and left subtrees (or the root node) is never more than one. The following figure 11.7 shows example of an AVL tree.

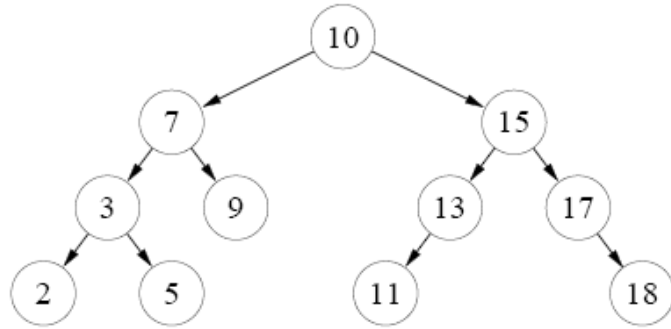


Figure 11.7: Example of an AVL tree

### 11.6.3 Heaps

A binary tree has the heap property iff

- it is empty or
- the key in the root is larger than that in either child and both subtrees have the heap property.

A heap can be used as a priority queue: the highest priority item is at the root and is trivially extracted. But if the root is deleted, we are left with two sub-trees and we must efficiently re-create a single tree with the heap property. The value of the heap structure is that we can both extract the highest priority item and insert a new one in  $O(\log n)$  time.

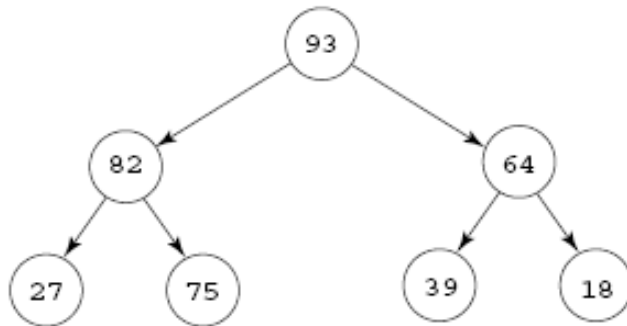


Figure 11.8: Example of a Heap

## 11.7 Traversing a binary tree

Traversing, as discussed in lesson 1, means visiting each node exactly once. There are three methods of traversing a binary tree.

**11.7.1 PreOrder:** If T is the root node, L is the left child node and R is the right child node of it, then the algorithm proceed as follows: (T L R)

1. Process the root T.

2. Traverse the left subtree L of T in preorder.
3. Traverse the right subtree R of T in preorder.

### 11.7.2 InOrder

1. Traverse the left subtree L of T in inorder.
2. Process the root T.
3. Traverse the right subtree R of T in inorder.

### 11.7.3 PostOrder

1. Traverse the left subtree L of T in postorder.
2. Traverse the right subtree R of T in postorder.
3. Process the root T.

The above three algorithms contain the same three steps and that the left subtree L of T is always traversed before right subtree R of T. The difference lies in the traversal time of root T. In preorder the root t is processed before subtrees, in inorder root T is processed between traversals of subtrees and in postorder root T is processed after traversing subtrees.

Consider the tree given in figure 11.9, the results of traversing the tree using the three techniques discussed above are as follows:

Preorder: a b d h e c f g

Inorder: h d b e a f c g

Postorder: h d e b f g c a

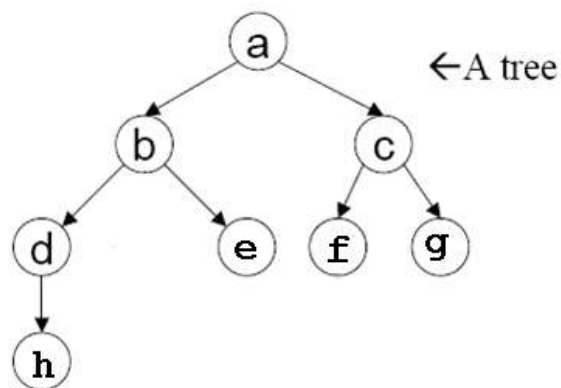


Figure 11.9: A binary tree with labeled nodes.

## 11.8 Summary

Tree is a data structure which has a root not from which all other nodes stem out. A tree with two children is called a binary tree and there are different types of binary trees like binary search tree, height balanced tree and heap. Tree depth is the number of levels in the tree. The size of a tree is the number of nodes in the tree. A tree T is said to be complete if all its levels, except possibly the last, have the maximum number of possible nodes and if all the nodes at the last level appear as far left as possible. Binary trees provide very efficient searching facility.

## 11.9 Questions

1. What is a tree?
2. What is the relationship between parent node and child nodes?
3. What are a key and a value?
4. What is a root node?
5. What is the purpose of a left child node and a right child node?
6. What is a leaf node?
7. What is the depth of a tree?
8. What is the size of a tree? How do you calculate the size of a tree?
9. What are the different types of binary trees? Briefly discuss each type.

## 11.10 Suggested readings

1. A. Tanenbaum, Y. Lanhgsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley.

**BINARY SEARCH TREE****12.1 Objectives****12.2 Introduction****12.3 Operations on Binary Search Tree****12.4 Types of BST****12.5 Uses of BST****12.6 Summary****12.7 Questions****12.8 Suggested readings****12.1 Objectives**

In this lesson we will discuss binary search tree (BST), which is a type of binary tree, operations performed on BST, type of BST and uses of BST.

**12.2 Introduction**

**In computer science, a binary search tree (BST) is a binary tree data structure which has the following properties:**

- Each node (item in the tree) has a distinct value.
- Both the left and right subtrees must also be binary search trees.
- The left subtree of a node contains only values less than the node's value.
- The right subtree of a node contains only values greater than or equal to the node's value.

The major advantage of binary search trees over other data structures is that the related sorting algorithms and search algorithms such as in-order traversal can be very efficient. Binary search trees can choose to allow or disallow duplicate values, depending on the implementation. Binary search trees are a fundamental data structure used to construct more abstract data structures such as sets, multisets and associative arrays.

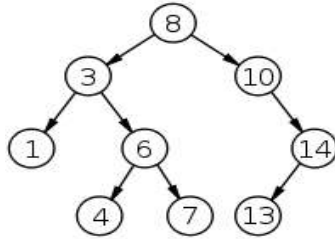


Figure 12.1: A binary search tree of size 9 and depth 3, with root 8 and leaves 1, 4, 7 and 13

An Example: Figure 12.2 shows a binary search tree. Notice that this tree is obtained by inserting the values 13, 3, 4, 12, 14, 10, 5, 1, 8, 2, 7, 9, 11, 6, 18 in that order, starting from an empty tree.

Note that inorder traversal of a binary search tree always gives a sorted sequence of the values. This is a direct consequence of the BST property. This provides a way of sorting a given sequence of keys: first, create a BST with these keys and then do an inorder traversal of the BST so created.

Note that the highest valued element in a BST can be found by traversing from the root in the right direction all along until a node with no right link is found (we can call that the rightmost element in the BST).

The lowest valued element in a BST can be found by traversing from the root in the left direction all along until a node with no left link is found (we can call that the leftmost element in the BST).

Search is straightforward in a BST. Start with the root and keep moving left or right using the BST property. If the key we are seeking is present, this search procedure will lead us to the key. If the key is not present, we end up in a null link.

Insertion in a BST is also a straightforward operation. If we need to insert an element  $x$ , we first search for  $x$ . If  $x$  is present, there is nothing to do. If  $x$  is not present, then our search procedure ends in a null link. It is at this position of this null link that  $x$  will be included.

If we repeatedly insert a sorted sequence of values to form a BST, we obtain a completely skewed BST. The height of such a tree is  $n - 1$  if the tree has  $n$  nodes. Thus, the worst case complexity of searching or inserting an element into a BST having  $n$  nodes is  $O(n)$ .



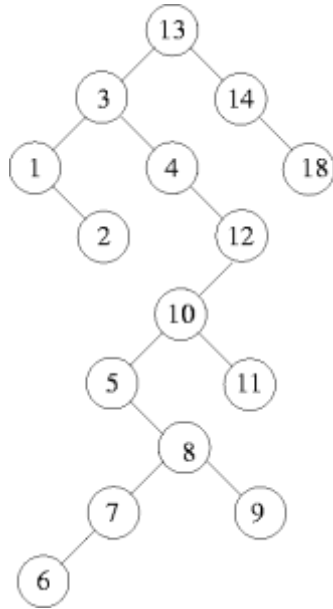


Figure 12.2: An example of a binary search tree

### 12.3 Operations on Binary Search Tree

Operations on a binary tree require comparisons between nodes. These comparisons are made with calls to a comparator, which is a subroutine that computes the total order (linear order) on any two values. This comparator can be explicitly or implicitly defined, depending on the language in which the BST is implemented.

#### Searching

Searching a binary tree for a specific value can be a recursive or iterative process. This explanation covers a recursive method.

We begin by examining the root node. If the tree is null, the value we are searching for does not exist in the tree. Otherwise, if the value equals the root, the search is successful. If the value is less than the root, search the left subtree. Similarly, if it is greater than the root, search the right subtree. This process is repeated until the value is found or the indicated subtree is null. If the searched value is not found before a null subtree is reached, then the item must not be present in the tree.

This operation requires  $O(\log n)$  time in the average case, but needs  $O(n)$  time in the worst-case, when the unbalanced tree resembles a linked list (degenerate tree).

#### Insertion

Insertion begins as a search would begin; if the root is not equal to the value, we search the left or right subtrees as before. Eventually, we will reach an external node and add the value as its right or left child, depending on the node's value. In other words, we

examine the root and recursively insert the new node to the left subtree if the new value is less than the root, or the right subtree if the new value is greater than or equal to the root.

Here's how a typical binary search tree insertion might be performed in C++:

```
/* Inserts the node pointed to by "newNode" into the subtree rooted at  
"treeNode" */  
void InsertNode(Node *&treeNode, Node *newNode)  
{  
    if (treeNode == NULL)  
        treeNode = newNode;  
    else if (newNode->key < treeNode->key)  
        InsertNode(treeNode->left, newNode);  
    else  
        InsertNode(treeNode->right, newNode);  
}
```

The above "destructive" procedural variant modifies the tree in place. It uses only constant space, but the previous version of the tree is lost.

The part that is rebuilt uses  $\Theta(\log n)$  space in the average case and  $\Omega(n)$  in the worst case (see big-O notation).

In either version, this operation requires time proportional to the height of the tree in the worst case, which is  $O(\log n)$  time in the average case over all trees, but  $\Omega(n)$  time in the worst case.

Another way to explain insertion is that in order to insert a new node in the tree, its value is first compared with the value of the root. If its value is less than the root's, it is then compared with the value of the root's left child. If its value is greater, it is compared with the root's right child. This process continues, until the new node is compared with a leaf node, and then it is added as this node's right or left child, depending on its value.

There are other ways of inserting nodes into a binary tree, but this is the only way of inserting nodes at the leaves and at the same time preserving the BST structure.

## Deletion

For the purpose of deletion there are several cases to be considered:

- **Deleting a leaf:** Deleting a node with no children is easy, as we can simply remove it from the tree.
- **Deleting a node with one child:** Delete it and replace it with its child.
- **Deleting a node with two children:** Suppose the node to be deleted is called N. We replace the value of N with either its in-order successor (the left-most child of the right subtree) or the in-order predecessor (the right-most child of the left subtree).

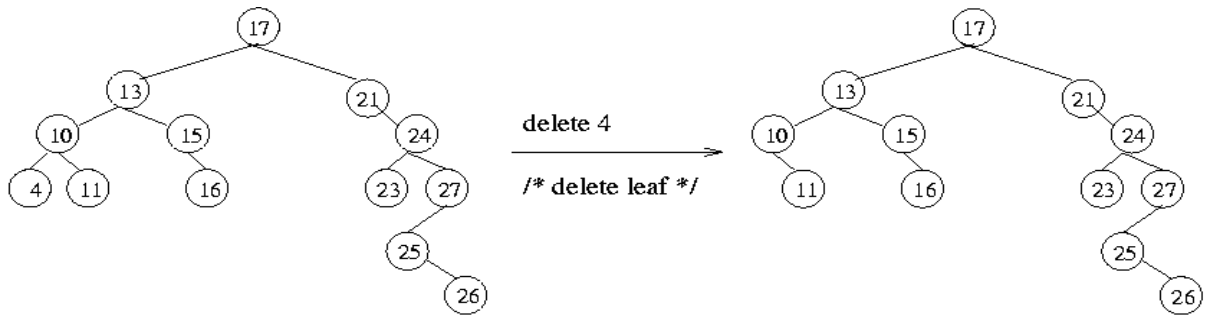


Figure 12.3: Deletion of 4 (a leaf node)

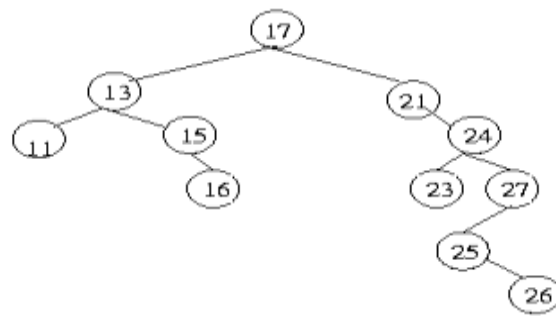


Figure 12.4 Delete 10 (delete a node with no left subtree)

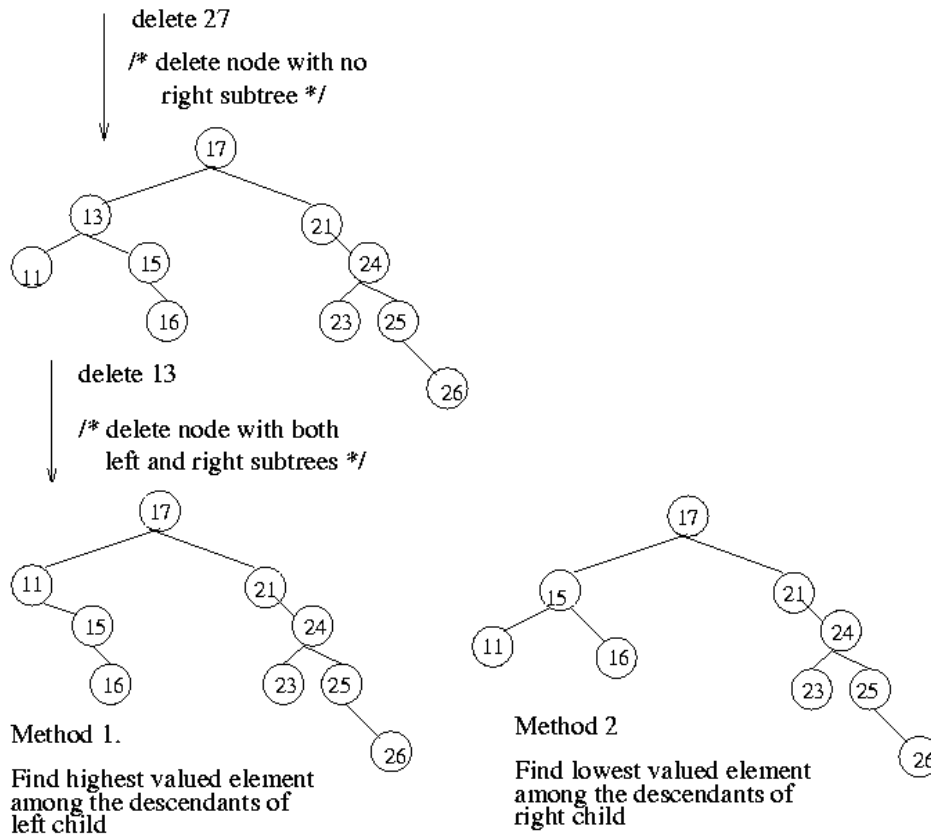


Figure 12.4 Delete 27 (with no right subtree), 13 (with both left and right subtree)

Once we find either the in-order successor or predecessor, swap it with  $N$ , and then delete it. Since both the successor and the predecessor must have fewer than two children, either one can be deleted using the previous two cases. A good implementation avoids consistently using one of these nodes, however, because this can unbalance the tree.

Here is C++ sample code for a recursive version of deletion.

```
struct BinarySearchTree* BstDelete(struct BinarySearchTree *node,int val)
{
    struct BinarySearchTree *successor,*node_delete;
    if(node)
    {
        if(node->val>val)
            node->left=BstDelete(node->left,val);
        else if(node->val<val)
            node->right=BstDelete(node->right,val);
        else
        {
            if(node->left==NULL)
            {
```

```

        node_delete=node;
        node=node->right;
        free(node_delete);
    }
    else if (node->right==NULL)
    {
        node_delete=node;
        node=node->left;
        free(node_delete);
    }
    else /*use inorder_predecessor every alternate delete
for better tree balancing*/
    {
        successor=inorder_successor_Bst (node->right);
        node->val=successor->val;
        node->right=Bst_Delete (node->right,successor->
val);
    }
    return node;
}
else
{
    printf("\nNode to be deleted not found\n");
    return NULL;
}
}

```

Although this operation does not always traverse the tree down to a leaf, this is always a possibility; thus in the worst case it requires time proportional to the height of the tree. It does not require more even when the node has two children, since it still follows a single path and does not visit any node twice.

## **Traversal**

There are several well-known ways to traverse, to travel throughout, a binary tree. We will look at three of them. The first is an inorder traversal. This consists of three overall steps: traverse the left subtree (recursively), visit the root node and traverse the right subtree (recursively). When we "visit" a node we typically do some processing on it, such as printing out the contents of the node.

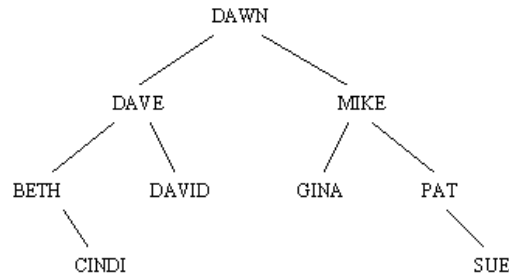


Figure 12.5 A sample BST for demonstrating traversal

For example, an inorder traversal of the binary search tree, given in figure 12.5, gives us the names in this order:

BETH

CINDI

DAVE

DAVID

DAWN

GINA

MIKE

PAT

SUE

Note that we got the data back in ascending order. This will always happen when doing an inorder traversal of a binary search tree. In fact, a sort can be done this way. One first inserts the data into a binary search tree and then does an inorder traversal to obtain the data in ascending order. Some people call this a tree sort.

Now, how exactly did we get the above list of data? Essentially, we did so by following the recursive definition for an inorder traversal. First we traverse the left subtree of the root, DAWN. That left subtree is the one rooted at DAVE. How do we traverse it? By using the same three-step process. We first traverse its left subtree, the one rooted at BETH. Of course, we then have to go through the three steps on the subtree rooted at BETH. We begin by traversing its left subtree, but it is empty, so we visit the root, BETH. That is the first data item printed. Then we traverse the right subtree, the one rooted at CINDI. We use the three-step process on it, but since its subtrees are empty, we simply print the root, CINDI, which is the second item printed. We then back up to where we left off with the subtree rooted at DAVE. We have now traversed its left subtree, so we go on to print the root, DAVE, and then traverse the right subtree. Since the right subtree itself has

empty subtrees, we end up just printing its root, DAVID. We continue in a similar fashion for the rest of this binary search tree.

The other two traversals that we will study are the preorder traversal and the postorder traversal. They are very similar to the inorder traversal in that they consist of the same three steps, but reordered slightly. The preorder traversal puts the step of visiting the root first. The postorder traversal puts the step of visiting the root last. Everything else stays the same. Here is an outline of the steps for all three of our traversals.

#### Preorder traversal

1. Visit the root
2. Traverse the left subtree
3. Traverse the right subtree

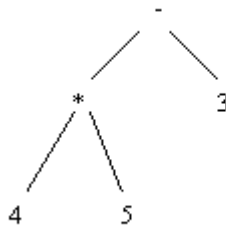
#### Inorder traversal

1. Traverse the left subtree
2. Visit the root
3. Traverse the right subtree

#### Postorder traversal

1. Traverse the left subtree
2. Traverse the right subtree
3. Visit the root

As an example, let's do a postorder traversal of the binary expression tree for  $4 * 5 - 3$  that we looked at earlier. The tree is shown again for convenience:

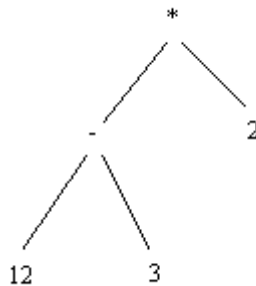


First we traverse the left subtree of the whole binary tree. This is the subtree rooted at  $*$ . To do so, we apply our three steps. We traverse its left subtree, which results in printing 4. Then we traverse the right subtree, which results in printing 5. Then we visit the root, printing  $*$ . Next, we back up to where we left off with the whole binary tree. We have now traversed the left subtree, so we traverse the right subtree, printing 3. Then we visit the root, printing  $-$ . Overall we end up printing  $4 * 5 - 3$ , the postfix form of the expression. A postfix expression is deciphered by looking at it left to right and using the fact that each operator (such as  $*$ ) applies to the two previous values. Note that the

traversal always works like this: a postorder traversal of a binary expression tree yields the postfix form of the expression. You may be familiar with postfix expressions in that some calculators use them. In ordinary mathematics we are used to using infix expressions, where operators such as + and \* come between the two values to which they apply.

For practice try a preorder traversal of the same binary expression tree for  $4 * 5 - 3$ . The result should be  $- * 4 5 3$ . This is the prefix form of the expression, that is, the form in which each binary operator precedes the two quantities to which it applies. A preorder traversal of a binary expression tree always gives the prefix form of the expression.

The natural conjecture, then, would be that the inorder traversal of a binary expression tree would produce the infix form of the expression, but that is not quite true. With the above expression it is true. However, try the infix expression  $(12 - 3) * 2$ . Here parentheses are used to indicate that the subtraction should be done before the multiplication. The binary expression tree looks like this:



As you can verify, an inorder traversal of this binary expression tree produces  $12 - 3 * 2$ , which is the infix form of a slightly different expression, one in which the multiplication is done before the subtraction. The problem is that we did not get the parentheses back. It is possible to modify the code for an inorder traversal so that it always parenthesizes things, but a plain inorder traversal does not give any parentheses.

Traversal requires  $\Omega(n)$  time, since it must visit every node. This algorithm is also  $O(n)$ , and so it is asymptotically optimal.

## Sort

A binary search tree can be used to implement a simple but efficient sorting algorithm. Similar to heapsort, we insert all the values we wish to sort into a new ordered data structure—in this case a binary search tree—and then traverse it in order, building our result:

The worst-case time of `build_binary_tree` is  $\Theta(n^2)$ —if you feed it a sorted list of values, it chains them into a linked list with no left subtrees. For example, `build_binary_tree([1, 2, 3, 4, 5])` yields the tree `(None, 1, (None, 2, (None, 3, (None, 4, (None, 5, None))))`.



There are several schemes for overcoming this flaw with simple binary trees; the most common is the self-balancing binary search tree. If this same procedure is done using such a tree, the overall worst-case time is  $O(n \log n)$ , which is asymptotically optimal for a comparison sort. In practice, the poor cache performance and added overhead in time and space for a tree-based sort (particularly for node allocation) make it inferior to other asymptotically optimal sorts such as heapsort for static list sorting. On the other hand, it is one of the most efficient methods of incremental sorting, adding items to a list over time while keeping the list sorted at all times...

## 12.4 Types of BST

There are many types of binary search trees. AVL trees and red-black trees are both forms of self-balancing binary search trees. A splay tree is a binary search tree that automatically moves frequently accessed elements nearer to the root. In a treap ("tree heap"), each node also holds a priority and the parent node has higher priority than its children.

Two other titles describing binary search trees are that of a complete and degenerate tree. A complete tree is a tree with  $n$  levels, where for each level  $d \leq n - 1$ , the number of existing nodes at level  $d$  is equal to  $2^d$ . This means all possible nodes exist at these levels. An additional requirement for a complete binary tree is that for the  $n$ th level, while every node does not have to exist, the nodes that do exist must fill from left to right. A degenerate tree is a tree where for each parent node, there is only one associated child node. What this means is that in a performance measurement, the tree will essentially behave like a linked list data structure.

## 12.5 Uses of BST

A binary search tree can be a very useful data structure. We have already seen that it can be used to create a sort routine. Such a sort routine is normally pretty fast. In fact, it is  $\Theta(n \cdot \lg(n))$  on the average. However, it does have a bad worst case, namely when the data is already in ascending or descending order. (Try it. Start with a list of data items in order and insert them one by one into a binary search tree. What does this tree look like? Why would this make it slow to access the data later when we do the inorder traversal?)

Another use of a binary search tree is in storing data items for fast lookup later. In the average case it is pretty fast to insert a new item into a binary tree, because in the average case the data is fairly random and the binary tree is reasonably "bushy". (In such a tree it is known that the height of the binary tree is  $\Theta(\lg(n))$ , so that insertion is a  $\Theta(\lg(n))$  operation.) Similarly, doing a lookup of an item already in the binary tree follows the same pattern as used when it was inserted. Thus lookup is  $\Theta(\lg(n))$  on average.

For example, to look up GINA in the binary tree above, one compares GINA to DAWN, the root. Since GINA is larger, move to the right child, MIKE. Now compare GINA to

MIKE. Since GINA is smaller, move to the left child GINA. Now compare GINA to the item in the node, also GINA and we see that we have a match. All lookups are like this. One starts at the root and follows a path from the root to the matching item (or to a leaf if no match is ever found).

## 12.6 Summary

In computer science, a binary search tree (BST) is a binary tree data structure whose nodes have distinct values and left and right subtrees must also be binary search trees. The left subtree of a node contains only values less than the node's value. The right subtree of a node contains only values greater than or equal to the node's value. The operations performed on binary search tree include insertion, deletion, searching, traversing and sorting. Traversing can be done in preorder, inorder or postorder. AVL (discussed in next lesson) and red-black are both forms of BST.

## 12.7 Questions

1. Define Binary Search Tree.
2. What are the properties of BST?
3. How sorting is performed using BST?
4. What are the different cases for deletion of a node from BST?
5. What are the different methods of traversing a BST? Write algorithm for each method.
6. What are the used of BST? Explain.
7. How many maximum nodes can be there in a BST of height  $N$ ?
8. Why duplicate keys are not permitted in BST? Explain.
9. WAP to implement BST.

## 12.8 Suggested readings

1. A. Tanenbaum, Y. Lanhsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley.

## HEIGHT BALANCED TREE

### 13.1 Objectives

### 13.2 Introduction

### 13.3 Inserting a new item

13.3.1 Case 1

13.3.2 Case 2

13.3.3 Case 3

### 13.4 Deletion form AVL tree

### 13.5 Performance of AVL tree

### 13.6 Summary

### 13.7 Questions

### 13.8 Suggested readings

### 13.1 Objectives

In this lesson, we shall discuss a binary tree with a special property that it is always partially balanced.

### 13.2 Introduction

These are self-adjusting, height-balanced binary search trees and are named after the inventors: Adelson-Velskii and Landis. A **balanced** binary search tree has  $\Theta(\lg n)$  height and hence  $\Theta(\lg n)$  worst case lookup and insertion times. However, ordinary binary search trees have a bad worst case. When sorted data is inserted, the binary search tree is very unbalanced, essentially more of a linear list, with  $\Theta(n)$  height and thus  $\Theta(n)$  worst case insertion and lookup times. AVL trees overcome this problem.

An AVL tree is a special type of binary tree that is always "partially" balanced. The criteria that is used to determine the "level" of "balanced-ness" is the difference between

the heights of subtrees of a root in the tree. The "height" of tree is the "number of levels" in the tree. Or to be more formal, the height of a tree is defined as follows:

1. The height of a tree with no elements is 0
2. The height of a tree with 1 element is 1
3. The height of a tree with > 1 element is equal to 1 + the height of its tallest subtree.

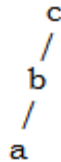
An AVL tree is a binary tree in which the difference between the height of the right and left subtrees (or the root node) is never more than one.

166

The idea behind maintaining the "AVL-ness" of an AVL tree is that whenever we insert or delete an item, if we have "violated" the "AVL-ness" of the tree in anyway, we must then restore it by performing a set of manipulations (called "rotations") on the tree. These rotations come in two flavors: single rotations and double rotations (and each flavor has its corresponding "left" and "right" versions). An example of a single rotation is as follows: Suppose I have a tree that looks like this:



Now I insert the item "a" and get the resulting binary tree:



Now, this resulting tree violates the "AVL criteria", the left subtree has a height of 2 but the right subtree has a height of 0 so the difference in the two heights is "2" (which is greater than 1). SO what we do is perform a "single rotation" (or RR for a single right rotation, or LL for a single left rotation) on the tree (by rotating the "c" element down clockwise to the right) to transform it into the following tree:

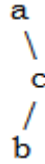


This tree is now balanced.

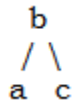
An example of a "double rotation" (or RL for a double right rotation or LR for a double left rotation) is the following: Suppose I have a tree that looks like this:



Now I insert the item "b" and get the resulting binary tree:



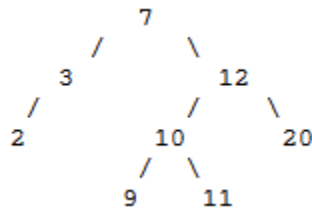
This resulting tree also violates the "AVL criteria" so we fix it by first rotating "c" down to the right (so we get "a-b-c"), and then rotating "a" down to the left so that the tree is transformed into this:



In order to detect when a "violation" of the AVL criteria occurs we need to have each node keep track of the difference in height between its right and left subtrees. We call this "difference" the "balance" factor and define it to be the height of the right subtree minus the height of the left subtree of a tree. So as long as the "balance" factor of each node is never  $>1$  or  $<-1$  we have an AVL tree. As soon as the balance factor of a node becomes 2 (or -2) we need to perform one or more rotations to ensure that the resultant tree satisfies the AVL criteria.

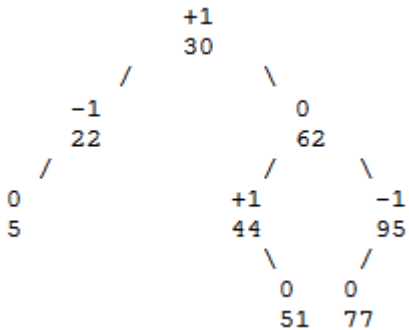
### Definitions

The **height** of a binary tree is the maximum path length from the root to a leaf. A single-node binary tree has height 0 and an empty binary tree has height -1. As another example, the following binary tree has height 3.

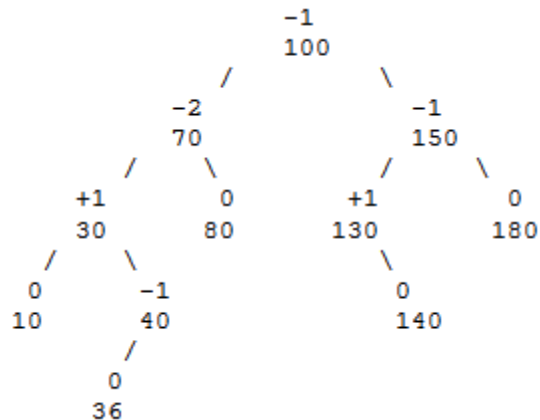


An **AVL tree** is a binary search tree in which every node is **height balanced**, that is, the difference in the heights of its two subtrees is at most 1. The balance factor of a node is the height of its right subtree minus the height of its left subtree. An equivalent definition, then, for an AVL tree is that it is a binary search tree in which each node has a balance factor of -1, 0, or +1. Note that a balance factor of -1 means that the subtree is left-heavy and a balance factor of +1 means that the subtree is right-heavy. For example, in the following AVL tree, note that the root node with balance factor +1 has a right

subtree of height 1 more than the height of the left subtree. (The balance factors are shown at the top of each node.)



The idea is that an AVL tree is close to being completely balanced. Hence it should have  $\Theta(\lg n)$  height (it does - always) and so have  $\Theta(\lg n)$  worst case insertion and lookup times. An AVL tree does not have a bad worst case, like a binary search tree which can become very unbalanced and give  $\Theta(n)$  worst case lookup and insertion times. The following binary search tree is **not** an AVL tree. Notice the balance factor of -2 at node 70.

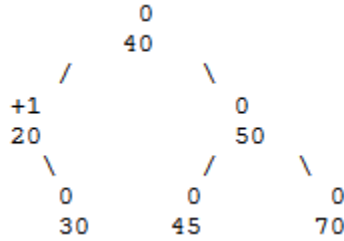


### 13.3 Inserting a new item

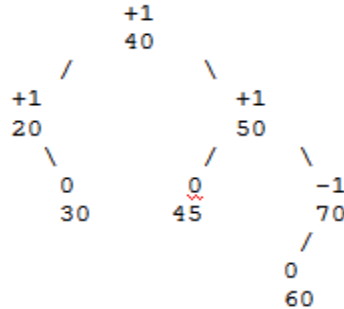
Initially, a new item is inserted just as in a binary search tree. Note that the item always goes into a new leaf. The tree is then readjusted as needed in order to maintain it as an AVL tree. There are three main cases to consider when inserting a new node.

#### 13.3.1 Case 1:

A node with balance factor 0 changes to +1 or -1 when a new node is inserted below it. No change is needed at this node. Consider the following example. Note that after an insertion one only needs to check the balances along the path from the new leaf to the root.

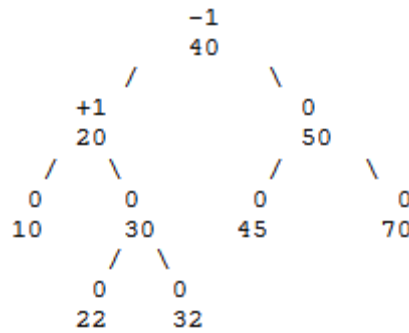


After inserting 60 we get:

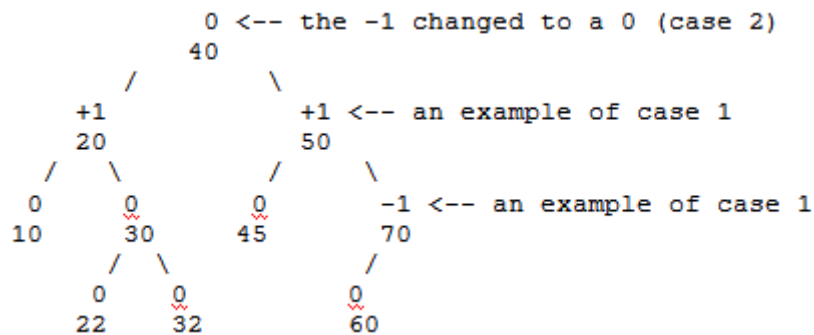


### 13.3.2 Case 2:

A node with balance factor -1 changes to 0 when a new node is inserted in its right subtree. (Similarly for +1 changing to 0 when inserting in the left subtree.) No change is needed at this node. Consider the following example.



After inserting 60 we get:

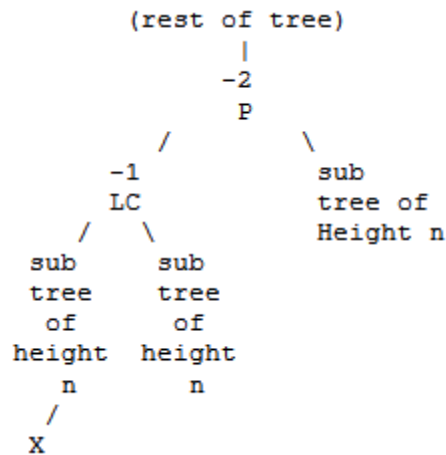


### 13.3.3 Case 3:

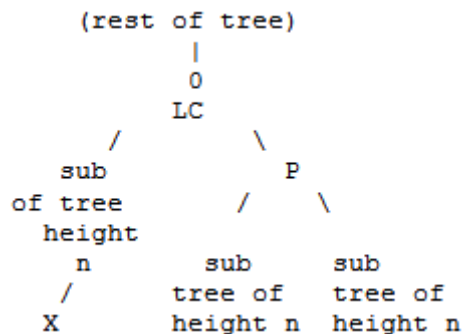
A node with balance factor -1 changes to -2 when a new node is inserted in its left subtree. (Similarly for +1 changing to +2 when inserting in the right subtree.) Change **is** needed at this node. The tree is restored to an AVL tree by using a rotation.

#### Subcase A:

This consists of the following situation, where P denotes the parent of the subtree being examined, LC is P's left child, and X is the new node added. Note that inserting X makes P have a balance factor of -2 and LC have a balance factor of -1. The -2 must be fixed. This is accomplished by doing a right rotation at P. Note that rotations do not mess up the order of the nodes given in an inorder traversal. This is very important since it means that we still have a legitimate binary search tree. (Note, too, that the mirror image situation is also included under subcase A.)

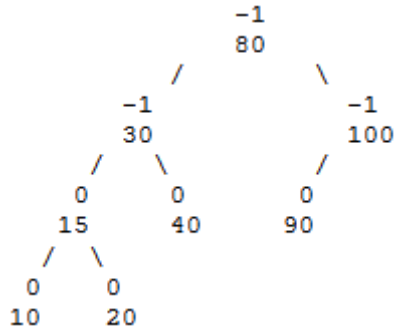


The fix is to use a single right rotation at node P. (In the mirror image case a single left rotation is used at P.) This gives the following picture.

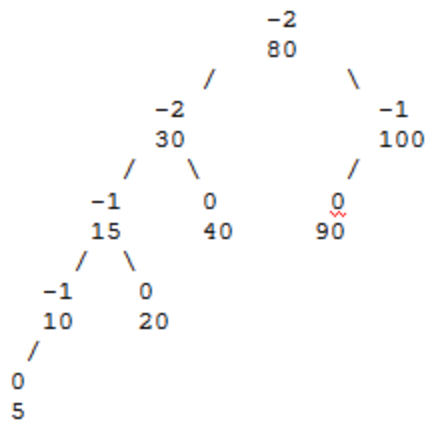


Consider the following more detailed example that illustrates subcase A.

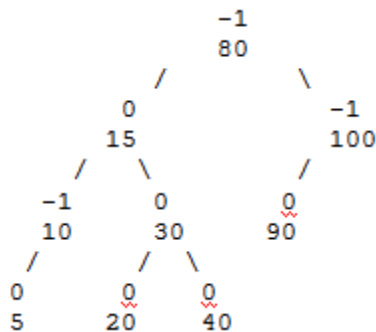




We then insert 5 and then check the balance factors from the new leaf up toward the root. (Always check from the bottom up.)

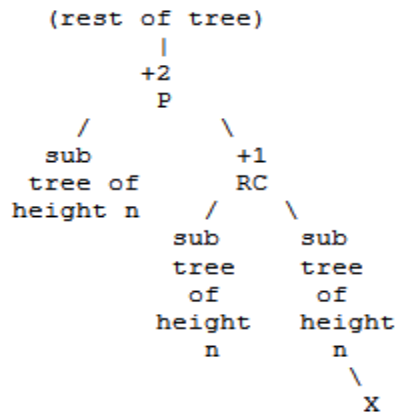


This reveals a balance factor of -2 at node 30 that must be fixed. (Since we work bottom up, we reach the -2 at 30 first. The other -2 problem will go away once we fix the problem at 30.) The fix is accomplished with a right rotation at node 30, leading to the following picture.



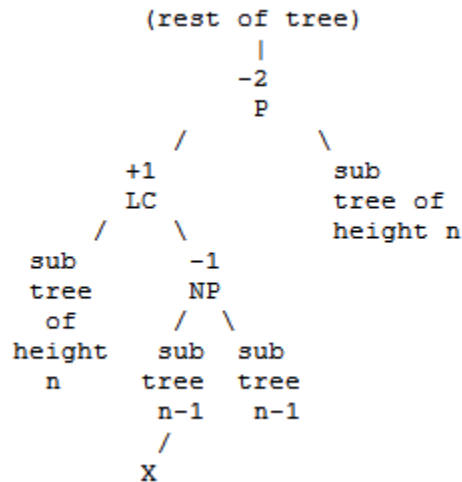
Recall that the mirror image situation is also included under subcase A. The following is a general illustration of this situation. The fix is to use a single left rotation at P. See if you can draw a picture of the following after the left rotation at P. Then draw a

picture of a particular example that fits our general picture below and fix it with a left rotation.

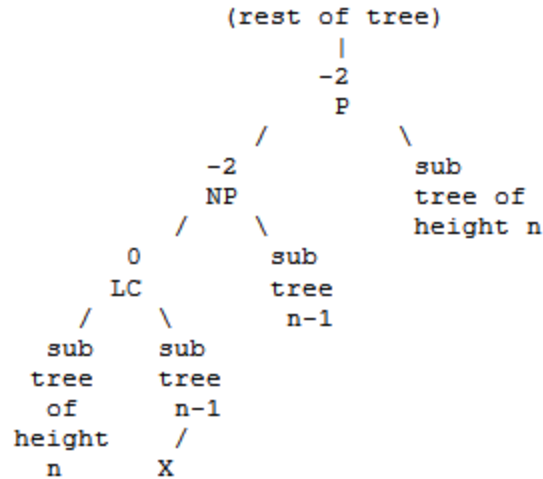


**Subcase B:**

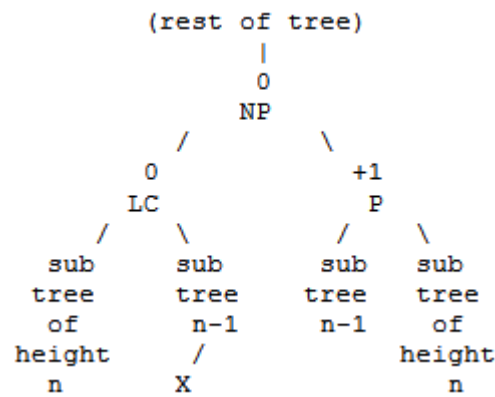
This consists of the following situation, where P denotes the parent of the subtree being examined, LC is P's left child, NP is the node that will be the new parent, and X is the new node added. X might be added to either of the subtrees of height n-1. Note that inserting X makes P have a balance factor of -2 and LC have a balance factor of +1. The -2 must be fixed. This is accomplished by doing a double rotation at P (explained below). (Note that the mirror image situation is also included under subcase B.)



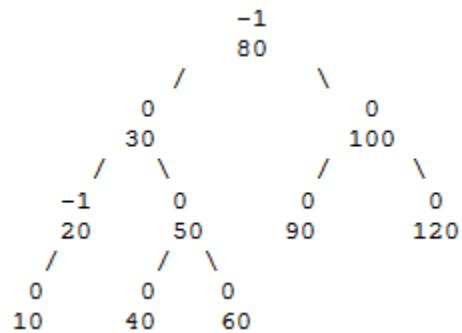
The fix is to use a double right rotation at node P. A double right rotation at P consists of a single **left** rotation at LC followed by a single right rotation at P. (In the mirror image case a double left rotation is used at P. This consists of a single right rotation at the right child RC followed by a single left rotation at P.) In the above picture, the double rotation gives the following (where we first show the result of the left rotation at LC, then a new picture for the result of the right rotation at P).



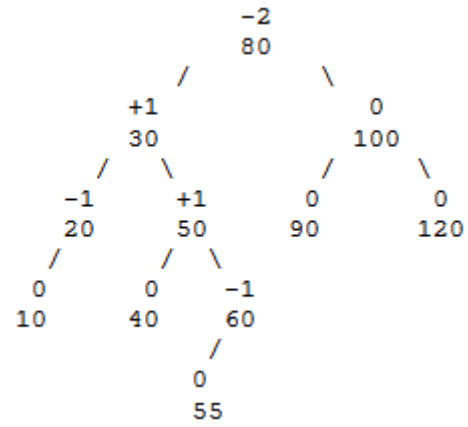
Finally we have the following picture after doing the right rotation at P.



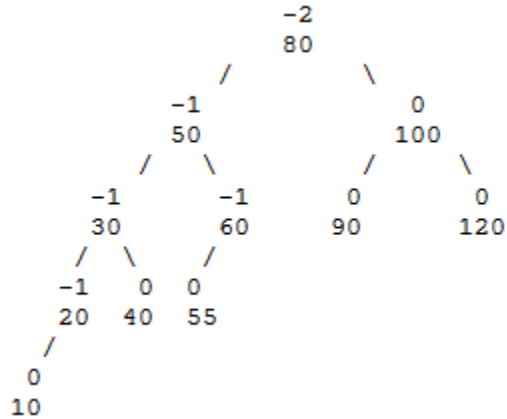
Consider the following concrete example of subcase B.



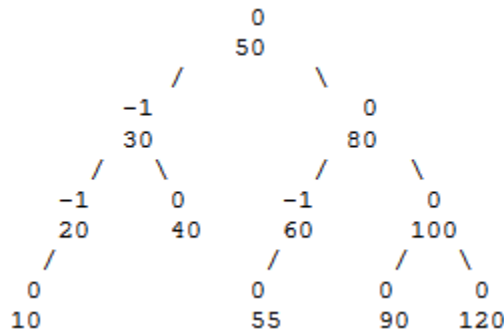
After inserting 55, we get a problem, a balance factor of -2 at the root node, as seen below.



As discussed above, this calls for a double rotation. First we do a single left rotation at 30. This gives the following picture.



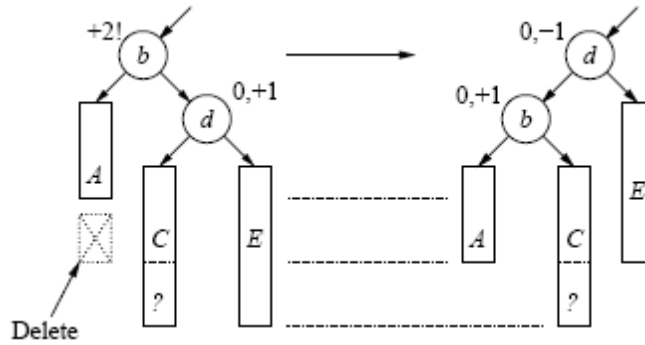
Finally, the right rotation at 80 restores the binary search tree to be an AVL tree. The resulting picture is shown below.



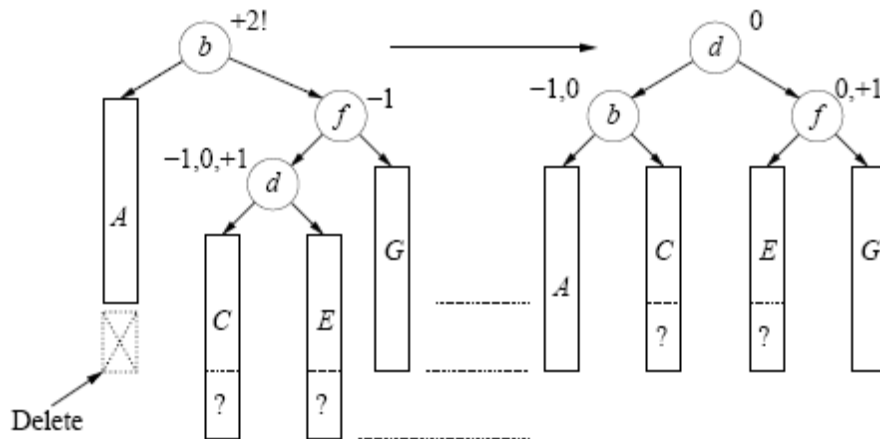
### 13.4 Deletion from AVL tree

Deletion is similar to insertion in that we start by applying the deletion algorithm for unbalanced binary trees. Recall that this breaks into three cases, leaf, single child, and two children. In the two children case we need to find a replacement key. Once the deletion is finished, we walk back up the tree (by returning from the recursive calls), updating the balance factors (or heights) as we go. Whenever we come to an unbalanced node, we apply an appropriate rotation to remedy the situation.

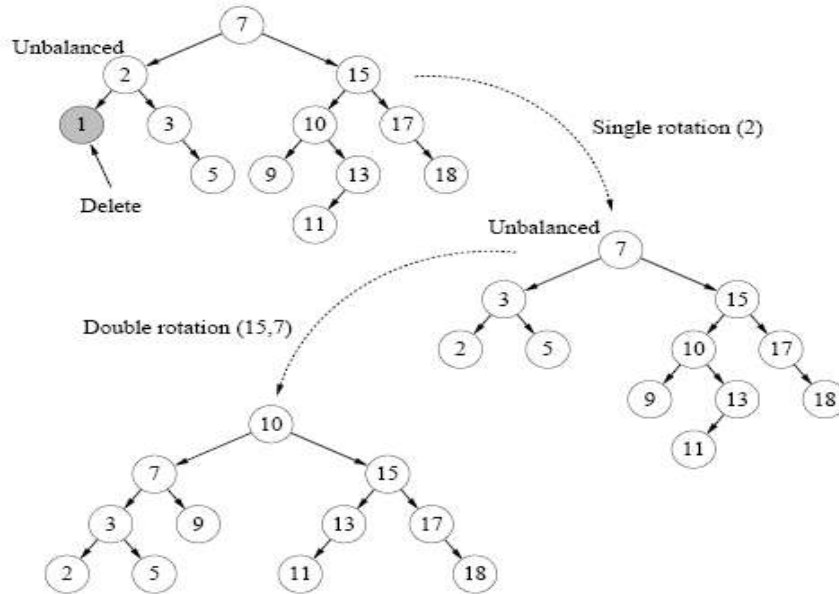
The deletion code is messier than the insertion code. (We will not present it.) But the idea is the same. Suppose that we have deleted a key from the left subtree and that as a result this subtree's height has decreased by one, and this has caused some ancestor to violate the height balance condition. There are two cases. First, if balance factor for the right child is either 0 or +1 (that is, it is not left heavy), we can perform a single rotation as shown in the figure below. We list multiple balance factors, because there are a number of possibilities. (We leave it as an exercise to figure out which balance factors apply to which cases.)



On the other hand, if the right child has a balance factor of  $-1$  (it is left heavy) then we need to perform a double rotation, as shown in the following figure. A more complete example of a deletion is shown in the figure below. We delete element 1. This causes node 2 to be unbalanced. We perform a single left rotation at 2. However, now the root is unbalanced. We perform a right-left double rotation to the root.



**Lazy Deletion:** The deletion code for AVL trees is almost twice as long as the insertion code. It is not all that more complicated conceptually, but the number of cases is significantly larger. Our text suggests an interesting alternative for avoiding the pain of coding up the deletion algorithm, called lazy deletion. The idea is to not go through the entire deletion process. For each node we maintain a boolean value indicating whether this element is alive or dead. When a key is deleted, we simply declare it to be dead, but leave it in the tree. If an attempt is made to insert the same key value again, we make the element alive again. Of course, after a long sequence of deletions and insertions, it is possible that the tree has many dead nodes. To fix this we periodically perform a garbage collection phase, which traverses the tree, selecting only the live elements, and then building a new AVL tree with these elements.



### 13.5 Performance of AVL tree

The maximum height of an AVL tree is  $1.44 * \lg n$ , which is an  $O(\lg n)$  function. This means that in the worst possible case, a lookup in a large AVL tree needs no more than 44% more comparisons than a lookup in a completely balanced tree. Even in the worst case, then, AVL trees are efficient; they still have  $O(\lg n)$  lookup times. On average, for large  $n$ , AVL trees have lookup times of  $(\lg n) + 0.25$ , which is even better than the above worst case figure, though still  $O(\lg n)$ . On average, a rotation (single or double) is required in 46.5% of insertions. Only one (single or double) rotation is needed to readjust an AVL tree after an insertion throws it out of balance.

An insertion into an AVL tree requires at most one rotation to rebalance a tree. A deletion may require  $\lg(N)$  rotations to rebalance the tree.

Simulation Results AVL

n	C(n)	E[ h(n) ]	R(n)
5	2.2	3.0	0.213
10	2.907	4.0	0.318
50	4.930	6.947	0.427
100	5.889	7.999	0.444

500	8.192	10.923	0.461
1000	9.202	11.998	0.463
5000	11.555	14.936	0.465
10000	12.568	15.996	0.465

$n$  = number of nodes in the tree

$C(n)$  = average number of comparisons to find a key

$E[ h(n) ]$  = expected height of tree

$R(n)$  = average number of rotations per insertion/deletion

### 13.6 Summary

A height balanced tree, also known as AVL tree, is a sophisticated self balancing tree. It can be thought of as the smarter, younger brother of the binary search tree. Unlike its older brother the AVL tree avoids worst case linear complexity runtimes for its operations. The AVL tree guarantees via the enforcement of balancing algorithms that the left and right subtrees differ in height by at most 1 which yields at most a logarithmic runtime complexity.

The criteria that is used to determine the "level" of "balanced-ness" is the difference between the heights of subtrees of a root in the tree. The "height" of tree is the "number of levels" in the tree. An AVL tree is a binary tree in which the difference between the height of the right and left subtrees (or the root node) is never more than one.

### 13.7 Questions

1. Define AVL trees. What are its properties?
2. How balancing factor of an AVL tree is calculated?
3. What are the different cases for insertion in an AVL tree?
4. Why deletion from an AVL tree is more complicated as compared to insertion?
5. Explain the concept of lazy deletion.

### 13.8 Suggested readings

1. A. Tanenbaum, Y. Lanhgsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.



3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley.
6. William Ford, William Topp, "Data Structures with C++", Prentice-Hall.

## HEAPS

### 14.1 Objectives

### 14.2 Introduction

### 14.3 Need of using Heap

### 14.4 Implementation of Heap

### 14.5 Storage of a Heap

### 14.6 Insertion in a Heap

### 14.7 Deletion from a Heap

### 14.8 Efficiency of a Heap

### 14.9 Heap Sort

### 14.10 Summary

### 14.11 Questions

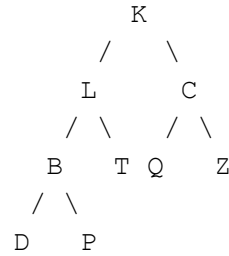
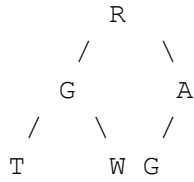
### 14.12 Suggested readings

### 14.1 Objectives

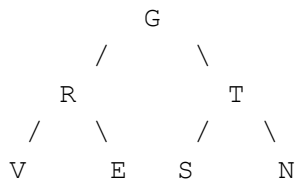
This lesson focuses on another type of tree called heap. We shall discuss the meaning, creation and operations on heap tree.

### 14.2 Introduction

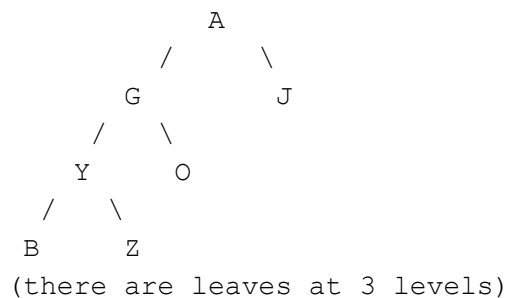
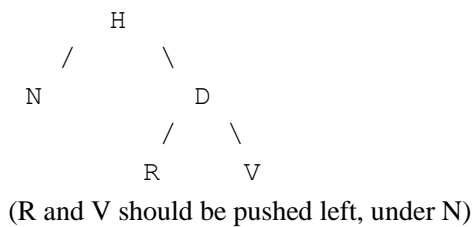
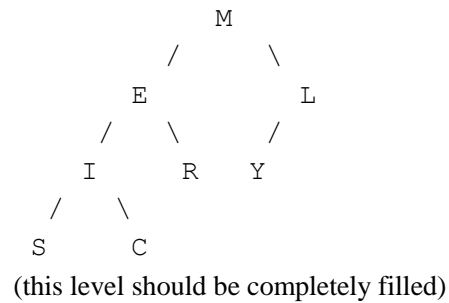
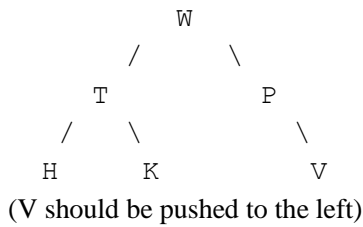
Recall that in a binary tree each node can have a left child node and/or a right child node. A leaf is a node with no children. An **almost complete** binary tree is a binary tree in which the following 3 conditions hold: all the leaves are at the bottom level or the bottom 2 levels, all the leaves are in the leftmost possible positions, and (except possibly for the bottom level) all levels are completely filled with nodes. Here are some examples of almost complete binary trees. Note that there is no particular ordering of the letters.



The following example is a complete binary tree. (Of course, it also fits the definition of almost complete. Any complete binary tree is automatically almost complete.)



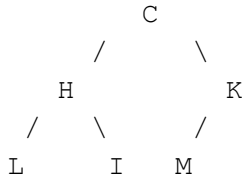
The following binary trees are *not* almost complete:



A **minimal heap** (descending heap) is an almost complete binary tree in which the value at each parent node is less than or equal to the values in its child nodes.

Obviously, the minimum value is in the root node. Note, too, that any path from a leaf to the root passes through the data in descending order.

Here is an example of a minimal heap:



A heap is a partially sorted binary tree. Although a heap is not completely in order, it conforms to a sorting principle: every node has a value less (for the sake of simplicity, we will assume that all orderings are from least to greatest) than either of its children. Additionally, a heap is a "complete tree" -- a complete tree is one in which there are no gaps between leaves. For instance, a tree with a root node that has only one child must have its child as the left node. More precisely, a complete tree is one that has every level filled in before adding a node to the next level, and one that has the nodes in a given level filled in from left to right, with no breaks.

### 14.3 Need of using Heap

A heap can be thought of as a priority queue; the most important node will always be at the top, and when removed, its replacement will be the most important. This can be useful when coding algorithms that require certain things to be processed in a complete order, but when you don't want to perform a full sort or need to know anything about the rest of the nodes. For instance, a well-known algorithm for finding the shortest distance between nodes in a graph, Dijkstra's Algorithm, can be optimized by using a priority queue.

Heaps can also be used to sort data. A heap sort is  $O(n \log n)$  efficiency, though it is not the fastest possible sorting algorithm. Check out a tutorial on heap sort for more information related to heap sort.

### 14.4 Implementation of Heap

Although the concept of a heap is simple, the actual implementation can appear tricky. How do you remove the root node and still ensure that it is eventually replaced by the correct node? How do you add a new node to a heap and ensure that it is moved into the proper spot?

The answers to these questions are more straightforward than meets the eye, but to understand the process, let's first take a look at two operations that are used for adding and removing nodes from a heap: upheaping and downheaping.

**Upheap:** The upheap process is used to add a node to a heap. When you upheap a node, you compare its value to its parent node; if its value is less than its parent node, then you switch the two nodes and continue the process. Otherwise the condition is met that the node is less than its parent node and so you can stop the process. As you know that in a heap, a node is always less than its parent node, you are assured that if the node you are

upheaping is less than its parent node, that that node is also less than the parents of the parent node.

**Downheap:** The downheap process is similar to the upheaping process. When you downheap a node, you compare its value with its two children. If the node is less than both of its children, it remains in place; otherwise, if it is greater than one or both of its children, then you switch it with the child of lowest value, thereby ensuring that of the three nodes being compared, the new parent node is lowest. Of course, you cannot be assured that the node being downheaped is in its proper position -- it may be greater than one or both of its children; the downheap process must be repeated until the node is less than both of its children.

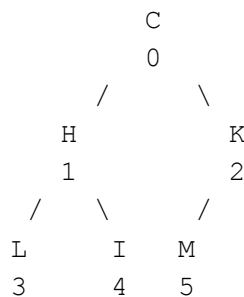
When you add a new node to a heap, you add it to the rightmost unoccupied leaf on the lowest level. Then you upheap that node until it has reached its proper position. In this way, the heap's order is maintained and the heap remains a complete tree.

Removing the root node from a heap is almost as simple: when you take the node out of the tree, you replace it with "last" node in the tree: the node on the last level and rightmost on that level.

Once the top node has been replaced, you downhead the node that was moved until it reaches its proper position. As usual, the result will be a proper heap, as it will be complete, and even if the node in the last position happens to be the greatest node in the entire heap, it will do no worse than end up back where it started.

### 14.5 Storage of a Heap

The typical storage method for a heap, or any almost complete binary tree, works as follows. Begin by numbering the nodes level by level from the top down, left to right. For example, consider the following heap. The numbering has been added below the nodes.



Then store the data in an array as shown below:

C	H	K	L	I	M
0	1	2	3	4	5

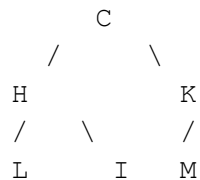
The advantage of this method over using the usual pointers and nodes is that there is no wasting of space due to storing two pointer fields in each node. Instead, starting with the current index, CI, one calculates the index to use as follows:

$$\begin{aligned} \text{Parent}(CI) &= (CI - 1) / 2 \\ \text{RightChild}(CI) &= 2 * (CI + 1) \\ \text{LeftChild}(CI) &= 2 * CI + 1 \end{aligned}$$

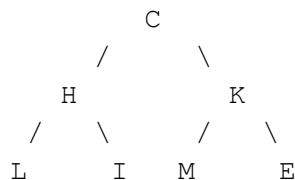
For example, if we start at node H (with index 1), the right child is at index  $2 * (1 + 1) = 4$ , that is, node I.

## 14.6 Insertion in a Heap

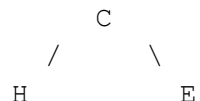
This is done by temporarily placing the new item at the end of the heap (array) and then calling a FilterUp routine to make any needed adjustments on the path from this leaf to the root. For example, let's insert E into the following heap:

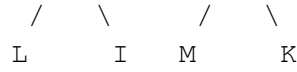


First, temporarily place E in the next available position:

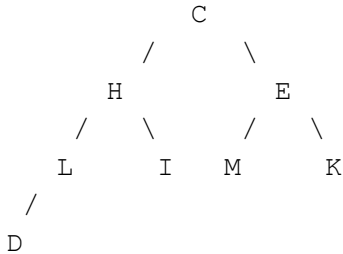


Of course, the new tree might not be a heap. The FilterUp routine now checks the parent, K, and sees that things would be out of order as they are. So K is moved down to where E was. Then the parent above that, C, is checked. It is in order relative to the target item E, so the C is not moved down. The hole left behind is filled with E, then, as this is the correct position for it.

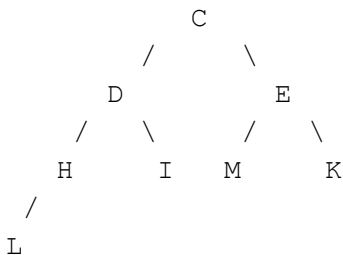




For practice, let's take the above heap and insert another item, D. First, place D temporarily in the next available position:



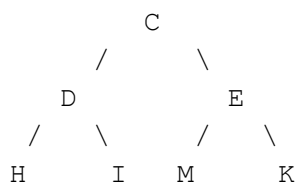
Then the FilterUp routine checks the parent, L, and discovers that L must be moved down. Then the parent above that, H, is checked. It too must be moved down. Finally C is checked, but it is OK where it is. The hole left where the H had been is where the target D is then inserted.



Things have now been adjusted so that we again have a heap!

### 14.7 Deletion from a Heap

We always remove the item from the root. That way we always get the smallest item. The problem is then how to adjust the binary tree so that we again have a heap (with one less item). The algorithm works like this: First, remove the root item and replace it temporarily with the item in the last position. Call this replacement the target. A FilterDown routine is then used to check the path from the root to a leaf for the correct position for this target. The path is chosen by always selecting the smaller child at each node. For example, let's remove the C from this heap:



```

  /
 L

```

First we remove the C and replace it with the last item (the target), L:

```

      L
     / \
    D   E
   / \ / \
  H  I M  K

```

The smaller child of L is D. Since D is out of order compared to the target L, we move D up. The smaller child under where D had been is H. When H is compared to L we see that the H too needs to be moved up. Since we are now at a leaf, this empty leaf is where the target L is put.

```

      D
     / \
    H   E
   / \ / \
  L  I M  K

```

For another example, let's remove the E from the following heap:

```

      E
     / \
    G   K
   / \ / \
  J  N K  X
 / \ /
X  Y P

```

First remove the E and replace it with the target P (the last item):

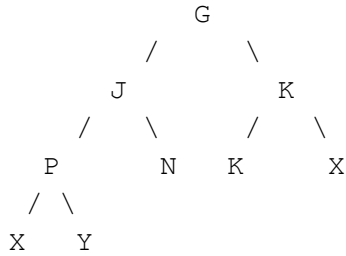
```

      P
     / \
    G   K
   / \ / \
  J  N K  X
 / \
X  Y

```

Now use the FilterDown routine to filter the P down to its correct position by checking the smaller child, G, which should be moved up and then the smaller child below that, J, which should also be moved up. Finally, the smaller child, X, under where J had been is checked, but it does not get moved since it is OK relative to the target P. The P is then placed in the empty node above X. We then have the following heap:





### 14.8 Efficiency of a Heap

Whenever you work with a heap, most of the time taken by the algorithm will be in upheaping and downheaping. As it happens, the maximum number of levels of a complete tree is  $\log(n)+1$ , where  $n$  is the number of nodes in the tree. Because upheap or downheap moves an element from one level to another, the order of adding to or removing from a heap is  $O(\log n)$ , as you can make switches only  $\log(n)$  times, or one less time than the number of levels in the tree (consider that a two level tree can have only one switch).

### 14.9 Heap Sort

Heapsort is performed by somehow creating a heap and then removing the data items one at a time. The heap could start as an empty heap, with items inserted one by one. However, there is a relatively easy routine to convert an array of items into a heap, so that method is often used. This routine is described below. Once the array is converted into a heap, we remove the root item (the smallest), readjust the remaining items into a heap, and place the removed item at the end of the heap (array). Then we remove the new item in the root (the second smallest), readjust the heap and place the removed item in the next to the last position, etc.

Heapsort is  $\Theta(n \cdot \lg(n))$ , either average case or worst case. This is great for a sorting algorithm! No appreciable extra storage space is needed either. On average, quicksort (which is also  $\Theta(n \cdot \lg(n))$  for the average case) is faster than heapsort. However, quicksort has that bad  $\Theta(n^2)$  worst case running time.

### Example Trace of Heapsort

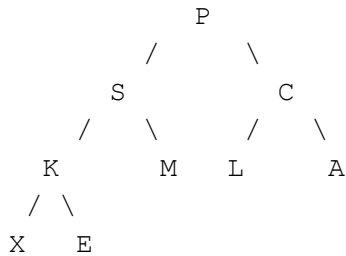
Let's heapsort the following array:

P	S	C	K	M	L	A	X	E
0	1	2	3	4	5	6	7	8

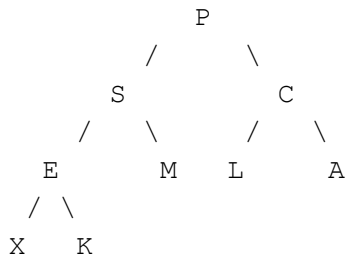
To convert this to a heap, first go to the index of the last parent node. This is given by  $(\text{HeapSize} - 2) / 2$ . In this case,  $(9 - 2) / 2 = 3$ . Thus K is the last parent in the tree. We then apply the FilterDown routine to each node from this index down to index 0. (Note that

this is *each* node from 3 down to 0, not just the nodes along the path from index 3 to index 0.)

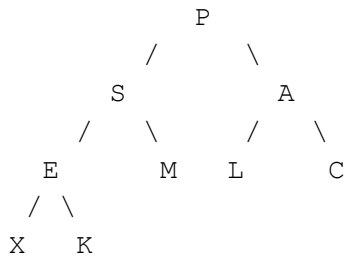
In our example, the array corresponds directly to the following binary tree. Note that this is not yet a heap.



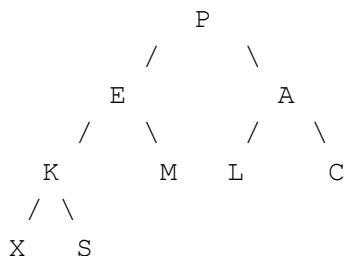
Applying FilterDown at K gives the following. (Note that E is the smaller child under K.)



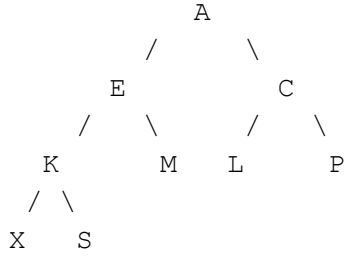
Now apply FilterDown at index 2, that is, at node C. (Under C, A is the smaller child.)



Next, apply FilterDown at index 1, that is, at node S. Check the smaller child, E, and then the smaller child under that, namely K. Both E and K get moved up.

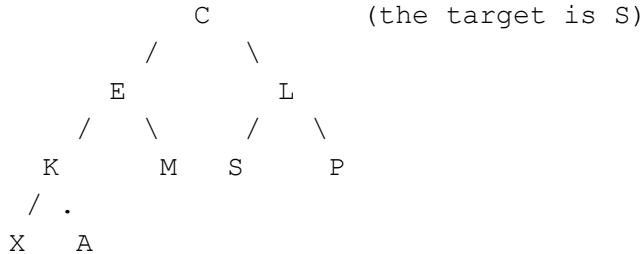


Finally, apply FilterDown at index 0, that is, at the root node. We check the smaller child, A, and then the smaller child, C, relative to the target P. Both A and C get moved up.



Now we have a heap! The first main step of heapsort has been completed. The other main component of heapsort was described earlier: to repeatedly remove the root item, adjust the heap, and put the removed item in the empty slot toward the end of the array (heap).

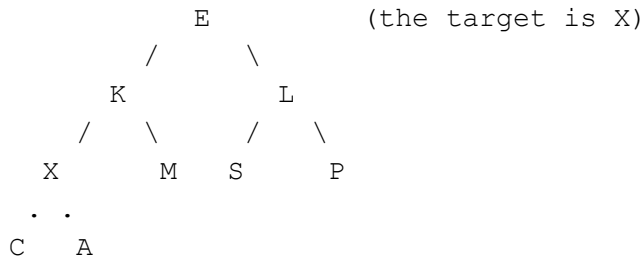
First we remove the A, adjust the heap by using FilterDown at the root node and place the A at the end of the heap (where it is not really part of the heap at all and so is not drawn below as connected to the tree).



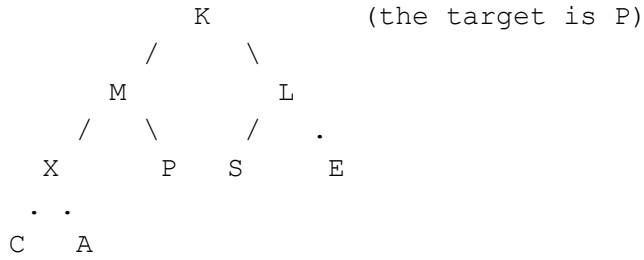
Of course, all of this is really taking place in the array that holds the heap. At this point it looks like the following. Note that the heap is stored from index 0 to index 7. The A is after the end of the heap.

C	E	L	K	M	S	P	X	A
0	1	2	3	4	5	6	7	8

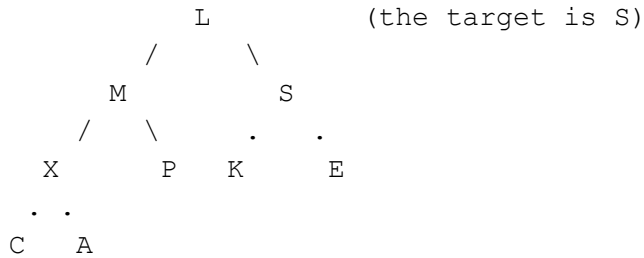
Next we remove the C, adjust the heap by using FilterDown at the root node and place the C at the end of the heap:



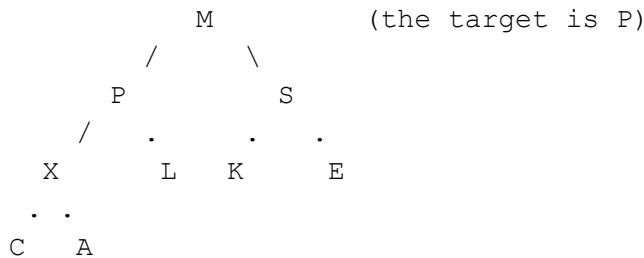
Next we remove the E, adjust the heap by using FilterDown at the root node and place the E at the end of the heap:



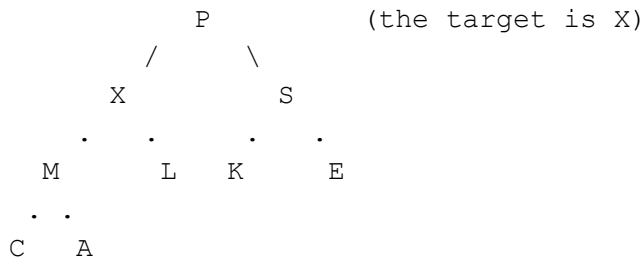
Next we remove the K, adjust the heap by using FilterDown at the root node, and place the K at the end of the heap:



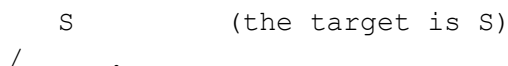
Next we remove the L, adjust the heap by using FilterDown at the root node and place the L at the end of the heap:

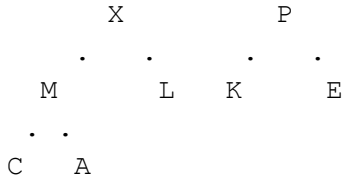


Next we remove the M, adjust the heap by using FilterDown at the root node and place the M at the end of the heap:

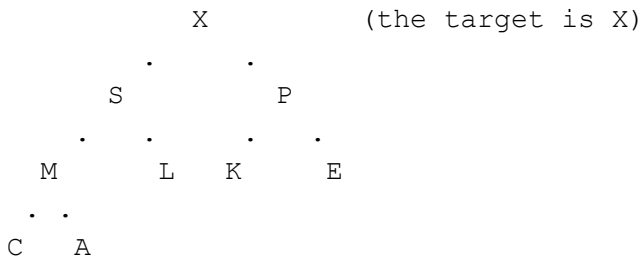


Next we remove the P, adjust the heap by using FilterDown at the root node and place the P at the end of the heap:





Next we remove the S, adjust the heap (now a trivial operation), and place the S at the end of the heap:



Since only the item X remains in the heap, and since we have removed the smallest item, then the second smallest, etc., the X must be the largest item and should be left where it is. If you now look at the array that holds the above items you will see that we have sorted the array in descending order:

X	S	P	M	L	K	E	C	A
0	1	2	3	4	5	6	7	8

### 14.10 Summary

A heap is a binary tree if it is empty or the key in the root is larger than the the key of both the children and both subtrees have the heap property. Insertion in heap is done by temporarily place the new element at the end of the heap and then making adjustments on the path from the newly inserted leaf node to the root. Deletion is always done from the root node. After deleting the root node it is replaced with the last node of the heap and then by checking the path from the root to a leaf for the correct position for this target. Heaps provide us with a method of sorting, known as heapsort.

### 14.11 Questions

1. Define a complete binary tree.
2. What are the properties of a heap?
3. What is the difference between minheap and maxheap?
4. How heap is stored in memory? Explain giving example.
5. Write the procedure of creating a heap?

6. Create a heap using the following data: 45 32 98 56 12 41 90 67 48 82 71.
7. Write and explain the heapsort algorithm.

#### **14.12 Suggested readings**

1. A. Tanenbaum, Y. Lanhgsam and A. J. Augenstein, "Data Structures Using C++", PHI.
2. M. A. Weiss, "Data Structures and Algorithm Analysis in C++", Pearson Education.
3. R. Sedgewick, "Algorithms in C++", Pearson Education.
4. S. Lipschutz, "Data Structures", Tata McGraw Hill.
5. Donald E. Knuth, "The Art Of Computer Programming", Vol 1-4. 3rd ed., Addison Wesley.

## Mandatory Student Feedback Form

<https://forms.gle/KS5CLhvpwrpgjwN98>

Note: Students, kindly click this google form link, and fill this feedback form once.