

CSCI2100 Data Structures

Ryan Chan

March 15, 2025

Abstract

This is a note for **CSCI2100 Data Structures** for self-revision purpose ONLY. Some contents are taken from lecture notes and reference book.

Mistakes might be found. So please feel free to point out any mistakes.

Contents are adapted from the lecture notes of CSCI2100, prepared by **Irwin King**, as well as some online resources.

Contents

1	Introduction	2
1.1	Overview	2
1.2	Algorithm	2
1.3	Study of Data	3
2	Analysis	4
2.1	Complexity	4
2.2	Recurrence Relations	5
3	ADT, List, Stack and Queue	6
3.1	Abstract Data Type (ADT)	6
3.2	List	6
3.3	Stack	8
3.4	Queue	9
4	Trees	11
4.1	General Tree	11
4.2	Binary Tree	12
4.3	Expression Tree	13
4.4	Binary Search Tree	13
4.5	AVL Tree	14
4.6	B-Tree	16
5	More on Tree	18
5.1	Tries	18
5.2	B-Tree	19
6	Hashing	21
6.1	Introduction	21
6.2	Hash Table	22
6.3	Hash Function	22
6.4	Collision Resolution	23

Chapter 1

Introduction

1.1 Overview

A **data structure** is a way to organize and store data in a computer program, allowing for efficient access and manipulation.

An **algorithm** is different from a **program**. An algorithm is a process or set of rules used for calculation or problem-solving. It is a step-by-step outline or flowchart showing how to solve a problem. A program, on the other hand, is a series of coded instructions that control the operation of a computer or other machines. It is the implemented code of an algorithm.

For example, to solve the greatest common divisor (GCD) problem, we can use the following algorithm.

Algorithm 1.1: Euclid's Algorithm

Data: $m, n \in \mathbb{Z}^+$

Result: $\text{GCD}(m, n)$

```
1 while  $m > 0$  do
2   if  $n > m$  then
3     swap  $m$  and  $n$ 
4   subtract  $n$  from  $m$ 
5 return  $n$ 
```

By using a mathematical method to prove this algorithm, we can show that it is correct, provided that it terminates.

Having proved the correctness, we also need to use different test cases to check if there is anything wrong with the coding or the proof. We should consider special cases, including large values, swapped values, etc.

We are also interested in the time and space (computer memory) it uses, which we call **time complexity** and **space complexity**. Typically, complexity is a function of the values of the inputs, and we would like to know which function. We can also consider the best case, average case, and worst-case scenarios.

For example, in the above algorithm, the best case would be $m = n$, with just one iteration. If $n = 1$, there are m iterations, which is the worst case. However, for the average case, it is difficult to analyze.

Also, for space complexity, it is constant since we only use space for the three integers: m , n , and t .

To improve the above algorithm, we can use mod, so we don't need to keep doing subtraction.

1.2 Algorithm

An algorithm is a finite set of instructions which, if followed, accomplishes a particular task. Every algorithm must satisfy the following criteria:

-
- Input: There are zero or more quantities that are externally supplied.
 - Output: At least one quantity is produced.
 - Definiteness: Each instruction must be clear and unambiguous.
 - Finiteness: If we trace out the instructions of an algorithm, then for all cases, the algorithm will terminate after a finite number of steps.
 - Effectiveness: Every instruction must be sufficiently basic that it can, in principle, be carried out by a person using only pencil and paper.

We also define an algorithm as any well-defined computational procedure that takes some value, or set of values, as input and produces some value, or set of values, as output. It is thus a sequence of computational steps that transform the input into output.

It can also be viewed as a tool for solving a well-specified computational problem. The problem statement specifies, in general terms, the desired input or output relationship, and the algorithm describes a specific computational procedure for achieving that input or output relationship.

An algorithm is said to be correct if, for every input instance, it halts with the correct output. It can solve the given computational problem. In contrast, an incorrect algorithm might not halt at all on some input instances, and sometimes it can even produce useful results.

1.3 Study of Data

A data structure is a particular way of storing and organizing data in a computer so that it can be used efficiently.

A data structure is a set of domains D , a designated domain $d \in D$, a set of functions F , and a set of axioms A .

An implementation of a data structure d is a mapping from d to a set of other data structures e .

Chapter 2

Analysis

2.1 Complexity

Before, we talked about the definition of an algorithm. In this part, we would like to know how we can estimate the time required for a program, how to reduce the running time of a program, what the storage complexity is, and how to deal with trade-offs.

We can analyze the runtime by comparing functions. For example, given two functions $f(N)$ and $g(N)$, we can compare their relative rates of growth. There are three types of comparisons that we can make: $f(n) = \Theta(g(n))$ represents the exact bound, $f(n) = O(g(n))$ represents the upper bound, and $f(n) = \Omega(g(n))$ represents the lower bound.

By using bounds, we can establish a relative order among functions. Here, we often use $O(n)$ to analyze time complexity.

For the definition of the upper bound, it says that there is some point n_0 past which $cf(N)$ is always at least as large as $T(N)$. Then we say that $T(N) = O(f(N))$, where $f(N)$ is the upper bound on $T(N)$.

Definition 2.1.1. We say that $f(n) = O(g(n))$ iff there exists a constant $c > 0$ and an $n_0 \geq 0$ such that

$$f(n) \leq cg(n) \quad \text{for all } n \geq n_0$$

Or we can use the following notation

$$\exists c > 0, n_0 \geq 0 \text{ such that } f(n) \leq cg(n) \forall n \geq n_0$$

There are some rules to follow:

- Transitivity

$$\text{If } f(n) = O(g(n)) \text{ and } g(n) = O(h(n)), \text{ then } f(n) = O(h(n))$$

- Rule of sums

$$f(n) + g(n) = O(\max\{f(n), g(n)\})$$

- Rule of products

$$\text{If } f_1(n) = O(g_1(n)), f_2(n) = O(g_2(n)), \text{ then } f_1(n)f_2(n) = O(g_1(n)g_2(n))$$

We do not include constants or lower-order terms inside Big-O notation.

If $f(n)$ is a polynomial in n with degree r , then $f(n) = O(n^r)$, but for $s < r$, $f(n) \neq O(n^s)$.

Also, any logarithm of n grows more slowly than any positive power of n as it increases. Hence, $\log n$ is $O(n^k)$ for any $k > 0$, but n^k is never $O(\log n)$ for any $k > 0$.

Order	Time
$O(1)$	constant time
$O(n)$	linear time
$O(n^2)$	quadratic time
$O(n^3)$	cubic time
$O(2^n)$	exponential time
$O(\log n)$	logarithmic time
$O(\log^2 n)$	log-squared time

On a list of length n , sequential search has a running time of $O(n)$.

On an ordered list of length n , binary search has a running time of $O(\log n)$.

The sum of the sums of integer indices of a loop from 1 to n is $O(n^2)$.

In summary, Big-O notation provides an upper bound of the complexity in the **worst-case**, helping to quantify performance as the input size becomes arbitrarily large. However, it doesn't measure the actual time, but represents the number of operations an algorithm will execute.

2.2 Recurrence Relations

Recurrence relations are useful in certain counting problems, for example, recursive algorithms. They relate the n -th element of a sequence to its predecessors.

By definition, a recurrence relation for the sequence a_0, a_1, \dots is an equation that relates a_n to certain of its predecessors a_0, a_1, \dots, a_{n-1} . Initial conditions for the sequence are explicitly given values for a finite number of the terms of the sequence.

To solve a recurrence relation, we can use iteration. We use the recurrence relation to write the n -th term a_n in terms of certain of its predecessors. We then successively use the recurrence relation to replace each of a_{n-1}, \dots by certain of their predecessors. We continue until an explicit formula is obtained.

For example, the Fibonacci sequence is also defined by the recurrence relation.

Example (Tower of Hanoi). Find an explicit formula for a_n , the minimum number of moves in which the n -disk Tower of Hanoi puzzle can be solved.

Given $a_n = 2a_{n-1} + 1$, $a_1 = 1$, by applying the iterative method, we obtain:

$$\begin{aligned}
 a_n &= 2a_{n-1} + 1 \\
 &= 2(2a_{n-2} + 1) + 1 \\
 &= 2^2a_{n-2} + 2 + 1 \\
 &= 2^2(2a_{n-3} + 1) + 2 + 1 \\
 &= 2^3a_{n-3} + 2^2 + 2 + 1 \\
 &= \dots \\
 &= 2^{n-1}a_1 + 2^{n-2} + 2^{n-3} + \dots + 2 + 1 \\
 &= 2^{n-1} + 2^{n-2} + 2^{n-3} + \dots + 2 + 1 \\
 &= 2^n - 1
 \end{aligned}$$

Chapter 3

ADT, List, Stack and Queue

3.1 Abstract Data Type (ADT)

We use data abstraction to simplify software development since it facilitates the decomposition of the complex task of developing a software system. To put it simply, data abstraction shows only the essential details of data, while the implementation details are hidden.

For example, as will be discussed later, List, Stack, and Queue (LSQ) are forms of data abstraction, or what we call abstract data types. We can use them to retrieve or store data, but we don't know how they are actually stored or indexed.

Data encapsulation, or information hiding, is the concealing of the implementation of a data object from the outside world. Through data abstraction, we can separate the specification of a data object from its implementation.

A data type is a collection of objects and a set of operations that act on those objects. An abstract data type (ADT) is a data type organized in such a way that we can separate the specification of the object and the specification of the operations on the object. Abstract data types are simply a set of operations, and they are mathematical abstractions.

Note that abstraction is like a functional description without knowing how to use it, while implementation, on the contrary, is something that can be used and executed.

In summary, an ADT is a high-level description of how data is organized and the operations that can be performed on it. It abstracts the details of its implementation and only exposes the operations that are allowed on data structures.

3.2 List

The first abstract data type in this chapter is List.

3.2.1 Definition

When dealing with a general list of the form a_1, a_2, \dots, a_n , we say that the size of this list is n . If the list is of size 0, we call it the **null list**. Except null list, we say that a_{i+1} follows/succeeds a_i ($i < n$) and that a_{i-1} precedes a_i ($i > 1$).

The first element of the list is a_1 , and the last element is a_n . The predecessor of a_1 and the successor of a_n is not defined.

3.2.2 Operations

A list of elements of type T is a finite sequence of elements of T together with the following operations:

- Create the list and make it empty.
- Determine whether the list is empty or not.
- Determine whether the list is full or not.
- Find the size of the list.
- Retrieve any entry from the list, provided that the list is not empty.
- Store a new entry, replacing the entry at any position in the list, provided that the list is not empty.
- Insert a new entry into the list at any position, provided that the list is not full.
- Delete any entry from the list, provided that the list is not empty.
- Clear the list to make it empty.

With these operations, we can perform various tasks on the list ADT.

3.2.3 Implementation

We can use an array to implement a list. Most of the operations follow linear time, for example, `print_list`, `make_null`, `find`, and `find_kth`. For insertion, there could be cases where the list is full, and that's why we need dynamically allocated space. For deletion, we need to find the element, perform the deletion, and reallocate space, which might require more time. Thus, we introduce the linked list.

3.2.4 Linked List

There are several types of linked lists:

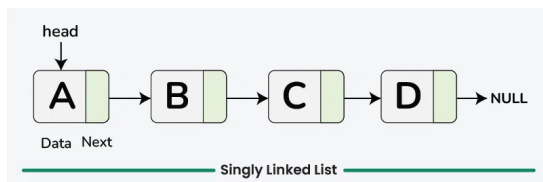


Figure 3.1: Singly Linked List

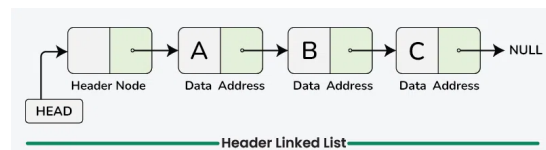


Figure 3.2: Singly Linked List with Header

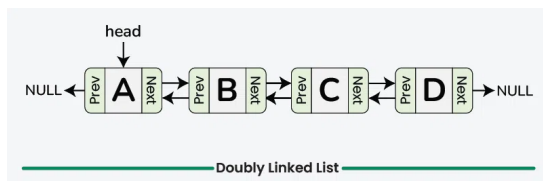


Figure 3.3: Doubly Linked List

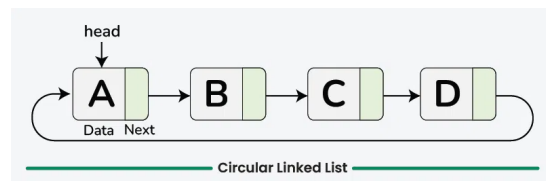


Figure 3.4: Circularly Linked List

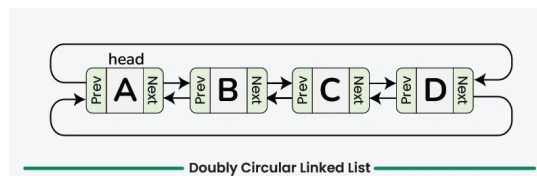


Figure 3.5: Circularly Doubly Linked List

A polynomial can be represented as

$$F(X) = \sum_{i=0}^N A_i X^i$$

For example, $F(X) = 4X^3 + 2X^2 + 5X + 1$. We may want to perform operations like addition, subtraction, multiplication, and differentiation. Using an array data structure, the time complexity may be larger due to the need to store all terms, including zero coefficients. However, with a linked list, we can efficiently perform these operations by traversing the linked list and processing only the non-zero terms.

Also, note that a circular list saves space but not time. It is useful for smaller datasets. However, for a larger number of students and courses, the use of such a circular list might be a waste of space.

In summary, a list abstract data type represents an ordered collection of elements. They can be added or removed at any position in the list. It provides methods to access elements by their position.

By using different types of linked lists, we can achieve various goals. For example, we can print all the elements in reverse using a doubly linked list.

3.3 Stack

3.3.1 Definition

A stack is an ordered list in which all insertions and deletions are made at one end, called the top. It follows the Last In, First Out (LIFO) rule.

3.3.2 Operations

A stack of elements of type T is a finite sequence of elements of T along with the following operations:

- Create the stack.
- Determine if the stack is empty or not.
- Determine if the stack is full or not.
- Determine the number of entries in the stack.
- Insert (Push) a new entry at one end of the stack, called its top, if the stack is not full.
- Retrieve the entry at the top of the stack, if the stack is not empty.
- Delete (Pop) the entry at the top of the stack, if the stack is not empty.
- Clear the stack to make it empty.

3.3.3 Implementation

We can use a doubly linked list to implement the stack, where both Push and Pop operations happen at the front of the list. We can use the Top operation to examine the element at the front of the list. In this case, the space complexity would be $O(3n)$, and the time complexity would be $O(c)$, where c is a constant.

Alternatively, we can use an array to implement the stack, since we only perform insertion and deletion at the top. However, we need to declare the size ahead of time and use TopOfStack as the counter to point to the top of the stack. If an array is used, the space complexity would be $O(n)$, and the time complexity would be $O(1)$.

3.3.4 Application

Balance Symbols

We can use a stack to balance symbols, which is commonly used in compilers to check for syntax errors. While compiling, an empty stack is created, and an opening symbol is pushed onto the stack. Then, when a closing symbol is encountered, it is popped from the stack. There could be four types of errors:

1. Stack Overflow: Too many brackets.

-
2. Mismatched Symbols: The opening and closing symbols don't match.
 3. Empty Stack: Attempting to pop from an empty stack.
 4. Non-Empty Stack: The stack isn't empty at the end of the process.

Reverse Polish Calculator

We can also use a stack to make a Reverse Polish Calculator. There are three forms of notation: prefix, postfix, and infix. For example, if the expression is $a \times b$, then:

1. Prefix: $\times ab$
2. Postfix: $ab\times$
3. Infix: $a \times b$

In Reverse Polish notation (postfix), parentheses are not needed. Using a stack, we can calculate the answer by evaluating the postfix expression. For example, when the character is a number, it is pushed onto the stack. If the character is an operator, two elements are popped from the stack, and the operation is performed on those two elements.

In summary, the stack abstract data type is a collection of elements with two main operations: push and pop. All operations happen at the top, and it follows the Last In, First Out (LIFO) approach.

3.4 Queue

3.4.1 Definition

A Queue is an ordered list in which all insertions take place at one end, the rear, while all deletions take place at the other end, the front. It follows the First In, First Out (FIFO) rule.

3.4.2 Operations

Several operations can be performed on a queue:

- Create the queue
- Determine if the queue is empty or not.
- Insert (Enqueue) a new entry at one end of the queue, called its rear, if the queue is not full.
- Delete (Dequeue) an entry at the other end of the queue, called its front, if the queue is not empty.
- Retrieve the entry at the front of the queue, if the queue is not empty.
- Clear the queue and make it empty.

3.4.3 Implementation

We can use both linear and circular arrays to implement a queue. For a linear array, we can have two indices that always increase. However, this might lead to overflow. Additionally, the array may need to be shifted forward or backward after each enqueue or dequeue operation. For a circular array, we have the following possibilities:

- Front and rear indices, with one position left vacant.
- Front and rear indices, with a Boolean variable indicating fullness or emptiness.
- Front and rear indices, with an integer variable counting entries.
- Front and rear indices taking special values to indicate emptiness.

3.4.4 Application

Queues are commonly used in various applications, such as in printer queues, airline control systems, and bank queues.

In summary, the physical model of a queue is a linear array, with the front always in the first position. All entries are moved up the array whenever the front is deleted. It follows the First In, First Out (FIFO) rule.

Chapter 4

Trees

In this chapter, we will introduce some fundamental concepts of trees. Trees are hierarchical data structures widely used for various applications such as representing hierarchical relationships, optimizing search operations, and organizing data efficiently.

4.1 General Tree

4.1.1 Nodes

A tree is a collection of nodes. The collection can be empty, which is sometimes denoted as A . Otherwise, a tree consists of a distinguished node r , called the root, and zero or more subtrees T_1, T_2, \dots, T_k , each of whose roots are connected by a directed edge to r . The root of each subtree is said to be a child of r , and r is the parent of each subtree root.

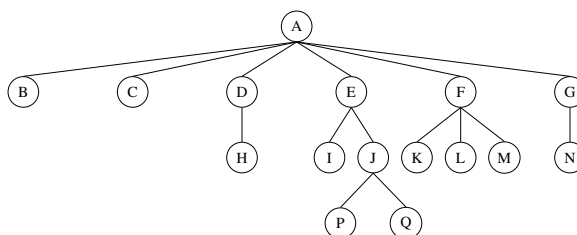


Figure 4.1: General Tree

Each node in a tree has a parent and may have an arbitrary number of children, possibly zero. Nodes with no children are known as leaves, and nodes with the same parent are called siblings. For example, in the above graph, A is the parent of D , and B , C , and D are siblings.

A path from node n_1 to n_k is defined as a sequence of nodes n_1, n_2, \dots, n_k such that n_i is the parent of n_{i+1} for $1 \leq i < k$. The length of this path is the number of edges on the path, namely $k - 1$. There is a path of length zero from every node to itself, and there is exactly one path from the root to each node.

Also, if there is a path from n_1 to n_2 , we call n_1 the ancestor of n_2 , while n_2 is the descendant of n_1 . If $n_1 \neq n_2$, we call them a proper ancestor or proper descendant.

For any node n_i , the depth of n_i is the length of the unique path from the root to n_i . Thus, the root is at depth 0. The height of n_i is the longest path from n_i to a leaf. Therefore, all leaves are at height 0. For example, in the graph above, E is at depth 1 and height 2.

The height of a tree is equal to the height of the root, and the depth of a tree is the depth of the deepest leaf. These two values are always equal, representing the longest path from the root to any leaf.

To implement a tree, we could store both the data and a pointer to each child in the node. However, this approach might not work well for a large number of children. We can solve this issue by keeping the children of each node in a linked list of tree nodes.

4.2.3 Implementation

Since a binary tree has at most two children, we can implement it by keeping direct pointers to them. A node can be represented as an element in a doubly linked list, storing key information along with two pointers to its left and right children.

We can also implement a node using two pointers: one pointing to its left child and the other pointing to its next sibling.

4.3 Expression Tree

An expression tree is also a binary tree, which is used to calculate the result of an expression. For example, we can express the expression $((9 - (2 + 3)) * (7 - 1))$ using the expression tree on the right.

In an expression tree, the leaves represent operands, and the internal nodes represent operators. By using inorder traversal, we can recover the original expression. From the binary tree, using postorder traversal, we can obtain the postfix notation. With the use of a stack, we can implement a calculator, as shown in the previous chapter.

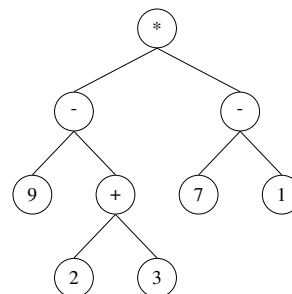


Figure 4.7: Expression Tree

4.4 Binary Search Tree

4.4.1 Definition

A binary search tree (BST) has the same physical property as a binary tree, meaning that nodes have at most two children. However, it also has an ordering property: for each node, all the nodes in its left subtree have smaller values, and all the nodes in its right subtree have larger values. This ordering property turns a binary tree into a binary search tree, and it implies that all the elements in the tree can be ordered in a consistent manner.

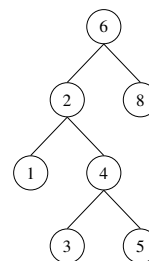


Figure 4.8: Binary Search Tree

4.4.2 Operations

There are some typical operations that can be done on a binary search tree, like make null, find, find max, find min, insertion and deletion.

For the **Find** operation, it generally requires returning a pointer to the node in tree T that has key x , or **null** if there is no such node. The structure of the tree makes this simple. If T is empty, then we can just return **null**. If the key stored at T is x , we can return T . Otherwise, we make a recursive call on a subtree of T , either left or right, depending on the relationship of x to the key stored in T .

For the **Find_min** and **Find_max** operations, these routines return the position of the smallest and largest elements in the tree, respectively. To perform **Find_min**, start at the root and go left as long as there is a left child. The stopping point is the smallest element. The **Find_max** routine is the same, except that branching is to the right child.

For insertion, we proceed down the tree. If x is found, we do nothing (or "update" something). Otherwise, we insert x at the last spot on the path traversed. Duplicates can be handled by keeping an extra field in the node record that indicates the frequency of occurrence.

However, for deletion, it is more difficult since we need to consider several possibilities. If the node is a leaf, it can be deleted immediately. If the node has one child, the node can be deleted after its parent adjusts a pointer to bypass the node. The complicated case is when we need to delete a node with two children. The general idea is to replace the key of the node with the smallest (leftmost) key of the right subtree and recursively delete the node.

4.4.3 Analysis

As mentioned before, the average depth of a binary search tree is $O(\log n)$. Therefore, intuitively, all operations, except `make_null`, should take $O(\log n)$ time. The running time of all the operations, except `make_null`, is $O(d)$, where d is the depth of the node containing the accessed key. However, how do we get the average depth of $O(\log n)$?

Proof. Let $D(n)$ be the internal path length for some tree T of n nodes. The internal path length is the sum of the depths of all nodes in a tree, and $D(1) = 0$. A n -node tree consists of an i -node left subtree and an $(n - i - 1)$ -node right subtree, plus a root at depth zero for $0 \leq i < n$.

Then we have $D(i)$ as the internal path length of the left subtree with respect to its root, and we obtain

$$D(n) = D(i) + D(n - i - 1) + n - 1$$

If all subtree sizes are equally likely, which is true for a binary search tree, then the average value of both $D(i)$ and $D(n - i - 1)$ is

$$\frac{1}{n} \sum_{j=0}^{n-1} D(j)$$

Which yields

$$D(n) = \frac{2}{n} \left[\sum_{j=0}^{n-1} D(j) \right] + n - 1$$

■

This recurrence gives an average value of $D(n) = O(n \log n)$. Thus, the expected depth of any node should be $O(\log n)$.

Another thing to notice is that for the deletion operation, it favors the left subtree. So, after many insertions and deletions, we may end up with an unbalanced binary tree, which would look like a degenerated tree. To ensure that all the nodes can be operated on in $O(\log n)$ time, we need to make the binary search tree balanced. This is why we have the AVL tree.

4.5 AVL Tree

4.5.1 Definition

An AVL (Adelson-Velskii and Landis) tree is a binary search tree with a balancing condition. It is identical to a binary search tree, having the same physical and ordering properties. However, for every node in the AVL tree, the heights of the left and right subtrees can differ by at most 1.

With an AVL tree, all tree operations, except insertion, can be performed in $O(\log n)$ time.

To construct the smallest AVL tree of height n , we can use the two smallest AVL subtrees of heights $n - 1$ and $n - 2$. By doing so recursively, we can find the smallest AVL tree.

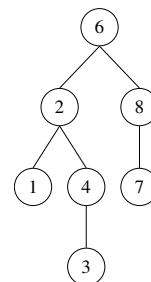


Figure 4.9: AVL Tree

The height of an empty tree is defined to be -1 . Height information is kept for each node. The height of an AVL tree is at most roughly $1.44 \log(n + 2) - 0.328$, but in practice, it is about $\log(n + 1) + 0.25$.

4.5.2 Operations

All tree operations can be performed in $O(\log n)$ time, except possibly insertion. Insertion and deletion operations need to update the balancing information since they might violate the AVL tree property. Therefore, we need to restore the property by means of rotations.

A single rotation involves only a few pointer changes and alters the structure of the tree while preserving the search tree property. Rotations happen from the bottom up, meaning we start checking balancing conditions from the lowest affected node and move upward.

For insertion into the right subtree of the right child, we perform a Left Rotation; for insertion into the left subtree of the left child, we perform a Right Rotation.

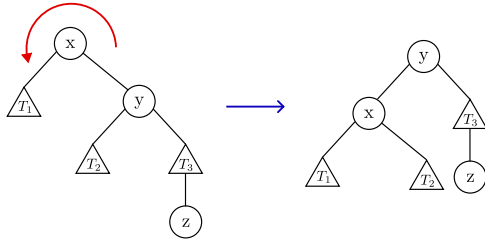


Figure 4.10: Left Rotation

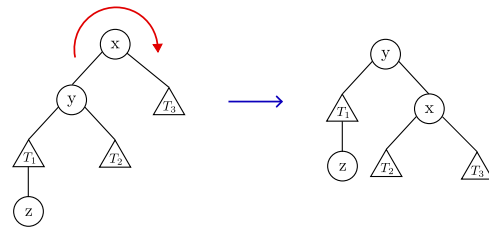


Figure 4.11: Right Rotation

However, single rotation might not work for more complex cases, where the height imbalance is caused by a node inserted into the tree containing the middle element, while the other subtrees have identical heights. In such cases, we use a double rotation, which involves four subtrees instead of three.

For insertion into the right subtree of the left child, we use a Left-Right Rotation. For insertion into the left subtree of the right child, we use a Right-Left Rotation. These double rotations help restore balance by first performing a single rotation on the child node and then performing a second rotation on the parent node, ensuring that the AVL tree's balance property is maintained.

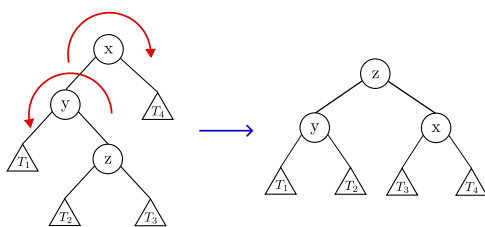


Figure 4.12: Left-Right Rotation

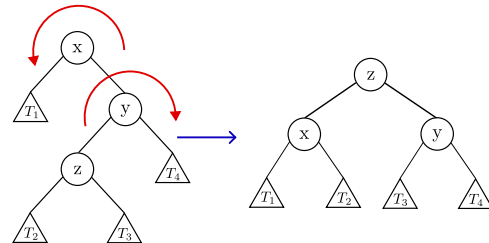
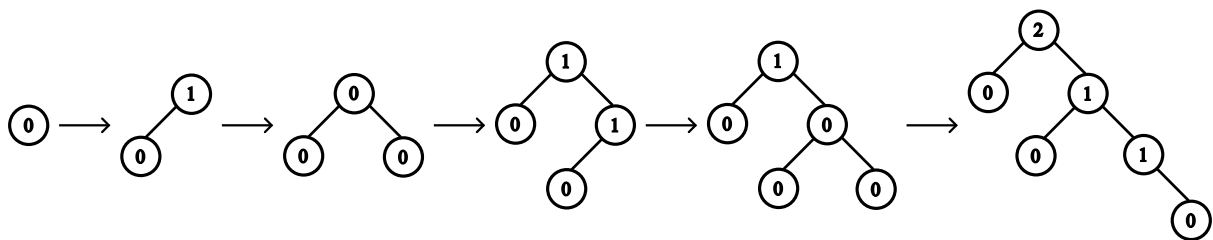


Figure 4.13: Right-Left Rotation

These four rotations help us perform all the necessary balancing operations.

To insert a new node with key x into an AVL tree T , we recursively insert x into the appropriate subtree of T_{lr} . If the height of T_{lr} does not change, then we are done. Otherwise, if a height imbalance occurs, we perform the appropriate single or double rotation depending on x and the keys in T and T_{lr} . Finally, we update the height. For example, the method for checking imbalance is demonstrated below:



The implementation of an AVL tree is similar to a binary search tree, with the addition of height information stored in each node. This height information allows the tree to maintain balance by ensuring that the difference in height between the left and right subtrees of any node is at most 1.

4.6 B-Tree

4.6.1 Definition

A B-Tree of order m is a tree with the following properties:

- All the leaf nodes must be at the same level (have the same depth).
- All non-leaf nodes except the root must have at least $\lceil \frac{m}{2} \rceil - 1$ keys and a maximum of $m - 1$ keys.
- All non-leaf nodes except the root (i.e., all internal nodes) must have children between $\lceil \frac{m}{2} \rceil$ and m .
- The root is either a leaf or has between 2 and m children.
- A non-leaf node with $n - 1$ keys must have n children.
- All the key values within a node must be in ascending order.

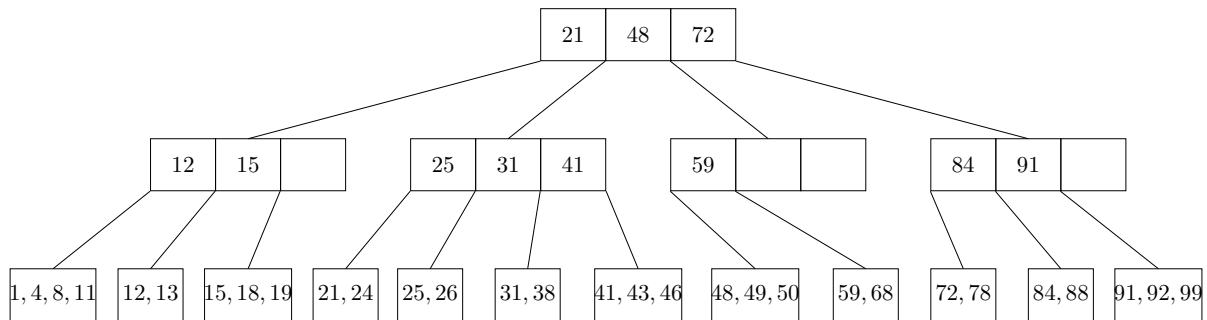


Figure 4.14: B-tree of order 4

A B-tree of order 4 is also known as a 2-3-4 tree; a B-tree of order 3 is also known as a 2-3 tree.

4.6.2 Operations

To insert a key into a node, we need to note the maximum number of values that the node can store. If the node isn't full, we can simply insert the key. However, there are some cases that need to be considered:

To insert 1, we find the location. However, since the leftmost node is full, we cannot insert it there. This can be solved by splitting the node into two nodes with two keys, then adjusting the information of the parent. Next, we try to insert 19. Since the rightmost node is full, we need to split it into two nodes with two children. Then, we continue splitting upwards to the root until we either reach the root node or find a node with fewer than two children.

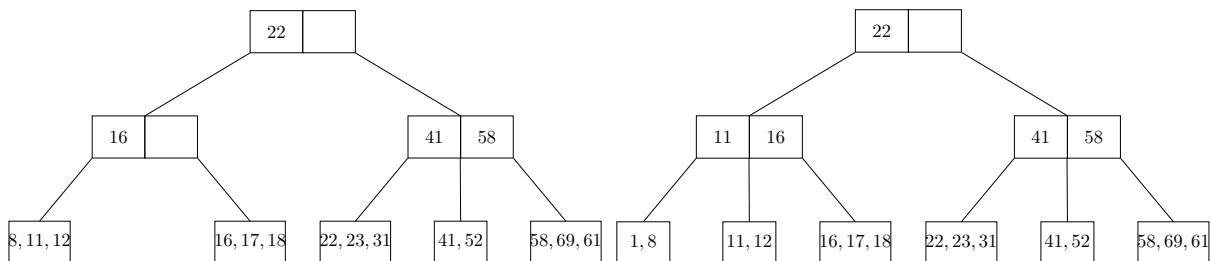


Figure 4.15: Original Tree

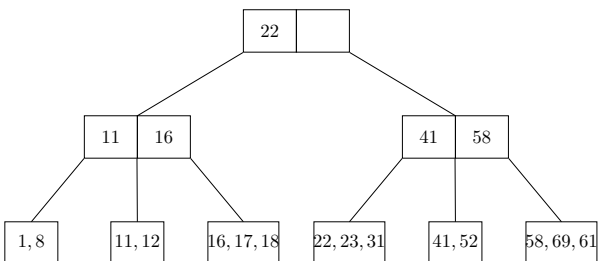


Figure 4.16: Insert 1

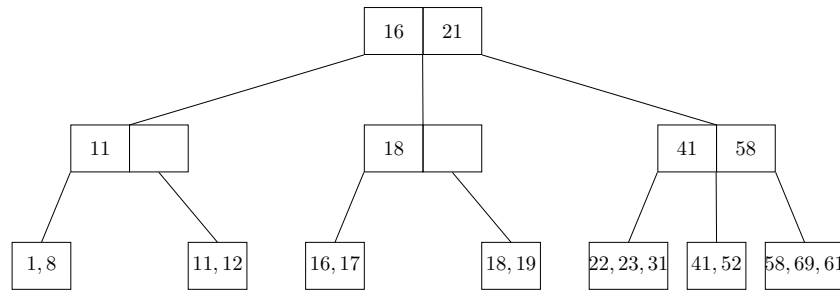


Figure 4.17: Insert 19

The depth of a B-tree is at most $\lceil \log_{\lceil \frac{m}{2} \rceil} n \rceil$. At each node along the path, we perform $O(\log m)$ work to determine which branch to take. An insertion or deletion could require $O(m)$ work to fix up all the information at the node. The worst case for insertion and deletion would be $O(m \log mn) = O\left(\frac{m}{\log m} \log n\right)$.

In summary, trees are used in operating systems, compiler design, and searching. In practice, all the balanced tree schemes are worse than the simple binary search tree, but this is acceptable.

Chapter 5

More on Tree

5.1 Tries

5.1.1 Definition

A trie, also called a digital tree, radix tree, or prefix tree, is a special type of tree used to store associative data structures. The name trie comes from its use for **retrieval**, because the trie can find a single word in a dictionary with only a prefix of the word.

For example, strings are stored in a top-to-bottom manner based on their prefixes in a trie. All prefixes of length 1 are stored at level 1, all prefixes of length 2 are stored at level 2, and so on. For the string set $S = \{\text{bear, bell, bid, bull, but, sell, stock, stop}\}$, we can have the tree structure as shown on the right.

This figure shows how the strings are stored in the trie. Inserting and deleting words in this tree is straightforward. For example, when typing "belt," it can be inserted in the subtree starting from $b \rightarrow e \rightarrow l$.

Unlike a binary search tree, no node in a trie stores the key associated with that node; instead, its position in the tree defines the key with which it is associated.

All the descendants of a node share a common prefix in the string associated with that node, and the root is associated with the empty string.

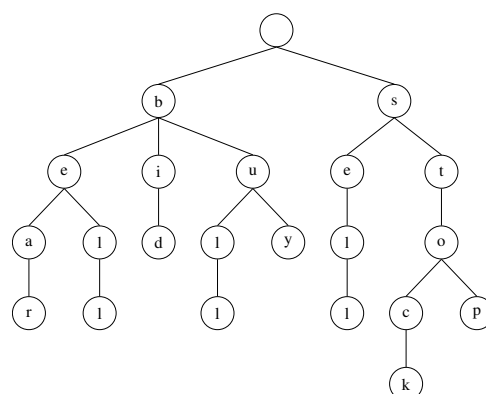


Figure 5.1: Tries

5.1.2 Applications

A trie can also be used to replace a hash table, offering the following advantages:

1. **Faster lookup:** Looking up data in a trie is faster than the worst case of a hash table. For a trie, the time complexity is $O(m)$, where m is the length of the search string. However, the time for an imperfect hash table would be $O(n)$, where n is the total number of strings.
2. **No collisions:** There are no collisions of different keys in a trie, unlike hash tables, where collisions can occur when two keys map to the same hash value.
3. **No need for a hash function:** In a trie, there is no need to provide a hash function or change it as more keys are added, unlike hash tables that require such adjustments.
4. **Alphabetical ordering:** A trie can provide an alphabetical ordering of the entries by key, making it useful for applications like autocomplete or lexicographical ordering.

A common application of a trie is storing a predictive text or autocomplete dictionary, such as those

found on mobile phones. It is also used in web browsers, which autocomplete your text or show possible completions for the text you are typing. Additionally, a trie can act as an orthographic corrector, checking whether every word you type exists in a dictionary.

5.1.3 Analysis

A standard trie uses $O(n)$ space and supports searches, insertions, and deletions in time $O(dm)$, where n is the total size of the strings in S , m is the size of the string parameter of the operation, and d is the size of the alphabet.

5.2 B-Tree

Here we continue the discussion on B-Trees.

A B-Tree is a self-balanced search tree with multiple keys in each node and can have more than two children per node. The properties of B-Trees have been discussed in the previous chapter.

Following these properties, for example, a B-Tree of order 4 contains a maximum of 3 key values in a node and a maximum of 4 children for a node.

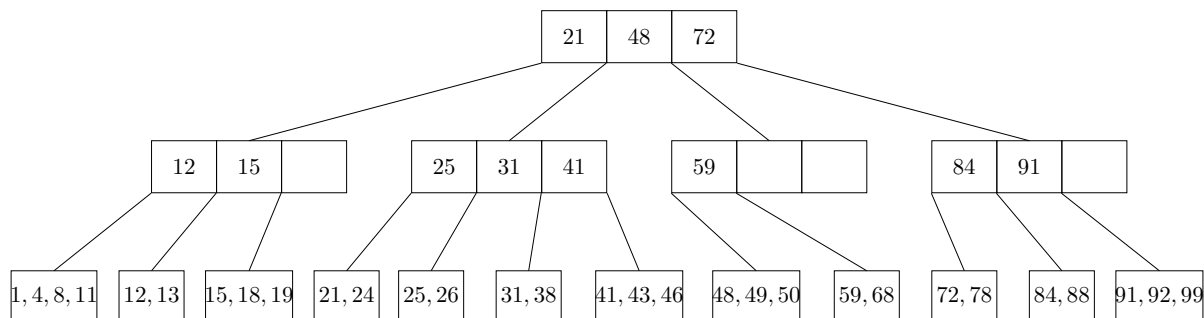


Figure 5.2: B-tree of order 4

5.2.1 Operations

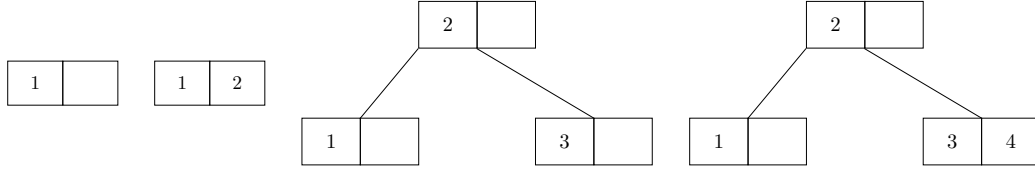
To search for a certain element in a B-Tree, we first read the value from the user. Then, we compare the search element with the first key value of the root node in the tree. If they match, we terminate the function immediately. If they do not match, we check whether the search element is smaller or larger, then move to the appropriate subtree where the element could be found. This process is done recursively until we either find the exact match or complete the comparisons with the last key value in a leaf node.

In a B-Tree, a new element must be added only at a leaf node. This means new key values are always attached to leaf nodes.

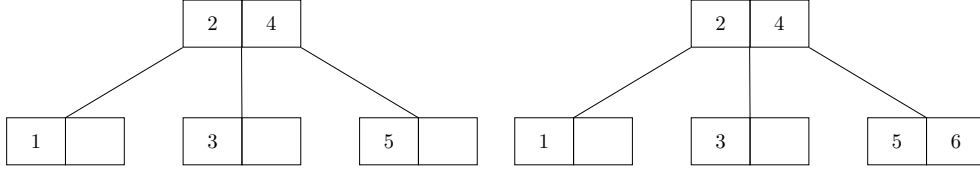
To perform an insertion, we first check if the tree is empty. If it is empty, we create a node with the new key value and insert it into the tree as the root node. Otherwise, we find the appropriate leaf node where the new key value can be added, using binary search tree logic. If the leaf node has an empty position, we add the new key value to the leaf node, maintaining the ascending order of key values within the node.

If the leaf node is already full, we split that leaf node by sending the middle value to its parent, and repeat this process until the value is placed into a node. If the split reaches the root, the middle value becomes the new root node, increasing the height of the tree by one.

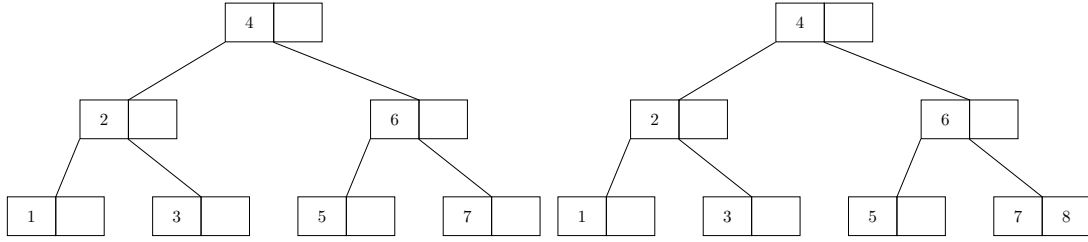
Following, we construct a B-Tree of order 3 by inserting numbers from 1 to 8.



For 1 and 2, we can directly do the insertion. For 3, since the node is full, we split the node by sending the middle value 2 to the parent node. Since there is an empty position, we can directly insert 4.



Since the rightmost node is full, we insert 5 by sending the middle value to the parent node and split the node. Then, we can directly insert 6.



Since the rightmost node is again full, we insert 7 by sending the middle value and splitting the node. Since the parent is also full, we further split it and send the middle value, in this case 4, to the parent node. To insert 8, it can be done directly.

To do deletion, if the key to be deleted is in a leaf and it contains more than the minimum number of keys, then this key can be deleted with no further action. If it is not a leaf, swap it with its successor under the natural order of the keys, then delete the key from the leaf.

If the node contains the minimum number of keys, consider the two immediate siblings of the node. If one of these siblings has more than the minimum number of keys, then redistribute one key from this sibling to the parent node, and one key from the parent to the deficient node. This is a rotation that balances the nodes.

If both immediate siblings have exactly the minimum number of keys, then merge the deficient node with one of the immediate sibling nodes and one key from the parent node. If this leaves the parent node with too few keys, then the process is propagated upward.

5.2.2 Analysis

The depth of a B-Tree is h , then

$$h \leq \log_t \frac{n+1}{2}, \quad t = \left\lceil \frac{m}{2} \right\rceil, \quad n = \text{the number of keys.}$$

The search, insertion, and deletion time is $O(t \log_t n)$.

Chapter 6

Hashing

From linear search, which takes $O(n)$ time, to binary search, which takes $O(\log n)$ time, the time is reduced. However, is there another data structure that allows better time? The answer is hashing.

6.1 Introduction

In some applications, fully utilizing the entire array is rare, leading to sparse arrays or sparse matrices. To avoid the multiplication operations required to calculate the index of an entry, we use an access table. This auxiliary table contains values like $0, n, 2n, 3n, \dots, (m-1)n$. When referencing the rectangular array, the index for entry $[i, j]$ is computed by accessing the value at position i in the auxiliary table, adding j , and using the resulting position.

A table with index set I and base type T can be viewed as a function from I to T , along with two main operations:

- Table access: Evaluate the function at any index in I .
- Table assignment: Modify the function by changing the value at a specified index in I to a new value.

Insertion involves adding a new element x to the index set I and defining the corresponding value for the function at x . Deletion removes an element x from the index set I , restricting the function to the resulting smaller domain.

However, array indices are not natural identifiers for the items we want to store, access, and retrieve. This limitation is overcome by using hashing.

6.1.1 Hashing

Hashing is a technique used to perform insertions, deletions, and lookups in constant average time. However, hashing does not efficiently support tree operations that require ordering information among elements. There are three important components of hashing:

1. Hash function: to generate a key;
2. Hash table: to store the elements;
3. Collision resolution: to resolve conflicts when two keys hash to the same index.

A hash table is an abstract data type that allows storing and retrieving elements in constant time, $O(1)$. This is true when the indices are known and the value at the target index of a store operation can be discarded. Without a complete set of items, we cannot determine the index of an element in a sorted list.

To solve this problem, we use the item as a key and convert it into unique integers that can be used as array indices. Since the index integer is not known at the entry of an item, it would be helpful if the item itself could be used as a key to index the cell where it will be stored. For example, to store names in an array, rather than assigning arbitrary indices, we can sum the ASCII values of each character in

the key (e.g., $a = 1, b = 2, \dots$).

6.1.2 General Idea

As defined previously, a function that performs the conversion is called a hash function. The conversion process is called hashing, and the storage structure used is called a hash table or scatter-storage.

In general, a hash table is an array of fixed size that contains keys. Each key is associated with a value and mapped to a number in the range 0 to $H_SIZE - 1$, placing it in the corresponding cell.

The mapping of keys to indices is done by the hash function, which ideally should be simple to compute and should ensure that distinct keys are placed in different cells. However, this is challenging because there are a finite number of cells and a virtually limitless number of potential keys. Therefore, the goal is to design a hash function that distributes the keys as evenly (uniformly) as possible among the cells.

There are several important issues that need to be addressed, which will be discussed in this chapter:

1. Choosing the hash function.
2. Handling collisions.
3. Handling deletions.

6.2 Hash Table

The idea of a hash table is to allow many of the different possible keys that might occur to be mapped to the same location in an array under the action of the index function. This is sometimes referred to as scatter-storage or key-transformation.

A hash function takes a key and maps it to an index in the array. However, two or more records may be mapped to the same location, leading to a collision. Therefore, a collision resolution procedure must be devised to handle this situation.

6.3 Hash Function

6.3.1 Analysis

The two principal criteria in selecting a hash function are that:

1. it should be easy and quick to compute,
2. it should achieve an even distribution of the keys that actually occur across the range of indices.

If the input keys are integers, then simply returning $\text{key} \bmod H_SIZE$ is generally a reasonable strategy. For example, $\text{student_ID} \bmod 10000$ would be a reasonable strategy. Also, it is usually a good idea to ensure that the table size is prime. When the input keys are random integers, this function is simple to compute and also distributes the keys evenly.

6.3.2 Truncation

Truncation can be done by ignoring part of the key and using the remaining part directly as the index (considering non-numeric fields as their numerical codes). For example, if the keys are eight-digit integers and the hash table has 1000 locations, then the first, second, and fifth digits from the right make the hash function, so that 62538194 maps to 394. Although this method is fast, it often fails to distribute the keys evenly through the table.

6.3.3 Folding

We can also partition the key into several parts and combine the parts in a convenient way (often using addition or multiplication) to obtain the index. For example, we can map 62538194 to $625 + 381 + 94 = 1100$, which is truncated to 100.

6.3.4 Modular Arithmetic

We can use modular arithmetic to convert the key into the index that we want. We can simply use the ASCII values of all the characters in the string, and then return the `result mod H_SIZE`.

For example, `'abcd' mod 100 = (64 + 65 + 66 + 67)%100 = 62`.

We can also take only 3 characters and use `hash_val = key[0] + 27*key[1] + 729*key[2]`, and then again return the `hash_val mod H_SIZE`, assuming that the key has at least two characters plus the NULL terminator. If the three characters are random, and the table size is 10007, then we would have a reasonably equitable distribution. However, since English is not random, the percentage of the table being hashed to could be less than expected.

We can improve the hash function by using `hash_val = hash_val << 5 + *key++`, then returning `hash_val mod H_SIZE`. This hash function involves all the characters in the key by computing:

$$\sum_{i=0}^{\text{keySize} - 1} \text{key}[\text{keySize} - i] \times 32^i,$$

which follows Horner's Rule.

It is common to not use all the characters in the key, since the length and properties of the key would influence the choice.

6.4 Collision Resolution

Though a hash table is an efficient data structure, there could be cases where different elements share the same index, causing a collision. Thus, we need to handle such collisions, and there are several resolutions.

6.4.1 Open Hashing

The first strategy, commonly known as open hashing or separate chaining, is to keep a list of all elements that hash to the same value.

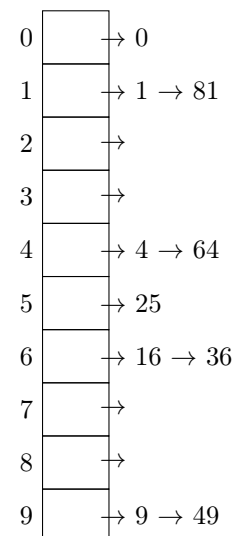
To illustrate, we assume that the keys are the first 10 perfect squares, where the hashing function is given by $\text{hash}(x) = x \bmod 10$.

In open hashing, each node contains a linked list, allowing multiple elements to be stored at the same index.

We can perform **find** in open hashing by using the hash function to determine which list to traverse. Then, we traverse this list in the normal manner, returning the position where the item is found.

To perform **insertion**, we first traverse the list to check whether the element is already present. If it is a new element, it is inserted either at the front or the end of the list. Often, new elements are inserted at the front, as this is convenient and because recently inserted elements are frequently the most likely to be accessed again soon.

The **deletion** routine is straightforward, similar to linked lists. After performing a find operation, we delete the item as we would in a linked list.



The use of linked storage can save a considerable amount of space. It allows for simple and efficient collision handling, and the size of the hash table no longer needs to exceed the number of records we

have. Additionally, deletion is faster compared to other methods.

However, all the pointers (links) require space, and even if the number of records is small, the space used is comparatively larger. Thus, it takes longer time to allocate the new cells, and it requires the implementation of another data structure. Thus, we introduce Closed Hashing.

6.4.2 Closed Hashing

In closed hashing, if a collision occurs, alternate cells are tried until an empty cell is found.

For example, cells $h_0(x), h_1(x), \dots$ are tried in succession where $h_i(x) = (\text{hash} + f(i)) \bmod \text{H_SIZE}$, with $f(0) = 0$.

The function $f(i)$ is called the collision resolution strategy. Because all the data goes inside the table, a bigger table is needed for closed hashing than for open hashing. Generally, the load factor should be below $l = 0.5$ for closed hashing.

Again, we want to perform insertion. An array must be declared that will hold the hash table. Then, all locations in the array need to be initialized to show that they are empty. To insert a record, the hash function for the key is first calculated. If the corresponding location is empty, then the record can be inserted. Otherwise, insertion of the new record would not be allowed, which is what we call a collision.

For the find operation, we again calculate the key first. If the desired record is in the corresponding location, then the retrieval has succeeded. Otherwise, we follow the same procedure as for collision resolution, examining all locations until we find the correct record. However, if the position is empty, then no record with the given key is in the table, and the search is thus unsuccessful.

What we haven't discussed yet is the way to resolve collisions. In general, there are three methods.

Linear Probing

This is the simplest method to resolve a collision. Starting with the hash address where the collision happens, do a sequential search for the desired key or an empty location.

However, data could become clustered in this method. That is, records start to appear in long strings of adjacent positions with gaps between the strings.

For example, given $x = 89, 18, 49, 58, 69$, with $\text{hash}(x) = x \bmod 10$, we can perform insertion as shown on the right.

For values 89 and 18, they can be inserted directly since the corresponding cell for the address is empty. However, when inserting 49, the key value 9 is repeated. Then it finds the next available cell, which in this case would be at position 0. The same applies to 58 and 69.

	Empty Table	89	18	49	58	69
0				49	49	49
1					58	58
2						69
3						
4						
5						
6						
7						
8			18	18	18	18
9		89	89	89	89	89

Although they can all be inserted, the time to search for an empty cell may be long. Since the records are not distributed uniformly and become progressively more unbalanced, this leads to the problem of primary clustering. This problem is essentially one of instability. If a few keys happen randomly to be near each other, then it becomes more and more likely that other keys will join the cluster.

Assuming a very large table and that each probe is independent of the previous probes, the number of probes for a successful search is equal to the number of probes required when the particular element is inserted. When an element is inserted, it is done as a result of an unsuccessful search. We can use the cost of an unsuccessful search to compute the average cost of a successful search.

Since the fraction of empty cells is $1 - \lambda$, the number of cells we expect to probe is $\frac{1}{1-\lambda}$.

Then, it can be shown that the expected number of probes using linear probing is roughly

$$\frac{1}{2} \left(1 + \frac{1}{(1 - \lambda)^2} \right)$$

for insertions and unsuccessful searches. It takes roughly

$$\frac{1}{2} \left(1 + \frac{1}{(1 - \lambda)} \right)$$

for successful searches. Here, λ is the ratio of the number of elements in the hash table to the table size.

This is inefficient since, intuitively, we need to keep probing, which takes longer. As λ increases, the expected number of probes also increases. This issue becomes more significant if the table is expected to be more than half-full.

Thus, we have quadratic probing.

Quadratic Probing

Quadratic probing avoids the primary clustering problem of linear probing. If there is a collision at hash address H , the method called quadratic probing looks in the table at locations $h+0, h+1, h+4, h+9, \dots$, that is, at location $h + i^2$ for $i = 0, 1, 2, \dots$.

This reduces clustering, but it is obvious that it will not probe all locations in the table. The idea is that if quadratic probing is used and the table size is prime, then a new element can always be inserted if the table is at least half empty.

Here we use the same example. Now we use quadratic probing instead. After 18 and 89 being inserted, for 49, it goes to $(9 + 1^2) \bmod 10 = 0$ since location 9 is occupied, for 58, it goes to $(8 + 2^2) \bmod 10 = 2$, and for 69 it goes to $(9 + 2^2) \bmod 10 = 3$.

Note that if the table is even more than half full, insertion could fail. It is also crucial that the table size is a prime number. Otherwise, the number of alternate locations can be severely reduced.

	Empty Table	89	18	49	58	69
0				49	49	49
1						
2					58	58
3						69
4						
5						
6						
7						
8			18	18	18	18
9		89	89	89	89	89

Standard deletion cannot be performed in a closed hash table since the cell might have caused a collision to pass it. Simply put, deleting an entry could break the probing chain, causing searches to miss other items. Thus, lazy deletion is required.

In lazy deletion, we delete an entry by placing a special marker or key in the deleted position. This marker indicates that the position is free for future insertions but should not be treated as empty when searching for other items — ensuring the probing sequence remains unbroken.

Random Probing

Rather than having the increment depend on the number of probes already made, we can let it be a simple function of the key itself. For example, we could truncate the key to a single character and use its code as the increment.

For random probing, we use a pseudo-random number generator to obtain the increment. The generator should always produce the same sequence when given the same seed. This method is excellent at avoiding clustering, but it is likely to be slower than others.

Double Hashing

Another closed hashing method is double hashing. For double hashing, one popular choice is:

$$f(i) = i \times h_2(x).$$

We apply a second hash function to x and probe at distances $h_2(x), 2h_2(x), \dots$, and so on. However, a poor choice of $h_2(x)$ could be disastrous. The function must never evaluate to zero, and it needs to ensure that all cells can be probed.

For instance, the obvious choice $h_2(x) = x \bmod 9$ would not help if 99 were inserted into the input in the previous example. A function such as:

$$h_2(x) = R - (x \bmod R),$$

with R being a prime number smaller than H_SIZE , works well. One may also perform triple hashing, and so on.

Again, we perform insertion for $x = 89, 18, 49, 58, 69$ using double hashing. Here, we have $h_2(x) = R - (x \bmod R)$, with $R = 7$.

After inserting 89 and 18, 49 causes a collision. Then we use $h_2(49) = 7 - (49 \bmod 7) = 7$, so it goes to $(7 + 49) \bmod 10 = 6$. Inserting 58 also causes a collision, then we have $h_2(58) = 7 - (58 \bmod 7) = 5$, so it goes to $(5 + 58) \bmod 10 = 3$.

	Empty Table	89	18	49	58	69
0						69
1						
2						
3					58	58
4						
5						
6				49	49	49
7						
8			18	18	18	18
9		89	89	89	89	89

6.4.3 Rehashing

There could be cases where the hash table becomes too full, causing the running time for operations to deteriorate, especially when there are too many deletions intermixed with insertions.

To solve this issue, we can build another table that is about twice as large with an associated new hash function. Then, we scan the entire original hash table, compute the new hash value for each element, and insert it into the new table. This process is known as rehashing.

For example, for a table of size 7, let $h(x) = x \bmod 7$. After inserting $x = 13, 15, 24, 6$, we insert 23, which makes the table now 70% full. Thus, we create a new table with size 17. The new hash function is then $h(x) = x \bmod 17$. The old table is scanned, and elements 6, 23, 24, 13, and 15 are then inserted sequentially.

The running time is $O(n)$, which is expensive. Therefore, it should not be done too frequently. This leads us to the concept of extendible hashing.

0	6
1	15
2	
3	24
4	
5	
6	13

0	6
1	15
2	23
3	24
4	
5	
6	13

0	
1	
2	
3	
4	
5	
6	6
7	23
8	24
9	
10	
11	
12	
13	13
14	
15	15
16	
17	

When the amount of data is too large to fit in main memory and must be stored on disk, we minimize disk access by using a tree.

In definition, the root of the tree is called the *directory*, and the leaves are called *buckets*. The number of bits used by the root is called the *global depth*, and the *bucket size* is the maximum number of elements that a leaf can contain.

Suppose that our data consists of several 6-bit integers. With global depth D equal to 2, the root of the tree contains 4 pointers determined by the leading two bits of the data. Each leaf can hold up to $m = 4$ elements.

For example, we have the following tree. To insert 36, we first convert it to binary 100100_2 . The hash function here is the D most significant bits (MSBs), which is the 2 MSBs in this case. Since the bucket of 00 is full, and the global depth is equal to the bucket depth, the directory is split, and the roots now have $D = 3$.

Then, we update the content by inserting 100100_2 into the directory 100. Note that for the same bucket, there could be more than one pointer pointing to it.

Next, to insert 000000_2 , we again encounter a collision. However, since the global depth is now larger than the bucket depth, the bucket itself is split instead of the directory. The element can then be inserted.

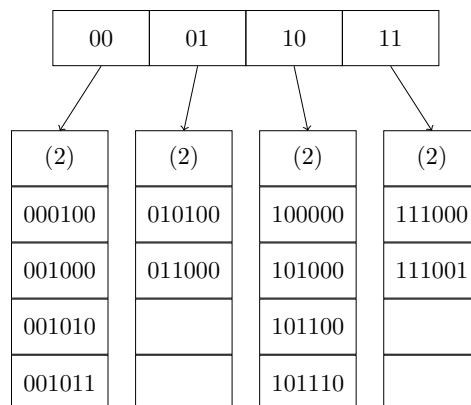


Figure 6.1: Original

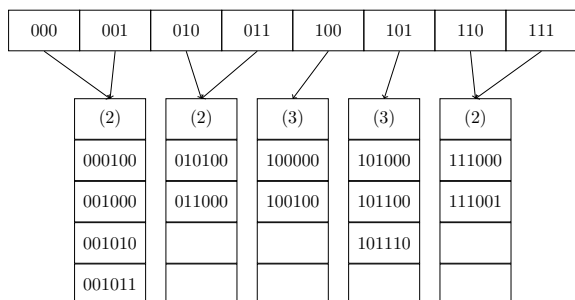


Figure 6.2: Insertion 100100

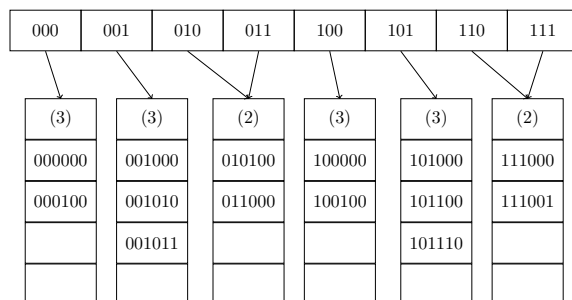


Figure 6.3: Insertion 000000

It is possible that several directory splits will be required if the elements in a leaf agree in more than $D + 1$ leading bits, i.e., overflow. For example, to insert $111010_2, 111011_2, 111100_2$, the directory size must be increased to 4.

This algorithm does not work when there are more than m duplicates. It is important for the bits to be fairly random.

In summary, hash tables can be used to implement the insertion and find operations in constant average time. It is especially important to pay attention to details such as the load factor when using hash tables. The choice of hash function is also crucial, especially when the key is not a short string or integer.

For open hashing, the load factor should be 1, meaning all the cells can be filled. For closed hashing, the load factor should not exceed 0.5, unless unavoidable. It is not possible to find the minimum or maximum since there is no sorting order.

In applications, hash tables can be used in compilers to keep track of declared variables in source code. They are useful for any graph theory problem where nodes have real names instead of numbers. Additionally, hash tables are commonly used in programs that play games or even online spelling checkers.