

Formal DSA in C++

Nathan Warner



Northern Illinois
University

Computer Science
Northern Illinois University
United States

Contents

1	Preface	9
1.1	Universe	9
1.2	Dynamic and static sets	9
2	Elementary complexity theory	10
3	Linked lists	17
3.1	Singly-linked lists	17
3.1.1	Structure of the node	17
3.1.2	The list class/struct	18
3.1.3	Interface of a singly linked list stack	18
3.1.4	Traversing	19
3.1.5	Printing	20
3.1.6	Printing in reverse	21
3.1.7	Getting the length	22
3.1.8	Clearing	23
3.1.9	Reversing	24
3.1.10	Pushing	25
3.1.11	Inserting	26
3.1.12	Popping	27
3.1.13	Erasing	28
3.1.14	Searching	29
3.2	Doubly-linked list	30
3.2.1	Node structure	30
3.2.2	The list class	30
3.2.3	push_back and push_front	31
3.2.4	pop_back and pop_frot	32

3.2.5	Searching	33
4	Recursion	34
4.1	Recursion vs iteration	34
4.2	Elementary recursion	35
4.3	Base cases	35
4.3.1	Factorials	35
4.3.2	Powers	36
4.4	Tail recursion	37
4.5	Indirect Recursion	39
4.6	Nested Recursion	40
4.7	Excessive Recursion	40
4.8	Backtracking	42
4.8.1	N-Queens	43
4.8.2	Greedy backtracking: Traveling salesman problem (TSP)	44
4.8.3	Combinations	48
4.8.4	Combinations that sum to m	49
4.8.5	Permutations	50
4.9	Recursion in singly linked lists	51
4.9.1	Traversing	51
4.9.2	Printing	51
4.9.3	Printing in reverse	51
4.9.4	Getting the length	52
4.9.5	Clearing	52
4.9.6	Reversing	53
4.9.7	Pushing	55
4.9.8	Inserting	56
4.9.9	Popping	57
4.9.10	Erasing	58
4.9.11	Searching	59
5	Dynamic programming	60
5.1	Key Concepts of DP	60

5.2	Bottom-up DP	60
5.2.1	Key Characteristics	60
5.2.2	Steps for Bottom-Up DP	60
5.3	Bottom-up DP: Fibonacci	61
5.4	Memoization (top down DP)	61
5.5	Top-down DP: Fibonacci	62

6 Binary trees 63

6.1	Terminology	63
6.2	Type of binary trees	64
6.3	Maximum height of a binary tree	65
6.3.1	Minimum height of a binary tree	65
6.3.2	Number of Leaves in a Binary Tree	65
6.3.3	Relationship Between Internal Nodes and Leaves:	65
6.3.4	Maximum Number of Nodes at Height h	65
6.3.5	Number of Edges in a Binary Tree:	66
6.4	Full trees	67
6.4.1	Number of leaves	67
6.4.2	Number of nodes	68
6.4.3	Number of internal nodes	68
6.5	Complete Binary Tree	69
6.5.1	Number of nodes	69
6.5.2	Height	69
6.5.3	Number of Leaf Nodes (L) in a Complete Binary Tree	69
6.5.4	Number of internal nodes	70
6.5.5	Parent and Child Relationships in a Complete Binary Tree	70
6.6	Perfect binary tree	71
6.6.1	Number of Nodes	71
6.6.2	Number of Leaf Nodes	71
6.6.3	Height of the Tree	71
6.6.4	Number of Internal Nodes	71
6.6.5	Depth	71

7	Applications of binary trees	72
7.1	Binary search trees	72
7.1.1	Interface	72
7.1.2	Traversals	73
7.1.2.1	Level order	73
7.1.2.2	Preorder	74
7.1.2.3	Inorder	74
7.1.2.4	Postorder	76
7.1.3	Successor of a node	76
7.1.4	Predecessor	77
7.1.5	The node	77
7.1.6	The class	77
7.1.7	Recursive Insertion	78
7.1.8	A better recursive insert	79
7.1.9	Iterative insert	80
7.1.10	Recursive removing	81
7.1.11	Clearing	83
7.1.12	Counting the height of the tree (root)	84
7.1.13	Counting the height of a node	84
7.1.14	Getting the depth of the node	85
7.1.15	Counting the number of nodes	86
7.1.16	Comparison traversals	86
7.1.17	Finding the smallest and largest values	88
7.1.18	Getting the widths of a bst	89
7.1.19	Degenerate Binary Search trees	91
7.1.20	Verifying a binary search tree	91
7.1.21	Complexities	92
7.2	Adelson-Velsky and Landis Trees (AVL trees)	93
7.2.1	Definition	93
7.2.2	AVL Nodes	94
7.2.3	Storing the height	94
7.2.4	Defining balance factors in C++ with enums	95
7.2.5	Defining balance factors with a height calculation	95
7.2.6	Interface	96

7.2.7	Balancing an AVL tree	96
7.2.8	Rotations: Right tree	97
7.2.8.1	Case 1: Right higher	97
7.2.8.2	Case 2: Left higher	100
7.2.9	C++ Rotations	102
7.2.10	Balancing	104
7.2.11	Insertions	105
7.2.12	Removing nodes	106
7.3	Red-black trees	108
7.3.1	Rotations	110
7.3.2	Inserting	112
7.3.3	Deletion	117
7.3.4	C++ implementation	121
7.3.4.1	The node structure	121
7.3.4.2	The tree class and defining nil	121
7.3.4.3	Rotation methods	122
7.3.4.4	Insert	123
7.3.4.5	Insert fixup	124
8	Heaps and Priority Queues (Zero based)	125
8.1	Max and Min heaps	127
8.2	Heapify an array	127
8.3	Min-heap in c++	128
8.4	Max-heap in c++	129
8.5	Percolating	130
8.5.1	Percolate up	130
8.5.2	Percolate down	131
8.6	Inserting into a heap	132
8.7	Removing the root	133
8.8	Removing an arbitrary node	134
8.9	Priority queues	135
8.9.1	Interface	135
8.9.2	Insert, pop, and top	136
8.9.3	Size and Empty	137
9	Sorting	138

9.1	Bubble, selection, insertion	138
9.1.1	Bubble sort	138
9.1.1.1	Complexity	139
9.1.2	Selection sort	140
9.1.2.1	Complexity	140
9.1.3	Insertion sort	141
9.1.3.1	Complexity	142
9.2	Heap sort	143
9.2.1	Complexity	143
9.3	BST Sort	144
9.3.1	Insert	144
9.3.2	The inorder traversal	144
9.3.3	The BST sort function	145
9.3.4	Inplace sorting	146
9.3.5	Complexity	146
9.4	Quick sort	147
9.4.1	Partitioning	147
9.4.2	The quick sort procedure	148
9.4.3	A better partition	150
9.4.4	Median of three partition	150
9.4.5	Quicksort with iterators	152
10	Multi-way (m-way) search trees	153
10.1	Multi-way Searching	153
10.2	2-4 (2-3-4) Trees	154
10.2.1	Insertion	154
10.2.2	Removal	155
10.2.3	Properties	157
10.3	B-trees	158
11	Hashing (hash tables)	159
11.1	Direct-address table	159
11.2	Hash tables	160
11.3	Independent uniform hashing	161
11.4	Collision resolution by chaining	161

11.5	Analysis of hashing with chaining	162
11.6	Hash functions	162
11.6.1	Static hashing	163
11.6.2	The division method	163
11.6.3	The multiplication method	164
11.7	Open addressing	164
11.7.1	Linear probing	164
11.7.2	Quadratic probing	165
11.7.3	Deletion problem in open addressing	165
11.7.4	Tombstoning	165
11.7.5	Lazy deletion	166
11.7.6	Rehashing	166
11.8	Load factor	167
11.9	Static hashing	167
11.10	Static linear hashing	167
11.11	Static hashing with linear probe implementation	168
11.11.1	Interface	168
11.11.2	The basics	168
11.11.3	The hash function	168
11.11.4	Inserts	169
11.11.5	Searching	170
11.11.6	Removing	170
11.11.7	Searching and removing with tombstones	170
11.12	Static hashing with quadratic probing implementation	173
11.12.1	Inserting	173
11.12.2	Searching	174
11.12.3	Removing	175
11.13	Static hashing with chaining implementation	176
11.13.1	The basics	176
11.13.2	Inserting	176
11.13.3	Searching	176
11.13.4	Removing	177
11.14	Injective hashing	178

11.15	Perfect hashing	178
12	Table indexing and row-major order	179
12.1	Row-major order	179
13	Graphs	180
13.1	Simple Graph	180
13.2	Undirected graph	180
13.3	Directed graph	180
13.4	Weighted Graph	180
13.5	Complete Graph	181
13.6	More terms	181
13.7	Graph representations	181
13.8	Interface of adjacency lists and matrices	182
13.9	Complexities for adjacency operations	182
13.10	Graph traversal/search	182
13.10.1	Breadth-first-traversal	183
13.10.2	Depth-first traversal	183
13.11	Shortest path problems	183
14	Math algorithms	184
14.1	Euclidean GCD Algorithm	184
14.2	Fibonacci numbers in constant time	184
14.3	Sterlings factorial approximation	184

Preface

1.1 Universe

In the context of data structures and algorithms, the term universe refers to the complete set of all possible elements or values that could be involved in a particular problem or data structure. It defines the range or domain from which data can be selected. We denote the universe with a capital U

1.2 Dynamic and static sets

A dynamic set is a data structure that supports not only membership queries but also allows the insertion, deletion, and sometimes modification of elements over time. Unlike static sets, which are fixed once defined, dynamic sets can change as elements are added or removed

Elementary complexity theory

- **Idea:** The same problem can frequently be solved with algorithms that differ in efficiency. The differences between the algorithms may be immaterial for processing a small number of data items, but these differences grow with the amount of data. To compare the efficiency of algorithms, a measure of the degree of difficulty of an algorithm called computational complexity was developed by Juris Hartmanis and Richard E. Stearns

Computational complexity indicates how much effort is needed to apply an algorithm or how costly it is. This cost can be measured in a variety of ways, and the particular context determines its meaning. This book concerns itself with the two efficiency criteria: time and space. The factor of time is usually more important than that of space, so efficiency considerations usually focus on the amount of time elapsed when processing data. However, the most inefficient algorithm run on a Cray computer can execute much faster than the most efficient algorithm run on a PC, so run time is always system-dependent. For example, to compare 100 algorithms, all of them would have to be run on the same machine. Furthermore, the results of run-time tests depend on the language in which a given algorithm is written, even if the tests are performed on the same machine. If programs are compiled, they execute much faster than when they are interpreted. A program written in C or Ada may be 20 times faster than the same program encoded in BASIC or LISP.

- **Units:** To evaluate an algorithm's efficiency, real-time units such as microseconds and nanoseconds should not be used. Rather, logical units that express a relationship between the size n of a file or an array and the amount of time t required to process the data should be used

If there is a linear relationship between the size n and time t , that is, $t_1 = cn_1$, then an increase of data by a factor of 5 results in the increase of the execution time by the same factor. If $n_2 = 5n_1$, then $t_2 = 5t_1$

Similarly, if $t_1 = \log_2 n$, then doubling n increases t by only one unit of time. Therefore, if $t_2 = \log_2(2n)$, then $t_2 = t_1 + 1$.

- **Eliminating insignificant terms:** A function expressing the relationship between n and t is usually much more complex, and calculating such a function is important only in regard to large bodies of data; any terms that do not substantially change the function's magnitude should be eliminated from the function. The resulting function gives only an approximate measure of efficiency of the original function. However, this approximation is sufficiently close to the original, especially for a function that processes large quantities of data.
- **Asymptotic complexity:** This measure of efficiency is called asymptotic complexity and is used when disregarding certain terms of a function to express the efficiency of an algorithm or when calculating a function is difficult or impossible and only approximations can be found
- **Big-O Notation:** The most commonly used notation for specifying asymptotic complexity—that is, for estimating the rate of function growth—is the big-O notation introduced in 1894 by Paul Bachmann.

Given two positive-valued functions f and g , consider the following definition:

$f(n)$ is $O(g(n))$ if there exist positive numbers c and N such that $f(n) \leq c \cdot g(n)$ for all $n \geq N$.

$$f(n) \text{ is } O(g(n)) \iff \exists c, N \in \mathbb{Z}^+ \mid f(n) \leq cg(n) \forall n \geq N.$$

Big-O notation says that for large enough n , the function $f(n)$ does not grow faster than a constant multiple of $g(n)$. So, $g(n)$ provides an upper bound on how fast $f(n)$ can grow as n increases.

In other words, f is big-O of g if there is a positive number c such that f is not larger than $c \cdot g$ for sufficiently large ns ; that is, for all ns larger than some number N . The relationship between f and g can be expressed by stating either that $g(n)$ is an upper bound on the value of $f(n)$ or that, in the long run, f grows at most as fast as g .

The problem with this definition is that, first, it states only that there must exist certain c and N , but it does not give any hint of how to calculate these constants. Second, it does not put any restrictions on these values and gives little guidance in situations when there are many candidates. In fact, there are usually infinitely many pairs of c 's and N 's that can be given for the same pair of functions f and g .

For example, suppose

$$f(n) = 2n^2 + 3n + 1 = O(n^2).$$

Where $g(n) = n^2$. Candidate values for c and N are

c	≥ 6	$\geq 3^{\frac{3}{4}}$	$\geq 3^{\frac{1}{9}}$	$\geq 2^{\frac{13}{16}}$	$\geq 2^{\frac{16}{25}}$	\dots	\rightarrow	2
N	1	2	3	4	5	\dots	\rightarrow	∞

We obtain these values by solving the inequality:

$$2n^2 + 3n + 1 \leq cn^2.$$

Or equivalently

$$2 + \frac{3}{n} + \frac{1}{n^2} \leq c.$$

For different n 's

For large n , the terms $\frac{3}{n}$ and $\frac{1}{n^2}$ get smaller. Let's find N such that for all $n \geq N$, the right-hand side stays bounded.

As n gets larger, $\frac{3}{n}$ and $\frac{1}{n^2}$ approach zero. To simplify the analysis, choose $N = 1$ initially and check how small $\frac{3}{n}$ and $\frac{1}{n^2}$ are:

$$2 + \frac{3}{1} + \frac{1}{1^2} = 2 + 3 + 1 = 6.$$

From the inequality, at $N = 1$, we have $6 \leq c$. Therefore, we can choose $c = 6$. This ensures that for all $n \geq 1$, the inequality holds:

$$2 + \frac{3}{n} + \frac{1}{n^2} \leq 6.$$

Thus, you can choose $c = 6$ and $N = 1$.

different pairs of constants c and N for the same function $g(= n^2)$ can be determined.

- **Choosing the best c, N :** To choose the best c and N , it should be determined for which N a certain term in f becomes the largest and stays the largest.

In the example above, The only candidates for the largest term are $2n^2$ and $3n$; these terms can be compared using the inequality $2n^2 > 3n$ that holds for $n > 1.5$. Thus, $N = 2$ and $c \geq \frac{15}{4} = 3.75$.

- **Significance:** What is the practical significance of the pairs of constants just listed? All of them are related to the same function $g(n) = n^2$ and to the same $f(n)$. For a fixed g , an infinite number of pairs of c 's and N 's can be identified. The point is that f and g grow at the same rate. The definition states, however, that g is almost always greater than or equal to f if it is multiplied by a constant c . "Almost always" means for all n 's not less than a constant N . The crux of the matter is that the value of c depends on which N is chosen, and vice versa.
- **Inherent imprecision: Choosing best $g(n)$:** The inherent imprecision of the big-O notation goes even further, because there can be infinitely many functions g for a given function f . For example, the f from Equation 2.2 is big-O not only of n^2 , but also of $n^3, n^4, \dots, n^k, \dots$ for any $k \geq 2$. To avoid this embarrassment of riches, the smallest function g is chosen, n^2 in this case.
- **Big-o as approximating terms:** The approximation of function f can be refined using big-O notation only for the part of the equation suppressing irrelevant information. For example, in the equation below, the contribution of the third and last terms to the value of the function can be omitted

$$\begin{aligned} f(n) &= n^2 + 100n + \log(n) + 1000 \\ \implies f(n) &= n^2 + 100n + O(\log(n)). \end{aligned}$$

Similarly,

$$\begin{aligned} f(n) &= 2n^2 + 3n + 1 \\ \implies f(n) &= 2n^2 + O(n). \end{aligned}$$

This equation says that for large values of n , the expression $2n^2 + 3n + 1$ behaves like $2n^2$ plus some terms that grow linearly or slower (captured by $O(n)$). The exact contributions of $3n$ and 1 are not important for asymptotic analysis; what matters is that their growth is slower compared to $2n^2$.

- **Algorithm analysis: Most common time complexities:** Ranked slowest to fastest growth
 - $O(1)$: Constant time
 - $O(\log(\log(n)))$: Logarithmic time
 - $O(\log(n))$: Logarithmic time
 - $O(n)$: Linear time
 - $O(n \log(n))$: Log-linear time
 - $O(n^k)$, $k > 1$: Polynomial time
 - $O(a^n)$, $a > 1$: Exponential time
 - $O(n!)$: Factorial time

- **Ranking complexities from slowest to fastest: Process:** Given

- (a) $O(25)$
- (b) $O(n^{\frac{1}{2}} + \log^2(n))$
- (c) $O(\log^{200}(n))$
- (d) $O(n^3 \log^4(n))$
- (e) $O(n^{200} + 3^n)$
- (f) $O(n \log^{40}(n))$
- (g) $O(4^n \log(n))$
- (h) $O(n^3 \log(\log(n)))$

How can we go about sorting these slowest to fastest. Well, to start, in the expressions with plus or minus, we can throw out the slower terms. Thus,

- (a) $O(n^{\frac{1}{2}})$
- (b) $O(25)$
- (c) $O(\log^{200}(n))$
- (d) $O(n^3 \log^4(n))$
- (e) $O(3^n)$
- (f) $O(n \log^{40}(n))$
- (g) $O(4^n \log(n))$
- (h) $O(n^3 \log(\log(n)))$

In product terms, we disregard the slower term unless there are complexities with the same dominant term. For example, $O(n^3 \log(\log(n)))$ grows slower than $O(n^3 \log^4(n))$ because although they have the same dominant term n^3 , $\log(\log(n))$ grows slower than $\log^4(n)$. Thus, the correct sequence is

- (b) $O(25)$
- (c) $O(\log^{200}(n))$
- (a) $O(n^{\frac{1}{2}} + \log^2(n))$
- (f) $O(n \log^{40}(n))$
- (h) $O(n^3 \log(\log(n)))$
- (d) $O(n^3 \log^4(n))$
- (e) $O(n^{200} + 3^n)$
- (g) $O(4^n \log(n))$

- **Properties of Big-O notation**

1. **Transitivity:** If $f(n)$ is $O(g(n))$ and $g(n)$ is $O(h(n))$, then $f(n)$ is $O(h(n))$.

Proof: According to the definition, $f(n)$ is $O(g(n))$ if there exist positive numbers c_1 and N_1 such that $f(n) \leq c_1 g(n)$ for all $n \geq N_1$, and $g(n)$ is $O(h(n))$ if there exist positive numbers c_2 and N_2 such that $g(n) \leq c_2 h(n)$ for all $n \geq N_2$. Hence, $c_1 g(n) \leq c_1 c_2 h(n)$ for $n \geq N$ where N is the larger of N_1 and N_2 . If we take $c = c_1 c_2$, then $f(n) \leq c h(n)$ for $n \geq N$, which means that f is $O(h(n))$.

2. **Addition:** If $f(n)$ is $O(h(n))$ and $g(n)$ is $O(h(n))$, then $f(n) + g(n)$ is $O(h(n))$.

Proof: If $f(n) \leq c_1 h(n)$, and $g(n) \leq c_2 h(n)$, then $f(n) + g(n) \leq c_1 h(n) + c_2 h(n) \leq (c_1 + c_2) h(n)$. Let $c = c_1 + c_2$, then $f(n) + g(n) \leq c h(n)$ and $f(n) + g(n)$ is $O(h(n))$

3. **Polynomial bounds:** The function an^k is $O(n^k)$

Proof: $an^k \leq cn^k$ for $c \geq a$. Since we can always find some constant $c \geq a$, an^k is $O(n^k)$

Observation: For $an^k \leq cn^k$ to hold, $c \geq a$ is necessary

4. **Domination of higher-degree polynomials:** n^k is $O(n^{k+j}) \forall j > 0$

This statement holds if $c = N = 1$

It follows from all these facts that every polynomial is big-O of n raised to the largest power, or

$$f(n) = a_k n^k + a_{k-1} n^{k-1} + \dots + a_1 n + a_0 \text{ is } O(n^k).$$

5. **Logs:** The function $\log_a n$ is $O(\log_b n)$ for any positive numbers a and $b \neq 1$.

This correspondence holds between logarithmic functions. The fact above states that regardless of their bases, logarithmic functions are big-O of each other; that is, all these functions have the same rate of growth.

Proof: Let $\log_a(n) = x$, and $\log_b(n) = y$, then $a^x = n$, $b^y = n$. Take the natural log of both sides

$$\begin{aligned} \ln(a^x) &= \ln(n) & \ln(b^y) &= \ln(n) \\ \implies x \ln(a) &= \ln(n) & y \ln(b) &= \ln(n) \\ \implies x \ln(a) &= y \ln(b). \end{aligned}$$

Since $x = \log_a(n)$, and $y = \log_b(n)$, then we have

$$\begin{aligned} \ln(a) \log_a(n) &= \ln(b) \log_b(n) \\ \log_a(n) &= \frac{\ln(b)}{\ln(a)} \log_b(n). \end{aligned}$$

Let $c = \frac{\ln(b)}{\ln(a)}$, then $\log_a(n) = c \log_b(n)$, which proves that $\log_a(n)$ and $\log_b(n)$ are multiples of each other. Thus, $\log_a(n)$ is $O(\log_b(n))$

Note: Because the base of the logarithm is irrelevant in the context of big-O notation, we can always use just one base.

$$\therefore \log_a(n) \text{ is } O(\lg n).$$

For any positive $a \neq 1$, where $\lg(n)$ is $\log_2(n)$

- **Big-Ω.** The function $f(n)$ is $\Omega(g(n))$ iff $\exists c, N \in \mathbb{R}^+ \mid f(n) \geq cg(n) \forall n \geq N$.

In other words, $cg(n)$ is a lower bound on the size of $f(n)$, or, in the long run, f grows at least at the rate of g

There is an interconnection between these two notations expressed by the equivalence

$$f(n) \text{ is } \Omega(g(n)) \text{ iff } g(n) \text{ is } O(f(n)).$$

There are an infinite number of possible lower bounds for the function f ; that is, there is an infinite set of g s such that $f(n)$ is $\Omega(g(n))$ as well as an unbounded number of possible upper bounds of f . This may be somewhat disquieting, so we restrict our attention to the smallest upper bounds and the largest lower bounds. Note that there is a common ground for big-O and Ω notations indicated by the equalities in the definitions of these notations: Big-O is defined in terms of " \leq " and Ω in terms of " \geq "; " $=$ " is included in both inequalities. This suggests a way of restricting the sets of possible lower and upper bounds.

- **Big- Θ :** $f(n)$ is $\Theta(g(n))$ iff $\exists c_1, c_2, N \in \mathbb{R}^+ \mid c_1 g(n) \leq f(n) \leq c_2 g(n) \forall n \geq N$

We see that $f(n)$ is $\Theta(g(n))$ if $f(n)$ is $O(g(n))$ and $f(n)$ is $\Omega(g(n))$.

When applying any of these notations, do not forget that they are approximations that hide some detail that in many cases may be considered important.

- **Double O notation:** f is $OO(g(n))$ if it is $O(g(n))$ and the constant c is too large to have practical significance. Thus, $10^8 n$ is $OO(n)$. However, the definition of "too large" depends on the particular application.
- **Using asymptotic complexity to estimate time:** If an algorithm is $O(n^2)$, the time to process n elements is proportional to n^2 .

Let $T(n)$ represent the time, so $T(n) = k \cdot n^2$ where k is a constant.

To find the time for 1 million elements ($n = 10^6$):

$$T(10^6) = k \cdot (10^6)^2 = k \cdot 10^{12}$$

For example, if processing 1000 elements takes 1 second, then:

$$T(1000) = k \cdot 1000^2 = k \cdot 10^6 \implies k = \frac{1}{10^6}$$

Now, for $n = 10^6$:

$$T(10^6) = \frac{1}{10^6} \cdot (10^6)^2 = 10^6 \text{ seconds} = 1,000,000 \text{ seconds} \approx 11.57 \text{ days}.$$

- **Finding asymptotic complexities:** Asymptotic bounds are used to estimate the efficiency of algorithms by assessing the amount of time and memory needed to accomplish the task for which the algorithms were designed. This section illustrates how this complexity can be determined. In most cases, we are interested in time complexity, which usually measures the number of assignments and comparisons performed during the execution of a program. For now let's focus on assignments

Consider a simple loop to calculate the sum of numbers in an array

```
1  for (i = sum = 0; i < n; i++)
2      sum += a[i];
```


First, two variables are initialized, then the for loop iterates n times, and during each iteration, it executes two assignments, one of which updates `sum` and the other of which updates `i`. Thus, there are $2 + 2n$ assignments for the complete run of this for loop; its asymptotic complexity is $O(n)$.

Complexity usually grows if nested loops are used, as in the following code, which outputs the sums of all the subarrays that begin with position 0:

```

1  for (i = 0; i < n; i++) {
2      for (j = 1, sum = a[0]; j <= i; j++)
3          sum += a[j];
4      cout<<"sum for subarray 0 through "<< i <<" is
   ↪      "<<sum<<endl;
5  }
```

Before the loops start, i is initialized. The outer loop is performed n times, executing in each iteration an inner for loop, print statement, and assignment statements for i , j , and `sum`. The inner loop is executed i times for each $i \in \{1, \dots, n-1\}$ with two assignments in each iteration: one for `sum` and one for j . Therefore, there are

$$1 + 3n + \sum_{i=1}^{n-1} 2i = 1 + 3n + 2(1 + 2 + \dots + n - 1) = 1 + 3n + n(n - 1)$$

$= O(n) + O(n^2) = O(n^2)$ assignments executed before the program is completed.

- **Amortized complexity:** amortized analysis can be used to find the average complexity of a worst case sequence of operations

Linked lists

3.1 Singly-linked lists

If a node contains a data member that is a pointer to another node, then many nodes can be strung together using only one variable to access the entire sequence of nodes. Such a sequence of nodes is the most frequently used implementation of a linked list, which is a data structure composed of nodes, each node holding some information and a pointer to another node in the list. If a node has a link only to its successor in this sequence, the list is called a singly linked list

Each node resides on the heap

Linked lists can easily grow and shrink in size without reallocating memory or moving elements. Adding or removing nodes (especially at the beginning or middle) is more efficient compared to arrays, as no shifting of elements is required. Memory is allocated as needed, avoiding wasted space typical in arrays with fixed sizes.

However, each node requires extra memory for the pointer to the next node. Accessing elements requires traversal from the head, making lookups slower ($O(n)$) compared to arrays, which offer $O(1)$ access via indexing. Nodes are scattered in memory, leading to poor cache performance compared to arrays, which have contiguous memory locations.

3.1.1 Structure of the node

The node structure is typically implemented in the following way

```
1  struct node {
2      node* next = nullptr;
3      T data = 0;
4
5      node() = default;
6      node(data) : data(data) {}
7      node(next, data) : next(next), data(data) {}
8  }
```

A node includes two data members: info and next. The info member is used to store information, and this member is important to the user. The next member is used to link nodes to form a linked list. It is an auxiliary data member used to maintain the linked list. It is indispensable for implementation of the linked list, but less important (if at all) from the user's perspective. Note that node is defined in terms of itself because one data member, next, is a pointer to a node of the same type that is just being defined. Objects that include such a data member are called self-referential objects.

3.1.2 The list class/struct

We also implement the list structure as a class or struct.

```
1  class single_list {  
2      node* head = nullptr;  
3  public:  
4      ...  
5  };
```

3.1.3 Interface of a singly linked list stack

The interface typically includes the following operations:

1. **Insert:** Add a node at the beginning, end, or a specific position in the list.
2. **Delete:** Remove a node from the beginning, end, or a specific position.
3. **Search:** Find a node with a given value.
4. **Traverse:** Iterate through the list to access or print each node's data.
5. **IsEmpty:** Check if the list is empty.
6. **Size:** Return the number of nodes in the list. The first node is called the head, and the last node points to nullptr (indicating the end of the list).

3.1.4 Traversing

Traversing a list is simple.

```
1  node* curr = head;  
2  
3  while (curr) {  
4      curr = curr->next;  
5      ...  
6  }
```

3.1.5 Printing

Now that we can traverse, we can print each node

```
1  node* curr = head;
2  while (curr) {
3      cout << curr->data;
4      curr=curr->next;
5  }
```

3.1.6 Printing in reverse

Printing in reverse requires creating a stack.

```
1  if (!head) return; // noop, dont even bother creating a vector.
2
3  vector<node*> stack;
4  node* curr = head;
5
6  while (curr) {
7      stack.push_back(curr);
8      curr=curr->next;
9  }
10
11  for (int i=(int)stack.size()-1; i>=0; --i) {
12      cout << stack[i]->data << " ";
13  }
14  cout << endl;
```

3.1.7 Getting the length

While we traverse, just increment a counter.

```
1  size_t len() {  
2      size_t len = 0;  
3      for (node* curr = head; curr; curr=curr->next, ++len);  
4      return len;  
5  }
```

3.1.8 Clearing

```
1  void clear() {  
2      node* curr=head, *prev=NULLPTR;  
3  
4      while (curr) {  
5          prev=curr;  
6          curr=curr->next;  
7          delete prev;  
8      }  
9      head = NULLPTR;  
10 }
```


3.1.9 Reversing

Reversing is pretty straight forward

```
1  void reverse() {  
2      node* prev=nullptr, *curr=head, *next=nullptr;  
3  
4      while(curr) {  
5          next=curr->next;  
6          curr->next = prev;  
7          prev = curr;  
8          curr=next;  
9      }  
10  
11     head = prev;  
12 }
```

In each iteration, next temporarily holds the next node so you don't lose track of it when reversing the link.

The curr->next pointer is set to prev, effectively reversing the link.

Prev is then updated to curr, and curr is updated to next to continue the process.

3.1.10 Pushing

```
1 void push(int element) {  
2     if (!head) {  
3         head = new node(element);  
4         return;  
5     }  
6  
7     node* curr = head;  
8     while (curr->next) {  
9         curr=curr->next;  
10    }  
11    curr->next = new node(element);  
12 }
```

3.1.11 Inserting

```
1 void insert(int pos, int element) {
2     if (!head || pos == 0) {
3         node* new_node = new node(element);
4         new_node->next = head;
5         head = new_node;
6         return;
7     }
8     node* curr = head;
9
10    int count=0;
11    while (count != pos-1 && curr->next) {
12        curr=curr->next;
13        ++count;
14    }
15    node* new_node = new node(element);
16
17    new_node->next = curr->next;
18    curr->next = new_node;
19 }
```

1. Check if the list is empty or inserting at the head (position 0):

- If head is nullptr (meaning the list is empty) or pos == 0 (you want to insert at the beginning), a new node is created with the given element.
- The new node's next pointer is set to the current head (which could be nullptr if the list is empty), and then head is updated to point to this new node.
- This handles the case where the new node becomes the first node in the list.

2. Traverse to the correct position:

- If you are inserting somewhere other than the head, the function uses a loop to find the node just before the desired position (pos - 1).
- It starts at the head and moves along the list until it reaches the node right before where the new node will be inserted.

3. Insert the new node:

- Once the loop finds the right place (curr points to the node before the insertion position), a new node is created.
- The new node's next pointer is set to curr->next (the node currently in the target position).
- Then, curr->next is updated to point to the new node, effectively inserting the new node into the list.

3.1.12 Popping

```
1  void pop() {
2      if (!head) return;
3      if (!head->next) {
4          delete head;
5          head=nullptr;
6          return;
7      }
8
9      node* prev=nullptr, *curr = head;
10     while (curr->next) {
11         prev=curr;
12         curr=curr->next;
13     }
14     delete curr;
15     prev->next=nullptr;
16 }
```

1. **Empty List Check:** If the list is empty (`head == nullptr`), it does nothing.
2. **Single Node Case:** If the list has only one node, it deletes the head and sets head to `nullptr`.
3. **Multiple Nodes:** It traverses to the last node using two pointers (`prev` and `curr`), deletes the last node (`curr`), and sets the second-to-last node's next pointer (`prev->next`) to `nullptr` to mark the new end of the list.

3.1.13 Erasing

```
1 void erase(int element) {
2     if (!head) return;
3
4     while (head->data == element) {
5         if (head->next && head->data == element) {
6             node* tmp = head;
7             head = head->next;
8             delete tmp;
9         }
10    }
11
12    node* prev=nullptr, *curr=head;
13
14    while (curr) {
15        if (curr->data == element) {
16            node* tmp = curr;
17            prev->next = curr->next;
18            curr=curr->next;
19            delete tmp;
20        } else {
21            prev=curr;
22            curr=curr->next;
23        }
24    }
25 }
```

This erase function removes all nodes with a specific value (element) from the list:

- **Empty List Check:** If the list is empty (`head == nullptr`), it returns immediately.
- **Head Node Deletion:** If the head contains the target value, it deletes the head and updates it to the next node. We keep doing this until the head node no longer contains the data we want to remove
- **Traverse and Delete:** It iterates through the list, and for each node with the target value, it removes the node by adjusting the next pointer of the previous node and deleting the current node.

3.1.14 Searching

```
1  node* search(int element) {  
2      node* curr = head;  
3      while (curr) {  
4          if (curr->data == element) {  
5              return curr;  
6          }  
7      }  
8      return nullptr;  
9  }
```

3.2 Doubly-linked list

A doubly linked list is a data structure consisting of nodes where each node contains three components:

1. **Data:** The value or information the node holds.
2. **Next pointer:** A reference to the next node in the list.
3. **Previous pointer:** A reference to the previous node in the list.

Each node links to both its predecessor and successor. Traversal is possible in both forward and backward directions. It requires more memory than a singly linked list due to the extra pointer.

3.2.1 Node structure

Instead of just a next pointer, a doubly linked list also has a prev pointer

```
1  struct node {  
2      int data = 0;  
3      node* next = nullptr, *prev = nullptr;  
4  
5      node(int data) : data(data) {}  
6  };
```

3.2.2 The list class

The list class implements the interface

```
1  class list {  
2      node* head=nullptr, *tail=nullptr;  
3      size_t n = 0;  
4  
5  public:  
6      void push_back(int element);  
7      void push_front(int element);  
8      void pop_back();  
9      void pop_front();  
10     bool find();  
11     size_t size();  
12     bool empty();  
13     void print();  
14     void clear();  
15 };
```

Size, empty, print and clear are trivial, so their implementation will not be discussed

3.2.3 push_back and push_front

```
1 void push_back(int element) {
2     if (!head) {
3         head = new node(element);
4         tail = head;
5         ++n;
6         return;
7     }
8     tail->next = new node(element);
9     tail->next->prev = tail;
10    tail = tail->next;
11    ++n;
12 }
13
14 void push_front(int element) {
15     if (!head) {
16         head = new node(element);
17         tail = head;
18         ++n;
19         return;
20     }
21     head->prev = new node(element);
22     head->prev->next = head;
23     head = head->prev;
24     ++n;
25 }
```

The `push_back` function adds a new element to the end of the doubly linked list. If the list is empty (indicated by a null head), it creates a new node as both the head and tail of the list and increments the size counter `n`. Otherwise, it creates a new node, links it as the next node of the current tail, sets the previous pointer of the new node to the old tail, updates the tail to the new node, and increments `n`.

The `push_front` function inserts a new element at the beginning of the doubly linked list. If the list is empty, it operates similarly to `push_back` by creating a new node as both the head and tail. For non-empty lists, it creates a new node, links it as the previous node of the current head, sets the next pointer of the new node to the old head, updates the head to the new node, and increments `n`. Both functions maintain the integrity of the doubly linked list by appropriately updating the pointers.

3.2.4 pop_back and pop_frot

```
1 void pop_back() {
2     node* save = tail->prev;
3     save->next = nullptr;
4     delete tail;
5     tail = save;
6     --n;
7 }
8
9 void pop_front() {
10    node* save = head->next;
11    save->prev = nullptr;
12    delete head;
13    head = save;
14    --n;
15 }
```

The `pop_back` function removes the last element from the doubly linked list. It starts by saving the previous node of the current tail into a temporary pointer (`save`). Then, it disconnects the tail by setting the next pointer of the saved node to `nullptr`. Afterward, it deletes the current tail node, updates the tail pointer to the saved node, and decrements the size counter `n`.

The `pop_front` function removes the first element from the doubly linked list. It saves the next node of the current head into a temporary pointer (`save`). It then disconnects the head by setting the prev pointer of the saved node to `nullptr`. The current head is deleted, the head pointer is updated to the saved node, and `n` is decremented. Both functions ensure proper memory management and maintain the structural integrity of the list.

3.2.5 Searching

- **Start from the head:** Traverse the list from the beginning since it is more straightforward and avoids redundant checks.
- **Compare each node's data with the target value:** If a match is found, return a pointer to the node.
- **Continue until the end:** If no match is found by the time you reach the end of the list (when current becomes nullptr), return nullptr to indicate the value is not in the list.

```
1  bool find(int element) {
2      node* curr = head;
3      while (curr) {
4          if (curr->data == element) return true;
5          curr=curr->next;
6      }
7      return false;
8  }
```

Recursion

4.1 Recursion vs iteration

In theory, any problem that can be solved recursively can be solved iteratively. This also means that any problem that can be solved iteratively can also be solved recursively.

The question is, for any problem that can be solved, which method can be used such that the problem is easier to solve.

4.2 Elementary recursion

A recursive definition consists of two parts. In the first part, called the anchor or the ground case, the basic elements that are the building blocks of all other elements of the set are listed. In the second part, rules are given that allow for the construction of new objects out of basic elements or objects that have already been constructed. These rules are applied again and again to generate new objects. For example, to construct the set of natural numbers, one basic element, 0, is singled out, and the operation of incrementing by 1 is given as:

1. $0 \in \mathbb{N}$
2. If $n \in \mathbb{N}$, then $(n + 1) \in \mathbb{N}$
3. There are no other objects in the set \mathbb{N}

It is more convenient to use the following definition, which encompasses the whole range of Arabic numeric heritage:

1. $0, 1, 2, 3, 4, 5, 6, 7, 8, 9 \in \mathbb{N}$
2. If $n \in \mathbb{N}$, then $n0, n1, n2, n3, n4, n5, n6, n7, n8, n9 \in \mathbb{N}$
3. These are the only natural numbers

Recursive definitions serve two purposes: generating new elements, as already indicated, and testing whether an element belongs to a set. In the case of testing, the problem is solved by reducing it to a simpler problem, and if the simpler problem is still too complex it is reduced to an even simpler problem, and so on, until it is reduced to a problem indicated in the anchor

4.3 Base cases

In recursion, a base case is a condition that stops further recursive calls and provides a direct answer without further recursion

If there were no base case, there would be nothing to stop the recursion. Thus, it would go on until the program crashes. For this reason, all recursive functions must have at least one base case.

If a base case in a recursive function returns a value, then every recursive call leading up to that base case should also return a value. This is necessary to ensure that the result of the recursion is propagated back up the call stack.

In a recursive function, the base case stops the recursion, and if the base case returns something (e.g., a node pointer, integer, etc.), the recursive calls that occur before reaching the base case need to return that result so it can propagate back to the original caller.

4.3.1 Factorials

$$n! = \begin{cases} 1 & \text{if } n = 0 \\ n(n-1)! & \text{if } n \neq 0 \end{cases}.$$

```

1  int factorial(int n) {
2      if (n == 0) return 1;
3      return n * factorial(n-1);
4
5      // Expands to
6      // n * n-1 * n-2 * ... * 1
7  }

```

4.3.2 Powers

Consider the recursive definition for a power of x

$$x^n = \begin{cases} 1 & \text{if } n = 0 \\ x \cdot x^{n-1} & \text{if } n > 0 \end{cases}.$$

```

1  constexpr int power(int x, int n) {
2      if (n == 0) return 1;
3      return x * power(x,n-1);
4  }

```

The function `power()` can be implemented differently, without using any recursion, as in the following loop:

```

1  int power2(int x, int n) {
2      int res = 1;
3
4      for (res = x; n > 1; --n) {
5          res*=x;
6      }
7      return res;
8  }

```

Do we gain anything by using recursion instead of a loop? The recursive version seems to be more intuitive because it is similar to the original definition of the power function. The definition is simply expressed in C++ without losing the original structure of the definition. The recursive version increases program readability, improves self-documentation, and simplifies coding. In our example, the code of the nonrecursive version is not substantially larger than in the recursive version, but for most recursive implementations, the code is shorter than it is in the nonrecursive implementations

4.4 Tail recursion

Tail recursion is a type of recursion where the recursive call is the last thing the function does before returning a result. This means there are no more computations or operations to perform after the recursive call.

Because of this, tail recursion can be optimized by some compilers or interpreters to avoid adding new frames to the call stack, making it more memory-efficient than regular recursion.

In simple terms, if a recursive function calls itself, and after that call there's nothing left to do, it's tail recursion. This allows the function to reuse the same memory space, preventing stack overflow in cases with deep recursion.

the recursive call is not only the last statement but there are no earlier recursive calls, direct or indirect. For example, the function `tail()` defined as

```
1 void tail(int i) {  
2     if (i > 0) {  
3         cout << i << ' ';  
4         tail(i-1);  
5     }  
6 }
```

Is an example of a function with tail recursion, whereas the function `nonTail()` defined as

```
1 void nonTail(int i) {  
2     if (i > 0) {  
3         nonTail(i-1);  
4         cout << i << ' ';  
5         nonTail(i-1);  
6     }  
7 }
```

Is not. Tail recursion is simply a glorified loop and can be easily replaced by one. In this example, it is replaced by substituting a loop for the if statement and decrementing the variable `i` in accordance with the level of recursive call. In this way, `tail()` can be expressed by an iterative function:

```
1 void iterativeEquivalentOfTail(int i) {  
2     for ( ; i > 0; i--)  
3         cout << i << ' ';  
4 }
```

Is there any advantage in using tail recursion over iteration? For languages such as C++, there may be no compelling advantage, but in a language such as Prolog, which has no explicit loop construct (loops are simulated by recursion), tail recursion acquires a much greater weight. In languages endowed with a loop or its equivalents, such as an if statement combined with a goto statement, tail recursion should not be used.

Another problem that can be implemented in recursion is printing an input line in reverse order. Here is a simple recursive implementation:

```
1 void reverse() {
2     char ch;
3     cin.get(ch);
4     if (ch != '\n') {
5         reverse();
6         cout.put(ch);
7     }
8 }
```

Compare the recursive implementation with a nonrecursive version of the same function:

```
1 void simpleIterativeReverse() {
2     char stack[80];
3     int top = 0;
4     cin.getline(stack,80);
5     for (top = strlen(stack) - 1; top >= 0;
6     ↪ cout.put(stack[top--]));
7 }
```

functions like `strlen()` and `getline()` from the standard C++ library can be used. If we are not supplied with such functions, then our iterative function has to be implemented differently:

```
1 void iterativeReverse() {
2     char stack[80];
3
4     register int top = 0;
5     cin.get(stack[top]);
6
7     while(stack[top]!='\n') {
8         cin.get(stack[++top]);
9     }
10    for (top -= 2; top >= 0; cout.put(stack[top--]));
11 }
```

4.5 Indirect Recursion

The preceding sections discussed only direct recursion, where a function $f()$ called itself. However, $f()$ can call itself indirectly via a chain of other calls. For example, $f()$ can call $g()$, and $g()$ can call $f()$. This is the simplest case of indirect recursion. The chain of intermediate calls can be of an arbitrary length, as in:

$$f() \rightarrow f_1() \rightarrow f_2() \rightarrow \dots \rightarrow f_n() \rightarrow f().$$

There is also the situation when $f()$ can call itself indirectly through different chains. Thus, in addition to the chain just given, another chain might also be possible. For instance

$$f() \rightarrow g_1() \rightarrow g_2() \rightarrow \dots \rightarrow g_m() \rightarrow f().$$

This situation can be exemplified by three functions used for decoding information. `receive()` stores the incoming information in a buffer, `decode()` converts it into legible form, and `store()` stores it in a file. `receive()` fills the buffer and calls `decode()`, which in turn, after finishing its job, submits the buffer with decoded information to `store()`. After `store()` accomplishes its tasks, it calls `receive()` to intercept more encoded information using the same buffer. Therefore, we have the chain of calls

$$\text{recieve}() \rightarrow \text{decode}() \rightarrow \text{store}() \rightarrow \text{recieve}() \rightarrow \text{decode}() \rightarrow \dots$$

4.6 Nested Recursion

A more complicated case of recursion is found in definitions in which a function is not only defined in terms of itself, but also is used as one of the parameters. The following definition is an example of such a nesting

$$h(n) = \begin{cases} 0 & \text{if } n = 0 \\ n & \text{if } n > 4 \\ h(2 + h(n)) & \text{if } n \leq 4 \end{cases}.$$

4.7 Excessive Recursion

Logical simplicity and readability are used as an argument supporting the use of recursion. The price for using recursion is slowing down execution time and storing on the run-time stack more things than required in a nonrecursive approach. If recursion is too deep (for example, computing $5.6^{100,000}$), then we can run out of space on the stack and our program crashes. But usually, the number of recursive calls is much smaller than 100,000, so the danger of overflowing the stack may not be imminent

However, if some recursive function repeats the computations for some parameters, the run time can be prohibitively long even for very simple cases

Consider Fibonacci numbers. A sequence of Fibonacci numbers is defined as follows:

$$\text{Fib}(n) = \begin{cases} n & \text{if } n < 2 \\ \text{Fib}(n-2) + \text{Fib}(n-1) & \text{otherwise} \end{cases}.$$

The definition states that if the first two numbers are 0 and 1, then any number in the sequence is the sum of its two predecessors. But these predecessors are in turn sums of their predecessors, and so on, to the beginning of the sequence.

How can this definition be implemented in C++? It takes almost term-by-term translation to have a recursive version, which is

```
1  constexpr unsigned long fib(int n) {  
2      if (n < 2) return n;  
3      return fib(n-2) + fib(n-1);  
4  }
```

The function is simple and easy to understand but extremely inefficient. To see it, compute $\text{Fib}(6)$, the seventh number of the sequence, which is 8. Based on the definition, the computation runs as follows:

$$\begin{aligned} \text{Fib}(6) &= \text{Fib}(4) + \text{Fib}(5) \\ &= \text{Fib}(2) + \text{Fib}(3) + \text{Fib}(5) \\ &= \text{Fib}(0) + \text{Fib}(1) + \text{Fib}(3) + \text{Fib}(5) \\ &= 0 + 1 + \text{Fib}(3) + \text{Fib}(5) \\ &= 1 + \text{Fib}(1) + \text{Fib}(2) + \text{Fib}(5) \\ &= 1 + \text{Fib}(1) + \text{Fib}(0) + \text{Fib}(1) + \text{Fib}(5). \end{aligned}$$

Etc... The source of this inefficiency is the repetition of the same calculations because the system forgets what has already been calculated. For example, Fib() is called eight times with parameter $n = 1$ to decide that 1 can be returned. For each number of the sequence, the function computes all its predecessors without taking into account that it suffices to do this only once.

It takes almost a quarter of a million calls to find the twenty-sixth Fibonacci number, and nearly 3 million calls to determine the thirty-first! This is too heavy a price for the simplicity of the recursive algorithm. As the number of calls and the run time grow exponentially with n , the algorithm has to be abandoned except for very small numbers

An iterative algorithm may be produced rather easily as follows:

```

1  unsigned long iterativeFib(unsigned long n) {
2      if (n < 2)
3          return n;
4      else {
5          register long i = 2, tmp, current = 1, last = 0;
6          for ( ; i <= n; ++i) {
7              tmp = current;
8              current += last;
9              last = tmp;
10         }
11         return current;
12     }
13 }
```

However, there is another, numerical method for computing Fib(n), using a formula discovered by Abraham de Moivre:

$$\text{Fib}(n) = \frac{\phi^n - \hat{\phi}^n}{\sqrt{5}}.$$

Where $\phi = \frac{1}{2}(1 + \sqrt{5})$, and $\hat{\phi} = 1 - \phi = \frac{1}{2}(1 - \sqrt{5})$. $\hat{\phi}$ becomes very small when n grows, thus it can be omitted.

$$\text{Fib}(n) = \frac{\phi^n}{\sqrt{5}}.$$

Approximated to the nearest integer

```

1  unsigned long deMoivreFib(unsigned long n) {
2      return ceil(exp(n*log(1.6180339897) - log(2.2360679775)) -
3      ↪ .5);
4  }
```

4.8 Backtracking

In solving some problems, a situation arises where there are different ways leading from a given position, none of them known to lead to a solution. After trying one path unsuccessfully, we return to this crossroads and try to find a solution using another path. However, we must ensure that such a return is possible and that all paths can be tried. This technique is called backtracking, and it allows us to systematically try all available avenues from a certain point after some of them lead to nowhere. Using backtracking, we can always return to a position that offers other possibilities for successfully solving the problem. This technique is used in artificial intelligence, and one of the problems in which backtracking is very useful is the eight queens problem.

The eight queens problem attempts to place eight queens on a chessboard in such a way that no queen is attacking any other. To solve this problem, we try to put the first queen on the board, then the second so that it cannot take the first, then the third so that it is not in conflict with the two already placed, and so on, until all of the queens are placed. What happens if, for instance, the sixth queen cannot be placed in a nonconflicting position? We choose another position for the fifth queen and try again with the sixth. If this does not work, the fifth queen is moved again. If all the possible positions for the fifth queen have been tried, the fourth queen is moved and then the process restarts. This process requires a great deal of effort, most of which is spent backtracking to the first crossroads offering some untried avenues. In terms of code, however, the process is rather simple due to the power of recursion, which is a natural implementation of backtracking

```
1  putQueen(row)
2      for every position col on the same row
3          if position col is available
4              place the next queen in position col;
5              if (row < 8)
6                  putQueen(row+1);
7              else success;
8              remove the queen from position col;
```

This algorithm finds all possible solutions without regard to the fact that some of them are symmetrical.

Backtracking is a systematic approach to exploring all possible solutions to a problem using recursion. It involves:

- **Recursive Exploration:** A recursive function is used to explore all potential solutions.
- **Pruning:** When a partial solution cannot lead to a valid answer, the algorithm "backtracks" by undoing the last decision and tries a different path.
- **State Restoration:** The state of the problem is restored before exploring a new path, ensuring no side effects from previous choices affect future decisions.

Backtracking often involves a loop because it systematically explores multiple options at each recursive step. The loop allows the algorithm to iterate over all possible choices available at a particular point. Without the loop, you would only be able to handle a single choice at each step, which would limit the ability to explore multiple branches of the problem space.

The loop ensures that each choice is explored one by one before moving to the next recursive call. After each recursive call, the algorithm backtracks and tries the next option.

Backtracking is essentially a tree traversal algorithm. The loop represents branching into different nodes of the tree. Without a loop, you'd only traverse a linear path, which isn't sufficient for exploring all solutions.

4.8.1 N-Queens

Assume we have already have an 8×8 array to represent the chess board. At each spot in the board, a zero will indicate no queen, and a one will indicate a queen. We also assume we have a function to check if a position is safe. The only function left to implement is a *place_queen* function.

```

1  bool place_queens(bool arr[][COLS], int row_n) {
2      // Base case: Placed a queen on all rows
3      if (row_n >= ROWS)
4          return true;
5
6      // For each column in a row
7      for (int j=0; j<=7; ++j) {
8          // If the spot is safe
9          if (safe(arr, row_n, j)) {
10             // Place the queen
11             arr[row_n][j] = true;
12             // Check the row beneath
13             if (place_queens(arr, row_n+1)) {
14                 return true;
15             }
16             // If the row beneath cannot place a queen at any
↪      column, backtrack one row, place that queen in a different
↪      column
17             arr[row_n][j] = false;
18         }
19     }
20     // Queen could not be placed at any column
21     return false;
22 }

```

4.8.2 Greedy backtracking: Traveling salesman problem (TSP)

Greedy backtracking is a variation of backtracking that incorporates greedy principles, attempting to make the "best" or "most promising" choice at each step rather than systematically exploring all possible choices. It combines the greedy approach (making locally optimal choices) with backtracking to handle situations where the greedy choice might fail.

- **Greedy Choice:** At each step, a heuristic ¹ is used to choose the most promising option that seems likely to lead to a solution.
- **Backtracking for Correction:** If the greedy choice leads to a dead end, the algorithm backtracks and tries the next best option.
- **Efficiency:** By attempting the most promising options first, it often reduces the search space compared to regular backtracking.

TSP: The goal is to find the shortest path that visits all cities exactly once and returns to the starting city.

- **Greedy Step:** At each step, visit the nearest unvisited city.
- **Backtracking:** If the current path does not lead to a valid solution, backtrack and try the next nearest city.

Assume we have a distance matrix

¹proceeding to a solution by trial and error or by rules that are only loosely defined.

```

1  vector<vector<double>> distances = {
2      {0, 10, 15, 20},
3      {10, 0, 35, 25},
4      {15, 35, 0, 30},
5      {20, 25, 30, 0}
6  };

```

The `vector<vector<double>>` `distances` represents a distance matrix for a graph where the nodes represent cities (or points of interest), and the entries in the matrix represent the distances between pairs of cities.

In the context of the Traveling Salesman Problem (TSP), this matrix provides the distances between any two cities in the graph.

- **`distances[i][j]`:** Represents the distance between city i and city j
- **Diagonal entries (`distances[i][i]`):** These are typically 0, because the distance from a city to itself is 0

Assume we also have a visited array, to mark the cities visited.

```

1  vector<bool> visited(distances.size(), false);

```

The solution function signature will be

```

1  double tspGreedyBacktracking(vector<vector<double>>& distances,
    ↪ vector<bool>& visited, int currentCity, int citiesVisited,
    ↪ double currentCost, double& bestCost);

```

- **`vector<vector<double>>& distances`:** Represents the distance matrix for the graph. Each entry `distances[i][j]` gives the distance from city i to city j .

Passed by reference to avoid copying the matrix for each function call, which would be computationally expensive.

- **`vector<bool>& visited`:** A boolean vector where `visited[i]` indicates whether city i has been visited. Ensures that each city is visited exactly once. Passed by reference to track visited cities across recursive calls.
- **`int currentCity`:** The city the algorithm is currently processing. Determines the row in the distance matrix to explore potential next cities. Helps calculate the distance for traveling to other unvisited cities.
- **`int citiesVisited`:** Tracks the number of cities visited so far in the current path. Used as a stopping condition: when `citiesVisited == distances.size()`, all cities have been visited.
- **`double currentCost`:** The cumulative cost (distance) of the path taken so far. This is updated at each recursive step to include the distance of the last move.
- **`double& bestCost`:** Tracks the minimum cost (shortest path) found so far across all recursive paths. Passed by reference so that all recursive calls update the same variable. Used to prune paths: If the `currentCost` exceeds `bestCost`, the algorithm skips further exploration of that path.

```

1  double tspGreedyBacktracking(vector<vector<double>>& distances,
    ↪ vector<bool>& visited, int currentCity, int citiesVisited,
    ↪ double currentCost, double& bestCost) {
2      if (citiesVisited == visited.size()) {
3          // Add the cost to return to the starting city
4          return currentCost + distances[currentCity][0];
5      }
6
7      for (int nextCity = 0; nextCity < distances.size();
    ↪ ++nextCity) {
8          if (!visited[nextCity]) {
9              visited[nextCity] = true;
10
11              // Calculate the cost for this path
12              double nextCost = currentCost +
    ↪ distances[currentCity][nextCity];
13              if (nextCost < bestCost) { // Only proceed if this
    ↪ path is promising
14                  bestCost = min(bestCost,
    ↪ tspGreedyBacktracking(distances, visited, nextCity,
    ↪ citiesVisited + 1, nextCost, bestCost));
15              }
16
17              visited[nextCity] = false; // Backtrack
18          }
19      }
20      return bestCost;
21  }
22  int main() {
23      vector<vector<double>> distances = {
24          {0, 10, 15, 20},
25          {10, 0, 35, 25},
26          {15, 35, 0, 30},
27          {20, 25, 30, 0}
28      };
29      vector<bool> visited(distances.size(), false);
30      visited[0] = true;
31
32      double bestCost = DBL_MAX;
33      tspGreedyBacktracking(distances, visited, 0, 1, 0, bestCost);
34
35      cout << "Minimum cost: " << bestCost << endl;
36      return 0;
37  }

```

This function implements a backtracking approach to solve the Traveling Salesman Problem (TSP) with an optimization that leverages a greedy principle to prune unnecessary paths. The algorithm starts from a given city (`currentCity`) and recursively explores all possible paths to visit every other city exactly once. It maintains a `visited` array to track which cities have already been visited, ensuring that no city is revisited during the current path.

The function uses a loop to evaluate each unvisited city as the next possible destination. For each potential move, it calculates the cumulative cost of the current path (`currentCost`)

and compares it with the best cost (`bestCost`) found so far. If the current path exceeds `bestCost`, it abandons further exploration of that branch. Otherwise, it marks the city as visited, proceeds recursively to explore the next cities, and updates the `bestCost` if a better solution is found. After processing, the function backtracks by unmarking the city, restoring the state for the next iteration.

When all cities have been visited (`citiesVisited == visited.size()`), the function completes the path by adding the cost to return to the starting city. This final cost is returned and compared with `bestCost`, ensuring that the shortest valid path is recorded. The combination of backtracking for completeness and greedy pruning for efficiency helps the algorithm systematically explore the search space while avoiding unnecessary computations.

It will begin by starting from the first city, then going $1 \rightarrow 2$, $2 \rightarrow 3$, $3 \rightarrow 4$. At this point, it will have a best cost for path one (65 in this case). It will then begin backtracking back to city one, enabling us to explore other paths. For example, $1 \rightarrow 3$, $3 \rightarrow 2, \dots$, only taking the path if it is promising (currentcost less than best cost). Once we have explored all options starting from city one, we return the best cost.

4.8.3 Combinations

To perform combinations using backtracking in C++, you can recursively explore subsets of elements. Here's an example of how to generate all combinations of size k from a given array or vector

```
1  void generateCombinations(const std::vector<int>& nums, int k,  
    ↪ int start, std::vector<int>& current,  
    ↪ std::vector<std::vector<int>>& result) {  
2      if (current.size() == k) {  
3          result.push_back(current);  
4          return;  
5      }  
6  
7      for (int i=start; i<nums.size(); ++i) {  
8          current.push_back(nums[i]);  
9          generateCombinations(nums, k, i+1, current, result);  
10         current.pop_back();  
11     }  
12 }
```

The backtracking logic in this function generates combinations by systematically exploring subsets of the input array. The key idea is to build each combination step-by-step, adding elements to the current subset (`current`) until its size reaches the desired length k . At that point, the subset is complete and added to the result.

The function iterates through the array starting from the given index `start`. For each element, it adds the element to the current combination and recursively calls itself with the next index (`i+1`) to explore subsequent elements. This ensures that no element is reused and combinations are formed in a non-repeating manner.

Once the recursive call completes, the function removes the last element added (backtracking) to explore other possible combinations. This backtracking step is crucial, as it resets the state of `current` to allow the exploration of different subsets without interference from previous paths. This approach ensures all unique combinations of size k are generated.

4.8.4 Combinations that sum to m

Expanding on the code above, we can add a simple check before we push the combination to the result vector, the check to see if the sum of the combination fits a requested sum.

```
1  void generateCombinations(const std::vector<int>& nums, int k,  
    ↪ int start, std::vector<int>& current,  
    ↪ std::vector<std::vector<int>>& result, const int& sum, int&  
    ↪ currentSum) {  
2      if (current.size() == k) {  
3          if (currentSum == sum) {  
4              result.push_back(current);  
5          }  
6          return;  
7      }  
8  
9      for (int i=start; i<nums.size(); ++i) {  
10         current.push_back(nums[i]);  
11         // Update the current sum  
12         currentSum += nums[i];  
13  
14         generateCombinations(nums, k, i+1, current, result, sum,  
    ↪ currentSum);  
15  
16         // Decrease the current sum  
17         currentSum-=nums[i];  
18         current.pop_back();  
19     }  
20 }
```

4.8.5 Permutations

To generate permutations, the key idea is to explore all possible arrangements of elements where the order matters. The typical approach involves swapping elements to produce different orders and backtracking to restore the original state before exploring further.

```
1  void generatePermutations(std::vector<int>& nums, int start,
   ↪  std::vector<std::vector<int>>& result) {
2      // Base case: if the starting index is at the end, we've
   ↪  formed a permutation
3      if (start == nums.size()) {
4          result.push_back(nums);
5          return;
6      }
7
8      // Loop through the array to swap the current index with all
   ↪  subsequent indices
9      for (int i = start; i < nums.size(); ++i) {
10         // Swap the current element with the one at index `i`
11         std::swap(nums[start], nums[i]);
12
13         // Recur to generate permutations for the next index
14         generatePermutations(nums, start + 1, result);
15
16         // Backtrack: restore the original order
17         std::swap(nums[start], nums[i]);
18     }
19 }
```

- The function starts with a start index, which determines the position being fixed for the current permutation.
- For every index from start to the end of the array, it swaps the element at start with the element at the current index (i) to create a new arrangement.
- It recursively generates permutations for the remaining portion of the array (start + 1 onward).
- After the recursive call, it swaps back to restore the original array (backtracking), ensuring other permutations can be explored without interference.

4.9 Recursion in singly linked lists

4.9.1 Traversing

To traverse a linked list using recursion, you need to define a recursive function that processes the current node and then calls itself with the next node until the list is fully traversed (i.e., until the current node is nullptr).

```
1 void TraverseList(node* head) {  
2     if (!head) {  
3         return;  
4     }  
5     TraverseList(head->next);  
6  
7     // ...  
8 }
```

4.9.2 Printing

We can use this, for example, to print each nodes data member

```
1 void PrintList(node* head) {  
2     if (!head) return;  
3  
4     cout << head->data << " ";  
5     PrintList(head->next);  
6 }
```

4.9.3 Printing in reverse

We a slight alter in the print example, we can reverse print the list.

```
1 void PrintListReverse(node* head) {  
2     if (!head) return;  
3  
4     PrintListReverse(head->next);  
5     cout << head->data << " ";  
6 }
```

4.9.4 Getting the length

4.9.5 Clearing

We can also use this to clear the list

```
1  void clear() {  
2      std::function<void(node*)> r_clear = [&] (node* p) = {  
3          if (!head) return;  
4  
5          r_clear(head->next);  
6          delete head;  
7      }  
8      r_clear(head);  
9      head=nullptr;  
10     size=0;  
11 }
```

4.9.6 Reversing

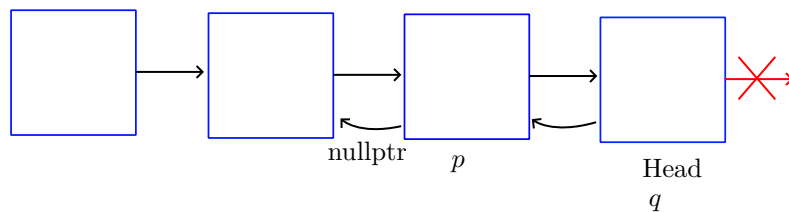
Let's first take a look at the reverse code

```
1 void reverse() {
2     std::function<void(node*)> r_reverse = [&] (node* p) -> void
    ↪ {
3         if (!p->next) {
4             head = p;
5             return;
6         }
7
8         r_reverse(p->next);
9         node* q = p->next;
10        q->next = p;
11        p->next = nullptr;
12
13    };
14    r_reverse(head);
15 }
```

The base case is that we are at the end, in this case we set head to this position. Head is now at the end of the list.

Once the base case is triggered and the head is set to the last node in the list, we will be sent back to the n-1 node call.

To get the intuition for linked list logic, we must examine a diagram of the list.



This figure shows the three operations done after each recursive call. In the figure above, we are at the node after the call that set the end to head. We

1. Get a pointer to the node ahead of the current (*q*).
2. This allows us to sever its old next pointer and reverse its direction.
3. Then, set *p* next to nullptr (set up for next return).

When the callstack returns to the first call, and does its operations, the list will be reversed

It is also a good idea to examine the iterative method.

```
1  void Itreverse() {
2      node* prev=nullptr, *curr=head, *next=nullptr;
3
4      while (curr) {
5          next=curr->next; // Move next to the next node
6          curr->next=prev; // Change the direction of current
7          ↪ nodes next pointer
8          prev=curr; // Advance prev
9          curr=next; // Advance curr
10     }
11     head=prev; // Prev is last node, set head to end
12 }
```

4.9.7 Pushing

```
1 void push(int data) {
2     if (!head) {
3         head = new node(nullptr, data);
4         return;
5     }
6     std::function<void(node*, int)> r_push = [&] (node* curr,
↪ int data) -> void {
7
8         if (!curr->next) {
9             curr->next = new node(nullptr, data);
10            ++size;
11            return;
12        }
13        r_push(curr->next, data);
14    };
15    r_push(head, data);
16 }
```

Base case:

1. **Empty list:** No recursion, make head the new node

Otherwise, recurse from the head until we get to the last node, simply set last nodes next pointer to new node and return

4.9.8 Inserting

```
1 void insert(unsigned pos, int element) {
2     std::function<void(node*&, unsigned)> r_insert = [&] (node*&
    ↪ p, unsigned curr_pos) {
3         if (curr_pos == 0) {
4             node* new_node = new node(nullptr, element);
5             new_node->next = p;
6             p = new_node;
7             return;
8         }
9         r_insert(p->next, curr_pos-1);
10    };
11    r_insert(head, pos);
12 }
```

Base case

1. **Recurse the same number of times as the pos arg:** In this case, make a new node, set next to current node in the recursive traversal, set current node to new node.

Otherwise, keep recursing, subtracting one from the curr_pos.

4.9.9 Popping

```
1  void pop() {
2      if (!head) return;
3
4      if (!head->next) {
5          delete head;
6          head=nullptr;
7          return;
8      }
9
10     std::function<void(node*)> r_pop = [&] (node* p) -> void {
11         if (!p->next->next) {
12             delete p->next;
13             p->next = nullptr;
14             --size;
15             return;
16         }
17         r_pop(p->next);
18     };
19     r_pop(head);
20 }
```

Base cases:

1. **Empty list:** Noop
2. **One node (head):** Delete then reset head

Otherwise, recurse until we are at the second to last node. Then, delete the second to last nodes next node, which is the last node. Set second to last nodes next pointer to nullptr.

4.9.10 Erasing

```
1 void erase(int element) {
2     std::function<void(node*&)> r_erase = [&] (node*& p) -> void
   ↪ {
3         if (p == nullptr) {
4             return;
5         }
6
7         r_erase(p->next);
8
9         if (p->data == element) {
10             node* tmp = p;
11             p = p->next;
12             delete tmp;
13         }
14     };
15     r_erase(head);
16 }
```

Base case:

1. **Reached the end:** Return, start unwinding

We traverse to the end of the list recursively, once we reach the end the recursion stops and we start unwinding the call stack, going backwards in the list.

For each node, we check if its data is equal to the element, if it is we set this node equal to its next node, then delete.

4.9.11 Searching

```
1  node* search(int element) {
2      std::function<node*(node*)> r_search = [&] (node* p) ->
    ↪ node* {
3          if (p == nullptr) {
4              return nullptr;
5          }
6          if (p->data == element) {
7              return p;
8          }
9          return r_search(p->next);
10     };
11     return r_search(head);
12 }
```

Base cases:

1. **Reached the end of the list:** Element is not in list, return nullptr
2. **Found the first node with the element:** Return the node

Otherwise, recurse through the nodes until we hit one of the base cases.

Dynamic programming

Dynamic Programming is a problem-solving paradigm used to solve optimization and combinatorial problems efficiently by breaking them into overlapping subproblems, solving each subproblem once, and storing its result for reuse. This avoids redundant calculations, making DP particularly effective for problems with overlapping subproblems and optimal substructure

5.1 Key Concepts of DP

- **Overlapping Subproblems:** Problems that can be broken down into smaller, overlapping subproblems. For example, the Fibonacci sequence involves repeatedly calculating the same terms, like $F(2)$, $F(3)$, etc., in a naive recursive approach.
- **Optimal Substructure:** A problem exhibits optimal substructure if the optimal solution of the problem can be composed from the optimal solutions of its subproblems. For example, in the shortest path problem, the shortest path from A to C passing through B is the sum of the shortest path from A to B and B to C .

5.2 Bottom-up DP

Bottom-up dynamic programming is a technique used to solve problems efficiently by breaking them into smaller subproblems, solving those subproblems iteratively (starting with the simplest cases), and building up solutions to the larger problems. It avoids the overhead of recursion by explicitly organizing computation in a tabular form, thereby reducing redundant calculations.

5.2.1 Key Characteristics

- **Iterative Approach:** Unlike top-down dynamic programming (which uses recursion with memoization), bottom-up DP explicitly fills a table iteratively.
- **Tabulation:** A table (usually an array or a matrix) is used to store solutions to subproblems, so they can be reused when solving larger problems.
- **Base Cases First:** The base cases are explicitly calculated and stored in the table before computing other values.
- **No Recursion:** Since it avoids recursion, it typically has better space efficiency because it does not use a call stack.

5.2.2 Steps for Bottom-Up DP

- **Define the State:** Determine what your subproblem represents. For example, $dp[i]$ might represent the solution to the problem for the first i elements.
- **Base Case Initialization:** Set the initial values in your DP table based on the simplest cases of the problem.
- **Transition Relation:** Determine how to compute the value of a subproblem (e.g., $dp[i]$) based on previously computed subproblems (e.g., $dp[j]$ where $j < i$).
- **Fill the Table:** Use loops to iteratively compute all the entries in the DP table up to the desired solution.
- **Extract the Solution:** The final entry in the table often represents the solution to the overall problem.

5.3 Bottom-up DP: Fibonacci

The Fibonacci sequence is a classic example where DP can be applied. The sequence is defined as

$$F(n) = F(n-1) + F(n-2), \quad F(0) = 0, F(1) = 1.$$

```

1  int fib(int n) {
2      if (n == 0) return 0;
3      if (n == 1) return 1;
4
5      int dp[n+1];
6      dp[0] = 0;
7      dp[1] = 1;
8
9      for (int i=2; i<=n; ++i) {
10         dp[i] = dp[i-1] + dp[i-2];
11     }
12     return dp[n];
13 }
```

5.4 Memoization (top down DP)

Top-Down Dynamic Programming, also known as memoization, is a technique where a problem is solved recursively, but results of subproblems are stored (memoized) to avoid redundant calculations. This is particularly effective for problems with overlapping subproblems.

- **Recursive Approach:** The problem is solved recursively by breaking it down into smaller subproblems.

The recursive function calls itself for these subproblems.

- **Memoization:** Subproblem solutions are stored in a data structure (e.g., an array, map, or dictionary).

If a subproblem has already been solved, its stored solution is reused, skipping further computation.

- **Base Cases:** Base cases are defined to stop recursion and provide known values.
- **Optimal Substructure:** The solution to a problem is built from solutions to its subproblems.

5.5 Top-down DP: Fibonacci

The Fibonacci sequence is defined as:

$$F(n) = \begin{cases} 0 & \text{if } n = 0 \\ 1 & \text{if } n = 1 . \\ F(n-1) + F(n-2) & \text{if } n > 1 \end{cases}$$

A simple recursive solution calculates $F(n)$ but does so inefficiently due to repeated calculations:

With memoization, we store already-computed Fibonacci values to avoid redundant work

```
1  int fibonacci(int n, int memo[]) {
2      // Base cases
3      if (n <= 1) return n;
4
5      // Check if the result is already computed
6      if (memo[n] != -1) return memo[n];
7
8      // Compute and store the result
9      memo[n] = fibonacci(n - 1, memo) + fibonacci(n - 2, memo);
10     return memo[n];
11 }
```

Binary trees

6.1 Terminology

- **Node:** The basic unit of a binary tree, containing data and references to left and right children.
- **Root:** The topmost node in a tree.
- **Child:** A node directly connected to another node when moving away from the root.
- **Descendants:** The descendants of a node are all nodes that come after a given node.
- **Parent:** The node directly above a child node.
- **Grandparents:** The grandparents of a node is all nodes above the parent up to the root.
- **Ancestors:** The ancestors of a node are all the nodes above a node up to the root
- **Leaf:** A node with no children.
- **Branch node:** A non-leaf node is called a branch node
- **Internal Node:** A branch node, a node with at least one child.
- **Subtree:** A tree consisting of a node and its descendants.
- **Height of a node:** The number of edges on the longest path from a node to a leaf.
- **Height of a tree:** The height of the tree is the height of the root
- **Depth:** The number of edges from the root to a node.
- **Depth of a tree:** The depth of a tree is the depth of the deepest node
- **Degree of a node:** The number of subtrees of a node is called the degree of the node. In a binary tree, all nodes have degree 0, 1, or 2.
- **Degree of a binary tree:** The degree of a tree is the maximum degree of a node in the tree. A binary tree is degree 2.

6.2 Type of binary trees

- **Full Binary Tree:** Every internal node has two children, all leaf nodes have zero children. Thus, all nodes are either zero or two, never one.
- **Complete Binary Tree:** All levels, except possibly the last, are fully filled, and all nodes are as far left as possible.
- **Perfect Binary Tree:** A binary tree where all internal nodes have exactly 2 children, and all leaf nodes are at the same level.
- **Balanced Binary Tree:** A binary tree where the height of the left and right subtrees of every node differs by at most one.
- **Degenerate (or pathological) Tree:** A tree where each parent node has only one child, essentially forming a linked list.
- **Skewed Tree:** A special case of a degenerate tree, where all nodes are skewed to the left or right, forming a linear structure.

6.3 Maximum height of a binary tree

The maximum height of a binary tree with n nodes can be as large as $n-1$ (in the case of a degenerate or skewed tree where each node has only one child). This is true for any binary tree:

$$h_{\max} = n - 1.$$

Which occurs for degenerate trees.

6.3.1 Minimum height of a binary tree

The minimum height (best case) for a binary tree with n nodes is achieved when the tree is perfectly balanced:

$$h_{\min} = \lfloor \log_2(n) \rfloor.$$

This is because the tree would need to spread nodes evenly across levels

6.3.2 Number of Leaves in a Binary Tree

For any binary tree with n nodes, the number of leaves l satisfies the following relationship:

$$l \leq \frac{n+1}{2}.$$

This formula gives the maximum number of leaves, assuming that the tree is full (every internal node has 2 children).

6.3.3 Relationship Between Internal Nodes and Leaves:

In any binary tree, the number of internal nodes i (nodes with at least one child) and the number of leaves l are related as follows:

$$i \leq l - 1.$$

6.3.4 Maximum Number of Nodes at Height h

The maximum number of nodes possible at a given height h (where the height is counted from the root as level 0) in a binary tree is:

$$\text{Max nodes at height } h = 2^h.$$

6.3.5 Number of Edges in a Binary Tree:

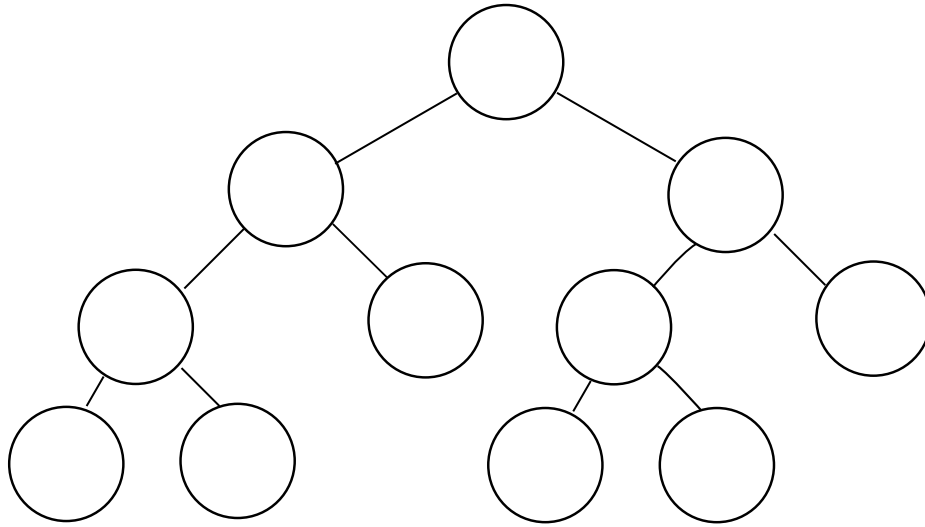
For any binary tree with n nodes, the number of edges e is always

$$e = n - 1.$$

This holds because every node (except the root) is connected to exactly one parent, so there are $n-1$ edges in the tree.

6.4 Full trees

A full tree is a tree where all internal nodes are degree two, and all leaf nodes are degree zero. Observe



The next three subsections refer to the *full binary tree theorem*, which states for a nonempty, full tree T

6.4.1 Number of leaves

If T has I internal nodes, the number of leaves is given by

$$L = I + 1.$$

If T has a total of N nodes, the number of leaves is

$$L = \frac{N + 1}{2}.$$

6.4.2 Number of nodes

If T has I internal nodes, the total number of nodes is

$$N = 2I + 1.$$

If T has L leaves, the total number of nodes is

$$N = 2L - 1.$$

6.4.3 Number of internal nodes

If T has a total of N nodes, the number of internal nodes is

$$I = \frac{N - 1}{2}.$$

If T has L leaves, the number of internal nodes is

$$I = L - 1.$$

6.5 Complete Binary Tree

A complete binary tree has a specific structure defined by how the nodes are filled level by level.

1. **All levels, except possibly the last, are fully filled:**

- In a complete binary tree, every level up to the second-to-last (penultimate) level must be completely filled with nodes.
- This means that if the tree has height h , levels 0 through $h - 1$ (from the root to the second-to-last level) will have the maximum possible number of nodes for that level.

2. **All nodes are as far left as possible:**

- On the last level, the nodes don't need to completely fill the level, but the nodes must be positioned as far to the left as possible.
- For example, if some nodes are missing from the last level, they will always be missing from the right side, not from the left.

Notes: The tree is balanced in terms of node distribution, with all the levels except possibly the last fully filled.

Nodes on the last level are always added from the leftmost position first.

6.5.1 Number of nodes

The height h of a complete binary tree is defined as the number of edges on the longest path from the root to a leaf node.

The total number of nodes in a complete binary tree is given by

$$n = 2^{h+1} - 1.$$

6.5.2 Height

The height h of a complete binary tree with n nodes can be derived as:

$$h = \lfloor \log_2(n) \rfloor.$$

6.5.3 Number of Leaf Nodes (L) in a Complete Binary Tree

The number of leaf nodes in a complete binary tree can be calculated based on the number of internal nodes or the height of the tree

$$L = \lceil \frac{n}{2} \rceil.$$

6.5.4 Number of internal nodes

The number of internal nodes (non-leaf nodes) in a complete binary tree can be calculated as:

$$I = N - L$$
$$I = \lfloor \frac{n}{2} \rfloor.$$

6.5.5 Parent and Child Relationships in a Complete Binary Tree

Parent of node at index i (1-based index):

$$\text{Parent}(i) = \left\lfloor \frac{i}{2} \right\rfloor$$

Left child of node at index i :

$$\text{Left child}(i) = 2i$$

Right child of node at index i :

$$\text{Right child}(i) = 2i + 1$$

These relationships assume a 1-based indexing system for the nodes in the tree (common in heaps or array-based representations).

6.6 Perfect binary tree

6.6.1 Number of Nodes

$$N = 2^{h+1} - 1.$$

6.6.2 Number of Leaf Nodes

$$L = 2^h.$$

6.6.3 Height of the Tree

$$h = \log_2(N + 1) - 1.$$

6.6.4 Number of Internal Nodes

$$I = N - L = 2^h - 1.$$

6.6.5 Depth

$$d = h.$$

Applications of binary trees

7.1 Binary search trees

A binary search tree (BST) is a binary tree in which each node has at most two children and follows these properties:

- **Left Subtree Property:** The value of each node in the left subtree is less than the value of the node itself.
- **Right Subtree Property:** The value of each node in the right subtree is greater than the value of the node itself.
- Both left and right subtrees must also be binary search trees.

7.1.1 Interface

The interface of a Binary Search Tree (BST) typically includes a set of operations for managing and accessing the tree's nodes.

- **Insert(value):** Inserts a new value into the BST while maintaining its properties.
- **Remove(value):** Removes a value from the BST, adjusting the structure to maintain its properties.
- **Predecessor(node):** Finds the predecessor of a node
- **Successor(node):** Finds the successor of a node
- **Find(value):** Searches for a value in the BST and returns the node containing it or null if not found.
- **FindMin():** Returns the node with the smallest value in the BST.
- **FindMax():** Returns the node with the largest value in the BST.
- **IsEmpty():** Checks if the BST is empty.
- **Traverse(order):** Traverses the tree in a specific order (e.g., in-order, pre-order, post-order).
- **Height():** Returns the height of the BST.
- **Clear():** Removes all nodes from the tree, making it empty

7.1.2 Traversals

We can traverse BST's in one of four ways

- Level order
- Preorder
- Inorder
- Postorder

7.1.2.1 Level order

Level-order traversal is a way of visiting all the nodes in a binary tree by levels, from top to bottom. It starts at the root and visits nodes level by level, left to right, for each level.

- Start with the root node (the topmost node).
- Visit all the nodes on the next level (children of the root) from left to right.
- Then, visit all nodes on the level below that (grandchildren of the root) from left to right, and so on.

A queue is often used to implement level-order traversal, as it helps keep track of nodes to visit in the correct order.

```
1  void levelorderPrint() {
2      if (!root) return; // noop for empty tree
3
4      queue<node*> q;
5      q.push(root);
6
7      while (!q.empty()) {
8          node* curr = q.front();
9          q.pop();
10
11         cout << curr->data << endl;
12         if (curr->left) {
13             q.push(curr->left);
14         }
15         if (curr->right) {
16             q.push(curr->right);
17         }
18     }
19 }
```

1. If the list is nonempty, construct a queue and push the root node.
2. While the queue is nonempty, grab the front, process the front, pop the front.
3. Push left and right nodes to queue, if they exist.

7.1.2.2 Preorder

Pre-order traversal is a way of visiting nodes in a binary tree where you:

1. Visit the root node first.
2. Recursively visit the left subtree.
3. Recursively visit the right subtree.

To explain simply:

1. Start with the root node.
2. Go as far left as possible, visiting each node along the way.
3. Once you've reached the end of the left subtree, backtrack and visit the right subtree.

```
1  void preorderPrint() {  
2      std::function<void(node*)> r_preorderPrint = [&] (node* p) {  
3          if (p == nullptr) return;  
4  
5          cout << p->data << endl;  
6          r_preorderPrint(p->left);  
7          r_preorderPrint(p->right);  
8      };  
9      r_preorderPrint(root);  
10 }
```

7.1.2.3 Inorder

in-order traversal is a way of visiting nodes in a binary tree where you:

1. Recursively visit the left subtree first.
2. Visit the root node.
3. Recursively visit the right subtree.

To explain simply:

1. Start by going all the way to the left, visiting nodes along the way.
2. Once you reach the leftmost node, visit it, then move up to its parent (the root).
3. After visiting the root, visit the right subtree.

For a BST, printing the tree with an inorder traversal yields a sorted sequence.

```

1  void inorderPrint() {
2      std::function<void(node*)> r_inorderPrint = [&] (node* p) ->
    ↪  void {
3          if (!p) return;
4
5          r_inorderPrint(p->left);
6          cout << p->data << endl;
7          r_inorderPrint(p->right);
8      };
9      r_inorderPrint(root);
10 }

```

7.1.2.4 Postorder

Post-order traversal is a way of visiting nodes in a binary tree where you:

1. Recursively visit the left subtree first.
2. Recursively visit the right subtree.
3. Finally, visit the root node.

To explain simply:

1. Start by going to the leftmost node, but don't visit it yet.
2. Then, go to the right subtree and process it.
3. After both subtrees have been visited, visit the root.

```
1 void postorderPrint() {  
2     std::function<void(node*)> r_postorderPrint = [&] (node* p)  
    ↪ -> void {  
3         if (!p) return;  
4  
5         r_postorderPrint(p->left);  
6         r_postorderPrint(p->right);  
7         cout << p->data << endl;  
8     };  
9     r_postorderPrint(root);  
10 }
```

7.1.3 Successor of a node

The successor of a node is defined mathematically as

$$\text{succ}(X) = \min\{A : A > X\}.$$

thus, we find the set of all nodes that have values greater than that of X , then find the minimum in that set.

By properties of binary search trees we find the successor of a node X by

1. **If X has a right child:** The successor is the leftmost node in the right subtree of X (the smallest node in the right subtree).
2. **If X has no right child:**
 - If the node is the left child of its parent, then the parent is its successor.
 - If the node is the right child of its parent, you move upward until you find a node that is the left child of its parent, and that parent is the successor.

7.1.4 Predecessor

The predecessor of a node X is defined as

$$\text{pred}(X) = \max\{A : A < X\}.$$

In other words it is the largest node that is less than X . To find the predecessor:

1. **If X has a left child:** The predecessor is the rightmost node in the left subtree
2. **If X has no left child:** The predecessor is the nearest ancestor for which the node is in the right subtree.

7.1.5 The node

The node is similar to a linked list node, but instead of a single next pointer, it has two. A left pointer and a right pointer.

```
1  struct node{
2      node* left = nullptr;
3      node* right = nullptr;
4      int data = 0;
5
6      node() = default;
7      node(int data) : data(data) {}
8      node(node* left, node* right, int data) : left(left),
    ↪ right(right), data(data) {}
9  };
```

7.1.6 The class

For simplicity, we often define the Binary Search Tree (BST) as a class. This allows each instance of the class to hold its own root node, along with other data members such as the size of the tree, that we may need.

If it were not a class, then each function would need to take the root node as an argument and return the (potentially modified) root node to maintain the structure.”

If it were not a class, then each function would have to take as an argument a root node, and return the root node to maintain the structure

```
1  class BST {
2  private:
3      node* root;
4      ...
5  public:
6      ...
7  };
```

7.1.7 Recursive Insertion

Because of the nature of BST's, we often use recursion to define the needed operations.

```
1 void insert(int element) {
2     // If the tree is empty, insert new element as root
3     if (!root) {
4         root = new node(element);
5         return;
6     }
7
8     std::function<void(node*)> r_insert = [&](node* p) -> void {
9
10        // If the element is less than current node, and p->left
11        ↪ exists, go left
12        if (element < p->data && p->left) {
13            r_insert(p->left);
14
15            // If the element is greater than current node, and
16            ↪ p->right exists, go right
17            } else if (element > p->data && p->right) {
18                r_insert(p->right);
19            }
20
21            // If the element is less than current node, and p->left
22            ↪ doesn't exist, insert node as current nodes left child
23            if (element < p->data && !p->left) {
24                p->left = new node(element);
25                return;
26
27                // If the element is greater than current node, and
28            ↪ p->right doesn't exist, insert node as current nodes right
29            ↪ child
30            } else if (element > p->data && !p->right) {
31                p->right = new node(element);
32                return;
33            }
34        };
35        // Start recursion from the root
36        r_insert(root);
37    }
```

If the tree is empty, it creates a new root node with the given element.

Otherwise, it uses a recursive lambda function (r_insert) to:

- Traverse the tree: going left if the element is smaller, or right if the element is larger.
- Once it finds an appropriate spot (where a left or right child doesn't exist), it inserts the new node as a left or right child accordingly.

The process starts from the root and recursively finds the right place to insert the new element.

7.1.8 A better recursive insert

```
1  node* r_insertC(node* p, int element) {
2      if (!p) return new node(element);
3
4      if (element < p->data) {
5          p->left = r_insertC(p->left, element);
6      } else if (element > p->data) {
7          p->right = r_insertC(p->right, element);
8      }
9      return p;
10 }
11
12 void insertC(int element) {
13     if (!root) root = new node(element);
14     r_insertC(root, element);
15 }
```

The main insertion function, `insertC`, initiates the process. If the tree's root is null, meaning the tree is empty, it creates a new root node with the given element. Otherwise, it calls the recursive helper function `r_insertC` on the root to handle the insertion process.

The `r_insertC` function operates recursively to find the correct location for the new element within the tree. Starting from the given node `p`, it checks whether `p` is null; if so, it creates and returns a new node with the specified element, making this node the new leaf of the tree at this position. If `p` is not null, the function compares the element to `p->data`. If the element is smaller, the function recursively calls `r_insertC` on `p->left` to continue searching in the left subtree, and if the element is larger, it calls `r_insertC` on `p->right` to search in the right subtree. After setting the appropriate child link, it returns the current node `p`, maintaining the correct structure of the tree at each level of recursion. This ensures the BST properties are preserved, with each node's left subtree containing values less than the node's data and the right subtree containing values greater.

Left as an exercise to the reader to see why we must return `p` to maintain the tree pointer chain.

7.1.9 Iterative insert

```
1 void insertB(int element) {
2     if (!root) {
3         root = new node(element);
4         return;
5     }
6
7     node* p = root, *trail = nullptr;
8     bool left;
9
10    while (p) {
11        trail = p;
12        if (element < p->data) {
13            p=p->left;
14            left=true;
15        } else if (element > p->data) {
16            p=p->right;
17            left=false;
18        } else {
19            return; // noop if already exists
20        }
21    }
22    if (left) {
23        trail->left = new node(element);
24    } else {
25        trail->right = new node(element);
26    }
27 }
```

If the tree is empty, it creates a new root node with the element.

It then iteratively traverses the tree starting from the root:

- Moves left if the element is smaller than the current node's data.
- Moves right if the element is larger.
- If the element already exists, it does nothing and returns.

Once it finds an empty spot (either left or right child is nullptr), it inserts the new node as the left or right child of the parent node (trail), depending on the comparison.

7.1.10 Recursive removing

To remove a node with a given value from a BST, there are three cases

1. Node has no children
2. Node has one child
3. Node has two children

For case I, we can simply set the nodes parent to nullptr, and then delete the node.

For case II, we must divert the connection from the nodes parent to the nodes child, and then free the node.

Case III is more involved, we first must find the successor of the node. Once we find the successor, we replace the nodes data value with its successor. Then, instead of deleting the node, we delete its successor. Since to be in this case the node must have exactly two children, the successor is found in the simple way.

1. Go right once
2. Go as far left as possible.

Once we have the successor node, it will either have no children, or exactly one child (a right child), if it were to have a left child, it would not be the true successor because we would have not gone as far left as possible.

```

1 void remove(int element) {
2     if (!root) return; // Noop for empty tree
3
4     std::function<void(node*&, node*&)> r_remove = [&] (node*&
↪ p, node*& last) -> void {
5         if (!p) return; // Not found in tree
6
7         if (element < p->data) {
8             r_remove(p->left, p);
9         } else if (element > p->data) {
10            r_remove(p->right, p);
11        } else { // Found
12            // Case I: Node has zero children
13            if (!p->left && !p->right) {
14                node* tmp = p;
15                p=nullptr;
16                delete tmp;
17                // Case II: Node has one child (note the use of
↪ xor)
18            } else if (!p->left ^ !p->right) {
19                node* tmp = p;
20                p = (p->left ? p->left : p->right);
21                delete tmp;
22                // Case III: Two children
23            } else {
24                node* successor = p->right;
25                node* successorParent = p;
26
27                // Find the in-order successor
28                while (successor->left) {
29                    successorParent = successor;
30                    successor = successor->left;
31                }
32
33                // Replace nodes value with successor value
34                p->data = successor->data;
35
36                // Now we need to delete the successor node
37                // The successor is a leaf or has a right child
38                if (successorParent->left == successor) {
39                    successorParent->left = successor->right;
40                } else {
41                    successorParent->right = successor->right;
42                }
43                delete successor;
44            }
45        }
46    };
47    r_remove(root,root);
48 }

```

7.1.11 Clearing

```
1  void clear() {  
2      if (!root) return;  
3  
4      std::function<void(node*)> r_clear = [&](node* p) -> void {  
5          if (!p) return;  
6  
7          r_clear(p->left);  
8          r_clear(p->right);  
9  
10         delete p;  
11     };  
12     r_clear(root);  
13     root = nullptr;  
14 }
```

This function deletes all nodes in a binary search tree. It recursively traverses the tree, deleting each node after its children have been deleted, and finally sets the root to `nullptr`, effectively clearing the entire tree

7.1.12 Counting the height of the tree (root)

```
1  size_t height() {
2      std::function<size_t(node*)> r_height = [&](node* p) ->
    ↪ size_t {
3          // Base case height of a nullptr is zero
4          if (!p) return 0;
5          return 1+std::max(r_height(p->left), r_height(p->right));
6      };
7      // Height is counting edges, so its number nodes in longest
    ↪ path from root to leaf - 1
8      return r_height(root) -1;
9  }
```

This code defines a `height()` function that calculates the height of a binary tree by counting the edges. It uses a recursive lambda function `r_height` to traverse the tree. For each node, it returns `1 + max(left subtree height, right subtree height)` to find the longest path from the root to any leaf. Since `r_height` counts nodes, the function subtracts 1 at the end to convert the node count to edge count, which is the definition of height.

7.1.13 Counting the height of a node

```
1  int getHeight(node* p) {
2      if (!p) return -1;
3      return 1 + std::max(getHeight(p->left), getHeight(p->right));
4  }
```

If height counts edges, then a `nullptr` nodes must return -1. If height counts vertices, then `nullptr` nodes must return 0.

7.1.14 Getting the depth of the node

```
1  int nodeDepth(node* p) {
2      if (p == root) return 1;
3      return r_nodeDepth(root, p, 1);
4  }
5
6  int r_nodeDepth(node* curr, node* p, int depth) {
7      if (curr == nullptr) {
8          return -1; // Node not found
9      }
10     if (p == curr) {
11         return depth; // Node found, return the depth
12     }
13
14     // Recursively search in the left subtree if p's data is
    ↪ smaller
15     if (p->data < curr->data) {
16         return r_nodeDepth(curr->left, p, depth + 1);
17     }
18     // Recursively search in the right subtree if p's data is
    ↪ greater
19     if (p->data > curr->data) {
20         return r_nodeDepth(curr->right, p, depth + 1);
21     }
22
23     return -1; // This should never be reached if the tree is
    ↪ valid
24 }
```

This code defines two functions that work together to calculate the depth of a given node p in a binary search tree (BST):

nodeDepth: This function is the entry point to calculate the depth of the node p .

It first checks if p is the root node. If so, it returns 1 since the root node is considered to have a depth of 1.

If p is not the root, it calls the helper function `r_nodeDepth` to recursively search for the node, starting from the root with an initial depth of 1.

r_nodeDepth: This is a recursive helper function that searches for the node p in the tree, while tracking the current depth.

It first checks if `curr` (the current node in the search) is `nullptr`, which indicates that the node p is not in the tree. In that case, it returns -1.

If `curr` matches p , it returns the current depth.

Otherwise, it recursively searches in the left or right subtree, depending on whether p 's data is less than or greater than the current node's data, and increments the depth by 1 at each recursive step.

7.1.15 Counting the number of nodes

```
1  template <typename NODE>
2  int count(NODE * root) {
3      if (!root) return 0;
4
5      return 1 + count(root->left) + count(root->right);
6  }
7
```

7.1.16 Comparison traversals

In this section we show a typical recursive algorithm to compare two binary search trees via an inorder traversal.

```
1  // Node compare
2  bool nc(node* p, node* q) {
3      return p->data == q->data;
4  }
5
6  bool r_inorderComp(node* p, node* q) {
7      if (!p && !q) return true;
8      if (!p || !q) return false;
9
10     if (!r_inorderComp(p->left, q->left)) return false;
11     if (!nc(p,q)) return false;
12     if (!r_inorderComp(p->right, q->right)) return false;
13
14     return true;
15 }
16
17 bool inorderComp(tree& t1, tree& t2) {
18
19     if (!t1.root && !t2.root) return true;
20     if (!t1.root || !t2.root) return false;
21
22     return r_inorderComp(t1.root, t2.root);
23 }
24
```

The core of the comparison is performed by two functions, `nc` and `r_inorderComp`, which are used to check if each corresponding node in the two trees holds the same data value and is structured identically.

The `nc` function is a helper that takes two nodes as parameters and simply compares their data values, returning true if the values are equal and false otherwise.

The function `r_inorderComp` performs a recursive, in-order traversal of two nodes, `p` and `q`. If both nodes are `nullptr`, it means that this position in both trees is empty, so they match and the function returns `true`. If one node is `nullptr` and the other is not, it indicates a structural mismatch, so the function returns `false`. The recursion first compares the left child nodes, then the current nodes themselves using `nc`, and finally the right child nodes. If any comparison fails, the function returns `false`; otherwise, it completes successfully, confirming that the structures and values match in this part of the trees.

The `inorderComp` function is the main entry point for the comparison. It takes two tree objects, `t1` and `t2`, and checks if both trees are empty, returning `true` if they are. If one tree is empty while the other is not, it returns `false` due to a structural difference. If both trees have a root node, `inorderComp` then calls `r_inorderComp` on the root nodes of `t1` and `t2`, performing the recursive in-order comparison on the entire structures.

7.1.17 Finding the smallest and largest values

By properties of BST's, to find the smallest value in a tree, we start from the root and go as far left as possible. To find the largest, we go as far right as possible

```
1  node* r_min(node* p) {
2      if (!p->left) return p;
3      return r_min(p->left);
4  }
5  node* min() {
6      if (!root) return nullptr;
7      return r_min(root);
8  }
9
10 node* r_max(node* p) {
11     if (!p->right) return p;
12     return r_max(p->right);
13 }
14 node* max() {
15     if (!root) return nullptr;
16     return r_max(root);
17 }
```

7.1.18 Getting the widths of a bst

```
1  std::vector<int> getWidths() {
2      if (!root) return;
3
4      std::vector<int> w(height());
5      std::queue<node*> q;
6      w[0] = 1;
7      q.push(root);
8      q.push(nullptr);
9
10     int level = 1, width=0;
11     while (!q.empty()) {
12         node* curr = q.front();
13         q.pop();
14
15         if (curr == nullptr) {
16             w[level++] = width;
17             width = 0;
18             q.push(nullptr);
19             continue;
20         }
21         if (level == height()) {
22             break;
23         }
24
25         if (curr->left) {
26             q.push(curr->left);
27             ++width;
28         }
29         if (curr->right) {
30             q.push(curr->right);
31             ++width;
32         }
33     }
34     return w;
35 }
```

This C++ code defines a `getWidths` function that calculates and returns a vector of integers representing the width of each level in a binary tree. The width of a level is defined as the number of nodes at that level.

The function first creates a vector `w` with a size equal to the height of the tree (obtained from a hypothetical `height()` function) and initializes the first element, `w[0]`, to 1, assuming the root level has one node. It also sets up a queue `q` to facilitate level-order traversal (breadth-first traversal) of the tree, starting by pushing the root node and a `nullptr` marker, which denotes the end of each level.

The main loop continues as long as the queue is not empty. It dequeues the front node in each iteration, checking if it's `nullptr`. When a `nullptr` is encountered, it means the current level has ended, so the function records the width for that level in the `w` vector, resets the width counter, and pushes another `nullptr` if there are more nodes to process at deeper levels. The loop then proceeds to the next level.

If the dequeued node is not nullptr, it examines its left and right children. If they exist, they are added to the queue, and the width counter is incremented for each child node. This process continues until all levels are processed or the traversal reaches the specified tree height, returning the vector `w` with each level's width at the end.

7.1.19 Degenerate Binary Search trees

A degenerate binary search tree is a tree where each parent node has only one child, causing the tree to resemble a linked list

Consider building a binary search tree, taking values from a sorted array from left to right. Since all subsequent entries will be greater than the previous, the insertions will only go in one direction, right. Thus, the final tree will resemble a linked list and thus we will not be able to use the $\lg(n)$ property of binary trees.

7.1.20 Verifying a binary search tree

To verify a binary search tree, we need to check that for a given node, all nodes to the left have values less than the node, and all nodes to the right have values greater than the node. This can be achieved using recursion and passing down minimum and maximum value constraints for each subtree.

```
1  bool r_is_bst(node* p, int min, int max) {
2      if (!p) return true;
3
4      if (p->data < min || p->data > max) return false;
5
6      return r_is_bst(p->left, min, p->data) && r_is_bst(p->right,
    ↪ p->data, max);
7  }
8
9  bool is_bst() {
10     if (!root) return true;
11     return r_is_bst(root, std::numeric_limits<int>::min(),
    ↪ std::numeric_limits<int>::max());
12 }
```

- We start with the root node, which has the range of INT_MIN to INT_MAX.
- For the left child, we update the maximum allowed value to the parent's value.
- For the right child, we update the minimum allowed value to the parent's value.

7.1.21 Complexities

Because BST have no guarantee of being well formed, (ie degenerate trees), the complexity of many operations in the worst case is $O(n)$.

Operation	Best Case	Average Case	Worst Case
Insertion	$O(\log n)$	$O(\log n)$	$O(n)$
Search	$O(1)$	$O(\log n)$	$O(n)$
Removal	$O(\log n)$	$O(\log n)$	$O(n)$
Height	$O(\log n)$	—	$O(n)$
Traversal	$O(n)$	$O(n)$	$O(n)$

Note: $\Omega(\lg(n))$ for BST operations like search, insert, and delete (best case).

$\Omega(\lg(n))$ for operations that require visiting all nodes, like traversals.

7.2 Adelson-Velsky and Landis Trees (AVL trees)

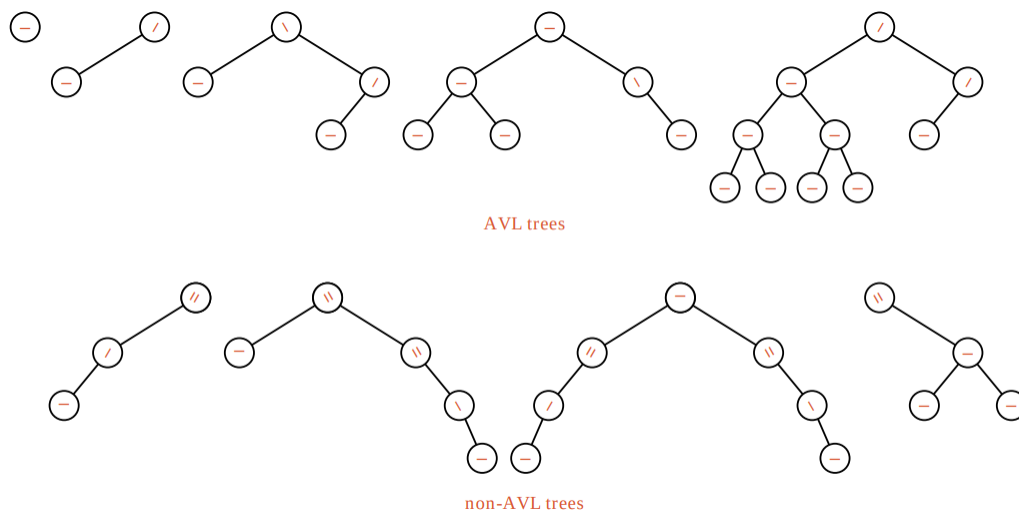
When dealing with binary search trees, insertions and removals occur continually, with no predictable order. In some of these applications, it is important to optimize search times by keeping the tree very nearly balanced at all times. The method in this section for achieving this goal was described in 1962 by two Russian mathematicians, G. M. ADEL'SON-VEL'SKI'I and E. M. LANDIS, and the resulting binary search trees are called AVL trees in their honor.

AVL trees achieve the goal that searches, insertions, and removals in a tree with n nodes can all be achieved in time that is $O(\log n)$, even in the worst case. The height of an AVL tree with n nodes, as we shall establish, can never exceed $\lg n$, and thus even in the worst case, the behavior of an AVL tree could not be much below that of a random binary search tree. In almost all cases, however, the actual length of a search is very nearly $\lg n$, and thus the behavior of AVL trees closely approximates that of the ideal, completely balanced binary search tree.

7.2.1 Definition

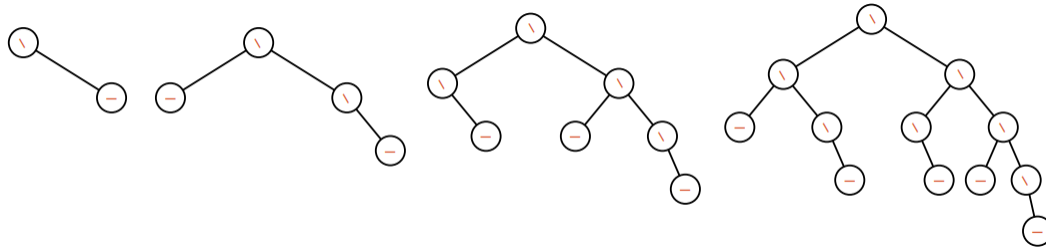
An AVL tree is a binary search tree in which the heights of the left and right subtrees of the root differ by at most 1 and in which the left and right subtrees are again AVL trees.

With each node of an AVL tree is associated a *balance factor* that is **lefthigher**, **equal-height**, or **right-higher** according, respectively, as the left subtree has height greater than, equal to, or less than that of the right subtree.



In drawing diagrams, we shall show a left-higher node by ‘/,’ a node whose balance

factor is equal by ‘-,’ and a right-higher node by ‘\’. The figure above shows several small AVL trees, as well as some binary trees that fail to satisfy the definition. Note that the definition does not require that all leaves be on the same or adjacent levels.



The figure above shows several AVL trees that are quite skewed, with right subtrees having greater height than left subtrees.

7.2.2 AVL Nodes

A typical node for an AVL tree is as follows

```

1  struct node {
2      node* left{nullptr}, *right{nullptr};
3      int data{0};
4      int height{0};
5      balance b{0};
6
7      // Constructors
8      ...
9  };

```

In AVL trees, it is common for the node structure to include a height member. This height field is essential for efficiently maintaining the balance of the tree, as the balance factor of a node (the difference in height between its left and right subtrees) is used to determine whether the tree needs rebalancing after insertions or deletions.

Without the height member, recalculating the height of each node during every operation would require traversing the entire subtree, significantly increasing the time complexity. By storing the height in each node, you can retrieve it in constant time, allowing rotations and rebalancing operations to remain efficient.

7.2.3 Storing the height

In AVL trees, we should store the height of a node as a data member, updating on insertions or removals into the tree.

7.2.4 Defining balance factors in C++ with enums

We employ an enumerated data type to record balance factors

```
1  enum Balance_factor { left_higher, equal_height, right_higher };
```

Balance factors must be included in all the nodes of an AVL tree, and we must adapt our former node specification accordingly.

7.2.5 Defining balance factors with a height calculation

Instead, we can define the balance factor of a node as $\text{height}(\text{left}) - \text{height}(\text{right})$. A negative balance implies a subtree is heavy on the right, a positive balance implies a subtree is heavy on the left. A balance of zero implies equal height subtrees. A balance factor ranges from -2 to 2. If $|\text{balance}| = 2$, we must rotate the tree to restore balance.

```
1  int getBalance(node* p) {  
2      if (!p) return -1;  
3      return (p->left ? p->left->height : -1) - (p->right ?  
    ↪ p->right->height : -1);  
4  }
```


7.2.6 Interface

The interface of a Binary Search Tree (BST) and an AVL Tree is generally the same. Both are types of binary trees, and they share similar operations such as:

- **Insert:** Insert a new element into the tree.
- **Remove/Delete:** Remove an element from the tree.
- **Search:** Find whether a particular element exists in the tree.
- **Traversal:** In-order, pre-order, post-order, and level-order traversals.

However, behind the scenes, the AVL tree performs additional work to maintain its balance property, but this doesn't typically change the public interface

7.2.7 Balancing an AVL tree

As we insert or remove nodes from the tree, it may happen that the resulting tree fails to satisfy the conditions imposed by AVL trees. To *rebalance* the tree, we have a set of operations, called rotations. We have

1. **Right Rotation (RR):** Applied when a left subtree is too deep. The subtree is rotated to the right, reducing the height of the left side.
2. **Left Rotation (LL):** Applied when a right subtree is too deep. The subtree is rotated to the left, reducing the height of the right side.
3. **Left-Right Rotation (LR):** Occurs when a left subtree has a deep right subtree. First, a left rotation is performed on the left subtree, followed by a right rotation.
4. **Right-Left Rotation (RL):** Occurs when a right subtree has a deep left subtree. First, a right rotation is performed on the right subtree, followed by a left rotation.

These rotations are applied based on the balance factor, ensuring the tree remains balanced with a height difference of at most 1 between subtrees.

Note: Left-right and Right-left rotations are also called double right and double left respectively.

When writing our rotation algorithms, we only need to define two. A left rotation algorithm and a right rotation algorithm

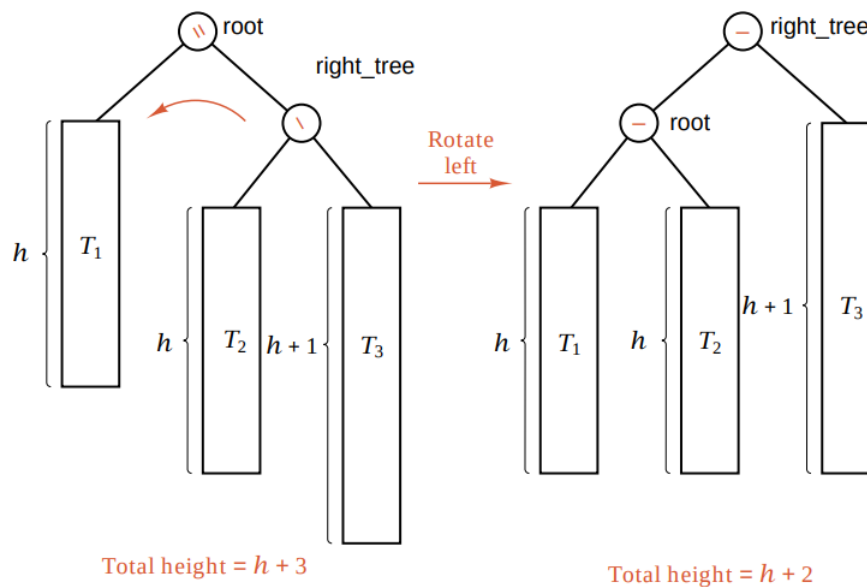
Note: Performing a rotation when the height difference is only 1 would be unnecessary and could actually disrupt the balancing of the tree. The tree is still balanced in this case because a height difference of 1 between subtrees is allowed in AVL trees.

7.2.8 Rotations: Right tree

Let us now consider the case when a new node has been inserted into the taller subtree of a root node and its height has increased, so that now one subtree has height 2 more than the other, and the tree no longer satisfies the AVL requirements. We must now rebuild part of the tree to restore its balance. To be definite, let us assume that we have inserted the new node into the right subtree, its height has increased, and the original tree was right higher. That is, we wish to consider the case covered by the function `right_balance`. Let `root` denote the root of the tree and `right_tree` the root of its right subtree.

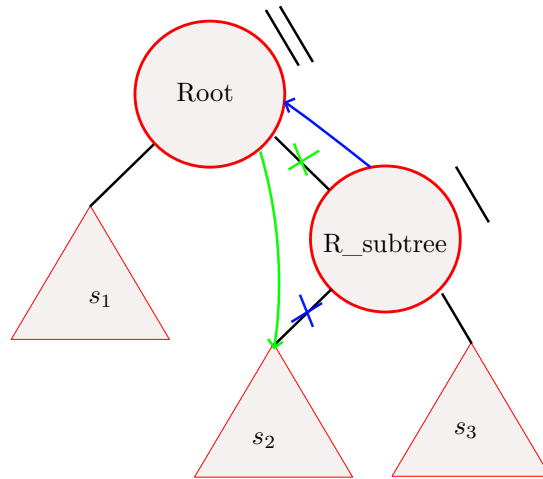
There are three cases to consider, depending on the balance factor of `right_tree`.

7.2.8.1 Case 1: Right higher The first case, when `right_tree` is right higher. The action needed in this case is called a left rotation. We have rotated the node `right_tree` upward to the root, dropping `root` down into the left subtree of `right_tree`; the subtree `T2` of nodes with keys between those of `root` and `right_tree` now becomes the right subtree of `root` rather than the left subtree of `right_tree`. A left rotation is succinctly described in the following C++ function. Note especially that, when done in the appropriate order, the steps constitute a rotation of the values in three pointer variables. Note also that, after the rotation, the height of the rotated tree has decreased by 1; it had previously increased because of the insertion; hence the height finishes where it began.



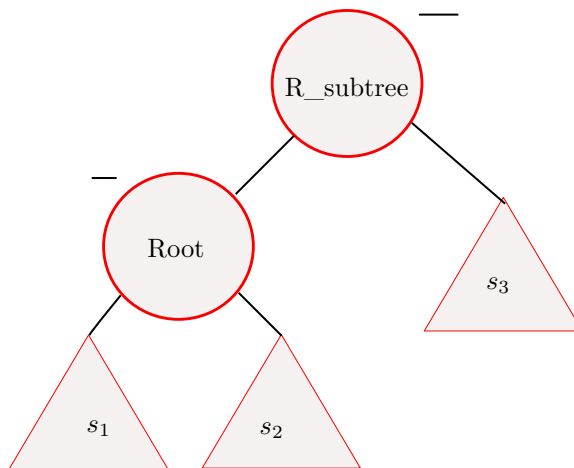
Thus, we come to the following implementation of a left rotation.

```
1 void left_rotate(node* root) {  
2     if (!root || !root->right) return;  
3  
4     right_subtree = root->right;  
5     root->right = right_subtree->left;  
6     right_subtree->left = root;  
7     right_subtree=root;  
8 }
```



1. **Step 1 (green):** Attach right_subtrees left subtree to the right of root
2. **Step 2 (blue):** Attach root to the left of right_subtree
3. **Step 3:** Make r_subtree the new root.

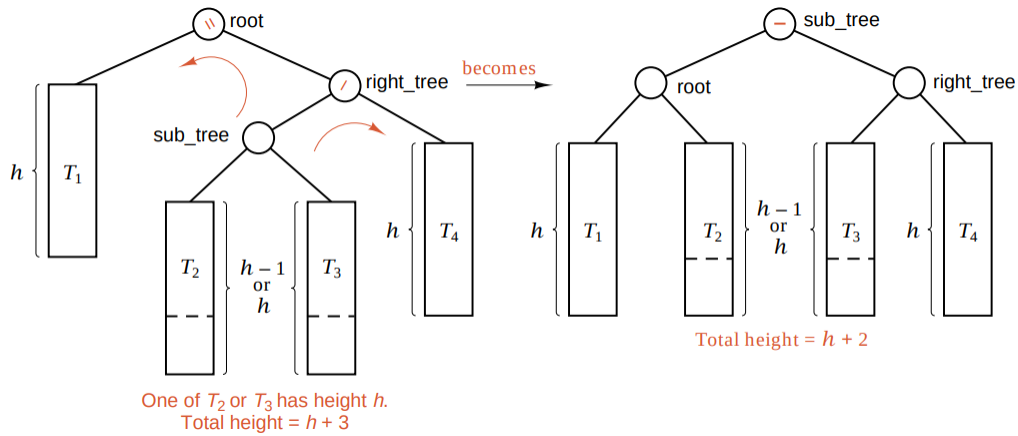
Yields the rotated tree



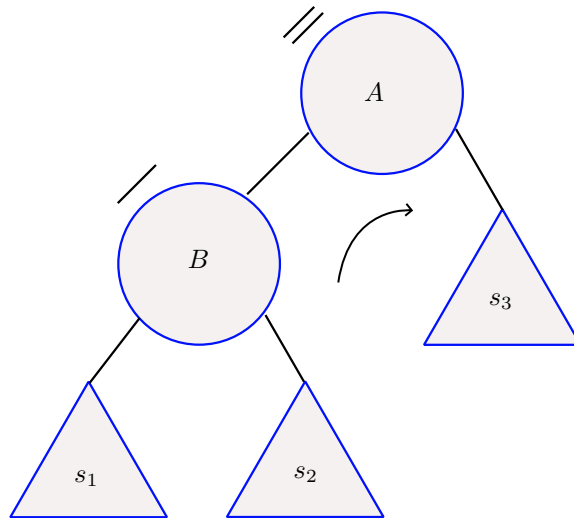
Note that we perform a left rotation when the roots right subtree has an insertion into its right subtree. Notice that the balance symbols point in the same direction. The only balance factors that change are the balance factors of root and right_subtree, they become even (balanced).

7.2.8.2 Case 2: Left higher

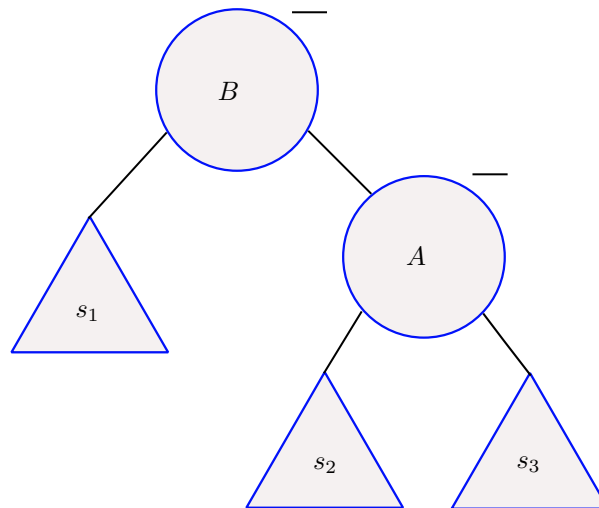
The second case, when the balance factor of `right_tree` is left higher, is slightly more complicated. It is necessary to move two levels, to the node `sub_tree` that roots the left subtree of `right_tree`, to find the new root. This process is shown in Figure 10.20 and is called a double rotation, because the transformation can be obtained in two steps by first rotating the subtree with root `right_tree` to the right (so that `sub_tree` becomes its root), and then rotating the tree pointed to by `root` to the left (moving `sub_tree` up to become the new root).



We see from the figure that we first perform a right rotation on `right_tree`, and `right_tree`'s left subtree, this shifts the subtree unbalance in `right_subtree` such that it becomes unbalanced on the right instead of on the left. This enables us to perform a left rotation on `root` and `right_tree`. Note that after the right rotation, the balance symbols point in the same direction. A right rotation is of the form



Which becomes



Thus, we see that like the left rotation, the only subtree that moves is s_2 . In the left rotation, the left subtree of B becomes the right subtree of A . In the right rotation, the right subtree of B becomes the left subtree of A . A right rotation occurs when a tree becomes unbalanced on the left. Whereas a left rotation occurs when a tree becomes unbalanced on the right.

Thus, we come to the following implementation of a left rotation.

```
1 void right_rotate(node* root) {
2     if (!root || !root->right) return;
3
4     left_subtree = root->left;
5     root->left = left_subtree->right;
6     left_subtree->right = root;
7     left_subtree=root;
8 }
```

7.2.9 C++ Rotations

```
1  node* leftRotate(node* p) {
2      node* subp = p->right;
3      p->right = subp->left;
4      subp->left = p;
5
6      p->height = 1 +
7          std::max(p->left ? p->left->height : -1,
8                  p->right ? p->right->height : -1);
9
10     subp->height = 1 +
11         std::max(subp->left ? subp->left->height : -1,
12                 subp->right ? subp->right->height : -1);
13
14     p->balance = getBalance(p);
15     subp->balance = getBalance(subp);
16     return subp;
17 }
18
19 node* rightRotate(node* p) {
20     node* subp = p->left;
21     p->left = subp->right;
22     subp->right = p;
23
24     p->height = 1 +
25         std::max(p->left ? p->left->height : -1,
26                 p->right ? p->right->height : -1);
27
28     subp->height = 1 +
29         std::max(subp->left ? subp->left->height : -1,
30                 subp->right ? subp->right->height : -1);
31
32     p->balance = getBalance(p);
33     subp->balance = getBalance(subp);
34     return subp;
35 }
```

In `leftRotate`, the function takes a node pointer `p` (the root of a subtree that is right-heavy) and performs a left rotation to reduce the height of its right subtree. First, it saves the right child of `p` (denoted as `subp`). Then, it adjusts the pointers so that the left child of `subp` becomes the new right child of `p`, and `subp` becomes the new root of this subtree with `p` as its left child. After adjusting the structure, it recalculates the height and balance of both `p` and `subp`, where the height of each node is determined by adding 1 to the maximum height of its left and right children. The function then returns `subp` as the new root of the subtree.

Similarly, `rightRotate` takes a node `p` (the root of a left-heavy subtree) and performs a right rotation to rebalance it. It stores the left child of `p` as `subp`, then updates pointers so that the right child of `subp` becomes the new left child of `p`, and `subp` becomes the new root with `p` as its right child. Heights and balances for `p` and `subp` are recalculated using the heights of their children, and `subp` is returned as the new root of this subtree. These rotations are essential for keeping the AVL tree balanced after insertions and deletions.

Note: To understand the return values of the rotateMethods, it is necessary to understand the insertion method below.

7.2.10 Balancing

```
1  node* balance(node* p) {
2      if (p->balance > 1) { // Left heavy
3          if (p->left && p->left->balance < 0) {
4              leftRotate(p->left); // Left-Right case
5          }
6          return rightRotate(p); // Left-Left case
7      }
8      if (p->balance < -1) { // Right heavy
9          if (p->right && p->right->balance > 0) {
10             rightRotate(p->right); // Right-Left case
11          }
12          return leftRotate(p); // Right-Right case
13      }
14      return p;
15  }
```

This function takes a pointer *p* to a node and checks its balance factor (the difference in height between its left and right subtrees) to determine if the node is unbalanced. If the balance factor is greater than 1, the tree is left-heavy, meaning the left subtree is taller than the right subtree. In this case, the function first checks if a "Left-Right" rotation is needed—this would happen if the left child of *p* is itself right-heavy, indicated by a negative balance factor. If so, a left rotation is applied to the left child of *p* to balance it before proceeding. Then, a right rotation is applied to *p* to correct the "Left-Left" imbalance.

Conversely, if the balance factor is less than -1, the node is right-heavy, indicating the right subtree is taller. For a "Right-Left" case, where the right child of *p* is left-heavy, a right rotation is applied to the right child of *p* first. This corrects the imbalance locally and prepares it for the final step: a left rotation applied to *p* to balance the "Right-Right" case. If *p* is already balanced, the function simply returns *p* without modification.

7.2.11 Insertions

```
1  node* r_insert(node* p, int data) {
2      if (!p) return new node(data);
3
4      if (data < p->data) {
5          p->left = r_insert(p->left, data);
6      } else if (data > p->data) {
7          p->right = r_insert(p->right, data);
8      }
9
10     p->height = 1 + std::max(p->left ? p->left->height : -1,
↪    p->right ? p->right->height : -1);
11     p->balance = getBalance(p);
12
13     return balance(p);
14 }
15
16 void insert(int data) {
17     root = r_insert(root, data);
18 }
```

The `r_insert` function takes two parameters: `p`, a pointer to the current node in the tree, and `data`, the integer value to be inserted. If `p` is `nullptr` (indicating an empty spot in the tree), a new node with the given data is created and returned, establishing the base case for recursion.

If the node already exists, `r_insert` decides where to place the new value by comparing `data` with `p->data`. If `data` is less than `p->data`, it recursively calls `r_insert` on the left child (`p->left`). Otherwise, if `data` is greater, it calls `r_insert` on the right child (`p->right`). This ensures that the AVL tree retains its binary search property.

After the recursive insertion, `r_insert` updates the height of the current node `p` by setting it to 1 plus the maximum height of its left and right children. It then calculates the node's balance factor using `getBalance(p)`, which is essential for determining if the subtree rooted at `p` has become unbalanced. Finally, `balance(p)` is called to apply any necessary rotations to maintain the AVL tree's balance, and the function returns the (possibly new) root of the subtree.

The `insert` function is a wrapper that starts the insertion process at the root node of the tree, calling `r_insert` with `root` and the specified data value. This encapsulates the recursive insertion logic and initiates it from the top of the tree.

Note: Worry about ancestry heights (and therefore balance factors) after rotations is not necessary, height and balance calculations occurs only when the insert method returns to that node in the recursion, and thus only after a rotation will an ancestry height and balance calculation be made.

7.2.12 Removing nodes

To implement removal from an AVL tree, you need to perform a few key steps: find and remove the node, balance the tree afterward, and adjust heights and balance factors as necessary

Start by locating the node with the target value. If found, delete it using standard BST deletion rules

After removing a node, retrace your steps back to the root, updating the height and balance factor of each ancestor. If a node is unbalanced, apply rotations to restore balance.

As with insertion, use recursive balancing and height adjustments after deletion.

```
1  node* succ(node* p) {
2      while (p->left) p = p->left;
3      return p;
4  }
5  node* r_remove(node* p, int data) {
6      if (!p) return nullptr; // Node not found
7      // Traverse the tree to find the node to delete
8      if (data < p->data) {
9          p->left = r_remove(p->left, data);
10     } else if (data > p->data) {
11         p->right = r_remove(p->right, data);
12     } else { // Node to delete is found
13         if (!p->left) {
14             node* temp = p->right;
15             delete p;
16             return temp;
17         } else if (!p->right) {
18             node* temp = p->left;
19             delete p;
20             return temp;
21         } else { // Node has two children
22             node* temp = succ(p->right);
23             p->data = temp->data; // Replace data
24             p->right = r_remove(p->right, temp->data);
25         }
26     }
27     // Update height and balance
28     p->height = 1 + std::max(p->left ? p->left->height : -1,
↪   p->right ? p->right->height : -1);
29     p->balance = getBalance(p);
30
31     // Balance the node if necessary
32     return balance(p);
33 }
34 void remove(int data) {
35     root = r_remove(root, data);
36 }
```

`r_remove`, performs a recursive search for the node to delete based on the value data. If data is less than `p->data`, it recurses into the left subtree; if data is greater, it recurses into the right subtree. When the node to delete is found, three cases are handled based on the structure of `p`:

1. If `p` has no left child, `p` is replaced by its right child, effectively bypassing it in the tree.
2. If `p` has no right child, `p` is replaced by its left child.
3. If `p` has two children, the in-order successor (smallest node in `p->right`) is found using `succ`. The data in `p` is replaced with this successor's data, and the successor node is then removed from `p->right` to avoid duplication.

After deletion, the function updates the height and balance of `p`, recalculating the height as one plus the maximum height of its children. The balance factor, which measures the difference in height between the left and right subtrees, is also recalculated. Finally, the function calls `balance(p)` to ensure that any imbalance introduced by the deletion is corrected. The `remove` function wraps this entire process, starting the deletion at the root. This ensures the AVL tree remains balanced after a node is removed.

7.3 Red-black trees

A red-black tree is a binary search tree with one extra bit of storage per node: its color, which can be either **RED** or **BLACK**. By constraining the node colors on any simple path from the root to a leaf, red-black trees ensure that no such path is more than twice as long as any other, so that the tree is approximately balanced.

as we're about to see, the height of a red-black tree with n keys is at most $2\lg(n+1)$, which is $O(\lg n)$

Each node of the tree now contains the attributes color, key, left, right, and p . If a child or the parent of a node does not exist, the corresponding pointer attribute of the node contains the value NIL. Think of these NILs as pointers to leaves (external nodes) of the binary search tree and the normal, key-bearing nodes as internal nodes of the tree.

A red-black tree is a binary search tree that satisfies the following red-black properties:

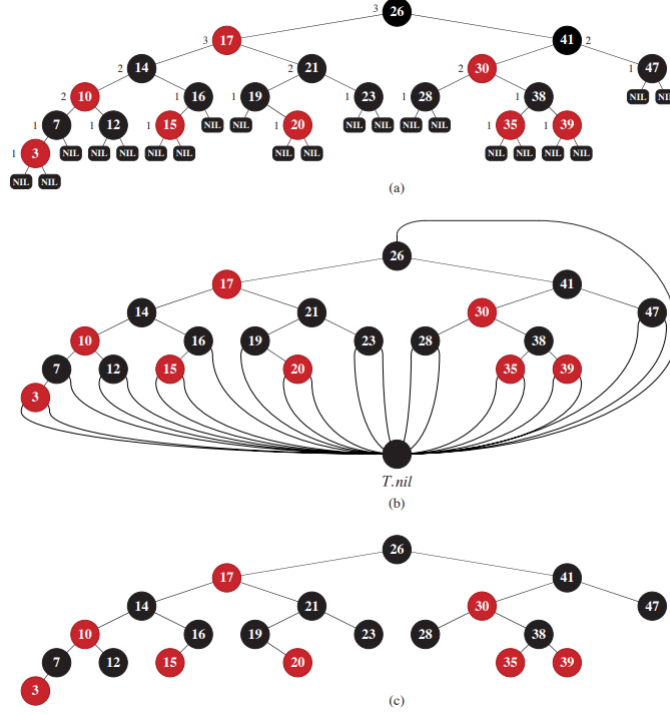
1. Every node is either red or black.
2. The root is black.
3. Every leaf (NIL) is black.
4. If a node is red, then both its children are black. (no double red)
5. For each node, all simple paths from the node to descendant leaves contain the same number of black nodes.

As a matter of convenience in dealing with boundary conditions in red-black tree code, we use a single sentinel to represent NIL. For a red-black tree T , the sentinel $T.\text{nil}$ is an object with the same attributes as an ordinary node in the tree. Its color attribute is BLACK, and its other attributes p , left, right, and key can take on arbitrary values

Why use the sentinel? The sentinel makes it possible to treat a NIL child of a node x as an ordinary node whose parent is x . An alternative design would use a distinct sentinel node for each NIL in the tree, so that the parent of each NIL is well defined. That approach needlessly wastes space, however. Instead, just the one sentinel $T.\text{nil}$ represents all the NILs, all leaves and the root's parent. The values of the attributes p , left, right, and key of the sentinel are immaterial. The red-black tree procedures can place whatever values in the sentinel that yield simpler code

We call the number of black nodes on any simple path from, but not including, a node x down to a leaf the black-height of the node, denoted $bh(x)$.

By property 5, the notion of black-height is well defined, since all descending simple paths from the node have the same number of black nodes. The black-height of a red-black tree is the black-height of its root.



We generally confine our interest to the internal nodes of a red-black tree, since they hold the key values. The remainder of this chapter omits the leaves in drawings of red-black trees, as shown in Figure 13.1(c).

Lemma. A red-black tree with n internal nodes has height at most $2\lg(n + 1)$

Proof. We start by showing that the subtree rooted at any node x contains at least $2^{bh(x)} - 1$ internal nodes. We prove this claim by induction on the height of x .

If the height of x is 0, then x must be a leaf ($T.nil$), and the subtree rooted at x indeed contains at least

$$2^{bh(x)} - 1 = 2^0 - 1 = 0$$

internal nodes.

For the inductive step, consider a node x that has positive height and is an internal node. Then node x has two children, either or both of which may be a leaf. If a child is black, then it contributes 1 to x 's black-height but not to its own. If a child is red, then it contributes to neither x 's black-height nor its own.

Therefore, each child has a black-height of either $bh(x) - 1$ (if it's black) or $bh(x)$ (if it's red). Since the height of a child of x is less than the height of x itself, we can apply the inductive hypothesis to conclude that each child has at least $2^{bh(x)-1} - 1$ internal nodes. Thus, the subtree rooted at x contains at least

$$(2^{bh(x)-1} - 1) + (2^{bh(x)-1} - 1) + 1 = 2^{bh(x)} - 1$$

internal nodes, which proves the claim.

To complete the proof of the lemma, let h be the height of the tree. According to property 4, at least half the nodes on any simple path from the root to a leaf, not including the root, must be black. Consequently, the black-height of the root must be at least $h/2$, and thus,

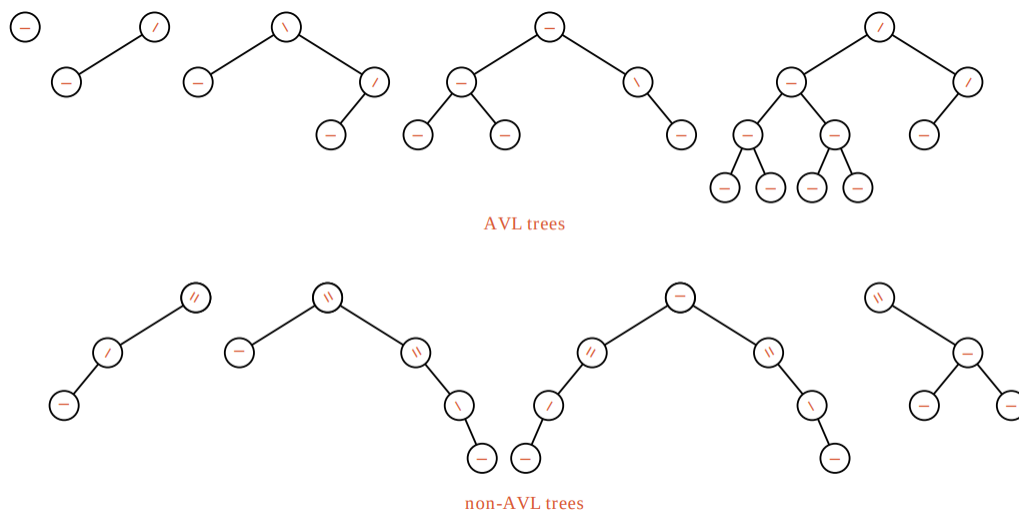
$$n \geq 2^{h/2} - 1.$$

Moving the 1 to the left-hand side and taking logarithms on both sides yields

$$\lg(n + 1) \geq h/2, \text{ or } h \leq 2 \lg(n + 1).$$

☺

As an immediate consequence of this lemma, each of the dynamic-set operations SEARCH, MINIMUM, MAXIMUM, SUCCESSOR, and PREDECESSOR runs in $O(\lg n)$ time on a red-black tree, since each can run in $O(h)$ time on a binary search tree of height h and any red-black tree on n nodes is a binary search tree with height $O(\lg n)$.



7.3.1 Rotations

red-black trees (RB trees) and AVL trees use the same basic rotation logic to maintain balance: left rotation and right rotation. However, they differ significantly in why and when they apply these rotations

The search-tree operations TREE-INSERT and TREE-DELETE, when run on a redblack tree with n keys, take $O(\lg n)$ time. Because they modify the tree, the result may violate the red-black properties. To restore these properties, colors and pointers within nodes need to change.

The pointer structure changes through *rotation*, which is a local operation in a search tree that preserves the binary-search-tree property. Figure 13.2 shows the two kinds of rotations: left rotations and right rotations. Let's look at a left rotation on a node x , which transforms the structure on the right side of the figure to the structure on the left. Node x has a right child y , which must not be $T.nil$. The left rotation changes the subtree originally rooted at x by "twisting" the link between x and y to the left. The new root of the subtree is node y , with x as y 's left child and y 's original left child (the subtree represented by β in the figure) as x 's right child.

Both LEFT-ROTATE and RIGHT-ROTATE run in $O(1)$ time. Only pointers are changed by a rotation, and all other attributes in a node remain the same.

```
1 left-rotate(T, x) {
2     y = x.right
3     // Start by attaching y's left to x's right
4     x.right = y.left
5
6     // If y's left existed, adjust its parent
7     if (y.left != nil)
8         y.left.p = x
9
10    y.p = x.p
11
12    // If x was the root, set y to the new root
13    if (x.p == nil)
14        root = y
15    // If x was not the root, check what x was to its parent
16    ↪ (left or right), change to y
17    else if (x.p.left == x)
18        x.p.left = y
19    else
20        x.p.right = y
21
22    // Finish up
23    y.left = x
24    x.p = y
25 }
```

A right rotate is symmetrical to left rotate

```
1 right-rotate(T, x) {
2     y = x.left
3     x.left = y.right
4
5     if (y.right != nil)
6         y.right.p = x
7
8     y.p = x.p
9     if (x.p == nil)
10        root = y
11    else if (x.p.left == x)
12        x.p.left = y
13    else
14        x.p.right = y
15
16    y.right = x
17    x.p = y
18 }
```


7.3.2 Inserting

The procedure RB-INSERT starts by inserting node z into the tree T as if it were an ordinary binary search tree, and then it colors z red.

Coloring z red ensures that:

1. The black-height of the tree remains unchanged.
2. No violations of the Red-Black Tree properties occur.

To guarantee that the red-black properties are preserved, an auxiliary procedure RB-INSERT-FIXUP on the facing page recolors nodes and performs rotations.

The call $\text{insert}(T, z)$ inserts node z , whose key is assumed to have already been filled in, into the red-black tree T .

```
1  insert(T,z) {
2      x=T.root
3      y=T.nil
4
5      while (x != T.nil) {
6          y = x;
7
8          if (z.key < x.key) {
9              x = x.left
10             } else x = x.right
11     }
12     z.p = y
13
14     // Tree was empty
15     if (y == T.nil) {
16         T.root = z
17     } else if (z.key < y.key) {
18         y.left = z
19     } else y.right = z
20
21     z.left = T.nil
22     z.right = T.nil
23     z.color = RED
24     insert-fixup(T,z)
25
26 }
```

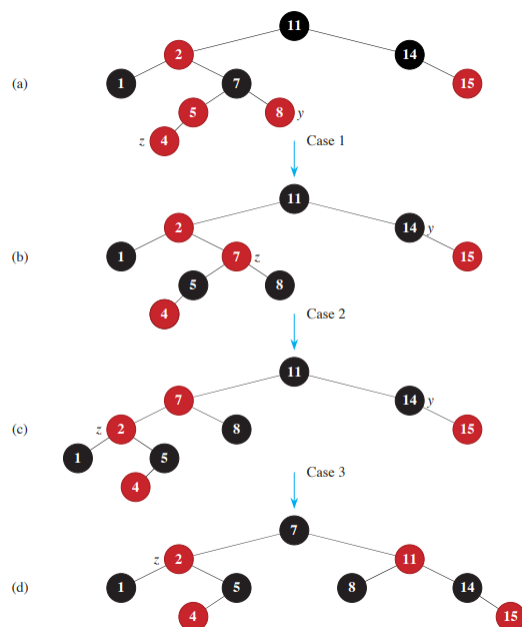
```

1  insert-fixup(T, z) {
2      while (z.p.color == red) {
3          if (z.p == z.p.p.left) { // Is z's parent a left child?
4              y = z.p.p.right // y is z's uncle
5              if (y.color == red) { // are z's parent and uncle
6                  ↪ both red?
7                      /* Case 1 */
8                      z.p.color = black
9                      y.color = black
10                     z.p.p.color = red
11                     z = z.p.p
12                     /* ----- */
13                 } else {
14                     if (z == z.p.right) {
15                         /* case 2 */
16                         z = z.p
17                         left-rotate(T,z)
18                         /* ----- */
19                     }
20                     /* case 3 */
21                     z.p.color = black
22                     z.p.p.color = red
23                     right-rotate(T,z.p.p)
24                     /* ----- */
25                 }
26                 // Right and left exchanged
27             } else {
28                 y = z.p.p.left
29                 if (y.color == red) {
30                     z.p.color = black
31                     y.color = black
32                     z.p.p.color = red
33                     z = z.p.p
34                 } else {
35                     if (z == z.p.left) {
36                         z = z.p
37                         right-rotate(T,z)
38                     }
39                     z.p.color = black
40                     z.p.p.color = red
41                     left-rotate(T,z.p.p)
42                 }
43             }
44         }
45         T.root.color = black

```

To understand how RB-INSERT-FIXUP works, let's examine the code in three major steps. First, we'll determine which violations of the red-black properties might arise in RB-INSERT upon inserting node z and coloring it red. Second, we'll consider the overall goal of the while loop in lines 13-29. Finally, we'll explore each of the three cases within the while loop's body (case 2 falls through into case 3, so these two cases are not mutually exclusive) and see how they accomplish the goal.

In describing the structure of a red-black tree, we'll often need to refer to the sibling of a node's parent. We use the term *uncle* for such a node.



What violations of the red-black properties might occur upon the call to **RB-INSERT-FIXUP**? Property 1 certainly continues to hold (every node is either red or black), as does property 3 (every leaf is black), since both children of the newly inserted red node are the sentinel $T.nil$. Property 5, which says that the number of black nodes is the same on every simple path from a given node, is satisfied as well, because node z replaces the (black) sentinel, and node z is red with sentinel children.

Thus, the only properties that might be violated are property 2, which requires the root to be black, and property 4, which says that a red node cannot have a red child. Both possible violations may arise because z is colored red. Property 2 is violated if z is the root, and property 4 is violated if z 's parent is red.

The **while** loop of lines 13–29 has two symmetric possibilities: lines 33–15 deal with the situation in which node z 's parent $z.p$ is a left child of z 's grandparent $z.p.p$, and lines 17–29 apply when z 's parent is a right child. Our proof will focus only on lines 33–15, relying on the symmetry in lines 17–29.

We'll show that the **while** loop maintains the following three-part invariant at the start of each iteration of the loop

1. Node z is red.
2. If $z.p$ is the root, then $z.p$ is black.
3. If the tree violates any of the red-black properties, then it violates at most one of them, and the violation is of either property 2 or property 4, but not both. If the tree violates property 2, it is because z is the root and is red. If the tree violates property 4, it is because both z and $z.p$ are red.

Part (c), which deals with violations of red-black properties, is more central to showing that **RB-INSERT-FIXUP** restores the red-black properties than parts (a) and (b), which we'll use along the way to understand situations in the code.

Because we'll be focusing on node z and nodes near it in the tree, it helps to know from part (a) that z is red. Part (b) will help show that z 's grandparent $z.p.p$ exists when it's referenced in lines 2, 3, 7, 8, 14, and 15 (recall that we're focusing only on lines 3–15).

To use a loop invariant, we need to show that the invariant is true upon entering the first iteration of the loop, that each iteration maintains it, that the loop terminates, and that the loop invariant gives us a useful property at loop termination. We'll see that each iteration of the loop has two possible outcomes: either the pointer z moves up the tree, or some rotations occur and then the loop terminates.

Before **RB-INSERT** is called, the red-black tree has no violations. **RB-INSERT** adds a red node z and calls **RB-INSERT-FIXUP**. We'll show that each part of the invariant holds at the time **RB-INSERT-FIXUP** is called:

- a. When **RB-INSERT-FIXUP** is called, z is the red node that was added.
- b. If $z.p$ is the root, then $z.p$ started out black and did not change before the call of **RB-INSERT-FIXUP**.
- c. We have already seen that properties 1, 3, and 5 hold when **RB-INSERT-FIXUP** is called.

If the tree violates property 2 (the root must be black), then the red root must be the newly added node z , which is the only internal node in the tree. Because the parent and both children of z are the sentinel, which is black, the tree does not also violate property 4 (both children of a red node are black). Thus this violation of property 2 is the only violation of red-black properties in the entire tree.

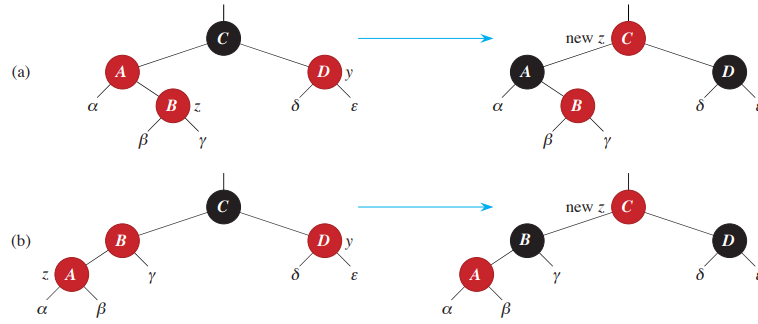
If the tree violates property 4, then, because the children of node z are black sentinels and the tree had no other violations prior to z being added, the violation must be because both z and $z.p$ are red. Moreover, the tree violates no other red-black properties.

Maintenance: There are six cases within the **while** loop, but we'll examine only the three cases in lines 3–15, when node z 's parent $z.p$ is a left child of z 's grandparent $z.p.p$. The proof for lines 17–29 is symmetric. The node $z.p.p$ exists, since by part (b) of the loop invariant, if $z.p$ is the root, then $z.p$ is black. Since **RB-INSERT-FIXUP** enters a loop iteration only if $z.p$ is red, we know that $z.p$ cannot be the root. Hence, $z.p.p$ exists.

Case 1 differs from cases 2 and 3 by the color of z 's uncle y . Line 3 makes y point to z 's uncle $z.p.p.right$, and line 4 tests y 's color. If y is red, then case 1 executes. Otherwise, control passes to cases 2 and 3. In all three cases, z 's grandparent $z.p.p$ is black, since its parent $z.p$ is red, and property 4 is violated only between z and $z.p$.

Case 1: z 's uncle y is red

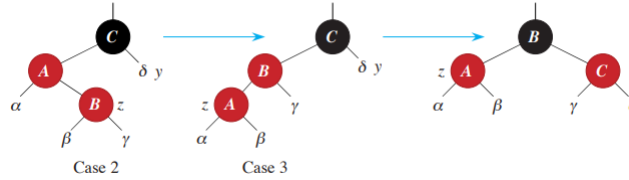
The figure below shows the situation for case 1 (lines 5–8), which occurs when both $z.p$ and y are red. Because z 's grandparent $z.p.p$ is black, its blackness can transfer down one level to both $z.p$ and y , thereby fixing the problem of z and $z.p$ both being red. Having had its blackness transferred down one level, z 's grandparent becomes red, thereby maintaining property 5. The **while** loop repeats with $z.p.p$ as the new node z , so that the pointer z moves up two levels in the tree.



Case 2: z 's uncle y is black and z is a right child

Case 3: z 's uncle y is black and z is a left child

In cases 2 and 3, the color of z 's uncle y is black. We distinguish the two cases, which assume that z 's parent $z.p$ is red and a left child, according to whether z is a right or left child of $z.p$. Lines 11–12 constitute case 2, which is shown in Figure 13.6 together with case 3. In case 2, node z is a right child of its parent. A left rotation immediately transforms the situation into case 3 (lines 13–15), in which node z is a left child. Because both z and $z.p$ are red, the rotation affects neither the black-heights of nodes nor property 5. Whether case 3 executes directly or through case 2, z 's uncle y is black, since otherwise case 1 would have run. Additionally, the node $z.p.p$ exists, since we have argued that this Node existed at the time that lines 2 and 3 were executed, and after moving z up one level in line 11 and then down one level in line 12, the identity of $z.p.p$ remains unchanged. Case 3 performs some color changes and a right rotation, which preserve property 5. At this point, there are no longer two red nodes in a row. The **while** loop terminates upon the next test in line 1, since $z.p$ is now black.



7.3.3 Deletion

Like the other basic operations on an n -node red-black tree, deletion of a node takes $O(\lg n)$ time. Deleting a node from a red-black tree is more complicated than inserting a node.

First, we discuss the *transplant* procedure. Red-Black Trees need a transplant function as part of their removal operation because it simplifies the process of replacing one subtree with another. This is particularly important for maintaining the binary search tree (BST) property and handling parent-child relationships efficiently during node deletion.

In a Red-Black Tree (and in a general BST), node removal often requires replacing the node being removed (u) with one of its subtrees (v). For example:

- If the node being removed has one child, that child replaces the node.
- If the node being removed has two children, its successor (or predecessor) replaces the node.

The transplant function encapsulates this replacement logic, making the process cleaner and more modular.

When replacing a node, the parent pointers of the replacement node and the replaced node's parent need to be updated. Manually handling these pointers can become error-prone, especially in a tree with many cases like a Red-Black Tree. The transplant function handles this pointer management consistently.

```
1  transplant(T,u,v) {
2      if (u.p == T.nil)
3          T.root = v
4      else if (u == u.p.left)
5          u.p.left = v
6      else u.p.right = v
7
8      v.p = u.p
9  }
```

The procedure **RB-DELETE** on the next page is like the **TREE-DELETE** procedure, but with additional lines of pseudocode. The additional lines deal with nodes x and y that may be involved in violations of the red-black properties. When the node z being deleted has at most one child, then y will be z . When z has two children, then, as in **TREE-DELETE**, y will be z 's successor, which has no left child and moves into z 's position in the tree. Additionally, y takes on z 's color. In either case, node y has at most one child: node x , which takes y 's place in the tree. (Node x will be the sentinel $T.nil$ if y has no children.) Since node y will be either removed from the tree or moved within the tree, the procedure needs to keep track of y 's original color. If the red-black properties might be violated after deleting node z , **RB-DELETE** calls the auxiliary procedure **RB-DELETE-FIXUP**, which changes colors and performs rotations to restore the red-black properties.

```

1  RB-DELETE(T,z) {
2      y = z
3      y-original-color = y.color
4      if (z.left == T.nil)
5          x = z.right
6          RB-TRANSPLANT(T,z,z.right) // replace z by its right
↪ child
7      else if (z.right == T.nil)
8          x = z.left
9          RB-TRANSPLANT(T,z,z.left) // replace z by its left child
10     else
11         y = MINIMUM(z.right) // y is z's successor
12         y-original-color = y.color
13         x = y.right
14         if (y != z.right) // is y farther down the tree?
15             RB-TRANSPLANT(T,y,y.right) // replace y by its
↪ right child
16             y.right = z.right // z's right child becomes
17             y.right.p = y // y's right child
18         else x.p = y // in case x is T.nil
19         RB-TRANSPLANT(T,z,y) // replace z by its successor y
20         y.left = z.left // and give z's left child to y,
21         y.left.p = y // which had no left child
22         y.color = z.color
23         if (y-original-color == BLACK) // if any red-black
↪ violations occurred,
24             RB-DELETE-FIXUP(T,x) // Correct them
25     }

```

Because node y 's color might change, the variable y -original-color stores y 's color before any changes occur. Lines 2 and 10 set this variable immediately after assignments to y . When node z has two children, then nodes y and z are Distinct. In this case, line 17 moves y into z 's original position in the tree (that is, z 's location in the tree at the time **RB-DELETE** was called), and line 20 gives y the same color as z . When node y was originally black, removing or moving it could cause violations of the red-black properties, which are corrected by the call of **RB-DELETE-FIXUP** in line 22.

As discussed, the procedure keeps track of the node x that moves into node y 's original position at the time of call. The assignments in lines 4, 7, and 11 set x to point to either y 's only child or, if y has no children, the sentinel $T.nil$.

Since node x moves into node y 's original position, the attribute $x.p$ must be set correctly. If node z has two children and y is z 's right child, then y just moves into z 's position, with x remaining a child of y . Line 12 checks for this case. Although you might think that setting $x.p$ to y in line 16 is unnecessary since x is a child of y , the call of **RB-DELETE-FIXUP** relies on $x.p$ being y even if x is $T.nil$. Thus, when z has two children and y is z 's right child, executing line 16 is necessary if y 's right child is $T.nil$, and otherwise it does not change anything.

Finally, if node y was black, one or more violations of the red-black properties might arise. The call of **RB-DELETE-FIXUP** in line 22 restores the red-black properties. If y was red, the red-black properties still hold when y is removed or moved, for the following reasons:

1. No black-heights in the tree have changed

2. No red nodes have been made adjacent. If z has at most one child, then y and z are the same node. That node is removed, with a child taking its place. If the removed node was red, then neither its parent nor its children can also be red, so moving a child to take its place cannot cause two red nodes to become adjacent. If, on the other hand, z has two children, then y takes z 's place in the tree, along with z 's color, so there cannot be two adjacent red nodes at y 's new position in the tree. In addition, if y was not z 's right child, then y 's original right child x replaces y in the tree. Since y is red, x must be black, and so replacing y by x cannot cause two red nodes to become adjacent.
3. Because y could not have been the root if it was red, the root remains black.

If node y was black, three problems may arise, which the call of **RB-DELETE-FIXUP** will remedy. First, if y was the root and a red child of y became the new root, property 2 is violated. Second, if both x and its new parent are red, then a violation of property 4 occurs. Third, moving y within the tree causes any simple path that previously contained y to have one less black node. Thus, property 5 is now violated by any ancestor of y in the tree. We can correct the violation of property 5 by saying that when the black node y is removed or moved, its blackness transfers to the node x that moves into y 's original position, giving x an “extra” black. That is, if we add 1 to the count of black nodes on any simple path that contains x , then under this interpretation, property 5 holds. But now another problem emerges: node x is neither red nor black, thereby violating property 1. Instead, node x is either “doubly black” or “red-and-black,” and it contributes either 2 or 1, respectively, to the count of black nodes on simple paths containing x . The color attribute of x will still be either **RED** (if x is red-and-black) or **BLACK** (if x is doubly black). In other words, the extra black on a node is reflected in x 's pointing to the node rather than in the color attribute.

The procedure **RB-DELETE-FIXUP** on the next page restores properties 1, 2, and 4

The goal of the while loop in lines 1-43 is to move the extra black up the tree until

1. x points to a red-and-black node, in which case line 44 colors x (singly) black;
2. x points to the root, in which case the extra black simply vanishes; or
3. having performed suitable rotations and recolorings, the loop exits.

Like **RB-INSERT-FIXUP**, the **RB-DELETE-FIXUP** procedure handles two symmetric situations: lines 3-22 for when node x is a left child, and lines 24-43 for when x is a right child. Our proof focuses on the four cases shown in lines 3-22


```
1  RB-DELETE-FIXUP(T, x)
```

7.3.4 C++ implementation

7.3.4.1 The node structure

```
1  struct node {
2      enum color {BLACK, RED};
3
4      int data=0;
5      bool color;
6      node* left, *right, *parent;
7
8      node() = default;
9      node(int data) : data(data) {}
10 };
```

A red-black bst node in this example has three pointers, left, right, and parent. It also has a color and a data field. The color enum defines the colors. Black is zero, red is one.

7.3.4.2 The tree class and defining nil

```
1  class tree {
2      node* root;
3      node* nil;
4
5      tree() {
6          nil = new node();
7          nil->color = node::BLACK;
8          nil->left = nil;
9          nil->right = nil;
10         nil->parent = nil;
11         root = nil;
12     }
13 };
```

We create two nodes, a root node and the nil node. In the constructor, we create the nil node, and set its attributes. An empty tree is one whose root is nil. Note that all nodes that would be nullptrs in a standard binary search tree now become the nil node in the red-black tree

7.3.4.3 Rotation methods

```
1  void left_rotate(node* x) {
2      node* y = x->right;
3      x->right = y->left;
4
5      if (y->left != nil) {
6          y->left->parent = x;
7      }
8
9      y->parent = x->parent;
10     if (x->parent == nil) {
11         root = y;
12     } else if (x == x->parent->left) {
13         x->parent->left = y;
14     } else {
15         x->parent->right = y;
16     }
17     y->left = x;
18     x->parent = y;
19 }
```

`right_rotate` is achieved by replacing `left` with `right`, and `right` with `left` in `left_rotate`.

7.3.4.4 Insert

```
1  void insert(int element) {
2      node* z = new node(element);
3      z->left = nil;
4      z->right = nil;
5      z->color = node::RED;
6
7      node* x = root;
8      node* y = nil;
9
10     while (x != nil) {
11         y = x;
12         if (z->data < x->data) {
13             x = x->left;
14         } else {
15             x = x->right;
16         }
17     }
18
19     z->parent = y;
20     if (y == nil) {
21         root = z;
22     } else if (z->data < y->data) {
23         y->left = z;
24     } else {
25         y->right = z;
26     }
27
28     insert_fixup(z);
29 }
```

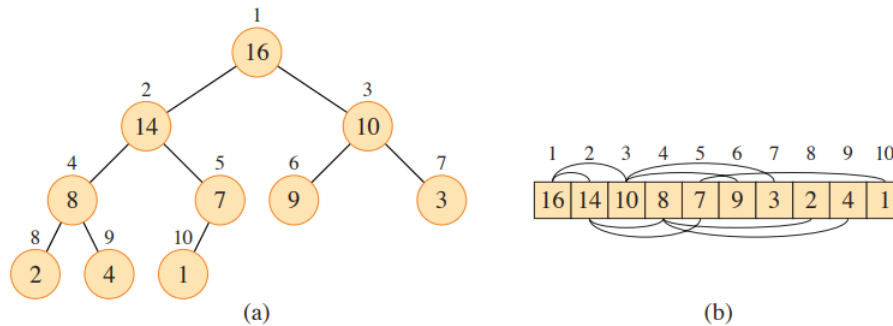
7.3.4.5 Insert fixup

```
1  void insert_fixup(node* z) {
2      while (z->parent->color == node::RED) {
3          if (z->parent == z->parent->parent->left) {
4              node* y = z->parent->parent->right;
5              if (y->color == node::RED) {
6                  z->parent->color = node::BLACK;
7                  y->color = node::BLACK;
8                  z->parent->parent->color = node::RED;
9                  z = z->parent->parent;
10             } else {
11                 if (z == z->parent->right) {
12                     z = z->parent;
13                     left_rotate(z);
14                 }
15                 z->parent->color = node::BLACK;
16                 z->parent->parent->color = node::RED;
17                 right_rotate(z->parent->parent);
18             }
19         } else {
20             node* y = z->parent->parent->left;
21             if (y->color == node::RED) {
22                 z->parent->color = node::BLACK;
23                 y->color = node::BLACK;
24                 z->parent->parent->color = node::RED;
25                 z = z->parent->parent;
26             } else {
27                 if (z == z->parent->left) {
28                     z = z->parent;
29                     right_rotate(z);
30                 }
31                 z->parent->color = node::BLACK;
32                 z->parent->parent->color = node::RED;
33                 left_rotate(z->parent->parent);
34             }
35         }
36     }
37     root->color = node::BLACK;
38 }
```

Heaps and Priority Queues (Zero based)

The (*binary*) *heap* data structure is an array object that we can view as a nearly complete binary tree (see Section B.5.3), as shown in Figure 6.1. Each node of the tree corresponds to an element of the array. The tree is completely filled on all levels except possibly the lowest, which is filled from the left up to a point. An array $A[0 : n - 1]$ that represents a heap is an object with an attribute $A.\text{heap-size}$, which represents how many elements in the heap are stored within array A . That is, although $A[0 : n - 1]$ may contain numbers, only the elements in $A[0 : A.\text{heap-size} - 1]$, where $0 \leq A.\text{heap-size} \leq n$, are valid elements of the heap. If $A.\text{heap-size} = 0$, then the heap is empty. The root of the tree is $A[0]$, and given the index i of a node,

The (*binary*) *heap* data structure is an array object that we can view as a nearly complete binary tree (see Section B.5.3), as shown in Figure 6.1. Each node of the tree corresponds to an element of the array. The tree is completely filled on all levels except possibly the lowest, which is filled from the left up to a point. An array $A[1 : n]$ that represents a heap is an object with an attribute $A.\text{heap-size}$, which represents how many elements in the heap are stored within array A . That is, although $A[1 : n]$ may contain numbers, only the elements in $A[1 : A.\text{heap-size}]$, where $0 \leq A.\text{heap-size} \leq n$, are valid elements of the heap. If $A.\text{heap-size} = 0$, then the heap is empty. The root of the tree is $A[1]$, and given the index i of a node,



there's a simple way to compute the indices of its parent, left child, and right child with the one-line procedures PARENT, LEFT, and RIGHT.

```
1  int parent(i) {
2      return (i+1)/2 // Floor division
3  }
4
5  int left(i) {
6      return 2i + 1
7  }
8
9  int right(i) {
10     return 2i + 2:
11 }
```

For one based indexing, we would have $\frac{i}{2}$, $2i$, and $2i + 1$

On most computers, the **LEFT** procedure can compute $2i$ in one instruction by simply shifting the binary representation of i left by one bit position. Similarly, the **RIGHT** procedure can quickly compute $2i + 1$ by shifting the binary representation of i left by one bit position and then adding 1. The **PARENT** procedure can compute $\lfloor \frac{i}{2} \rfloor$ by shifting i right one bit position. Good implementations of heapsort often implement these procedures as macros or inline procedures.

8.1 Max and Min heaps

There are two kinds of binary heaps: max-heaps and min-heaps. In both kinds, the values in the nodes satisfy a [heap property](#), the specifics of which depend on the kind of heap. In a [max-heap](#), the [max-heap property](#) is that for every node i other than the root,

$$A[\text{parent}(i)] \geq A[i].$$

That is, the value of a node is at most the value of its parent. Thus, the largest element in a max-heap is stored at the root, and the subtree rooted at a node contains values no larger than that contained at the node itself. A [min-heap](#) is organized in the opposite way: the [min-heap property](#) is that for every node i other than the root,

$$A[\text{parent}(i)] \leq A[i].$$

The smallest element in a min-heap is at the root.

8.2 Heapify an array

Heapify is a process used to maintain the heap property in a binary heap, either a max-heap or a min-heap. The heap property ensures that for every node in the heap:

- In a max-heap, the value of each node is greater than or equal to the values of its children.
- In a min-heap, the value of each node is less than or equal to the values of its children.

Heapify is typically used when you have an unsorted array and you want to turn it into a valid heap, or after inserting or deleting an element in a heap to restore the heap property.

The bottom-up approach to turning an array into a heap is also called the heapify process. This method is efficient for building a heap from an unordered array and runs in $O(n)$ time, which is faster than building the heap using successive insertions ($O(n \lg n)$).

We start from the first non-leaf node, which is at position $\frac{n-1}{2}$, where n is the size of the array.

Traverse all non-leaf nodes from the last one to the root (from right to left in the array).

For each non-leaf node, apply the sift down operation (also called percolate down), which ensures that the subtree rooted at this node satisfies the heap property.

Percolate Down Operation:

1. Compare the node with its children.
2. If the heap property is violated (for example, in a max-heap, the node is smaller than one of its children), swap the node with the largest child (for max-heap) or smallest child (for min-heap).
3. Repeat this process down the subtree until the heap property is restored or the node becomes a leaf.

8.3 Min-heap in c++

```
1 void min_heapify(int arr[], int n, int i) {
2     int smallest = i;    // Initialize smallest as root
3     int left = 2 * i + 1; // Left child index
4     int right = 2 * i + 2; // Right child index
5
6     // If left child is smaller than the root
7     if (left < n && arr[left] < arr[smallest])
8         smallest = left;
9
10    // If right child is smaller than the smallest so far
11    if (right < n && arr[right] < arr[smallest])
12        smallest = right;
13
14    // If the smallest is not the root
15    if (smallest != i) {
16        std::swap(arr[i], arr[smallest]);
17
18        // Recursively heapify the affected subtree
19        min_heapify(arr, n, smallest);
20    }
21 }
22
23 void build_heap(int arr[], int n) {
24     // Start from the last non-leaf node and heapify each node
25     for (int i = (n - 1) / 2; i >= 0; --i) {
26         min_heapify(arr, n, i);
27     }
28 }
```

The `min_heapify` function is designed to maintain the min-heap property for a subtree rooted at a given index i . The process begins by identifying the left and right children of the node at index i . The function then compares the current node with its children to find the smallest value. If one of the children is smaller than the current node, a swap occurs between the node and the smallest child. After the swap, the function is recursively called on the affected child to ensure that the min-heap property is maintained further down the subtree.

8.4 Max-heap in c++

```
1 void heapify(int arr[], int n, int i) {
2     int largest = i;
3     int left = 2 * i + 1;
4     int right = 2 * i + 2;
5
6     if (left < n && arr[left] > arr[largest]) {
7         largest = left;
8     }
9
10    if (right < n && arr[right] > arr[largest]) {
11        largest = right;
12    }
13
14    if (largest != i) {
15        std::swap(arr[i], arr[largest]);
16        heapify(arr, n, largest);
17    }
18 }
19
20
21 void build_heap(int arr[], int n) {
22     for (int i=n-1/2; i>=0; --i) {
23         heapify(arr, n, i);
24     }
25 }
```

8.5 Percolating

Percolating in heaps refers to adjusting a node's position to maintain the heap property, which can occur in two directions: percolate up or percolate down.

Note that the percolate direction refers to the direction in which we move a node to get its correct locating in the heap.

If we are comparing a given node to its children, its percolate down. If we compare with its parent, its percolate up.

8.5.1 Percolate up

Used when adding a new element to the heap (usually at the end). The new element is moved up the heap by comparing it with its parent node. If it violates the heap property (e.g., it's smaller than its parent in a min-heap), it swaps with the parent. This continues until the node is correctly positioned.

```
1 void percUp(vector<int>& v, int i) {  
2     while (i>0) {  
3         int parent = (i-1)/2;  
4         if (v[i] > v[parent]) {  
5             swap(v[i], v[parent]);  
6             i = parent;  
7         } else break; // In the correct place  
8     }  
9 }
```

8.5.2 Percolate down

Used when removing the root element (like in a heap deletion). The last element is moved to the root and then "percolates down" by swapping with the smaller child (in a min-heap) if it violates the heap property. This process repeats until the node is correctly placed.

```
1 void percdwn(vector<int>& v, int n, int i) {
2     // Assume the largest is the current node
3     int largest = i;
4     // Get index of both children
5     int left = 2*i+1;
6     int right = 2*i+2;
7
8     // If the left child exists, and is larger, update largest
9     if (left < n && v[left] > v[largest]) {
10         largest = left;
11     }
12     // If the right child exists, and is larger, update largest
13     if (right < n && v[right] > v[largest]) {
14         largest = right;
15     }
16
17     // If the largest was changed, swap. Then call percdwn on the
    ↪ node that had the largest value.
18     if (largest != i) {
19         swap(v[i], v[largest]);
20         percdwn(v,n,largest);
21     }
22 }
```

Because this function takes a node and compares it with its children, the value in the node moves down. Hence, percolate down.

8.6 Inserting into a heap

In a top-down approach to inserting into a min-heap, you maintain the heap property while inserting a new element. Specifically, you start by placing the new element at the bottom of the heap (in the next available position), and then you "bubble up" (also known as "heapify up" or "percolate up") to restore the heap property.

- **Insert the new element at the bottom:** Add the new element in the first available position (the next available spot in the array representation, which maintains a complete binary tree structure)
- **Bubble up (heapify up):** Compare the newly inserted element with its parent.

If the new element is smaller than its parent (which violates the min-heap property), swap the two.

Continue this process (i.e., compare with the next parent) until:

- The new element is larger than or equal to its parent, or
 - The element reaches the root of the heap.
- The process terminates when the new element is in a position where it is larger than its parent or becomes the root.

The code for a min-heap insert is as follows

```
1  void percUp(vector<int>& v, int i) {
2      int parent = (i-1)/2;
3
4      while (i>0 && v[i] < v[parent]) {
5          std::swap(v[i], v[parent]);
6
7          i = parent;
8      }
9  }
10
11 void heapInsert(vector<int>& v, int element) {
12     v.push_back(element);
13     int index = v.size() - 1;
14     percUp(v, index);
15 }
```

8.7 Removing the root

Removing an element from a min-heap (typically the root, i.e., the minimum element) involves maintaining the min-heap property after removal. The most common removal operation is to delete the root (the smallest element in a min-heap). The process of removal involves the following steps:

1. **Replace the root with the last element:** The root (at index 0) is the element to be removed, so we replace it with the last element in the heap (this ensures the complete binary tree property is maintained).
2. **Remove the last element:** After swapping, the last element is removed from the heap (usually by reducing the heap's size).
3. **Heapify down (percolate down):** Starting from the root, compare the new root element with its children.

If the root is larger than any of its children, swap it with the smaller child (to maintain the min-heap property).

Continue this process until the heap property is restored (i.e., the element is smaller than both children or it reaches a leaf node).

```
1  void removeRoot(vector<int>& v) {  
2      std::swap(v[0], v[v.size()-1]);  
3      v.pop_back();  
4      min_heapify(v,v.size(), 0);  
5  }
```

8.8 Removing an arbitrary node

To remove an arbitrary node from a min-heap (not just the root), the process is slightly more complex than simply removing the root

1. **Replace the node to be deleted with the last element:** Swap the node to be deleted with the last element in the heap. This ensures that the complete binary tree property is maintained.
2. **Remove the last element:** After the swap, the last element (now at the position of the deleted node) is removed from the heap.
3. **Restore the heap property:** After the swap, the heap property might be violated both upwards and downwards. So, depending on the value of the swapped element, either heapify up or heapify down from the position where the node was deleted.
 - Heapify up if the swapped element is smaller than its parent.
 - Heapify down if the swapped element is larger than its children

We continue this process until there are no more nodes that match the passed element. The c++ code for a min-heap erase function is as follows.

```
1 void erase(vector<int>& v, int element) {
2     bool found = false;
3     while (!found) {
4         found = false;
5         int i=0;
6         for (;i<(int)v.size(); ++i) {
7             if (v[i] == element) {
8                 found = true;
9                 std::swap(v[i], v[v.size()-1]);
10                v.pop_back();
11
12                if (v[i] > v[(i-1)/2]) {
13                    min_heapify(v, v.size(), i);
14                } else {
15                    percUp(v, i);
16                }
17            }
18        }
19        if (i == v.size()) break;
20    }
21 }
```

8.9 Priority queues

Heaps are used to implement priority queues because they efficiently maintain the highest (or lowest) priority element at the root. In a max-heap, the largest element is always at the top, while in a min-heap, the smallest element is. This allows for quick access to the highest priority element in $O(1)$ time. Insertion and removal (reordering) of elements take $O(\lg n)$ time, making heaps an optimal choice for priority queues where these operations need to be fast and frequent.

8.9.1 Interface

- **Insert (or Enqueue):** Adds an element to the priority queue based on its priority.
- **Pop:** Retrieve and remove the item with the highest priority (root of the heap)
- **Top:** Retrieve the item with the highest priority
- **Size:** Get the number of items in the queue
- **Empty:** Checks if the priority queue has any elements.

```
1  class priority_queue {
2      vector<int> heap; // or any random access container
3
4      void percDown(int n, int index);
5      void PercUp(int index)
6
7  public:
8      void insert(int element);
9      int pop();
10     int top();
11     size_t size();
12     bool empty();
13 };
```

Percolate down, up, insert, and pop can be found in the heap examples above.

8.9.2 Insert, pop, and top

These are simple methods to implement.

```
1  // Insert into the heap
2  void insert(int element) {
3      heap.push_back(element);
4      int index = heap.size() - 1;
5      percUp(index);
6  }
7
8  // Retrieve and remove the root
9  int pop() {
10     if (heap.empty()) return -1;
11
12     int ret = heap[0];
13     std::swap(heap[0], heap[heap.size()-1]);
14     heap.pop_back();
15     percDown((int)heap.size(), 0);
16     return ret;
17 }
18
19 // Retrieve the root
20 int top() {
21     return heap[0];
22 }
```

8.9.3 Size and Empty

```
1  size_t size() {  
2      return (size_t)heap.size();  
3  }  
4  
5  bool empty() {  
6      return heap.empty();  
7  }
```

Sorting

Sorting is a fundamental concept in data structures and algorithms where elements in a collection are arranged in a specific order, typically ascending or descending. Sorting enables efficient data access, searching, and data organization in applications. There are various sorting algorithms, each with different characteristics and efficiency, such as time and space complexity.

9.1 Bubble, selection, insertion

The three simplest sorting algorithms are bubble sort, selection sort, and insertion sort. They are all $O(n^2)$ comparison based sorting algorithms.

9.1.1 Bubble sort

Bubble Sort repeatedly compares adjacent elements and swaps them if they are in the wrong order, "bubbling" the largest unsorted element to the end of the list on each pass.

```
1  void bubbleSort(vector<int>& v, int n) {
2      bool swapped = true;
3      while (swapped) {
4          swapped = false;
5          for (int i=0; i<n-1; ++i) {
6              if (v[i+1] < v[i]) {
7                  swap(v[i], v[i+1]);
8                  swapped=true;
9              }
10         }
11         --n;
12     }
13 }
```

By decrementing n after each complete pass, you avoid unnecessary comparisons in already sorted parts of the vector, making the algorithm slightly more efficient without affecting the $O(n^2)$ time complexity.

- **Outer Loop:** The `while` loop continues as long as swaps are happening, so it terminates early if the list becomes sorted before completing all n passes.
- **Inner Loop:** The `for` loop compares adjacent elements and swaps them if they're out of order. After each pass, the largest unsorted element "bubbles up" to the correct position, so the next pass can ignore the last element.

9.1.1.1 Complexity

Worst Case: $O(n^2)$

In the worst case (a reverse-sorted list), Bubble Sort requires n passes, and each pass involves up to $n - 1$ comparisons.

Average Case: $O(n^2)$

On average, Bubble Sort still performs $O(n^2)$ comparisons due to repeated passes through the entire array.

Best Case: $O(n)$

If the list is already sorted, Bubble Sort can stop early if no swaps occur during a pass (often implemented with a flag to detect this). In this case, it only takes one pass to confirm that the list is sorted.

9.1.2 Selection sort

Selection Sort sorts an array by repeatedly finding the minimum element from the unsorted part and moving it to the beginning.

1. Start at the beginning of the array.
2. For each position i , find the smallest element in the unsorted portion (from i to the end).
3. Swap this minimum element with the element at position i
4. Move to the next position and repeat until the array is sorted.

```
1  void selectionSort(vector<int>& v, int n) {  
2      for (int i=0; i<n-1; ++i) {  
3          int min = i;  
4          for (int j=i+1; j<n; ++j) {  
5              if (v[j] < v[min]) {  
6                  min = j;  
7              }  
8          }  
9          std::swap(v[min], v[i]);  
10     }  
11 }
```

9.1.2.1 Complexity

Worst Case: $O(n^2)$

Selection Sort always performs n passes, each requiring a search through the unsorted portion to find the minimum element, resulting in $O(n^2)$ comparisons.

Average Case: $O(n^2)$

Selection Sort performs the same number of comparisons regardless of the initial order of elements, as it always looks for the minimum in the unsorted portion.

Best Case: $O(n^2)$

Even if the array is already sorted, Selection Sort will still go through n passes and $O(n^2)$ comparisons, making it inefficient for nearly sorted data.

9.1.3 Insertion sort

Insertion Sort sorts an array by building a sorted portion one element at a time

1. Start with the second element, assuming the first element is trivially sorted.
2. For each element, compare it with the elements in the sorted portion to its left.
3. Shift larger elements one position right until you find the correct spot for the current element.
4. Insert the current element in its correct position.
5. Repeat until all elements are sorted.

```
1      void insertionSort(vector<int>& v, int n) {  
2          for (int i=1; i<n; ++i) {  
3              int key=v[i];  
4              int j=i-1;  
5  
6              while (j>=0 && v[j] > key) {  
7                  v[j+1] = v[j];  
8                  --j;  
9              }  
10             v[j+1] = key;  
11         }  
12     }
```

Imagine sorting a hand of playing cards:

- You pick up cards one by one.
- For each new card, you compare it with the cards in your hand from right to left, sliding them over if they're larger than the new card, until you find the correct spot to insert it.
- Each card you add to your hand stays sorted with the previous ones, just like key gets inserted in the sorted sub-array.

This method is efficient when the list is nearly sorted since fewer shifts are needed, but it becomes less efficient for randomly ordered arrays due to repeated comparisons and shifts.

9.1.3.1 Complexity

Worst Case: $O(n^2)$

In the worst case (reverse-sorted list), each element is compared with all previous elements, resulting in $O(n^2)$ comparisons and shifts.

Average Case: $O(n^2)$

Insertion Sort generally performs $O(n^2)$ operations in the average case due to repeated comparisons and shifts.

Best Case: $O(n)$

If the array is already sorted, Insertion Sort only needs to confirm each element is in place, requiring just $O(n)$ comparisons and no shifts. This makes it very efficient for nearly sorted arrays.

9.2 Heap sort

Heap Sort is a comparison-based sorting algorithm that uses a binary heap data structure, specifically a max-heap, to sort elements in ascending order (or a min-heap for descending order).

1. **Build a Max-Heap:** Convert the array into a max-heap, where each parent node is greater than its children. This ensures the largest element is at the root.
2. **Extract Maximum and Swap:** Swap the root (largest element) with the last element in the heap. This moves the largest element to its correct position in the sorted array.
3. **Heapify:** After the swap, the heap may no longer satisfy the heap property. Restore the max-heap by performing a "heapify" operation on the root to maintain the max-heap structure.
4. **Repeat:** Continue extracting the maximum element, swapping, and heapifying, reducing the heap size each time. This process sorts the array in-place.

```
1 void heapSort(vector<int>& v, int n) {  
2     heapify(v, n);  
3     int last = n-1;  
4     while (last > 0) {  
5         swap(v[0], v[last--]);  
6         percDown(v, last, 0);  
7     }  
8 }
```

The functions *heapify* and *percDown* can be found in the *max-heap in c++* section above, where they are called *build-heap* and *heapify* respectively.

9.2.1 Complexity

Heap Sort has a time complexity of $O(n \lg n)$ in the best, worst, and average cases, making it more efficient than $O(n^2)$ sorting algorithms for larger datasets

Building the max-heap from an arbitrary container takes $O(n)$ time. Re-heapify the remaining elements (restore the max-heap property) by "sifting down" the new root takes $O(\lg n)$ time. Since we do this sift operation n times, the complexity of *heapSort* is therefore $O(n \lg n)$, which makes the entire algorithm $O(n) + O(n \lg n) = O(n \lg n)$.

9.3 BST Sort

BST Sort is a sorting algorithm that leverages a Binary Search Tree (BST) to sort elements. It works by inserting all elements into a BST and then performing an in-order traversal of the tree to retrieve the elements in sorted order.

9.3.1 Insert

The insert for bst follows the standard bst insert method described in a previous section.

```
1  node* insert(node* p, int data) {
2      if (!p) return new node(data);
3
4      if (data < p->data) {
5          p->left = insert(p->left, data);
6      } else {
7          p->right = insert(p->right, data);
8      }
9
10     return p;
11 }
```

9.3.2 The inorder traversal

The inorder follows the same logic as before, but when processing each node, we insert the data member into a passed vector.

```
1  void inorder(node* p, vector<int>& v) {
2      if (!p) return;
3      inorder(p->left, v);
4      v.push_back(p->data);
5      inorder(p->right, v);
6  }
```

9.3.3 The BST sort function

```
1  void bstsort(vector<int>& v) {  
2      vector<int> sorted;  
3      node* root = nullptr;  
4  
5      for (const auto& item : v) {  
6          root = insert(root, item);  
7      }  
8      inorder(root, sorted);  
9      clear(root);  
10  
11     v = sorted;  
12 }
```

The `bstsort` function sorts a vector of integers using a Binary Search Tree (BST) approach. It starts by creating an empty vector, `sorted`, which will eventually hold the sorted elements, and initializes `root` as `nullptr`, representing an initially empty BST.

For each item in the input vector `v`, the function calls `insert` to add the item to the BST, progressively building the tree in a way that respects the BST property (left child nodes are less than the parent node, and right child nodes are greater). After constructing the BST, the function performs an in-order traversal using `inorder`, which appends each node's data to the sorted vector in ascending order.

Finally, `clear` is called to delete all nodes from the BST, freeing up dynamically allocated memory. The function then assigns `sorted` back to the input vector `v`, so `v` now contains the elements in sorted order.

BST Sort is more of a conceptual exercise in data structures than a practical sorting method because balanced tree structures (e.g., AVL trees, red-black trees) are required to ensure $O(n \lg n)$ performance consistently.

9.3.4 Inplace sorting

To perform the sort "in-place" in `btsort`, we can eliminate the extra sorted vector by writing the sorted values directly back into the input vector `v` during the in-order traversal. This approach leverages an index to keep track of the position in `v` where each next sorted value should be placed

```
1 void inorder2(node* p, vector<int>& v, int& index) {
2     if (!p) return;
3     inorder2(p->left, v, index);
4     v[index++] = p->data;
5     inorder2(p->right, v, index);
6 }
7
8 void btsort(vector<int>& v) {
9     node* root = nullptr;
10
11     for (const auto& item : v) {
12         root = insert(root, item);
13     }
14     int index = 0;
15     inorder2(root, v, index);
16     clear(root);
17 }
```

9.3.5 Complexity

- **Average Case:** $O(n \lg n)$, where n is the number of elements. This is because inserting each element in a balanced BST takes $O(\lg n)$ time
- **Worst Case:** $O(n^2)$, if the tree becomes unbalanced, like when inserting elements in sorted order into a simple BST, resulting in a degenerate tree (like a linked list).
- **Space Complexity:** $O(n)$, for storing the BST nodes.

9.4 Quick sort

Quicksort, like merge sort, applies the divide-and-conquer method

- **Divide** by partitioning (rearranging) the array $A[p : r]$ into two (possibly empty) subarrays $A[p : q - 1]$ (*the low side*) and $A[q + 1 : r]$ (*the high side*) such that each element in the low side of the partition is less than or equal to the *pivot* $A[q]$, which is, in turn, less than or equal to each element in the high side. Compute the index q of the pivot as part of this partitioning procedure.
- **Conquer** by calling quicksort recursively to sort each of the subarrays $A[p : q - 1]$ and $A[q + 1 : r]$.
- **Combine** by doing nothing: because the two subarrays are already sorted, no work is needed to combine them. All elements in $A[p : q - 1]$ are sorted and less than or equal to $A[q]$, and all elements in $A[q + 1 : r]$ are sorted and greater than or equal to the pivot $A[q]$. The entire subarray $A[p : r]$ cannot help but be sorted!

The QUICKSORT procedure implements quicksort. To sort an entire n -element array $A[1 : n]$, the initial call is QUICKSORT($A, 1, n$).

9.4.1 Partitioning

The key to the algorithm is the PARTITION procedure, which rearranges the subarray $A[p : r]$ in place, returning the index of the dividing point between the two sides of the partition

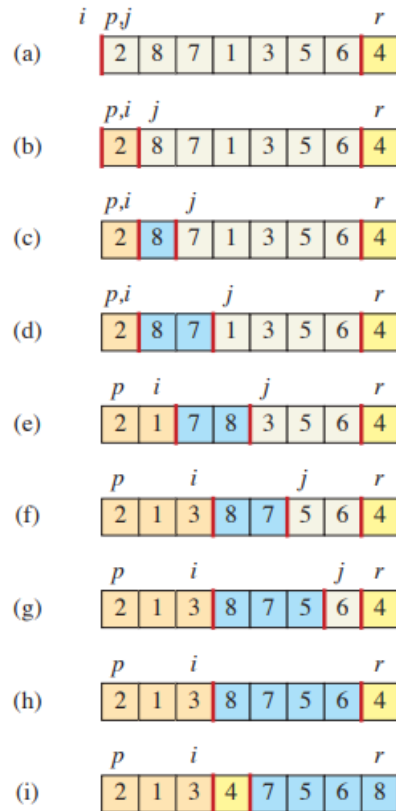
PARTITION always selects the element $A[r]$ (last element) as the pivot.

```
1  PARTITION(A, p,r)
2      x = A[r]                // the pivot
3      i = p-1                 // highest index into the low side
4      for j = p to r - 1     // process each element other than
        ↪ the pivot
5          if A[j] <= x        // does this element belong on the
        ↪ low side?
6              ++i             // index of a new slot in the
        ↪ low side
7          swap(A[i], A[j])    // put this element there
8      swap(A[i+1], A[r])      // pivot goes just to the right of
        ↪ the low side
9      return i + 1           // new index of the pivot
```

At the beginning of each iteration of the loop of lines 3–6, for any array index k , the following conditions hold:

1. if $p \leq k \leq i$, then $A[k] \leq x$ (the tan region of Figure 7.2);
2. if $i + 1 \leq k \leq j - 1$, then $A[k] > x$ (the blue region);
3. if $k = r$, then $A[k] = x$ (the yellow region).

Note: an invariant is a property or condition that remains true throughout the execution of a program, particularly during specific phases or loops. In other words, an invariant is something that holds steady or stays consistent at key points in an algorithm, typically at the beginning or end of each loop iteration



```

1  int partition(vector<int>& v, int p, int r) {
2      int x = v[r];
3      int i=p-1;
4
5      for (int j=p; j<r; ++j) {
6          if (v[j] <= x) {
7              ++i;
8              swap(v[j], v[i]);
9          }
10     }
11     swap(v[i+1], v[r]);
12     return i+1;
13 }

```

9.4.2 The quick sort procedure

```

1 void qsort(vector<int>& v, int start, int end) {
2     if (start < end) {
3         int pivot = partition(v, start, end);
4         qsort(v, start, pivot-1);
5         qsort(v, pivot+1, end);
6     }
7 }

```

In each recursive call to `qsort`, the function first checks if the start index is less than the end index. This condition serves as the base case for the recursion, ensuring that the function will terminate when the segment of the array being processed is either empty or has only one element (both of which are already sorted by default). If start is not less than end, the function simply returns without making further recursive calls, signaling that this segment is already sorted.

When start is less than end, the function calls `partition`, passing in the vector and the current segment boundaries (start and end). The partition function typically selects a pivot element, rearranges elements in the segment so that all elements less than or equal to the pivot appear on the left side, and all elements greater than the pivot appear on the right side. `partition` then returns the index of the pivot element after rearranging. This index is crucial because it marks the point in the array where the pivot element is correctly positioned in the sorted order.

After the array is partitioned around the pivot, the `qsort` function calls itself recursively on the two subarrays divided by the pivot. The first recursive call processes the left subarray, from start to pivot - 1, containing elements less than or equal to the pivot. The second recursive call processes the right subarray, from pivot + 1 to end, containing elements greater than the pivot. Each of these recursive calls will, in turn, partition their respective subarrays and further divide and sort them, following the same logic.

Through this recursive process, each segment of the array is progressively divided and sorted around pivot elements, until every element is in its correct position. This "divide and conquer" approach ensures that, by the end of the process, the entire array is sorted in ascending order. The function achieves efficiency by reducing the problem size with each recursive call, ultimately sorting the array in an average time complexity of $O(n \lg n)$, making QuickSort one of the fastest sorting algorithms in practice.

9.4.3 A better partition

In the example above, we always selected the last element as the pivot. We could instead select the middle element.

```
1  int partition2(std::vector<int>& v, int p, int r) {
2      int middle = p + (r - p) / 2; // Calculate the middle index
3      std::swap(v[middle], v[r]);   // Move the middle element to
    ↪ the end temporarily
4
5      int x = v[r];                  // Use the last element
    ↪ (original middle) as pivot
6      int i = p - 1;
7
8      for (int j = p; j < r; ++j) {
9          if (v[j] <= x) {
10             ++i;
11             std::swap(v[j], v[i]);
12         }
13     }
14     std::swap(v[i + 1], v[r]);      // Place the pivot in its
    ↪ correct position
15     return i + 1;                  // Return the pivot index
16 }
```

9.4.4 Median of three partition

To use the median-of-three approach for selecting the pivot in QuickSort, we need to compare three elements in the array — typically the first, middle, and last elements of the current subarray — and select the median (the middle value of these three) as the pivot. This approach helps improve the performance of QuickSort by reducing the chances of poor pivot choices in partially sorted or skewed data.

- **Identify the Three Elements:** We find the indices of the first element, the middle element, and the last element of the current subarray.
- **Calculate the Median of Three:** We choose the median value among these three elements and use it as the pivot.
- **Swap the Pivot to the End:** To keep the partition logic the same, we swap the median-of-three pivot element with the last element, so we can proceed with the same partitioning logic as before.

```

1  int partition3(std::vector<int>& v, int p, int r) {
2      int middle = p + (r-p) / 2;
3
4      // Determine the median of v[p], v[middle], and v[r]
5      if ((v[p] > v[middle]) != (v[p] > v[r])) {
6          std::swap(v[p], v[r]); // Median is at v[p]
7      } else if ((v[middle] > v[p]) != (v[middle] > v[r])) {
8          std::swap(v[middle], v[r]); // Median is at v[middle]
9      }
10     // If v[r] is the median, no swap is needed since we want it
    ↪ as the pivot
11
12     int x = v[r]; // The pivot is now the median of the three
    ↪ values
13     int i = p - 1;
14
15     // Standard partition logic
16     for (int j = p; j < r; ++j) {
17         if (v[j] <= x) {
18             ++i;
19             std::swap(v[j], v[i]);
20         }
21     }
22     std::swap(v[i + 1], v[r]); // Place the pivot in its correct
    ↪ position
23     return i + 1; // Return the pivot index
24 }

```


9.4.5 Quicksort with iterators

```
1  template<typename ForwardIt>
2  ForwardIt partition4(ForwardIt p, ForwardIt r) {
3      auto middle = p + std::distance(p, r) / 2; // Calculate the
    ↪ middle iterator
4      std::iter_swap(middle, r - 1); // Move the middle element
    ↪ to the end temporarily
5
6      auto pivot = *(r - 1); // Use the last element (original
    ↪ middle) as pivot
7      auto i = p; // Boundary for elements less than or equal to
    ↪ the pivot
8
9      for (auto j = p; j != r - 1; ++j) { // Iterate up to the
    ↪ pivot position
10         if (*j <= pivot) {
11             std::iter_swap(j, i);
12             ++i;
13         }
14     }
15     std::iter_swap(i, r - 1); // Place the pivot in its correct
    ↪ position
16     return i; // Return the pivot iterator
17 }
18
19 template <typename ForwardIt>
20 void qsort2(ForwardIt start, ForwardIt end) {
21     if (std::distance(start, end) > 1) { // Correct base case:
    ↪ at least two elements
22         auto pivot = partition4(start, end);
23         qsort2(start, pivot); // Sort elements before the pivot
24         qsort2(pivot + 1, end); // Sort elements after the pivot
25     }
26 }
```

Multi-way (m-way) search trees

A multiway search tree is a type of search tree where each node can have more than two children, unlike binary search trees, which are limited to two children per node. Multiway search trees generalize binary search trees by allowing nodes to hold multiple keys and have multiple pointers to child nodes, making them well-suited for managing large amounts of data in a balanced structure.

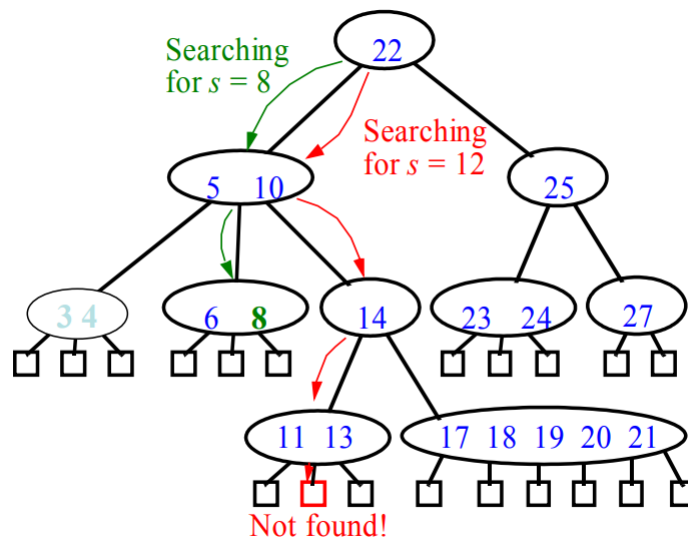
An m -Way tree of order m , each node contains a maximum of $m-1$ elements and m children.

The goal of an m -Way search tree of height h is to achieve $O(h)$ accesses for an insert, delete, or retrieval operation. Therefore, it ensures that the height h is close to $\log_m(n+1)$.

Each internal node of a multi-way search tree T :

1. Has at least two children
2. stores a collection of items of the form (k, x) , where k is a key and x is an element
3. contains $d-1$ items, where d is the number of children

Children of each internal node are "between" items. All keys in the subtree rooted at the child fall between keys of those items



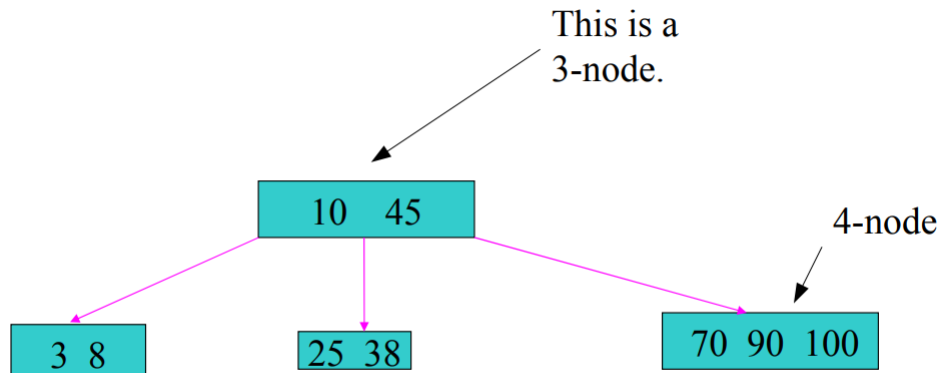
10.1 Multi-way Searching

Similar to binary searching, where if $s < k_1$, search the leftmost child. If $s > k_{d-1}$, search the rightmost child. But what if $d > 2$? Simply Find two keys k_{i-1} and k_i between which s falls, and search the child v_i .

10.2 2-4 (2-3-4) Trees

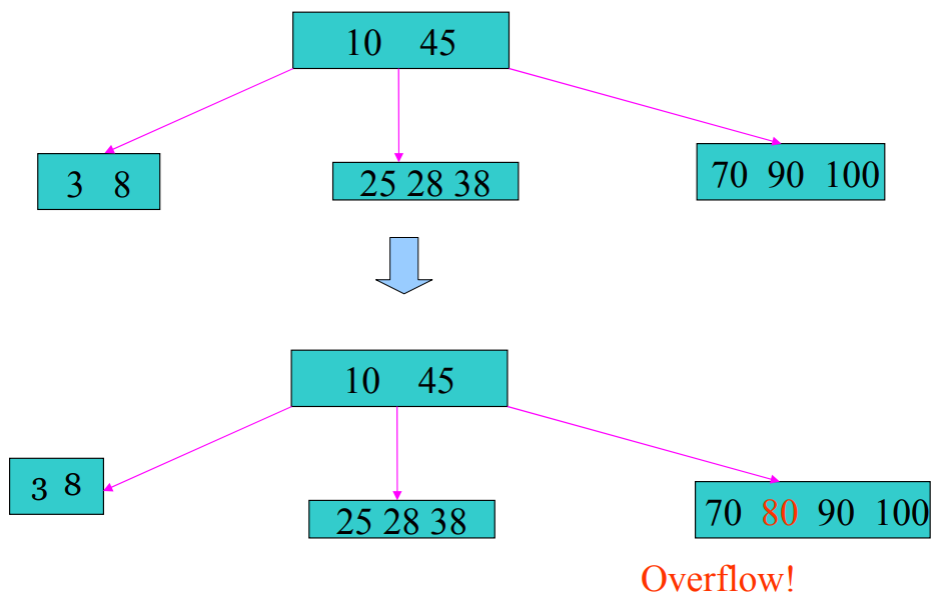
A 2-4 tree, also known as a 2-3-4 tree, is a multi-way search tree the following properties

1. Nodes may contain 1, 2 or 3 items
2. A node with k items has $k + 1$ children, except for leaf nodes. Such a node is called $(k + 1)$ -node
3. All leaves are on the same level

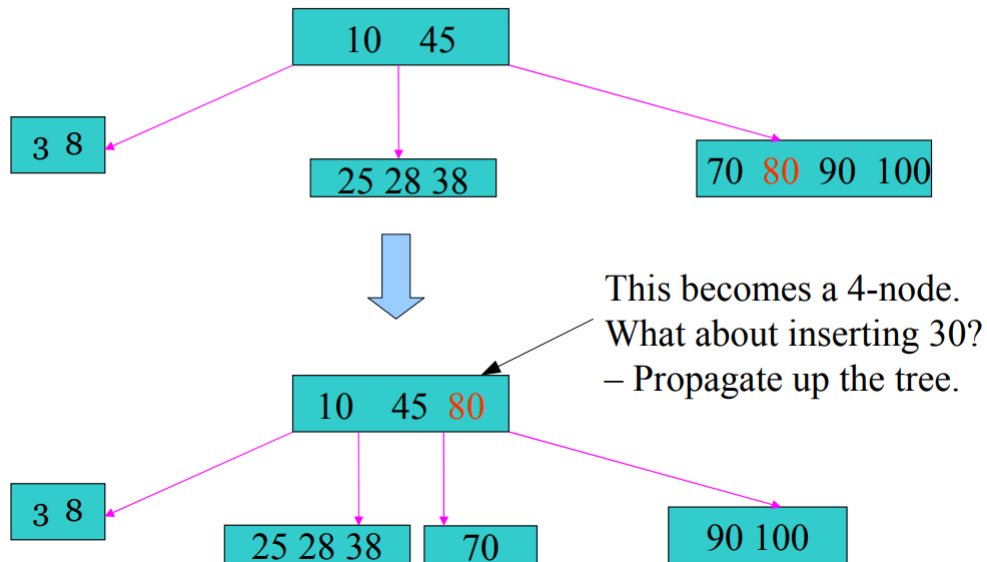


10.2.1 Insertion

To insert into a 2-4 tree, first, find the appropriate leaf. If there is room, just add the element to the leaf. If there is no room, move the middle item to parent and split remaining items among two children.



Move middle element to parent and split



10.2.2 Removal

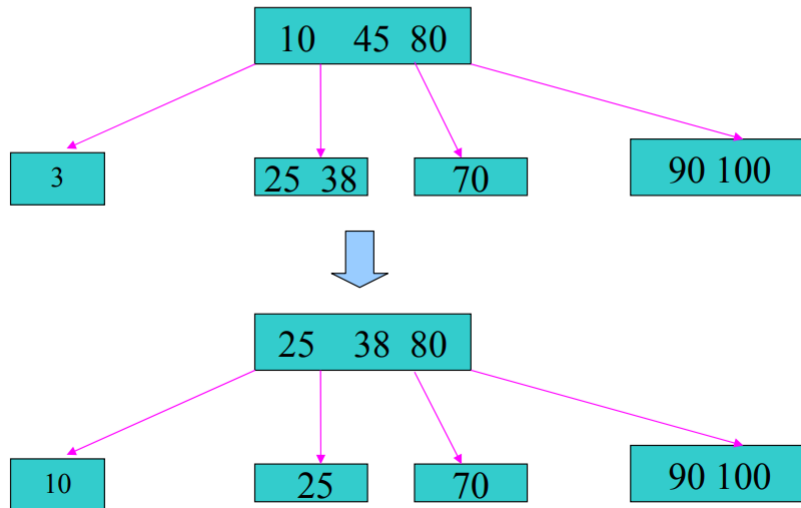
First, we find the key with a simple multi-way search. There are two cases

1. **Case 1:** It may be on the leaf
2. **Case 2:** It may be in internal node

If the item to delete is in internal node, we can reduce to case 1 by first finding its immediate predecessor, swapping them, and then removing the item.

✚ Remove 45

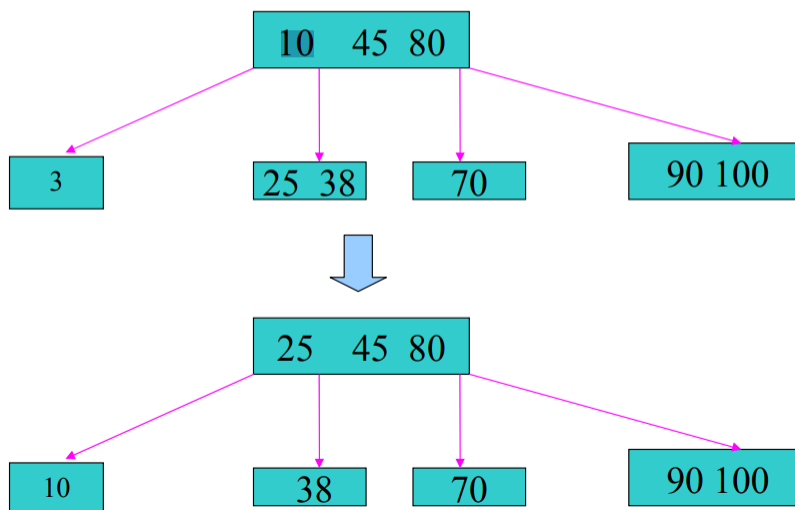
- ❖ Swap 38 and 45
- ❖ Remove 45 at leaf



But what if there are not enough items in the node after removal? this is known as an *underflow*. In this case, pull an item from the parent, replace it with an item from a sibling - transfer

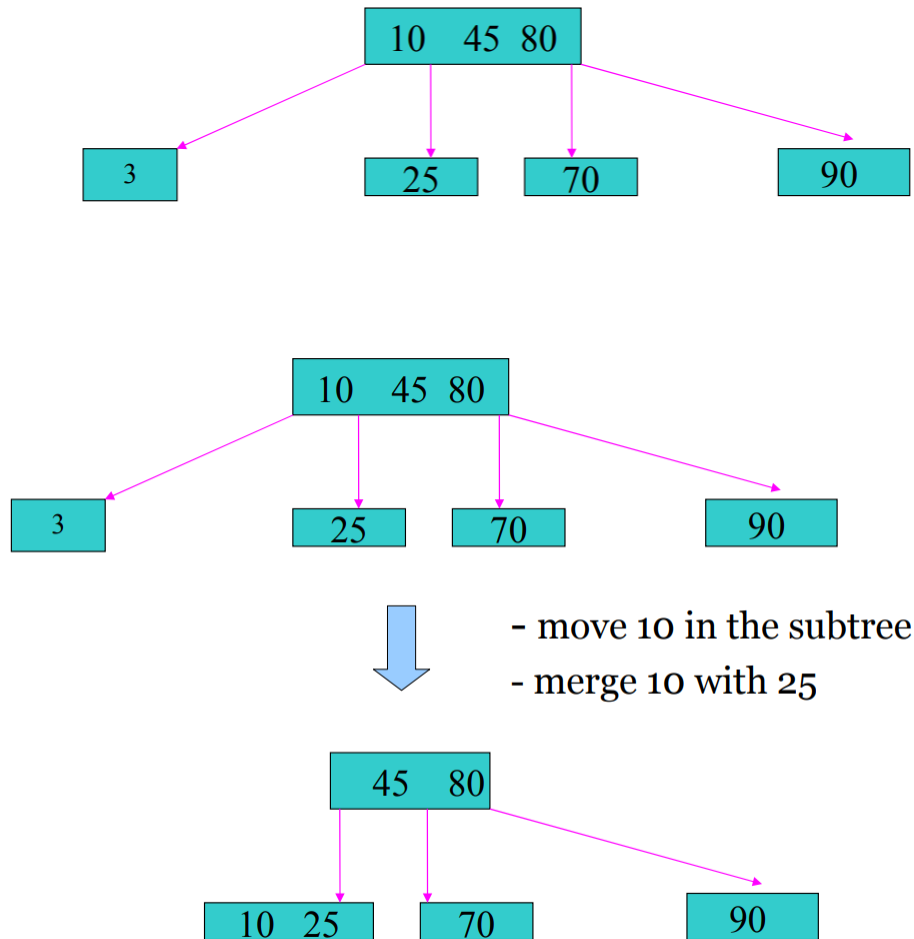
✚ Remove 3

- ❖ move 10 into the subtree
- ❖ move 25 into the parent



But what happens if the the siblings are 2-nodes (one element), we will not be able to steal from them. In this case, we perform *node merging*.

✦ Remove 3



10.2.3 Properties

- 2-4 trees are easy to maintain
- Insertion and deletion take $O(\log n)$
- Balanced trees

10.3 B-trees

Up to now, all data that has been stored in the tree has been in memory. If data gets too big for main memory, what do we do? If we keep a pointer to the tree in main memory, we could bring in just the nodes that we need.

For instance, to do an insert with a BST, if we need the left child, we do a disk access and retrieve the left child. If the left child is NIL, then we can do the insert, and store the child node on the disk

But storing the data requires disk accesses, which is expensive, compared to execution of machine instructions. If we can reduce the number of disk accesses, then the procedures run faster

The only way to reduce the number of disk accesses is to increase the number of keys in a node, we see the problem in using this technique with binary search trees... The BST allows only one key per leaf/node.

If we increase the number of keys in the nodes, how will we do any tree operations effectively?

A B-tree is a self-balancing search tree data structure that maintains sorted data and allows searches, insertions, deletions, and sequential access in logarithmic time. B-trees are especially useful for managing large blocks of data in systems like databases and file systems, where data is stored on disk or other slow-access storage and needs to be accessed efficiently.

Unlike binary trees, where each node has at most two children, B-trees are multi-way trees where each node can have multiple children. The number of children a node can have is determined by the order of the tree.

B-trees maintain balance by ensuring that every path from the root to a leaf node has the same length, which guarantees logarithmic height and thus efficient operations.

Each node can contain a range of keys (from a minimum to a maximum number). This enables efficient storage and retrieval by storing more keys in fewer nodes, which minimizes the number of disk reads needed to access data.

The order m of a B-tree defines the maximum number of children each node can have. An order- m B-tree node can have up to $m - 1$ keys and m children.

The minimum degree t (often used instead of order) specifies the minimum number of children a non-root node must have, which is t or more.

- **Insertion:** When adding a key, nodes split if they reach their maximum capacity, ensuring the B-tree remains balanced.
- **Deletion:** When deleting a key, nodes may need to merge with their siblings if they go below their minimum capacity.

Because each node contains multiple keys, B-trees have a low height relative to the number of keys they store. This makes B-trees very efficient for disk-based storage, as it minimizes the number of disk accesses needed.

B-trees are widely used in database indexing, file systems, and other applications that involve managing large volumes of data that cannot fit into main memory.

Hashing (hash tables)

Many applications require a dynamic set that supports only the dictionary operations INSERT, SEARCH, and DELETE. For example, a compiler that translates a programming language maintains a symbol table, in which the keys of elements are arbitrary character strings corresponding to identifiers in the language. A hash table is an effective data structure for implementing dictionaries. Although searching for an element in a hash table can take as long as searching for an element in a linked list - $\Theta(n)$ time in the worst case. In practice, hashing performs extremely well. Under reasonable assumptions, the average time to search for an element in a hash table is $O(1)$

A hash table generalizes the simpler notion of an ordinary array. Directly addressing into an ordinary array takes advantage of the $O(1)$ access time for any array element.

To use direct addressing, you must be able to allocate an array that contains a position for every possible key

When the number of keys actually stored is small relative to the total number of possible keys, hash tables become an effective alternative to directly addressing an array, since a hash table typically uses an array of size proportional to the number of keys actually stored. Instead of using the key as an array index directly, we compute the array index from the key.

11.1 Direct-address table

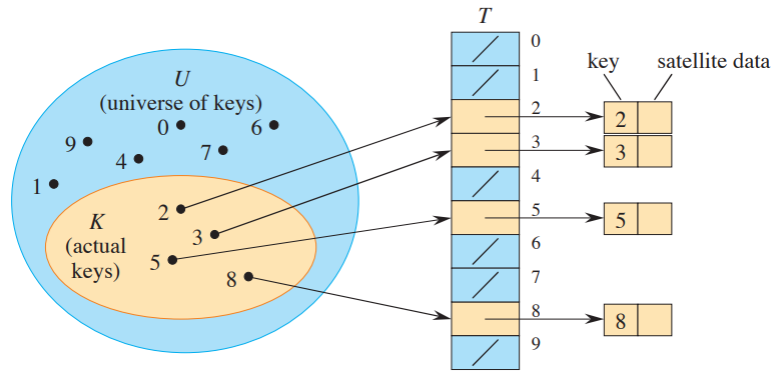
Direct addressing is a simple technique that works well when the universe U of keys is reasonably small. Suppose that an application needs a dynamic set in which each element has a distinct key drawn from the universe $U = \{0, 1, \dots, m-1\}$, where m is not too large.

To represent the dynamic set, you can use an array, or *direct-address table*, denoted by $T[0 : m-1]$, in which each position, or *slot*, corresponds to a key in the universe U . Figure 11.1 illustrates this approach. Slot k points to an element in the set with key k . If the set contains no element with key k , then $T[k] = \text{NIL}$.

The dictionary operations DIRECT-ADDRESS-SEARCH, DIRECT-ADDRESS-INSERT, and DIRECT-ADDRESS-DELETE on the following page are trivial to implement. Each takes only $O(1)$ time.

The dictionary operations DIRECT-ADDRESS-SEARCH, DIRECT-ADDRESS INSERT, and DIRECT-ADDRESS-DELETE on the following page are trivial to implement. Each takes only $O(1)$ time

For some applications, the direct-address table itself can hold the elements in the dynamic set. That is, rather than storing an element's key and satellite data in an object external to the direct-address table, with a pointer from a slot in the table to the object, save space by storing the object directly in the slot. To indicate an empty slot, use a special key. Then again, why store the key of the object at all? The index of the object is its key! Of course, then you'd need some way to tell whether slots are empty



```

1  direct-address-search(T,k)
2      return T[k]
3  direct-address-insert(T,x)
4      T[x.key] = x
5  direct-address-delete(T,x)
6      T[x.key] = nil

```

11.2 Hash tables

The downside of direct addressing is apparent: if the universe U is large or infinite, storing a table T of size $|U|$ may be impractical, or even impossible, given the memory available on a typical computer. Furthermore, the set K of keys actually stored may be so small relative to U that most of the space allocated for T would be wasted.

When the set K of keys stored in a dictionary is much smaller than the universe U of all possible keys, a hash table requires much less storage than a direct-address table. Specifically, the storage requirement reduces to $\Theta(|K|)$ while maintaining the benefit that searching for an element in the hash table still requires only $O(1)$ time. The catch is that this bound is for the *average-case time*¹, whereas for direct addressing it holds for the *worst-case time*.

With direct addressing, an element with key k is stored in slot k , but with hashing, we use a *hash function* h to compute the slot number from the key k , so that the element goes into slot $h(k)$. The hash function h maps the universe U of keys into the slots of a *hash table* $T[0 : m - 1]$:

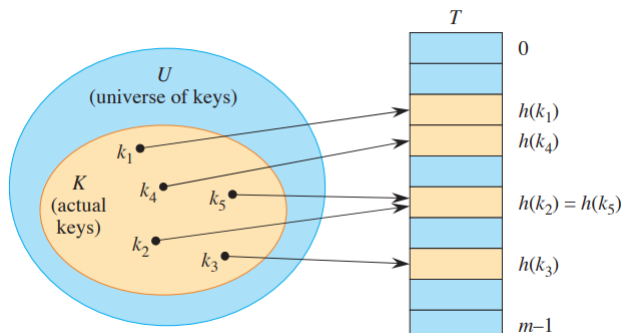
$$h : U \rightarrow \{0, 1, \dots, m - 1\}.$$

where the size m of the hash table is typically much less than $|U|$. We say that an element with key k *hashes to* slot $h(k)$, and we also say that $h(k)$ is the *hash value* of key k . The hash function reduces the range of array indices and hence the size of the array. Instead of a size of $|U|$, the array can have size m .

An example of a simple, but not particularly good, hash function is $h(k) = k \bmod m$.

There is one hitch, namely that two keys may hash to the same slot. We call this situation a *collision*. Fortunately, there are effective techniques for resolving the conflict created by collisions.

Of course, the ideal solution is to avoid collisions altogether. We might try to achieve this goal by choosing a suitable hash function h . One idea is to make h appear to be <random,= thus avoiding collisions or at least minimizing their number.



although a well-designed, random looking hash function can reduce the number of collisions, we still need a method for resolving the collisions that do occur

11.3 Independent uniform hashing

An "ideal" hashing function h would have, for each possible input k in the domain U , an output $h(k)$ that is an element randomly and independently chosen uniformly from the range $\{0, 1, \dots, m-1\}$. Once a value $h(k)$ is randomly chosen, each subsequent call to h with the same input k yields the same output $h(k)$.

We call such an ideal hash function an independent uniform hash function. Such a function is also often called a random oracle

When hash tables are implemented with an independent uniform hash function, we say we are using independent uniform hashing

Independent uniform hashing is an ideal theoretical abstraction, but it is not something that can reasonably be implemented in practice.

11.4 Collision resolution by chaining

At a high level, you can think of hashing with chaining as a nonrecursive form of divide-and-conquer: the input set of n elements is divided randomly into m subsets, each of approximate size $\frac{n}{m}$. A hash function determines which subset an element belongs to. Each subset is managed independently as a list.

each nonempty slot points to a linked list, and all the elements that hash to the same slot go into that slot's linked list. Slot j contains a pointer to the head of the list of all stored elements with hash value j . If there are no such elements, then slot j contains NIL

The worst-case running time for insertion is $O(1)$ The insertion procedure is fast in part because it assumes that the element x being inserted is not already present in the table. To enforce this assumption, you can search (at additional cost) for an element whose key is $x.key$ before inserting

For searching, the worst-case running time is proportional to the length of the list

Deletion takes $O(1)$ time if the lists are doubly linked. If the hash table supports deletion, then its linked lists should be doubly linked in order to delete an item quickly.

11.5 Analysis of hashing with chaining

Given a hash table T with m slots that stores n elements, we define the load factor α for T as $\alpha = \frac{n}{m}$, that is, the average number of elements stored in a chain. Our analysis will be in terms of α , which can be less than, equal to, or greater than 1.

The worst-case behavior of hashing with chaining is terrible: all n keys hash to the same slot, creating a list of length n . The worst-case time for searching is thus $\Theta(n)$ plus the time to compute the hash function—no better than using one linked list for all the elements. We clearly don't use hash tables for their worst-case performance.

The average-case performance of hashing depends on how well the hash function h distributes the set of keys to be stored among the m slots, on the average

for now we assume that any given element is equally likely to hash into any of the m slots. That is, the hash function is uniform. We further assume that where a given element hashes to is independent of where any other elements hash to. In other words, we assume that we are using independent uniform hashing

Because hashes of distinct keys are assumed to be independent, independent uniform hashing is universal: the chance that any two distinct keys k_1 and k_2 collide is at most $\frac{1}{m}$.

11.6 Hash functions

For hashing to work well, it needs a good hash function. Along with being efficiently computable, what properties does a good hash function have? How do you design good hash functions?

This section first attempts to answer these questions based on two ad hoc approaches for creating hash functions: hashing by division and hashing by multiplication.

Although these methods work well for some sets of input keys, they are limited because they try to provide a single fixed hash function that works well on any data, an approach called static hashing.

We then see that provably good average-case performance for any data can be obtained by designing a suitable family of hash functions and choosing a hash function at random from this family at runtime, independent of the data to be hashed. The approach we examine is called random hashing

A good hash function satisfies (approximately) the assumption of independent uniform hashing: each key is equally likely to hash to any of the m slots, independently of where any other keys have hashed to. What does "equally likely" mean here? If the hash function is fixed, any probabilities would have to be based on the probability distribution of the input keys

Unfortunately, you typically have no way to check this condition, unless you happen to know the probability distribution from which the keys are drawn. Moreover, the keys might not be drawn independently

Occasionally you might know the distribution. For example, if you know that the keys are random real numbers k independently and uniformly distributed in the range $0 < k < 1$, then the hash function

$$h(k) = \lfloor km \rfloor.$$

satisfies the condition of independent uniform hashing.

A good static hashing approach derives the hash value in a way that you expect to be independent of any patterns that might exist in the data. For example, the "division method" computes the hash value as the remainder when the key is divided by a specified prime number. This method may give good results, if you (somehow) choose a prime number that is unrelated to any patterns in the distribution of keys.

Random hashing picks the hash function to be used at random from a suitable family of hashing functions. This approach removes any need to know anything about the probability distribution of the input keys, as the randomization necessary for good average-case behavior then comes from the (known) random process used to pick the hash function from the family of hash functions, rather than from the (unknown) process used to create the input keys. We recommend that you use random hashing.

In practice, a hash function is designed to handle keys that are one of the following two types:

- A short nonnegative integer that fits in a w -bit machine word. Typical values for w would be 32 or 64.
- A short vector of nonnegative integers, each of bounded size. For example, each element might be an 8-bit byte, in which case the vector is often called a (byte) string. The vector might be of variable length.

To begin, we assume that keys are short nonnegative integers.

11.6.1 Static hashing

Static hashing uses a single, fixed hash function. The only randomization available is through the (usually unknown) distribution of input keys. This section discusses two standard approaches for static hashing: the division method and the multiplication method. Although static hashing is no longer recommended, the multiplication method also provides a good foundation for "nonstatic hashing", better known as random hashing, where the hash function is chosen at random from a suitable family of hash functions.

11.6.2 The division method

The division method for creating hash functions maps a key k into one of m slots by taking the remainder of k divided by m . That is, the hash function is

$$h(k) = k \bmod m.$$

For example, if the hash table has size $m = 12$ and the key is $k = 100$, then $h(k) = 4$. Since it requires only a single division operation, hashing by division is quite fast

The division method may work well when m is a prime not too close to an exact power of 2. There is no guarantee that this method provides good average-case performance, however, and it may complicate applications since it constrains the size of the hash tables to be prime.

It is best to ensure that the table size is a prime number, because keys are typically not randomly distributed, and usually have some pattern.

11.6.3 The multiplication method

The general multiplication method for creating hash functions operates in two steps. First, multiply the key k by a constant A in the range $0 < A < 1$ and extract the fractional part of kA . Then, multiply this value by m and take the floor of the result. That is, the hash function is

$$h(k) = \lfloor m(kA \bmod 1) \rfloor.$$

where " $kA \bmod 1$ " means the fractional part of kA , that is, $kA - \lfloor kA \rfloor$

The general multiplication method has the advantage that the value of m is not critical and you can choose it independently of how you choose the multiplicative constant A

11.7 Open addressing

Open addressing hashing is a collision resolution technique used in hash tables. When a hash collision occurs (i.e., two keys are hashed to the same index), open addressing searches for the next available slot in the array to store the colliding key. This eliminates the need for linked lists or separate chaining.

The two methods of open addressing we will discuss are

1. linear probing
2. quadratic probing

11.7.1 Linear probing

Linear probing is a simple and commonly used method for handling collisions in hash tables. Collisions occur when two different keys hash to the same index in the table. Linear probing resolves these collisions by searching sequentially through the table to find the next available slot.

The formula for linear probing is

$$h'(k, i) = (h(k) + i) \bmod m$$

- $h(k)$: Original hash function.
- i : Probe sequence (0, 1, 2, ...).
- m : Table size.

11.7.2 Quadratic probing

Quadratic probing is a collision resolution technique used in hash tables as part of open addressing. When a collision occurs, instead of probing linearly (as in linear probing), quadratic probing uses a quadratic function to calculate the next index to probe. This helps reduce clustering and improves performance.

The formula for quadratic probing is

$$h'(k, i) = (h(k) + c_1 i + c_2 i^2) \mod m$$

Where:

- $h(k)$: Original hash function.
- i : Probe number (starting from 0).
- c_1, c_2 : Constants to control the probing step size.
- m : Size of the hash table (typically a prime number).

Simple quadratic probing is when $c_1 = 0$, $c_2 = 1$

11.7.3 Deletion problem in open addressing

In open addressing, when an element is removed from a hash table, simply marking the corresponding slot as unoccupied (`occupied[index] = false`) can break the probe sequence. This makes it impossible to find keys that were inserted after the deleted element, even though they still exist in the table.

- Insert 10,20,30 into a hash table with linear probing.
- Suppose 10 and 20 are hashed to the same index due to collisions.
- Remove 10 by marking its slot as unoccupied.
- Now, if you search for 20, the probe sequence will terminate prematurely because it will encounter the unoccupied slot where 10 was stored.

The next few sections introduce solutions to this problem.

11.7.4 Tombstoning

One way to address the deletion problem in open addressing is tombstoning. This technique introduces a new boolean array called *tombstone*. When we remove an item, we mark its slot unoccupied, but we also mark the spot with a tombstone. This way, when we probe for an element, we can continue searching if we encounter a hole, as long as there is a tombstone. We will see this method more in depth in the upcoming C++ example.

11.7.5 Lazy deletion

Instead of marking a slot as a tombstone, simply leave the slot as occupied but ignore its value in future insertions.

Mark the slot as "deleted" logically but do not clear the occupied flag.

Searches and insertions treat these "deleted" slots as valid for probing but ignore their values.

11.7.6 Rehashing

Instead of handling deletions within the existing table structure, rebuild the entire hash table when a deletion occurs

When a key is deleted, construct a new hash table and reinsert all active keys from the old table into the new one.

This eliminates any empty or deleted slots and ensures clean probe sequences.

This simplifies search, insertion, and deletion logic (no need for tombstones or special flags).

It also avoids performance degradation from tombstones or empty clusters.

However, this method is expensive in terms of time complexity for frequent deletions, as rehashing involves reinserting all active keys.

```
1 void rehashTable() {
2     hashtable<TABLE_SIZE> newTable;
3     for (int i=0; i<TABLE_SIZE; ++i) {
4         if (occupied[i]) {
5             newTable.insert(table[i]);
6         }
7     }
8     *this = newTable;
9 }
```

Then we we remove an element, we mark its spot unoccupied, then rehash the table.

```
1 while (occupied[index]) {
2     if (table[index] == key) {
3         // Key found, mark as removed
4         occupied[index] = false;
5         rehashTable();
6         return true;
7     }
8     ...
}
```

11.8 Load factor

The load factor of a hash table is a measure of how full the table is, expressed as a ratio of the number of elements in the table to its total capacity. It helps determine when a hash table might need to be resized to maintain efficient performance.

The load factor α is calculated as:

$$\alpha = \frac{\text{number of elements in the table}}{\text{total number of slots in the table}}.$$

11.9 Static hashing

When a hash table uses a fixed-size array as its storage, the approach is generally called static hashing. In static hashing, the table's size is determined at the time of its creation and does not change as elements are added or removed. This is in contrast to dynamic hashing, where the table can grow or shrink based on the number of entries.

11.10 Static linear hashing

Static linear hashing is a hashing technique that combines aspects of static hashing (fixed-size table) with linear probing for collision resolution. In this approach:

- **Fixed-Size Table:** The hash table is created with a fixed number of slots, meaning its size does not grow or shrink based on the number of elements. This fixed size makes it suitable for scenarios where the maximum number of entries is known in advance or memory is constrained.
- **Hash Function:** A hash function maps each key to a specific index within the table. Since the table is of fixed size, the hash function ensures that the resulting index is within the range of available slots, typically using a modulus operation, such as

$$\text{hash}(\text{key}) \bmod \text{table size}.$$

- **Linear Probing for Collision Resolution:** When two keys hash to the same index (a collision), linear probing is used to find the next available slot. This means that if a collision occurs at index i , the algorithm checks $i + 1$, $i + 2$, and so on, wrapping around to the start of the table if necessary, until an empty slot is found.
- **Load Factor Constraints:** Because the table size is static, as the number of entries increases, so does the load factor. Higher load factors increase the likelihood of collisions, which can slow down insertion and search operations. Generally, static linear hashing performs best when the load factor is kept below a certain threshold (often around 0.7 to 0.8).

11.11 Static hashing with linear probe implementation

11.11.1 Interface

We support three operations

1. Inserting
2. Searching
3. Removing

11.11.2 The basics

```
1  template<size_t TABLE_SIZE>
2  class hashtable {
3      int table[TABLE_SIZE];
4      bool occupied[TABLE_SIZE];
5  public:
6      hashtable() {
7          fill(std::begin(occupied), std::end(occupied), 0);
8      }
9  };
```

By including `template<size_t TABLE_SIZE>` at the start, the class becomes a template, which means the table size is set by the user when they create an instance of the class. `size_t` is used as the type for the `TABLE_SIZE` parameter because it's ideal for holding non-negative values and represents sizes and counts. With this approach, users can create hash tables of varying sizes without needing dynamic memory allocation.

Within the class definition, there's an array named `table` that has a size of `TABLE_SIZE`. Since `TABLE_SIZE` is defined as a template parameter, it's treated as a compile-time constant, allowing the table array to be statically allocated. This fixed-size array avoids the need for dynamic memory management, which is often desirable in systems where memory usage needs to be predictable.

We also create a boolean array of the same size to track the spots in the hash table that have been filled. In the constructor, we fill the `occupied` array with 0 to mark all spots empty.

11.11.3 The hash function

For this example, we will use a simple hashing function

$$h(x) = x \bmod n.$$

Where n is the size of the table.

```
1  int hash(int element) {
2      return element % TABLE_SIZE;
3  }
```

11.11.4 Inserts

```
1  bool insert(int key) {
2      int index = hash(key);
3      int originalIndex = index;
4
5      // Linear probing to find an open slot
6      while (occupied[index]) {
7          if (table[index] == key) {
8              // Key already exists, insertion fails to avoid
↪      duplicates
9              return false;
10         }
11         // Increment index, looping around if necessary
12         index = (index + 1) % TABLE_SIZE;
13         if (index == originalIndex) {
14             // Table is full, insertion fails
15             return false;
16         }
17     }
18
19     table[index] = key;
20     occupied[index] = true;
21     return true;
22 }
```

11.11.5 Searching

```
1  bool search(int key) const {
2      int index = hash(key);
3      int originalIndex = index;
4
5      // Linear probing to search for the key
6      while (occupied[index]) {
7          if (table[index] == key) {
8              return true; // Key found
9          }
10         index = (index + 1) % TABLE_SIZE;
11         if (index == originalIndex) {
12             // Key not found, returned to original index
13             return false;
14         }
15     }
16     return false; // Key not found
17 }
```

11.11.6 Removing

```
1  bool remove(int key) {
2      int index = hash(key);
3      int originalIndex = index;
4
5      // Linear probing to find the key
6      while (occupied[index]) {
7          if (table[index] == key) {
8              // Key found, mark as removed
9              table[index] = -1;
10             occupied[index] = false;
11             return true;
12         }
13         index = (index + 1) % TABLE_SIZE;
14         if (index == originalIndex) {
15             return false;
16         }
17     }
18
19     return false;
20 }
```

11.11.7 Searching and removing with tombstones

Consider the hash table

0 :	142
1 :	27
2 :	15
3 :	empty
4 :	empty
5 :	83
6 :	18
7 :	empty add 109
8 :	47
9 :	empty
10 :	10
11 :	empty
12 :	90.

Suppose now we add key 109, which hashes to slot 5. Since 5 is taken, and we are using open addressing with a linear probe, 109 gets added to slot 7 instead.

Since the while loop depends on not finding any unoccupied slots in the path to the key, it could happen that on the path an element that previously existed was removed, thereby creating an unoccupied slot in the linear path to the requested element. In this case, the search function would incorrectly assume that the requested key does not exist! The while loop detected an empty slot before finding the key, therefore the key must not exist. To solve this problem, we need a way to indicate that slots that are empty have had data previously in them. We introduce a new array to track what we call tombstones, slots that are empty but previously held data.

In the example above, suppose we remove key 18, doing so creates a hole between 83 and 109, which both hash to 5. Searching for 109 in this case would give false, when it is clearly a member of the table.

Since the remove function also relies on searching for the key, it too needs the additional tombstone logic.

Thus, we simply create a boolean tombstone array, and make slight adjustments to the search and remove logic.

```

1  bool tombstone[TABLE_INDEX]; // Make sure to initially fill with
   ↪  zeros.
2
3  // Adjustments to search
4  while (occupied[index] || tombstone[index]) {
5      if (occupied[index] && table[index] == key) {
6          return true; // Key found
7      }
8      ...
9
10 // Adjustments to remove
11 while (occupied[index] || tombstone[index]) {
12     if (occupied[index] && table[index] == key) {
13         // Key found, mark as removed
14         occupied[index] = false;
15         tombstone[index] = true;
16         return true;
17     }
18     ...

```

11.12 Static hashing with quadratic probing implementation

The code for quadratic probing is similar to linear probing, but with a slight adjustment to the probe increment. For this example, we use

$$h'(x) = h(x) + i^2 \text{ for } i = 1, 2, 3, \dots, m.$$

Where m is the size of the table.

11.12.1 Inserting

```
1  bool insert(int key) {
2      int index = hash(key);
3      int originalIndex = index;
4      int i=1;
5
6      // quadratic probing to find an open slot
7      while (occupied[index]) {
8          if (table[index] == key) {
9              // Key already exists, insertion fails to avoid
10             ↪ duplicates
11                 return false;
12             }
13             index = (index + (i*i)) % TABLE_SIZE;
14             ++i; if (i>TABLE_SIZE) return false;
15         }
16         table[index] = key;
17         occupied[index] = true;
18         return true;
19     }
```

11.12.2 Searching

```
1  bool search(int key) const {
2      int index = hash(key);
3      int originalIndex = index;
4      int i=1;
5
6      // Linear probing to search for the key
7      while (occupied[index]) {
8          if (table[index] == key) {
9              return true; // Key found
10         }
11         index = (index + (i*i)) % TABLE_SIZE;
12         ++i; if (i>TABLE_SIZE) return false;
13     }
14     return false; // Key not found
15 }
```

11.12.3 Removing

```
1  bool remove(int key) {
2      int index = hash(key);
3      int originalIndex = index;
4      int i=1;
5
6      // Linear probing to find the key
7      while (occupied[index]) {
8          if (table[index] == key) {
9              // Key found, mark as removed
10             occupied[index] = false;
11             rehash();
12             return true;
13         }
14         index = (index + (i*i)) % TABLE_SIZE;
15         ++i; if (i > TABLE_SIZE) return false;
16     }
17
18     return false;
19 }
```


11.13 Static hashing with chaining implementation

An improvement on the collision resolution technique above would be to instead use linked lists as the hash table elements. This way, if we have collisions, we can just add the element that hashes to the same slot to that slot's linked list front. The hash function will remain the same

11.13.1 The basics

```
1  template<size_t TABLE_SIZE>
2  class hashtable {
3      list<int> table[TABLE_SIZE];
4      bool occupied[TABLE_SIZE];
5
6  public:
7      hashtable() {
8          fill(std::begin(occupied), std::end(occupied), 0);
9      }
10 };
```

Instead of a table of integers, we now use a table of doubly linked lists of integers. The occupied array remains the same.

11.13.2 Inserting

The insert function becomes much easier, we simply add the key to the slot's linked list.

```
1  void insert(int key) {
2      int hx = hash(key);
3      table[hx].push_front(key);
4      occupied[hx] = 1;
5  }
```

11.13.3 Searching

```
1  bool search(int key) const {
2      int hx = hash(key);
3      if (!occupied[hx])
4          return false;
5
6      if (std::find(table[hx].begin(), table[hx].end(), key) !=
    ↪ table[hx].end())
7          return true;
8      return false;
9  }
```

11.13.4 Removing

```
1  bool remove(int key) {  
2      int hx = hash(key);  
3  
4      if (!occupied[hx])  
5          return false;  
6  
7      table[hx].remove(key);  
8      if (table[hx].empty())  
9          occupied[hx] = 0;  
10  
11     return true;  
12 }
```

11.14 Injective hashing

A hash function is injective if it maps every distinct input to a distinct output, i.e., it has no collisions. In other words,

$$h(k_1) = h(k_2) \implies k_1 = k_2.$$

However, this is often impractical for large or infinite domains because it requires the output space to be at least as large as the input space.

11.15 Perfect hashing

A hash function is perfect if it is collision-free for a specific, finite set of inputs. For a given set S , a perfect hash function h maps every element in S to a unique value (injectivity is guaranteed within S).

If the output range of h matches the size of S , it is a minimal perfect hash function

It applies only to a known, finite set of inputs. Does not guarantee injectivity outside the specific set S . Computationally feasible and often used when the set of inputs is fixed, such as in compiler keyword lookups.

Table indexing and row-major order

Table indexing in multidimensional arrays is a method to access elements in an array by calculating their position in memory based on their indices. Since multidimensional arrays are often stored in contiguous memory as a single, linear array, table indexing is used to convert the array's multidimensional coordinates (like $A[i][j][k]$) into a single-dimensional index. This process allows programming languages to efficiently store and retrieve elements without the overhead of creating separate memory locations for each dimension.

12.1 Row-major order

In row-major order (used by C, C++, and many other languages), the elements are stored row by row. For a 2D array A with dimensions $[m][n]$, the element $A[i][j]$ is stored in a 1D array at position:

$$\text{index} = i \times n + j.$$

For a 3D array $A[x][y][z]$, the element $A[i][j][k]$ is at

$$\text{index} = i \times (y \times z) + j \times z + k.$$

This approach generalizes for any n -dimensional array by multiplying indices by the product of the sizes of subsequent dimensions.

For an n -dimensional array $A[d_1][d_2] \dots [d_n]$, where each d_k is the size of the k -th dimension, the index of element $A[i_1][i_2] \dots [i_n]$ can be calculated as:

$$\text{index} = i_1 \times (d_2 \times d_3 \dots \times d_n) + i_2 \times (d_3 \dots \times d_n) + \dots + i_n$$

Table indexing is essential for optimizing memory usage and access speed in multidimensional arrays. It enables direct memory access with a single calculation, avoiding nested loops or additional data structures, which is especially useful in low-level programming and performance-critical applications.

Graphs

A graph is a collection of nodes and the connections between them. There are many different types of graphs

- Simple graph
- Directed graph
- Multi-graph
- Weighted graph
- Complete graph

etc...

13.1 Simple Graph

A simple graph $G = (V, E)$ consists of a nonempty set V of vertices and a possibly empty set E of edges. Each edge connects two vertices from V :

A simple graph is a graph that satisfies the following properties:

1. **Undirected:** The edges have no direction.
2. **No Loops:** No vertex has an edge to itself.
3. **No Multiple Edges:** There is at most one edge between any pair of vertices.

In an undirected graph, $\{v_i, v_j\} = \{v_j, v_i\}$. We denote the number of vertices in a graph by $|V|$, and the number of edges in a graph by $|E|$

13.2 Undirected graph

An undirected graph is a graph where edges have no direction. If there is an edge between u and v , it is the same as the edge between v and u

13.3 Directed graph

A directed graph (digraph) $G = (V, E)$ consists of a nonempty set V of vertices and a possibly empty set E of edges (or archs). Each edge connects two vertices from V

$$\{v_i, v_j\} \neq \{v_j, v_i\}.$$

13.4 Weighted Graph

A weighted graph is a graph where edges have assigned numbers. The numbers could be distances values, lengths, costs, ..., etc

13.5 Complete Graph

For each pair of vertices, there is exactly one edge connecting them.

13.6 More terms

- **Subgraph:** G' of $G = (V, E)$ is a graph (V', E') such that $V' \subseteq V$, and $E' \subseteq E$
- **Adjacent vertices:** Two vertices v_i and v_j are adjacent if edge $(v_i, v_j) \subseteq E$

Such edge is called incident with vertices v_i and v_j

- **Degree:** The degree of a vertex v , $\deg(v)$, is the number of edges incident with v

If $\deg(v) = 0$, v is called isolated vertex.

- **Path:** A sequence of edges $(v_1, v_2), (v_2, v_3), \dots, (v_{n-1}, v_n)$, denoted as the path:

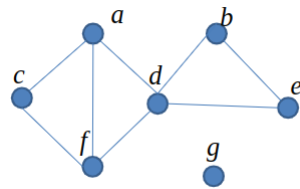
$$v_1 v_2 v_3 \dots v_n$$

- **Circuit:** There exists a path v_1, v_2, \dots, v_n where $v_1 = v_n$, and no edge is repeated.
- **Cycle:** If all vertices in a circuit are different.

13.7 Graph representations

One way to represent a graph is an *adjacency list*. Vertices are stored as records or objects, and every vertex stores a list of adjacent vertices.

This data structure allows the storage of additional data on the vertices.



a, c, d, f
 b, d, e
 c, a, f
 d, a, b, e, f
 e, b, d
 f, a, c, d
 $g.$

We can also store the same information as a matrix, called an *adjacency matrix*. Rows represent source vertices, columns represent destination vertices. Data on edges and vertices are stored externally

$$\begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Adjacency matrices have rows and columns in order, in the example above, row one column one is a,a.

Adjacency matrices are symmetric for simple graphs, non-symmetric for digraphs

In case of weighted graphs, values in matrix indicate weights of edges, but only if we don't allow weights to be zero.

13.8 Interface of adjacency lists and matrices

- Insert node
- Delete node
- Insert edge
- Delete edge
- Check existence of edge
- Get/set edge weight
- Get neighbors of node

13.9 Complexities for adjacency operations

Operation	Adjacency Matrix	Adjacency List
Add a node	$O(V ^2)$	$O(1)$
Remove a node	$O(V ^2)$	$O(E)$
Add an edge	$O(1)$	$O(1)$
Remove an edge	$O(1)$	$O(V)$
Get neighbors of a node	$O(V)$	$O(V)$
Test an edge	$O(1)$	$O(V)$
Get/set edge	$O(1)$	$O(V)$
Storage	$O(V ^2)$	$O(V + E)$

13.10 Graph traversal/search

We want to visit all vertices once and only once. There may be cycles, which can cause infinite loops. There may be isolated vertices or isolated sub- graphs

13.10.1 Breadth-first-traversal

Visit siblings first before visiting children.

```
1  BFS(Graph, startNode):
2      1. Create an empty queue `Q` to store nodes to be visited.
3      2. Create a set `visited` to track visited nodes.
4      3. Enqueue the `startNode` into `Q`.
5      4. Add `startNode` to the `visited` set.
6
7      5. While `Q` is not empty:
8          a. Dequeue a node `current` from `Q`.
9          b. Process `current` (e.g., print it or record it).
10
11         c. For each neighbor `neighbor` of `current`:
12             i. If `neighbor` is not in `visited`:
13                 - Enqueue `neighbor` into `Q`.
14                 - Add `neighbor` to the `visited` set
```

13.10.2 Depth-first traversal

Visit children before siblings, recursive.

```
1  DFS(Graph, currentNode, visited):
2      1. If `currentNode` is not in `visited`:
3          a. Add `currentNode` to `visited`.
4          b. Process `currentNode` (e.g., print it or record it).
5
6      c. For each neighbor `neighbor` of `currentNode`:
7          i. Call DFS(Graph, neighbor, visited).
```

13.11 Shortest path problems

- **Single pair shortest path problem:** Finding a path between two vertices that the sum of weights of its edges is minimal
- **Single source shortest path problem:** Finding shortest paths from one vertices to all others
- **All pair shortest path problem:** Finding shortest paths for all pairs of vertices

Math algorithms

14.1 Euclidean GCD Algorithm

The GCD of two integers a and b (with $a \leq b$) is the largest integer that divides both a and b . The Euclidean algorithm is based on the principle that

$$\gcd(a, b) = \gcd(b, a \bmod b).$$

This means that the GCD of two numbers doesn't change if the larger number is replaced by its remainder when divided by the smaller number. You keep repeating this until the remainder is 0, and the GCD will be the last non-zero remainder

```
1  int gcd(int a, int b) {  
2      if (!b) return a;  
3  
4      return gcd(b, a%b);  
5  }
```

14.2 Fibonacci numbers in constant time

The formula for the n th Fibonacci number $F(n)$ can be derived using Binet's formula, which expresses the Fibonacci sequence in terms of powers of the golden ratio.

$$F(n) = \frac{\phi^n - \psi^n}{\sqrt{5}}.$$

Where ϕ is the golden ration $\frac{1+\sqrt{5}}{2} \approx 1.618$, and ψ is the conjugate of the golden ration $\psi = \frac{1-\sqrt{5}}{2} \approx -0.618$

14.3 Sterlings factorial approximation

Stirling's Approximation provides an approximation for factorials, particularly useful for large values of n . The formula is:

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n.$$