

Algorithms and Data Structures

Notebook

Document produced with Jupyter Lab

Christoph Brauer | SE 02 Spring | Semester 2023

Table of Contents

- 1. Algorithm Time Complexity and Asymptotic Notations
 - 1.1. Counting Steps and RAM Model Assumptions
 - 1.1.1. Time Complexity
 - 1.1.2. Space Complexity
 - 1.2. Asymptotic Notations
 - 1.2.1. Big O
 - 1.2.2. Big Ω
 - 1.2.3. Big Θ
 - 1.3. Mathematical Functions
 - 1.3.1. Common Functions in Algorithm Analysis
 - 1.3.2. Growth Rates
 - 1.4. Analyzing Loops
 - 1.4.1. Nested Loops
 - 1.4.2. Recursive Loops
 - 1.5. Time Complexity of Searching and Sorting Algorithms
 - 1.5.1. Linear Search
 - 1.5.2. Binary Search
 - 1.5.3. Bubble Sort
 - 1.5.4. QuickSort
- 2. Mathematical Arguing and Reasoning
 - 2.1. Revision of Algorithm Analysis
 - 2.1.1. Time Complexity Examples
 - 2.1.2. Space Complexity Examples
 - 2.2. Mathematical Arguing and Reasoning
 - 2.2.1. Proof Techniques
 - 2.2.2. Inductive Proofs
 - 2.3. Logarithms Review
 - 2.3.1. Logarithm Basics and Properties
 - 2.3.2. Logarithms in Computer Science
- 3. Computational Machine Architectures

- 3.1. Von Neumann Architecture
 - 3.1.1. Components and Functions
 - 3.1.2. Memory Hierarchy
- 3.2. CPU Architecture
 - 3.2.1. Pipelining
 - 3.2.2. Cache Memory
- 3.3. CISC vs RISC Architectures
 - 3.3.1. Advantages and Disadvantages
 - 3.3.2. Evolution and Current Trends
- 3.4. GPU: Graphical Process Unit
 - 3.4.1. GPU Architecture
 - 3.4.2. Applications in Computing
- 3.5. Building Parallel Algorithms
 - 3.5.1. Parallelism Techniques
 - 3.5.2. Multithreading and Concurrency
- 4. Graph Data Structure and Search Algorithms
 - 4.1. Graph Data Structure Overview
 - 4.1.1. Terminology and Representation
 - 4.1.2. Types of Graphs
 - 4.2. Graph Applications in Real Life
 - 4.2.1. Social Networks
 - 4.2.2. Navigation Systems
 - 4.3. Search Algorithms on Graph Data Structure
 - 4.3.1. Breadth-first Search (BFS)
 - 4.3.2. Depth-first Search (DFS)
 - 4.3.3. Dijkstra's Algorithm
 - 4.3.4. A* Algorithm
- 5. Sources

Algorithm Time Complexity and Asymptotic Notations

Algorithm time complexity and asymptotic notations are fundamental concepts in computer science and algorithm analysis. Time complexity is a measure of the amount of time an algorithm takes to complete its tasks as a function of the input size. Asymptotic notations, on the other hand, provide a way to describe the growth of functions and compare their efficiency. Understanding these concepts allows developers to design, implement, and optimize algorithms that can efficiently solve computational problems. By studying algorithm time complexity and asymptotic notations, one can gain insights into how well an algorithm will perform under various circumstances and make informed decisions when choosing or designing algorithms for different applications.

1.1. Counting Steps and RAM Model Assumptions

In order to analyze an algorithm's time complexity, we need to count the number of elementary operations or steps that the algorithm performs. Counting steps allows us to quantify the resources an algorithm consumes, which can help us compare its efficiency with that of other algorithms. However, counting steps accurately can be challenging, as the number of steps may vary depending on the hardware and software environment.

To simplify this process, we can use the Random Access Machine (RAM) model of computation, which makes several assumptions to provide a consistent and easily understandable framework for algorithm analysis. The RAM model assumes that:

- Each elementary operation takes a constant amount of time.
- Memory access times are constant and uniform, regardless of the memory location.
- Instructions can be executed sequentially, one after another.

By making these assumptions, the RAM model allows us to focus on the key aspects of an algorithm's performance, abstracting away the hardware and software details that can make counting steps more complex. With the RAM model, we can concentrate on the high-level structure of the algorithm and count the number of elementary operations as a function of input size. This function forms the basis for determining the algorithm's time complexity, which will be discussed in the following sections.

1.1.1. Time Complexity

Time complexity is a measure of the efficiency of an algorithm in terms of the time it takes to execute as a function of the input size. It reflects the number of elementary operations or steps an algorithm performs to complete its tasks, and provides a means to compare the performance of different algorithms. Understanding the time complexity of an algorithm is crucial, as it helps developers optimize their code and choose the most efficient algorithm for a given problem.

When analyzing the time complexity of an algorithm, we focus on its growth rate rather than the absolute number of steps. This is because the growth rate indicates how the algorithm's performance scales as the input size increases, which is more important for determining its efficiency. We express the time complexity of an algorithm using asymptotic notations, which will be covered in the next section.

In general, there are three main types of time complexity:

1. Best-case time complexity: The minimum number of steps an algorithm takes to solve a problem, given the most favorable input. While this provides a lower bound on the algorithm's performance, it is not always a useful indicator, as real-world inputs might not be optimal.

2. Worst-case time complexity: The maximum number of steps an algorithm takes to solve a problem, given the least favorable input. This provides an upper bound on the algorithm's performance, and is often the most useful metric for comparing algorithms, as it guarantees a certain level of efficiency regardless of the input.
3. Average-case time complexity: The average number of steps an algorithm takes to solve a problem, considering all possible inputs. This provides a more realistic estimate of the algorithm's performance in practice but can be more challenging to compute.

In algorithm analysis, we usually focus on the worst-case time complexity, as it ensures a consistent level of performance and allows for meaningful comparisons between different algorithms.

1.1.2. Space Complexity

Space complexity is another crucial aspect of algorithm efficiency that measures the amount of memory an algorithm consumes as a function of the input size. It takes into account the storage required for the input data, temporary variables, and auxiliary data structures used by the algorithm. Similar to time complexity, understanding the space complexity of an algorithm is essential for optimizing code and selecting the most efficient algorithm for a given problem, especially when memory resources are limited.

When analyzing the space complexity of an algorithm, we focus on its growth rate in relation to the input size. This allows us to determine how the algorithm's memory consumption scales as the input size increases, which is crucial for assessing its efficiency in various scenarios.

There are two main components of space complexity:

1. Fixed space: This refers to the memory required for the input data, constants, and simple variables that do not depend on the input size. Fixed space usually does not contribute significantly to the overall space complexity, as it remains constant regardless of the input size.
2. Variable space: This refers to the memory required for data structures, such as arrays or linked lists, that depend on the input size. Variable space is the primary factor contributing to the space complexity of an algorithm, as it directly correlates with the input size.

When calculating the space complexity of an algorithm, we typically focus on the variable space, as it represents the memory consumption that grows with the input size. Similar to time complexity, we express the space complexity of an algorithm using asymptotic notations, which provide a concise and standardized way to describe the growth of functions and compare the efficiency of different algorithms.

1.2. Asymptotic Notations

Asymptotic notations are mathematical tools used to describe the growth of functions, particularly the time and space complexity of algorithms. They provide a concise and standardized way to express how the performance of an algorithm changes as the input size increases, allowing for meaningful comparisons between different algorithms. Asymptotic notations focus on the growth rate of a function and ignore constant factors and lower-order terms, which makes them more suitable for evaluating the efficiency of algorithms in the context of large inputs.

There are three primary asymptotic notations used in algorithm analysis:

- Big O (O),
- Big Omega (Ω)
- Big Theta (Θ).

These asymptotic notations allow us to analyze and compare the efficiency of different algorithms in a standardized manner, focusing on their growth rates rather than the absolute number of steps or memory consumption. By using asymptotic notations, we can gain insights into how well an algorithm will perform under various circumstances and make informed decisions when choosing or designing algorithms for different applications.

1.2.1. Big O

Big O notation describes the upper limit of an algorithm's time or space complexity, giving us a worst-case estimate. In simpler terms, Big O notation tells us the maximum amount of time or memory an algorithm could take based on the input size. It helps us understand the highest possible growth rate of an algorithm's time or space complexity as the input size increases.

1.2.2. Big Ω

Big Omega notation describes the lower limit of an algorithm's time or space complexity, giving us a best-case estimate. In simpler terms, Big Omega notation tells us the minimum amount of time or memory an algorithm could take based on the input size. It helps us understand the lowest possible growth rate of an algorithm's time or space complexity as the input size increases.

1.2.3. Big Θ

Big Theta notation describes the exact growth rate of an algorithm's time or space complexity, giving us an average-case estimate. In simpler terms, Big Theta notation tells us the typical amount of time or memory an algorithm would take based on the input size. It helps us understand the actual growth rate of an algorithm's time or

space complexity as the input size increases, when both the upper and lower limits are the same.

1.3. Mathematical Functions

Mathematical functions play a significant role in algorithm analysis, as they are used to represent the time and space complexity of algorithms. In this context, a function maps the input size of a problem (usually denoted as n) to the number of elementary operations or memory units required by the algorithm. By examining various mathematical functions and their properties, we can better understand the efficiency of different algorithms and how their performance scales with the input size.

1.3.1. Common Functions in Algorithm Analysis

There are several common mathematical functions frequently encountered in algorithm analysis:

1. Constant ($O(1)$): A constant function represents an algorithm whose time or space complexity does not depend on the input size. The performance of such an algorithm remains constant regardless of the input size.
2. Linear ($O(n)$): A linear function represents an algorithm whose time or space complexity increases proportionally with the input size. The performance of such an algorithm degrades linearly as the input size grows.
3. Quadratic ($O(n^2)$): A quadratic function represents an algorithm whose time or space complexity increases with the square of the input size. The performance of such an algorithm degrades rapidly as the input size grows, making it inefficient for large inputs.
4. Cubic ($O(n^3)$): A cubic function represents an algorithm whose time or space complexity increases with the cube of the input size. Similar to quadratic functions, the performance of such an algorithm degrades rapidly with increasing input size, making it inefficient for large inputs.
5. Logarithmic ($O(\log n)$): A logarithmic function represents an algorithm whose time or space complexity increases logarithmically with the input size. The performance of such an algorithm degrades slowly as the input size grows, making it efficient for large inputs.
6. Exponential ($O(2^n)$): An exponential function represents an algorithm whose time or space complexity increases exponentially with the input size. The performance of such an algorithm degrades extremely rapidly as the input size grows, making it impractical for all but the smallest inputs.
7. Factorial ($O(n!)$): A factorial function represents an algorithm whose time or space complexity increases with the factorial of the input size. The performance

of such an algorithm degrades extremely rapidly, making it impractical for all but the smallest inputs.

1.3.2. Growth Rates

Understanding the growth rates of these mathematical functions is essential for comparing the efficiency of different algorithms. The order of growth of a function determines how quickly its value increases as the input size grows. Functions with lower growth rates are generally more efficient than those with higher growth rates, especially for large input sizes.

Here is a ranking of the common functions mentioned above in terms of their growth rates, from the slowest to the fastest:

1. Constant ($O(1)$)
2. Logarithmic ($O(\log n)$)
3. Linear ($O(n)$)
4. Linearithmic ($O(n \log n)$)
5. Quadratic ($O(n^2)$)
6. Cubic ($O(n^3)$)
7. Exponential ($O(2^n)$)
8. Factorial ($O(n!)$)

By understanding the growth rates of these functions, we can make informed decisions when choosing or designing algorithms for different applications, aiming for the most efficient algorithm possible given the problem's constraints.

1.4. Analyzing Loops

Analyzing loops is a crucial aspect of algorithm analysis, as they often dictate the overall time complexity of an algorithm. Loops are used to perform repetitive tasks and can greatly influence the efficiency of an algorithm. Understanding the time complexity of loops allows us to optimize our algorithms and make informed decisions when choosing or designing algorithms for different applications.

1.4.1. Nested Loops

Nested loops occur when a loop is placed inside another loop. The time complexity of nested loops is generally the product of the time complexities of the individual loops.

```
In [12]: def nested_loops_example(n, m):  
        count = 0  
        for i in range(n):  
            for j in range(m):  
                count += 1  
        return count
```

```
n = 4
m = 3
result = nested_loops_example(n, m)
print(f"The operation inside the nested loops was executed {result} times")
```

The operation inside the nested loops was executed 12 times.

In this example, the outer loop runs n times, and the inner loop runs m times. The operation inside the inner loop will be executed $n * m$ times, resulting in a time complexity of $O(nm)$. If both n and m are the same, the time complexity would be $O(n^2)$.

1.4.2. Recursive Loops

Recursive loops occur when a function calls itself to perform a task repeatedly. The time complexity of recursive loops depends on the number of recursive calls made and the complexity of the operations performed inside the recursive function.

Here's a simple example of a recursive function in Python implementing the Fibonacci sequence:

```
In [13]: def fibonacci(n):
          if n == 0:
              return 0
          elif n == 1:
              return 1
          else:
              return fibonacci(n - 1) + fibonacci(n - 2)

          n = 10
          result = fibonacci(n)
          print(f"The {n}-th number in the Fibonacci sequence is {result}.")
```

The 10-th number in the Fibonacci sequence is 55.

In this example, the fibonacci function calls itself twice in each step (except for the base cases when n is 0 or 1). The time complexity of this algorithm is $O(2^n)$, which makes it highly inefficient for large input sizes.

Analyzing loops, whether they are nested or recursive, is important for understanding and optimizing the performance of algorithms. By identifying the time complexity of the loops in our code, we can make better decisions when designing and implementing algorithms and ultimately create more efficient solutions.

1.5. Time Complexity of Searching and Sorting Algorithms

Searching and sorting algorithms are fundamental operations in computer science, and their efficiency can greatly impact the performance of various applications.

Understanding the time complexity of these algorithms allows us to choose the best algorithm for a specific task and optimize our code accordingly. In this section, we will discuss the time complexity of some common searching and sorting algorithms.

These examples of searching and sorting algorithms illustrate the importance of understanding the time complexity of various algorithms. By knowing the time complexity, we can make better decisions when selecting the most suitable algorithm for a specific task or optimizing our code.

1.5.1. Linear Search

Linear search is a simple searching algorithm that iterates through each element in a list until the target element is found or the list is exhausted. The time complexity of linear search is $O(n)$, where n is the number of elements in the list.

```
In [14]: def linear_search(arr, target):
          for i, element in enumerate(arr):
              if element == target:
                  return i
          return -1

arr = [1, 5, 8, 12, 16, 20, 25]
target = 12
result = linear_search(arr, target)
print(f"Linear search found the target {target} at index {result}.")
```

Linear search found the target 12 at index 3.

1.5.2. Binary Search

Binary search is an efficient searching algorithm that works on sorted lists. It repeatedly divides the list in half, comparing the middle element with the target value. If the middle element is equal to the target, the search is successful. If the middle element is greater than the target, the search continues on the left half of the list; if it's less than the target, it continues on the right half. The time complexity of binary search is $O(\log n)$, where n is the number of elements in the list.

```
In [15]: def binary_search(arr, target):
          left, right = 0, len(arr) - 1

          while left <= right:
              mid = (left + right) // 2

              if arr[mid] == target:
                  return mid
              elif arr[mid] < target:
                  left = mid + 1
              else:
                  right = mid - 1

          return -1

arr = [1, 5, 8, 12, 16, 20, 25]
target = 12
result = binary_search(arr, target)
print(f"Binary search found the target {target} at index {result}.")
```

Binary search found the target 12 at index 3.

1.5.3. Bubble Sort

Bubble sort is a simple sorting algorithm that repeatedly steps through the list, comparing adjacent elements and swapping them if they are in the wrong order. The algorithm continues until no more swaps are needed. The time complexity of bubble sort is $O(n^2)$, where n is the number of elements in the list.

```
In [16]: def bubble_sort(arr):
          n = len(arr)

          for i in range(n - 1):
              for j in range(n - 1 - i):
                  if arr[j] > arr[j + 1]:
                      arr[j], arr[j + 1] = arr[j + 1], arr[j]

          example_array = [64, 34, 25, 12, 22, 11, 90]
          print("Original array:", example_array)

          bubble_sort(example_array)
          print("Sorted array:", example_array)
```

```
Original array: [64, 34, 25, 12, 22, 11, 90]
Sorted array: [11, 12, 22, 25, 34, 64, 90]
```

1.5.4. QuickSort

Quicksort is an efficient sorting algorithm that works by selecting a 'pivot' element from the list and partitioning the other elements into two groups: those less than the pivot and those greater than the pivot. The algorithm then recursively sorts the two groups. The average time complexity of quicksort is $O(n \log n)$, where n is the number of elements in the list. However, the worst-case time complexity is $O(n^2)$, which occurs when the pivot selection consistently results in an unbalanced partition.

```
In [17]: def partition(arr, low, high):
          pivot = arr[high]
          i = low - 1

          for j in range(low, high):
              if arr[j] <= pivot:
                  i += 1
                  arr[i], arr[j] = arr[j], arr[i]

          arr[i + 1], arr[high] = arr[high], arr[i + 1]
          return i + 1

          def quick_sort(arr, low, high):
              if low < high:
                  pivot_index = partition(arr, low, high)
                  quick_sort(arr, low, pivot_index - 1)
                  quick_sort(arr, pivot_index + 1, high)
              return arr

          arr = [64, 34, 25, 12, 22, 11, 90]
```

```
result = quick_sort(arr, 0, len(arr) - 1)
print(f"Quicksort sorted the list as: {result}")
```

Quicksort sorted the list as: [11, 12, 22, 25, 34, 64, 90]

2. Mathematical Arguing and Reasoning

Mathematical arguing and reasoning are essential skills in computer science and algorithm analysis, as they allow us to establish the correctness and efficiency of algorithms. In this section, we will explore some common proof techniques and mathematical concepts used in the analysis of algorithms.

2.1. Revision of Algorithm Analysis

Before diving into mathematical arguing and reasoning, let's briefly review the key concepts in algorithm analysis that we discussed earlier. Time complexity is a measure of the number of operations an algorithm takes to complete its tasks as a function of the input size, while space complexity measures the memory consumed by an algorithm. Asymptotic notations, such as Big O, Big Ω , and Big Θ , help us describe the growth of functions and compare the efficiency of different algorithms.

2.1.1. Time Complexity Examples

In this section, we will look at a few simple examples of time complexity to reinforce our understanding of the concept. Lets recall that time complexity is a measure of the number of operations an algorithm takes to complete its tasks as a function of the input size.

Example 1: Sum of an array

Suppose we have an array of integers, and we want to calculate the sum of its elements. We can write a simple algorithm that iterates through the array and adds each element to a running total:

```
In [18]: def array_sum(arr):
        total = 0
        for num in arr:
            total += num
        return total
```

For an array of length n , the algorithm performs a single addition operation for each element, resulting in a total of n operations. Therefore, the time complexity of this algorithm is $O(n)$.

Example 2: Matrix multiplication

Lets consider the problem of multiplying two square matrices of size $n \times n$. A simple algorithm for matrix multiplication involves three nested loops:

```
In [19]: def matrix_multiply(A, B):
          n = len(A)
          result = [[0] * n for _ in range(n)]

          for i in range(n):
              for j in range(n):
                  for k in range(n):
                      result[i][j] += A[i][k] * B[k][j]

          return result
```

In this algorithm, each loop iterates n times, and there is a constant number of operations inside the innermost loop. Therefore, the time complexity of this algorithm is $O(n^3)$.

2.1.2. Space Complexity Examples

Now, let's look at some simple examples of space complexity. Space complexity measures the memory consumed by an algorithm as a function of the input size.

Example 1: Reversing a string

Suppose we want to reverse a string. One straightforward approach is to create a new empty string and append each character from the original string in reverse order:

```
In [20]: def reverse_string(s):
          reversed_s = ""
          for i in range(len(s)-1, -1, -1):
              reversed_s += s[i]
          return reversed_s
```

In this example, the space complexity is determined by the additional memory used for the new reversed string. Since the length of the reversed string is equal to the length of the input string, the space complexity is $O(n)$, where n is the length of the input string.

Example 2: Finding the sum of an array

Let's consider another example where we want to find the sum of all elements in an array:

```
In [21]: def sum_array(arr):
          total = 0
          for element in arr:
              total += element
          return total
```

In this case, we only use a single variable `total` to store the sum of the elements in the array. No matter how large the input array is, the memory usage for this variable remains constant. Thus, the space complexity of this algorithm is $O(1)$, which means it uses constant space and does not depend on the input size.

By examining these examples, you can better understand how time complexity and space complexity are used to analyze the efficiency of algorithms, enabling you to make informed decisions when selecting or designing algorithms for different applications.

2.2. Mathematical Arguing and Reasoning

Mathematical arguing and reasoning involve the use of logical statements, precise definitions, and mathematical techniques to establish the validity of various claims, such as the correctness and efficiency of algorithms.

2.2.1. Proof Techniques

There are several common proof techniques used in computer science:

1. **Direct proof:** In a direct proof, we demonstrate the truth of a statement by a series of logical deductions based on known facts or previously proven statements.
2. **Proof by contradiction:** In a proof by contradiction, we assume the opposite of the statement we want to prove and then show that this assumption leads to a contradiction.
3. **Proof by induction:** Proof by induction is a technique used to prove statements involving a positive integer variable. It consists of two steps: the base case and the inductive step. In the base case, we prove that the statement is true for a specific value, usually the smallest value of the variable. In the inductive step, we assume that the statement is true for a certain value and then prove that it is also true for the next value.

2.2.2. Inductive Proofs

Inductive proofs are particularly useful in the analysis of algorithms, as they often involve proving properties of algorithms that hold for all input sizes. Let's take a look at a simple example using the concept of mathematical induction.

Example: Prove that the sum of the first n odd numbers is equal to n^2 .

Base case ($n = 1$):

The sum of the first odd number (1) is 1, which is equal to 1^2 .

The statement is true for $n = 1$.

Inductive step:

Assume that the statement is true for $n = k$, i.e., the sum of the first k odd numbers is equal to k^2 .

We need to prove that the statement is also true for $n = k + 1$.

The sum of the first $k + 1$ odd numbers can be written as:

$$\rightarrow (1 + 3 + 5 + \dots + (2k - 1)) + (2k + 1)$$

By our induction assumption, the sum of the first k odd numbers is k^2 , so the expression becomes:

$$\rightarrow k^2 + (2k + 1)$$

Now, let's simplify the expression:

$$\rightarrow k^2 + (2k + 1) = k^2 + 2k + 1 = (k + 1)^2$$

This shows that the sum of the first $k + 1$ odd numbers is equal to $(k + 1)^2$.

Thus, by the principle of mathematical induction, the statement is true for all positive integers n .

2.3. Logarithms Review

Logarithms are an essential concept in computer science and the analysis of algorithms. They are particularly useful for describing the growth rates of functions, especially in the context of time complexity. In this section, we will provide a brief review of logarithms and their properties, followed by an example of how logarithms are applied in computer science.

2.3.1. Logarithm Basics and Properties

A logarithm is the inverse operation of exponentiation. In other words, logarithms tell us the power to which a given base must be raised to produce a specific number. The logarithm of a number x with base b is denoted as $\log_b(x)$ and can be defined by the equation:

$$b^{(\log_b(x))} = x$$

For example, $\log_2(8) = 3$ because $2^3 = 8$. Here, the base is 2, and the logarithm tells us that we need to raise 2 to the power of 3 to obtain the number 8.

Some essential properties of logarithms are:

1. $\log_b(1) = 0$, because any number raised to the power of 0 is 1.
2. $\log_b(x * y) = \log_b(x) + \log_b(y)$, which shows that logarithms convert multiplication to addition.
3. $\log_b(x / y) = \log_b(x) - \log_b(y)$, which shows that logarithms convert division to subtraction.
4. $\log_b(x^p) = p * \log_b(x)$, which shows that logarithms can move exponents to coefficients.

2.3.2. Logarithms in Computer Science

Logarithms frequently appear in the analysis of algorithms, especially in cases where a problem is repeatedly divided into smaller subproblems. One common example is the binary search algorithm.

Binary search is an efficient algorithm for finding a specific value in a sorted list. The algorithm works by repeatedly dividing the list in half until the desired value is found or it's determined that the value is not in the list. At each step, the algorithm eliminates half of the remaining elements, so the number of elements to search is reduced by a factor of 2.

Here's a simple Python implementation of binary search:

```
In [22]: def binary_search(arr, target):
        low, high = 0, len(arr) - 1

        while low <= high:
            mid = (low + high) // 2
            mid_val = arr[mid]

            if mid_val == target:
                return mid
            elif mid_val < target:
                low = mid + 1
            else:
                high = mid - 1

        return -1
```

For a list of length n , the algorithm takes at most $\log_2(n)$ steps to find the target value or determine that it's not in the list. Therefore, the time complexity of binary search is $O(\log n)$.

Understanding logarithms and their properties is crucial in computer science, as it helps to analyze and optimize algorithms that involve repeated division of problems or data structures. Mastering logarithms will enable you to design more efficient solutions to complex problems and make informed decisions when selecting algorithms for various applications.

3. Computational Machine Architectures

Computational machine architectures define how computers are organized and how they execute instructions. Understanding these architectures helps programmers design efficient algorithms and optimize software performance. In this section, we will discuss the Von Neumann architecture, CPU architecture, CISC vs. RISC architectures, and the role of the Graphics Processing Unit (GPU) in modern computing.

3.1. Von Neumann Architecture

The Von Neumann architecture is a widely used model for computer design, named after mathematician and computer scientist John von Neumann. This architecture consists of several key components:

1. Memory: Stores both instructions and data.
2. Central Processing Unit (CPU): Executes instructions and processes data.
3. Input/Output (I/O) devices: Allow communication between the computer and external devices, such as keyboards, mice, and monitors.

The Von Neumann architecture operates on the stored-program concept, which means that programs and data are stored together in memory. This allows the CPU to fetch, decode, and execute instructions one after the other in a sequential manner.

3.1.1. Components and Functions

Memory: The memory component of the Von Neumann architecture stores both instructions and data. It is typically divided into an addressable array of storage cells, each of which can hold a fixed amount of information (usually one byte).

Central Processing Unit (CPU): The CPU is responsible for executing instructions and processing data. It contains the Arithmetic Logic Unit (ALU), which performs arithmetic and logical operations, and the Control Unit (CU), which manages the flow of data between the memory, ALU, and I/O devices.

Input/Output (I/O) devices: These devices enable communication between the computer and external components. Input devices, such as keyboards and mice, allow users to enter data, while output devices, like monitors and printers, display or produce the results of computations.

3.1.2. Memory Hierarchy

The memory hierarchy in the Von Neumann architecture is organized into multiple layers, each with different performance characteristics and capacities. The hierarchy typically includes:

1. Registers: Small, fast storage locations within the CPU.
2. Cache memory: Faster, more expensive memory located closer to the CPU. It is used to store frequently accessed data and instructions.
3. Main memory (RAM): Stores program instructions and data while the computer is running.
4. Secondary storage: Non-volatile storage devices, such as hard drives and solid-state drives, used for long-term data storage.

The memory hierarchy aims to provide an optimal balance between speed and storage capacity, allowing for efficient data access and processing.

3.2. CPU Architecture

The CPU is the heart of the computer, responsible for processing instructions and data. Modern CPUs are designed with several features that enhance their performance, such as pipelining and cache memory.

3.2.1. Pipelining

Pipelining is a technique that improves the CPU's throughput by allowing it to process multiple instructions simultaneously. In a pipelined CPU, the execution of instructions is divided into several stages, with each stage handling a different part of the instruction. This allows the CPU to work on several instructions at the same time, increasing the overall speed of execution.

3.2.2. Cache Memory

Cache memory is a small, fast memory located close to the CPU. It is used to store frequently accessed data and instructions, reducing the time needed for the CPU to access them. Cache memory can be organized into different levels (L1, L2, and L3) based on their proximity to the CPU and their capacity. L1 cache is the smallest and fastest, while L3 cache is larger and slower but still faster than main memory (RAM).

Cache memory uses various strategies to optimize its performance, such as:

- Spatial locality: Data that is close together in memory is likely to be accessed together.
- Temporal locality: Data that has been accessed recently is likely to be accessed again in the near future.

These strategies help predict which data and instructions are likely to be needed next, allowing the cache to prefetch them and minimize the time the CPU spends waiting for data.

3.3. CISC vs RISC Architectures

CISC (Complex Instruction Set Computer) and RISC (Reduced Instruction Set Computer) are two contrasting approaches to CPU design. They differ in their instruction sets, complexity, and execution strategies.

3.3.1. Advantages and Disadvantages

CISC architectures feature a large number of complex instructions that can perform multiple operations in a single instruction. This can lead to more efficient use of memory and a smaller program size. However, the complexity of CISC instructions can result in longer instruction execution times and increased power consumption.

RISC architectures, on the other hand, use a smaller set of simple instructions that execute quickly. Each instruction typically performs a single operation, making them easier to optimize for performance. RISC architectures often have a higher instruction throughput, resulting in faster overall execution. However, RISC programs can be larger in size, as more instructions are needed to perform the same tasks as a CISC program.

3.3.2. Evolution and Current Trends

Over the years, the distinction between CISC and RISC architectures has become less clear, as modern CPUs incorporate features from both designs to optimize performance. For example, many CISC processors now include RISC-like execution units to improve their instruction throughput. Similarly, RISC processors have adopted some CISC features to increase their capabilities and efficiency.

3.4. GPU: Graphical Process Unit

A Graphics Processing Unit (GPU) is a specialized processor designed to handle the complex calculations required for rendering graphics, particularly in 3D environments. GPUs have become increasingly important in modern computing due to their ability to perform parallel computations efficiently, making them suitable for a wide range of applications beyond graphics.

3.4.1. GPU Architecture

GPU architecture is designed to handle parallel processing, with many simple cores working together to perform calculations simultaneously. This makes GPUs particularly well-suited for tasks that can be divided into smaller, independent operations, such as rendering graphics or performing mathematical simulations.

3.4.2. Applications in Computing

While GPUs were initially developed for graphics processing, they have since found applications in a variety of computing tasks, including:

- Scientific simulations: GPUs can perform complex calculations quickly, making them ideal for simulating physical systems, weather patterns, and more.
- Machine learning and artificial intelligence: GPUs are widely used in training neural networks and other machine learning models due to their parallel processing capabilities.
- Cryptocurrency mining: GPUs are often used to perform the complex calculations required for mining cryptocurrencies like Bitcoin and Ethereum.

3.5. Building Parallel Algorithms

Parallel algorithms are designed to take advantage of multiple processing units working simultaneously to solve a problem. These algorithms can be executed on architectures such as multi-core CPUs and GPUs, which are capable of performing parallel computations. By dividing tasks into smaller, independent operations and distributing them across multiple processing units, parallel algorithms can significantly reduce computation time and improve overall performance.

3.5.1. Parallelism Techniques

There are several types of parallelism used in algorithm design:

- Data parallelism: Involves dividing the input data into smaller chunks and processing each chunk simultaneously. This is particularly suitable for tasks where the same operation is applied to different parts of the input data, such as image processing or matrix multiplication.
- Task parallelism: Focuses on breaking down a problem into smaller, independent tasks that can be executed concurrently. This approach is useful when different tasks within a problem can be performed independently of each other, such as solving a set of unrelated equations.
- Pipeline parallelism: Similar to pipelining in CPUs, this type of parallelism involves dividing a problem into stages, with each stage being executed by a different processing unit. As each stage completes, it passes its results to the next stage for further processing. This approach is useful for tasks that involve a sequence of operations, such as video encoding or speech recognition.

3.5.2. Multithreading and Concurrency

When designing parallel algorithms, there are several factors to consider:

- Scalability: A parallel algorithm should be able to efficiently utilize increasing numbers of processing units, resulting in improved performance as more resources are added.
- Load balancing: To achieve optimal performance, work should be evenly distributed among the available processing units, avoiding situations where some units are idle while others are overloaded.

- **Communication and synchronization:** Parallel algorithms often require communication and synchronization between processing units to share data and coordinate tasks. It's essential to minimize the overhead associated with these operations to avoid negatively impacting performance.
- **Granularity:** The granularity of a parallel algorithm refers to the size of the tasks it assigns to processing units. Fine-grained algorithms divide tasks into many small operations, while coarse-grained algorithms use fewer, larger tasks. The choice of granularity can significantly impact the efficiency and performance of a parallel algorithm.

By considering these factors and carefully designing parallel algorithms, developers can create efficient, high-performance solutions that leverage the full potential of modern computational architectures.

4. Graph Data Structure and Search Algorithms

Graphs are a versatile and powerful data structure used to represent relationships between objects, making them an essential tool in computer science, mathematics, and many other fields. They consist of nodes (vertices) and edges connecting these nodes, reflecting the connections between the elements in the graph. In this section, we will explore the basic concepts of graph data structures, their real-life applications, and search algorithms commonly used to traverse and analyze them.

4.1. Graph Data Structure Overview

4.1.1. Terminology and Representation

A graph G is defined as an ordered pair $G(V, E)$, where V is the set of vertices and E is the set of edges connecting the vertices. Each edge can be represented as a tuple (u, v) , where u and v are the vertices it connects.

Graphs can be classified as directed or undirected. In undirected graphs, edges have no direction, meaning the relationship between nodes is bidirectional. In directed graphs (also called digraphs), edges have a direction, indicating a one-way relationship between nodes.

There are several ways to represent graphs in computer programs, with two of the most common being adjacency lists and adjacency matrices:

- **Adjacency list:** For each vertex, a list of its adjacent vertices is maintained. In Python, this can be implemented using dictionaries, with keys representing vertices and values being lists of adjacent vertices.

```
In [23]: graph = {
    'A': ['B', 'C'],
    'B': ['A', 'D'],
    'C': ['A', 'E'],
    'D': ['B', 'E'],
    'E': ['C', 'D']
}
```

- Adjacency matrix: A two-dimensional matrix is used to represent the graph, with the rows and columns representing vertices, and the value at `matrix[i][j]` indicating the presence (1) or absence (0) of an edge between vertices `i` and `j`.

```
In [24]: graph_matrix = [
    [0, 1, 1, 0, 0],
    [1, 0, 0, 1, 0],
    [1, 0, 0, 0, 1],
    [0, 1, 0, 0, 1],
    [0, 0, 1, 1, 0]
]
```

4.1.2. Types of Graphs

There are many types of graphs, some of which include:

- Weighted graphs: Edges are assigned weights to represent attributes like distance, cost, or time. The adjacency matrix representation is commonly used for weighted graphs, with the value at `matrix[i][j]` representing the weight of the edge between vertices `i` and `j`.
- Trees: A special type of graph where there are no cycles, and each node has at most one parent.
- Bipartite graphs: Vertices can be partitioned into two disjoint sets, with every edge connecting a vertex from one set to a vertex from the other set.
- Cyclic and acyclic graphs: A graph is cyclic if it contains at least one cycle, and acyclic if it contains no cycles.

4.2. Graph Applications in Real Life

Graphs have numerous applications across various domains, some of which include:

4.2.1. Social Networks

Graphs are used to model social networks, with vertices representing people and edges representing connections or relationships between them. Analyzing these graphs can provide valuable insights into the structure and dynamics of social networks, such as identifying influential individuals or detecting communities.

4.2.2. Navigation Systems

Graphs can be employed to model transportation networks, with vertices representing locations and edges representing roads or connections between locations. Weighted graphs, in particular, are useful here, as edge weights can represent distances, travel times, or costs. Algorithms such as Dijkstra's or A* can be applied to find the shortest or most efficient path between two points on the graph.

4.3. Search Algorithms on Graph Data Structure

There are several search algorithms that can be used to traverse and analyze graph data structures. In this section, we will discuss some of the most common ones, including Breadth-first Search (BFS), Depth-first Search (DFS), Dijkstra's Algorithm, and A* Algorithm.

4.3.1. Breadth-first Search (BFS)

BFS is a graph traversal algorithm that explores all the vertices of a graph in breadth-first order, meaning that it visits all the neighbors of a vertex before moving on to their neighbors. BFS can be implemented using a queue data structure.

The following is a simple Python implementation of BFS on an undirected graph:

```
In [25]: from collections import deque

def bfs(graph, start):
    visited = set()
    queue = deque([start])

    while queue:
        vertex = queue.popleft()
        if vertex not in visited:
            visited.add(vertex)
            queue.extend(neighbor for neighbor in graph[vertex] if neighbor
                        not in visited)

    return visited

graph = {
    'A': ['B', 'C'],
    'B': ['A', 'D'],
    'C': ['A', 'E'],
    'D': ['B', 'E'],
    'E': ['C', 'D']
}

print(bfs(graph, 'A'))

{'A', 'D', 'B', 'E', 'C'}
```

4.3.2. Depth-first Search (DFS)

DFS is another graph traversal algorithm that explores vertices in a depth-first manner, visiting a vertex and recursively exploring its neighbors as deeply as

possible before backtracking. DFS can be implemented using recursion or a stack data structure.

The following is a simple Python implementation of DFS on an undirected graph using recursion:

```
In [26]: def dfs(graph, vertex, visited=None):
        if visited is None:
            visited = set()
            visited.add(vertex)

        for neighbor in graph[vertex]:
            if neighbor not in visited:
                dfs(graph, neighbor, visited)

        return visited

graph = {
    'A': ['B', 'C'],
    'B': ['A', 'D'],
    'C': ['A', 'E'],
    'D': ['B', 'E'],
    'E': ['C', 'D']
}

print(dfs(graph, 'A'))

{'A', 'D', 'B', 'E', 'C'}
```

In the following sections, we will discuss more advanced search algorithms like Dijkstra's and A* that can be used to find the shortest path between two points in a weighted graph.

4.3.3. Dijkstra's Algorithm

Dijkstra's Algorithm is a popular graph search algorithm used to find the shortest path between two vertices in a weighted graph with non-negative edge weights. The algorithm maintains a set of unvisited vertices, and for each vertex, it keeps track of the shortest known distance from the starting vertex. The algorithm repeatedly selects the unvisited vertex with the smallest known distance, updates the distances of its neighbors, and marks the selected vertex as visited.

The following is a simple Python implementation of Dijkstra's Algorithm on a directed graph:

```
In [27]: import heapq

def dijkstra(graph, start):
    distances = {vertex: float('infinity') for vertex in graph}
    distances[start] = 0
    queue = [(0, start)]

    while queue:
        current_distance, current_vertex = heapq.heappop(queue)
```

```

        if current_distance > distances[current_vertex]:
            continue

        for neighbor, weight in graph[current_vertex].items():
            distance = current_distance + weight

            if distance < distances[neighbor]:
                distances[neighbor] = distance
                heapq.heappush(queue, (distance, neighbor))

    return distances

graph = {
    'A': {'B': 1, 'C': 4},
    'B': {'A': 1, 'C': 2, 'D': 5},
    'C': {'A': 4, 'B': 2, 'D': 1},
    'D': {'B': 5, 'C': 1}
}

print(dijkstra(graph, 'A'))

{'A': 0, 'B': 1, 'C': 3, 'D': 4}

```

4.3.4. A* Algorithm

The A* Algorithm is a search algorithm that finds the shortest path between two vertices in a weighted graph by combining the benefits of Dijkstra's Algorithm and a heuristic function. The heuristic function estimates the cost of the shortest path from the current vertex to the goal vertex, allowing the algorithm to explore more promising paths first. A* is widely used in applications such as pathfinding in games and navigation systems.

The following is a simple Python implementation of the A* Algorithm on a grid with Manhattan distance as the heuristic:

```

In [28]: import heapq

def a_star(grid, start, end):
    def heuristic(a, b):
        return abs(a[0] - b[0]) + abs(a[1] - b[1])

    def neighbors(pos):
        for dx, dy in ((1, 0), (-1, 0), (0, 1), (0, -1)):
            x, y = pos[0] + dx, pos[1] + dy
            if 0 <= x < len(grid) and 0 <= y < len(grid[0]) and grid[x][y] != '#':
                yield (x, y)

    visited = set()
    queue = [(heuristic(start, end), 0, start)]
    costs = {start: 0}

    while queue:
        _, cost, current = heapq.heappop(queue)

        if current == end:
            return cost

```



```

        if current in visited:
            continue

        visited.add(current)

        for neighbor in neighbors(current):
            new_cost = cost + 1
            if neighbor not in costs or new_cost < costs[neighbor]:
                costs[neighbor] = new_cost
                heapq.heappush(queue, (new_cost + heuristic(neighbor, end)

        return None

grid = [
    ['S', '.', '.', '.', 'X'],
    ['.', 'X', '.', 'X', '.'],
    ['.', '.', '.', '.', '.'],
    ['X', 'X', '.', 'X', '.'],
    ['.', '.', '.', 'E', 'X']
]

start = (0, 0)
end = (4, 3)

print(a_star(grid, start, end))

```

7

In this example, we have a 5x5 grid where 'S' is the start position, 'E' is the end position, and 'X' represents obstacles. The code will return the shortest path cost from the start position to the end position, avoiding obstacles, using the A* Algorithm with the Manhattan distance heuristic.

5. Sources

This notebook on "Algorithms and Data Structures" has been created using the following resources:

- Cormen, T. H., Leiserson, C. E., Rivest, R. L., & Stein, C. (2009). Introduction to Algorithms (3rd ed.). MIT Press.
- Sedgewick, R., & Wayne, K. (2011). Algorithms (4th ed.). Addison-Wesley Professional.
- Skiena, S. S. (2020). The Algorithm Design Manual (3rd ed.). Springer.
- Grossman, R., & Tuma, J. (2019). A Programmer's Guide to Computer Science: A Virtual Degree for the Self-Taught Developer (Vol. 1). CreateSpace Independent Publishing Platform.
- Laakmann McDowell, G. (2015). Cracking the Coding Interview: 189 Programming Questions and Solutions (6th ed.). CareerCup.
- Khan Academy. (n.d.). Algorithms. Retrieved from <https://www.khanacademy.org/computing/computer-science/algorithms>
- Brilliant. (n.d.). Algorithms. Retrieved from <https://brilliant.org/courses/algorithms/>

- LeetCode. (n.d.). LeetCode. Retrieved from <https://leetcode.com/>
- GeeksforGeeks. (n.d.). Data Structures. Retrieved from <https://www.geeksforgeeks.org/data-structures/>
- GeeksforGeeks. (n.d.). Algorithms. Retrieved from <https://www.geeksforgeeks.org/fundamentals-of-algorithms/>