# CASE STUDY REPORT

## GROUP: 1

**Topic:** Brief description of 10-15 String matching algorithm and comparisons view of all and also 5 research papers published in string matching algorithm.

Under The Supervision of

Dr. Sukant Kishoro Bisoy

(Department of Computer Science & Engineering)

# Submitted by:

| NAME | Purnima  Pattnaik |
|---|---|
| REGD NO | 2201020960 |
| ROLL NO | CSE22292 |
| BRANCH | CSE |
| SEMESTER | 4th |
| GROUP | 4 'A' |

# **<u>ACKNOWLEDGEMENT</u>**

I would like to extend sincere and heartfelt thanks towards all those who have helped me in making this project. With their active guidance, help, cooperation and encouragement, I wouldn't  have been able to present the project on time.

 I extend my sincere gratitude to my teacher Mr. Dr. Sukant Kishoro Bisoy for the moral support and guidance during the tenure of my project. I also acknowledge with a deep sense of reverence, our gratitude towards my parents and other faculty members of the college for their valuable suggestions given to me in completing the project.

<u>Your's Sincerely ,</u>

 Purnima  Pattnaik

## **Abstract**

A string matching algorithm is a computational method used to find occurrences of a pattern within a longer text or string. Abstractly, these algorithms aim to efficiently determine if and where a smaller string (pattern) appears within a larger string (text). Common algorithms include brute-force, Knuth-Morris-Pratt, Boyer-Moore, and Rabin-Karp. These algorithms employ different strategies, such as comparing characters sequentially or utilizing pre-processing techniques to achieve faster matching times. The choice of algorithm often depends on factors such as the size of the text, the length of the pattern, and the desired efficiency of the search process.

 String matching algorithms are fundamental in computer science, utilized in various applications such as text processing, bioinformatics, and data mining. These algorithms aim to locate occurrences of a pattern within a larger text, a task crucial for tasks like search engines, plagiarism detection, and genome sequencing. string matching algorithms play a crucial role in numerous applications, ranging from text processing to bioinformatics and beyond. The diverse range of algorithms and techniques available cater to different requirements and scenarios, offering efficient solutions for tasks such as pattern searching, substring matching, and multiple pattern searching. Understanding the principles and characteristics of these algorithms enables developers to choose the most suitable approach for their specific applications, balancing factors such as time complexity, space complexity, and scalability.

# Contents

## Introduction to String Matching Algorithms:

String matching algorithms are fundamental tools in computer science used to efficiently locate occurrences of a pattern within a larger text or string. These algorithms play a crucial role in various applications such as text editing, data compression, bioinformatics, and information retrieval. Each algorithm employs unique techniques and approaches to achieve efficient pattern matching, considering factors like time complexity, space complexity, and practical performance.

## Why do we need String Matching Algorithms?

String matching algorithms are indispensable tools in computer science and various other fields due to several compelling reasons:

**Information Retrieval**: In applications such as search engines, databases, and information retrieval systems, string matching algorithms are crucial for efficiently locating relevant documents, records, or data based on user queries or specified patterns.

**Text Processing**: String matching algorithms enable tasks such as text editing, parsing, and manipulation. They facilitate operations like finding and replacing substrings, extracting information, and formatting text according to specified patterns or rules.

**Pattern Recognition**: In fields like bioinformatics, image processing, and signal analysis, string matching algorithms are essential for recognizing and analyzing patterns within sequences of data. This includes tasks such as DNA sequence alignment, image recognition, and speech processing.

**Data Compression**: String matching algorithms play a key role in data compression techniques such as dictionary-based and pattern-based compression. They help identify repetitive patterns within data, enabling more efficient storage and transmission of information.

**Security and Cryptography**: In cybersecurity and cryptography, string matching algorithms are used for tasks such as intrusion detection, malware detection, and password cracking. They help identify known patterns or signatures associated with malicious activities.

**Natural Language Processing**: In applications involving text analysis and language processing, string matching algorithms are utilized for tasks like named entity recognition, sentiment analysis, and language translation. They enable the identification and processing of specific words, phrases, or linguistic patterns.

# RABIN-KARP

Rabin-Karp is a string-searching algorithm that finds all occurrences of a given pattern in a text. It's based on hashing and is particularly useful for searching for patterns in large texts efficiently. The algorithm works by comparing hash values of the pattern and substrings of the text, and if the hash values match, it performs a full comparison to confirm the match. This algorithm has a time complexity of $O(n+m)$ in the average case, where n is the length of the text and m is the length of the pattern. It's named after its inventors, Michael O. Rabin and Richard M. Karp, and was introduced in 1987.

## AlgorithmDescription

The Rabin-Karp algorithm is a string-searching algorithm that uses hashing to efficiently find occurrences of a pattern within a text. Here's a high-level description of the algorithm:

**Preprocessing**: Calculate the hash value of the pattern and the hash values of all possible substrings of the text with the same length as the pattern. This step is typically done using a rolling hash function to efficiently calculate hash values for substrings.

**Matching***:* Slide the pattern over the text, comparing the hash value of the pattern with the hash values of substrings of the text. If the hash values match, perform a full comparison of the pattern with the substring to confirm the match. If there's a match, record the position of the match.

**RollingHash**: As you slide the pattern over the text, update the hash value of the current substring efficiently using a rolling hash function. This allows for constant-time updates to the hash value as you move from one substring to the next.

**HandlingHash Collisions**: Since hash collisions are possible, if the hash values of the pattern and the substring match, perform a full comparison of the pattern with the substring to confirm the match. This step ensures that false positives due to hash collisions are eliminated.

**Multiple Matches**: Continue the process until you've checked all possible substrings of the text. Record the positions of all matches found.

**TimeComplexity**: The average-case time complexity of the Rabin-Karp algorithm is $O(n + m)$, where n is the length of the text and m is the length of the pattern. This makes it efficient for searching for patterns in large texts.

Overall, the Rabin-Karp algorithm is particularly useful when you need to search for multiple occurrences of a pattern within a text efficiently. It's commonly used in applications such as plagiarism detection, DNA sequence matching, and text processing.

# *Code*

```c
#include <stdio.h>
#include <string.h>
#define d 256
voidsearch(char pattern[], char text[], int q) {
    intM=strlen(pattern);
    intN=strlen(text);
    inti, j;
    intp=0;
    intt=0;
    inth=1;
    for (i=0; i< M -1; i++)
        h = (h * d) % q;
    for (i=0; i< M; i++) {
        p = (d * p + pattern[i]) % q;
        t = (d * t + text[i]) % q;
    }
    for (i=0; i<= N - M; i++) {
        if (p == t) {
            for (j =0; j < M; j++) {
                if (text[i+ j] != pattern[j])
                    break;
            }
            if (j == M)
                printf("Pattern found at index %d\n", i);
        }
        if (i< N - M) {
            t = (d * (t - text[i] * h) + text[i+ M]) % q;
            if (t <0)
                t = (t + q);
        }
    }
}
intmain() {
    chartext[] ="AABAACAADAABAABA";
    charpattern[] ="AABA";
    intq=101;
    search(pattern, text, q);
    return0;
}
```

## *Output*

```
Pattern found at index 0
Pattern found at index 9
Pattern found at index 12
```

**7 |** P a g e

Text : A A B A A C A A D A A B A A B A

Pattern : A A B A

A A B A                    A A B A

A A B A A C A A D A A B A A B A
0   1   2   3   4   5   6   7   8   9   10  11  12  13  14  15

                              A A B A

Pattern Found at 0, 9 and 12

These resources provide valuable insights into the Rabin-Karp algorithm, its theoretical foundations, practical implementations, and applications in various domains. You can access these papers through academic databases such as IEEE Xplore, ACM Digital Library, or Google Scholar

**TimeComplexity**:

The average-case time complexity of the Rabin-Karp algorithm is $O(n + m)$, where n is the length of the text and m is the length of the pattern. This makes it efficient for searching for patterns in large texts.

Overall, the Rabin-Karp algorithm is particularly useful when you need to search for multiple occurrences of a pattern within a text efficiently. It's commonly used in applications such as plagiarism detection, DNA sequence matching, and text processing.

# Knuth-Morris-Pratt (KMP)

The Knuth-Morris-Pratt (KMP) algorithm is a string searching algorithm that efficiently finds occurrences of a pattern within a larger text. It was invented by Donald Knuth and Vaughan Pratt, and independently by James H. Morris. The algorithm improves on the basic brute-force approach by utilizing information about the pattern itself to avoid unnecessary comparisons.

## Algorithm Description:

Preprocessing: The KMP algorithm preprocesses the pattern to construct an auxiliary array, often referred to as the "failure function" or "prefix function". This array helps to determine the maximum length of proper prefix of the pattern that is also a suffix.

Searching: During the searching phase, the algorithm compares the characters of the pattern with the characters of the text, but instead of backtracking on mismatch, it uses the information from the failure function to shift the pattern efficiently

## How it works:

The preprocessing phase calculates the "failure function", which indicates the number of characters to skip ahead in the pattern whenever a mismatch occurs.

In the searching phase, when a mismatch occurs, the algorithm uses the information from the failure function to determine the new position in the pattern from which to resume comparisons, avoiding redundant comparisons.

## Code for the algorithm

```java
public class KMP {
    public static void main(String[] args) {
        String text="ababcabcabababd";
        String pattern="ababd";
        search(text, pattern);
    }

    public static void search(String text, String pattern) {
        int[] lps=computeLPSArray(pattern);
        int n=text.length();
        int m=pattern.length();
        int i=0; // index for text[]
        int j=0; // index for pattern[]
        while (i<n) {
            if (pattern.charAt(j) ==text.charAt(i)) {
                i++;
                j++;
            }
            if (j==m) {
                System.out.println("Pattern found at index "+ (i-j));
                j=lps[j-1];
            } else if (i<n&&pattern.charAt(j) !=text.charAt(i)) {
                if (j!=0)
```

```
                j=lps[j-1];
            else
                i++;
        }
    }
}

    publicstaticint[] computeLPSArray(Stringpattern) {
        intm=pattern.length();
        int[] lps=newint[m];
        intlen=0;
        inti=1;
        lps[0] =0; // lps[0] is always 0
        while (i<m) {
            if (pattern.charAt(i) ==pattern.charAt(len)) {
                len++;
                lps[i] =len;
                i++;
            } else {
                if (len!=0) {
                    len=lps[len-1];
                } else {
                    lps[i] =0;
                    i++;
                }
            }
        }
        returnlps;
    }
}
```
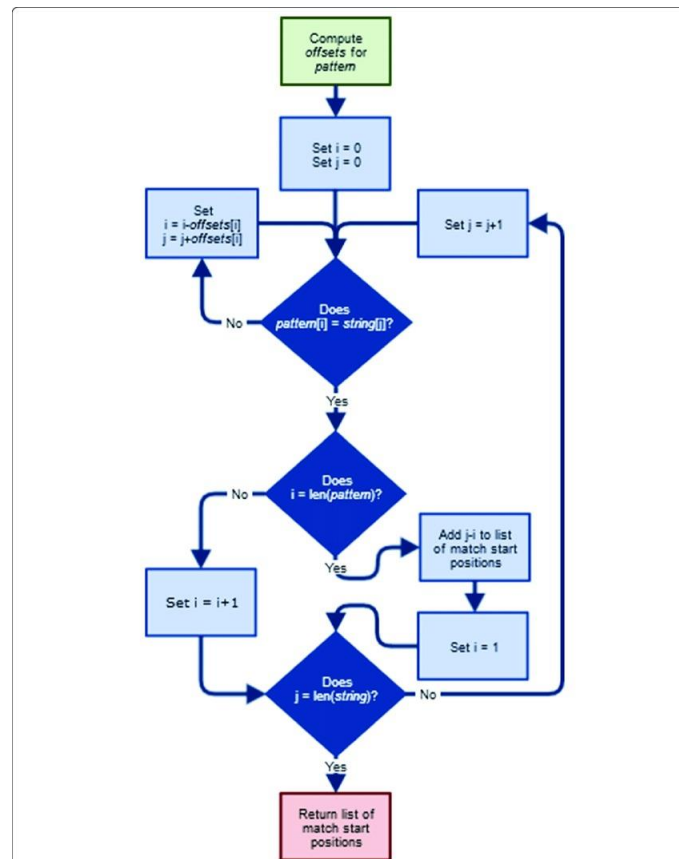
## Time Complexity:

The preprocessing phase of KMP runs in O(m) time, where m is the length of the pattern.

The searching phase runs in O(n) time, where n is the length of the text.

Thus, overall time complexity is O(m + n).

## Result of the Algorithm:

The KMP algorithm returns the indices of all occurrences of the pattern within the text.

## Efficiency:

KMP is efficient for searching patterns in text, especially when the pattern is long or when the pattern contains repetitive elements.

Its time complexity is better than the brute-force approach, particularly for large texts and patterns.

# AHO-CORASICK ALGORITHM

The Aho-Corasick algorithm is a string searching algorithm that efficiently finds all occurrences of a set of keywords within a given text. It was developed by Alfred V. Aho and Margaret J. Corasick in 1975. This algorithm is particularly useful when you need to search for multiple patterns simultaneously.
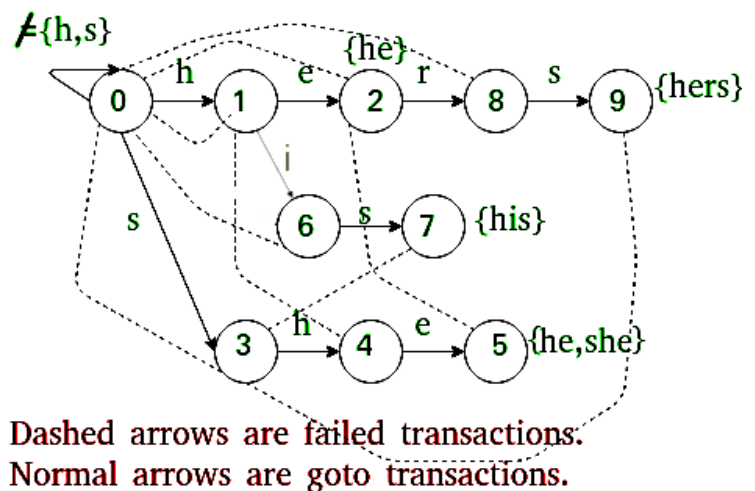
## ALGORITHM DESCRIPTION

**1.Construct Trie:** The algorithm constructs a finite state machine called a Trie (prefix tree) from the set of keywords. Each node in the Trie represents a prefix of one of the keywords.

**2.Add Failure Links:** After constructing the Trie, the algorithm adds additional edges called "failure links." These links help efficiently navigate the Trie during the search process. A failure link from a node points to the longest proper suffix of the current node that is also a prefix of one of the keywords. This allows the algorithm to quickly backtrack when a mismatch occurs.

**3.Search:** To search for keywords in a given text, the algorithm starts at the root of the Trie and processes each character of the text. It follows the edges of the Trie corresponding to the characters in the text.

**4.Output:** The algorithm outputs all occurrences of the keywords in the text along with their positions.

## TIME COMPLEXITY OF ALGORITHM :

- ➢ **The time complexity of the Aho-Corasick algorithm is O(N + L + Z), where**

- ➢ **N is the length of the text being searched**
- ➢ **L is the total length of all keywords in the dictionary (combined length of all patterns)**
- ➢ **Z is the number of occurrences of all keywords found in the text**

### CODE:

```
#include <stdio.h>
#include <string.h>
#include <stdbool.h>
#include <stdlib.h>
```

```c
#include <limits.h>

#define MAXS 500
#define MAXC 26

int out[MAXS];
int f[MAXS];
int g[MAXS][MAXC];

intbuildMatchingMachine(chararr[][MAXS], int k) {
    memset(out, 0, sizeof(out));
    memset(g, -1, sizeof(g));
    intstates=1;

    for (inti=0; i< k; ++i) {
        constchar*word =arr[i];
        intcurrentState=0;

        for (intj=0; word[j]; ++j) {
            intch= word[j] -'a';
            if (g[currentState][ch] ==-1)
                g[currentState][ch] = states++;
            currentState= g[currentState][ch];
        }

        out[currentState] |= (1<<i);
    }

    for (intch=0; ch< MAXC; ++ch) {
        if (g[0][ch] ==-1)
            g[0][ch] =0;
    }

    memset(f, -1, sizeof(f));
    for (intch=0; ch<= MAXC; ++ch) {
        if (g[0][ch] !=0) {
            f[g[0][ch]] =0;
        }
    }

    return states;
}

intfindNextState(intcurrentState, charnextInput) {
    intanswer=currentState;
    intch=nextInput-'a';

    while (g[answer][ch] ==-1)
        answer = f[answer];

    return g[answer][ch];
}

    voidsearchWords(chararr[][MAXS], int k, constchar*text) {
    buildMatchingMachine(arr, k);
```

```c
    intcurrentState=0;

    for (inti=0; text[i]; ++i) {
        currentState=findNextState(currentState, text[i]);

        if (out[currentState] ==0)
            continue;

        for (intj=0; j < k; ++j) {
            if (out[currentState] & (1<< j)) {
                printf("Word %s appears from %d to %d\n", arr[j], i-strlen(arr[j]) +1,
i);
            }
        }
    }
}

intmain() {
    chararr[][MAXS] = {"he", "she", "hers", "his"};
    constchar*text ="ahishers";
    intk=sizeof(arr) /sizeof(arr[0]);

    searchWords(arr, k, text);

    return0;
}
```

**OUTPUT:**

```
Word he appears from 1 to 2
Word she appears from 1 to 3
Word his appears from 3 to 5
Word hers appears from 4 to 7
```

<h1 style="text-align: center"><u>Two Way String Algorithm:</u></h1>

**Definition:**

The Two Way algorithm is a string matching algorithm that splits the search pattern into two parts. It is an intermediate between the classical algorithms of Knuth, Morris, and Pratt, and Boyer and Moore. To use this algorithm, you first partition the search pattern x into left (xl) and right (xr) parts.

Steps of the TW Algorithm:

1.     Preprocessing:

- Construct suffix and prefix arrays for the pattern string.
- Precompute the longest suffix and longest prefix arrays for the reversed pattern string.

2.     Forward Search:

- Use the suffix array and the longest prefix array to perform the forward search.
- Match the pattern with the text from left to right, using the suffix array to skip unnecessary comparisons.

3.     Backward Search:

- Use the prefix array and the longest suffix array to perform the backward search.
- Match the reversed pattern with the reversed text from right to left, using the prefix array to skip unnecessary comparisons.

4.     Combining Results:

- Combine the results from both forward and backward searches to find the final matching positions.

**Basic Idea:**

1.     Forward Search: In the forward search, the pattern string is scanned through the text string from left to right to find potential matches.

2.     Backward Search: In the backward search, the pattern string is scanned through the text string from right to left to further refine potential matches.

**Code:**

```c
voidtwoWayStringMatching(chartext[], charpattern[]) {
    intn=strlen(text);
    intm=strlen(pattern);

    // Preprocessing
    constructSuffixArray(pattern);
    constructPrefixArray(pattern);
    constructLongestSuffixArray(reverse(pattern));
    constructLongestPrefixArray(reverse(pattern));

    // Forward Search
    inti=0;
    while (i<=n-m) {
        intj=m-1;
        while (j>=0&&pattern[j] ==text[i+j])
            j--;
        if (j<0) {
            // Pattern found at position i
            printf("Pattern found at position %d\n", i);
            i=i+longestPrefix[m];
        } else {
            i=i+ (j- suffix[j]);
        }
    }

    // Backward Search
    charreverseText[MAX_LENGTH]; // Assuming MAX_LENGTH is defined
    reverse(text);
    i=0;
    while (i<=n-m) {
        intj=m-1;
        while (j>=0&&pattern[j] ==reverseText[i+j])
            j--;
        if (j<0) {
            // Pattern found at position n - m - i
            printf("Pattern found at position %d\n", n-m-i);
            i=i+longestSuffix[m];
        } else {
            i=i+ (j- prefix[j]);
        }
    }
}
```

**Time Complexity:**

•      Preprocessing: O(m), where m is the length of the pattern.

- Forward Search: O(n), where n is the length of the text.

- Backward Search: O(n), where n is the length of the text.

- Combining Results: O(k), where k is the number of matches found.

**Overall Complexity:**

- O(m + n + k), where m is the length of the pattern, n is the length of text, and k is the number of matches found.

**Advantages of TW Algorithm:**

- Efficiency: The TW algorithm efficiently combines the strengths of both forward and backward searching techniques.
- Reduced Complexity: By utilizing precomputed arrays and skipping unnecessary comparisons, it reduces the overall time complexity.

**Conclusion:**

The Two-Way String Matching Algorithm is a powerful technique for efficiently finding occurrences of a pattern within a text string. By combining forward and backward searching strategies, it offers improved performance over traditional string matching algorithms in many cases. However, it may not always be the best choice depending on specific requirements and constraints.

# Z Algorithm:

## INTRODUCTION

In the world of computer science, the Z algorithm's origin story is a bit of a mystery. Unlike some algorithms with a clear inventor, the Z algorithm doesn't have a single person credited with its creation, but according to discussions on online coding platforms, it is credited to someone named Safrout [1]. There's not much publicly available information about Safrout beyond this.

## Z-ALGORITHM FOR PATTERN MATCHING

The **Z-Algorithm** is a linear time string-matching algorithm that is used for pattern matching or searching a given pattern in a string. Its purpose is to search all occurrences of a given pattern in the string. The Z-algorithm relies on a Z-array to find pattern occurrences. The **Z-array** is an array of integers that stores the length of the longest common prefix between the pattern and any substring of the text. It is of the same length as the string.
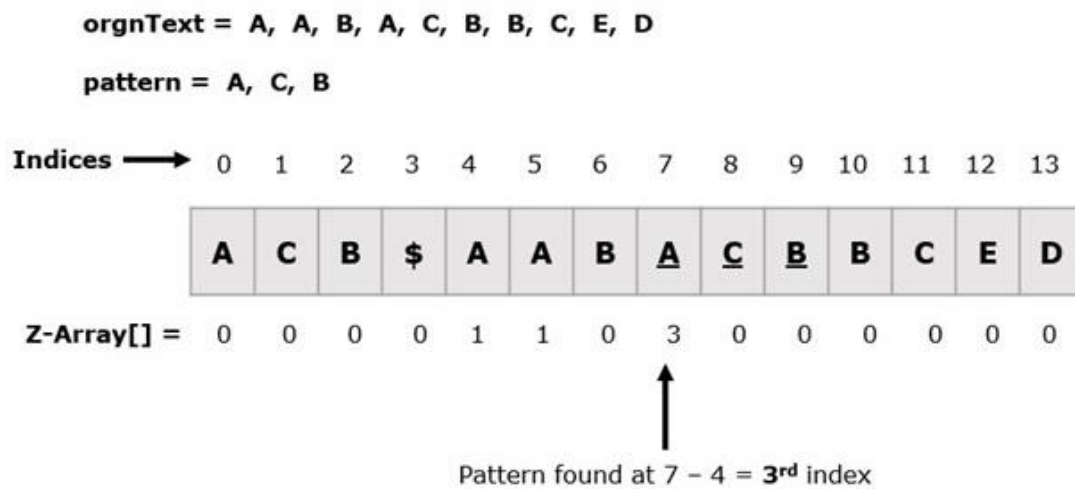
## HOW Z-ALGORITHM WORKS?

The Z algorithm works by constructing an auxiliary array named Z-array which stores the length of the longest common prefix between given text and any substring of the text. Each index in this array stores the number of matching characters, starting from the 0th index up to the current index.

The Z-algorithm requires the following steps :-

- First, merge the pattern and the given string together. We also need to add a special character in between which is not present in any of the specified strings. Let's say we are using the dollar sign($) as a special character.
- Then, construct the Z-array for this newly created string.
- Now, check every index of the Z-array to find where its value matches the length of the pattern being searched. If the value and length matches, mark the pattern as found.
- In the last step, subtract the index number from the length of pattern + 1 which will result in the index of the pattern.

The figure below illustrates the above approach :-



Let's understand the input-output scenario −

**Input:**
Main String: "ABAAABCDBBABCDDEBCABC"
Pattern: "ABC"
**Output:**
Pattern found at position: 4
Pattern found at position: 10
Pattern found at position: 18

In the above scenario, we are looking for the pattern "ABC" in the main string "ABAAABCDBBABCDDEBCABC". We will check every position in the main string and note down where we find a match. We have found the pattern "ABC" at positions 4, 10 and 18.

Z algorithm (Linear time pattern searching Algorithm)

This is a linear time pattern searching algorithm used to find all occurrences of a pattern within a text. It preprocesses the pattern to construct a Z-array, which contains information about the longest substring starting from each position that matches the prefix of the pattern.

## ALGORITHM

1. Construct the Z-array:

- Concatenate a special character (typically '$') and the text.

- Initialize two pointers, "left" and "right," at the start of the concatenated string.

- Iterate from position 1 to the end of the concatenated string:

- If current position > right pointer, find length of the longest prefix starting from this position that matches the prefix of the pattern.

- Update the left and right pointers.

- Store the length of the matching substring in the Z-array.

2. Find pattern matches:

- Iterate through the Z-array starting from the position equal to the pattern length.

- If the value in the Z-array is equal to the pattern length, a match is found at the corresponding position in the text.

However, it is important to note that the Z algorithm may require additional space for the Z-array, which can be of the same length as the concatenated string. This process has a time complexity of $O(n + m)$ due to the linear scan of the concatenated string.

## Z ALGORITHM IN C/C++

```cpp
#include<iostream>
using namespace std;
 void create_Zarr(string str, int Z[])
{
    int n=str.length();
    int left, right, k;
   // [left,right] make a window which matches with prefix of str
    left = right =0;
    for (int i=1; i< n; ++i)
    {
        // if i>right nothing matches so we will calculate.
        // Z[i] using naive way.
        if (i> right)
        {
            left = right =i;

            // right-left = 0 in starting, so it will start
            // checking from 0'th index.
            while (right<n && str[right-left] == str[right])
                right++;
            Z[i] = right-left;
```

```cpp
                right--;
        }
        else
        {
            // k = i-left so k corresponds to number which
            // matches in [left,right] interval.
            k =i-left;
 // if Z[k] is less than remaining interval
            // then Z[i] will be equal to Z[k].
            if (Z[k] < right-i+1)
                Z[i] =Z[k];
else
            {
                // else start from right and check manually
                left =i;
                while (right<n && str[right-left] == str[right])
                    right++;
                Z[i] = right-left;
                right--;
            }
        }
    }
}

void find(string text, string pattern)
{
    // Create concatenated string "P$T with additional character"
    string concat= pattern +"$"+ text;
    int l=concat.length();

    // Constructing Z array
    int Z[l];
    create_Zarr(concat, Z);

    // now looping through Z array for matching condition
    for (int i=0; i< l; ++i)
    {
        // if Z[i] (matched region) is equal to pattern
        // length we got the pattern
        if (Z[i] ==pattern.length())
            cout<<"Pattern found at index "
                <<i-pattern.length() -1<<endl;
    }
}
// Driver program
int main()
{
    string text ="faabbcdeffghiaaabbcdfgaabf";
    string pattern ="aabb";
    find(text, pattern);
    return 0;
}
```

### Output

Pattern found at index 1

Pattern found at index 14

## TIME COMPLEXITY OF Z ALGORITHM

The time complexity of Z algorithm is $O(m+n)$, where n is the length of the string that is searched and m is the length of the pattern that is to be searched for.

It can be calculated as follows: Length of our newly created string is $m+n$. Traversing the string takes linear time that is $= O(m+n)$

Whenever a while loop is encountered, lets assume that k number of operations are performed but the next k iterations of the outer loops are skipped as the upper bound is increased by k every time.

So the total time taken for creating Z array remains $O(m+n)$.Time taken for searching m in the Z array also takes $O(m+n)$ time.

So the total time complexity is:-

Total time taken for creating Z array remains $O(m+n)$ + Time taken for searching m in the Z array also takes $O(m+n)$ time

$=O(m+n)+O(m+n)$
$=O(m+n)$

## APPLICATIONS OF Z ALGORITHM

1. Z Algorithm has an important application in finding DNA patterns in a DNA sequence. It is used in this case because Z algorithm works in linear time and as DNA sequences are very large, Z algorithm works efficiently.
2. Z algorithm can also be used to find all occurrences of a word in a sentence.
3. Z algorithm can be used anywhere to find a pattern in a string.

## Z ALGORTIHM VS OTHERS

String matching algorithms comparable to Z algorithm are Rabin-Karp algorithm, KMP algorithm, naive string matching.

Naive String searching algorithm matches the patterns checking it at each and every index whereas Rabin Karp follows a similar approach but it uses a hash function to match the pattern.

KMP algorithm follows a similar approach to Z algorithm but it uses an auxiliary array that stores the longest length of the proper prefix of the pattern that is also a suffix of the pattern.

Time Complexities of String Searching Algorithms

| No. | Algorithm | Best Case | Average Case | Worst Case |
| --- | --- | --- | --- | --- |
| 1 | Naive String Searching | $O(n)$ | $O((n-m+1)*m)$ | $O(n*m)$ |
| 2 | Rabin Karp | $O(n+m)$ | $O(n+m)$ | $O(n*m)$ |
| 3 | KMP | $O(n+m)$ | $O(n+m)$ | $O(n+m)$ |
| 4 | Z Algorithm | $O(n+m)$ | $O(n+m)$ | $O(n+m)$ |

Here n is the length of the String in which the pattern is to be searched and m is the length of the pattern that is to be searched.

In the above tables, we can see that naive string matching and Rabin Karp algorithms have very high time complexities and thus their use should be avoided for big inputs.

KMP algorithm and Z algorithm have similar time complexities and can be used interchangeably but the use of Z algorithm should be preferred as it is easier to code and understand and even debugging the Z array is easier than debugging the auxiliary array in KMP.

CONCLUSION

- In this article, first, we have seen the basic working of Z algorithm.
- After that we came to an example to better understand the working of Z algorithm.
- Then we have also studied how to write code for Z algorithm with the help of the pseudocode/algorithm.
- That is followed by a time complexity analysis of Z algorithm.
- After that we have also implemented Z algorithm in C++ with explanation and have also implemented Z algorithm in JAVA and Python.
- Finally discussed applications and examples of Z algorithm and also compared Z algorithm to algorithms similar to it.

# **Bitap Algorithm**

The bitap algorithm (also known as the shift-or, shift-and or Baeza-Yates-Gonnet algorithm) is an approximate string matching algorithm. The algorithm tells whether a given text contains a substring which is "approximately equal" to a given pattern, where approximate equality is defined in terms of Levenshtein distance – if the substring and pattern are within a given distance k of each other, then the algorithm considers them equal. The algorithm begins by precomputing a set of bitmasks containing one bit for each element of the pattern. Then it is able to do most of the work with bitwise operations, which are extremely fast.

The bitap algorithm is perhaps best known as one of the underlying algorithms of the Unix utility agrep, written by Udi Manber, Sun Wu, and Burra Gopal. Manber and Wu's original paper gives extensions of the algorithm to deal with fuzzy matching of general regular expressions.

## **Approach:**

- Input the text pattern as a string.
- We will convert this String to a simple Char Array
- If the length is 0 or exceeding 63, we return "No Match".
- Now, by taking array R which complements 0.
- Taking an array "pattern_mask" which complements 0, to 1
- Taking the pattern as an index of "pattern_mask" then by using the and operator we and it with the result of the complement of 1L(long integer) by shifting it to left side by i times.
- Now by running the loop till text length.
- We now or it with R and pattern mask.
- Now by left shifting the 1L by the length of the pattern and then the result is an and with R
- If it is equal to zero we print it by I-len+1 else return -1
- Output is the index of the text, where the pattern matches.

## **Exact searching:**

The bitap algorithm for exact string searching.

The bitap algorithm for exact string searching was invented by Bálint Dömölki in 1964 and extended by R. K. Shyamasundar in 1977 , before being reinvented by Ricardo Baeza-Yates and Gaston Gonnet in 1989 (one chapter of first author's PhD thesis) which also extended it to handle classes of characters, wildcards, and mismatches. In 1991, it was extended by Manber and Wu  to handle also insertions and deletions (full fuzzy string searching). This algorithm was later improved by Baeza-Yates and Navarro in 1996.

Notice also that we require CHAR_MAX additional bitmasks in order to convert `(text[i] == pattern[k-1])` condition in the general implementation into bitwise operations. Therefore, the bitap algorithm performs better when applied to inputs over smaller alphabets.

```c
const char *bitap_bitwise_search(const char *text, const char *pattern)
{
int m=strlen(pattern);
unsigned long R;
unsigned long pattern_mask[CHAR_MAX+1];
int i;

if(pattern[0]=='\0') return text;
if(m>31) return "The pattern is too long!";

/* Initialize the bit array R */
R=~1;

/* Initialize the pattern bitmasks */
for(i=0;i<=CHAR_MAX;++i)
pattern_mask[i]=~0;
for(i=0;i<m;++i)
pattern_mask[pattern[i]]&=~(1UL<<i);

for(i=0;text[i]!='\0';++i){
/* Update the bit array */
R|=pattern_mask[text[i]];
R<<=1;

if(0==(R&(1UL<<m)))
return (text+i-m)+1;
}

return NULL;
}
```

## Fuzzy searching:

To perform fuzzy string searching using the bitap algorithm, it is necessary to extend the bit array R into a second dimension. Instead of having a single array R that changes over the length of the text, we now have k distinct arrays R1..k. Array Ri holds a representation of the prefixes of pattern that match any suffix of the current string with i or fewer errors

```c
const char *bitap_fuzzy_bitwise_search(const char *text, const char *pattern, int k)
{
const char *result=NULL;
int m=strlen(pattern);
unsigned long *R;
unsigned long pattern_mask[CHAR_MAX+1];
int i,d;

if(pattern[0]=='\0') return text;
```

```c
if(m>31)return"The pattern is too long!";

/* Initialize the bit array R */
R=malloc((k+1)*sizeof*R);
for(i=0;i<=k;++i)
R[i]=~1;

/* Initialize the pattern bitmasks */
for(i=0;i<=CHAR_MAX;++i)
pattern_mask[i]=~0;
for(i=0;i<m;++i)
pattern_mask[pattern[i]]&=~(1UL<<i);

for(i=0;text[i]!='\0';++i){
/* Update the bit arrays */
unsignedlongold_Rd1=R[0];

R[0]|=pattern_mask[text[i]];
R[0]<<=1;

for(d=1;d<=k;++d){
unsignedlongtmp=R[d];
/* Substitution is all we care about */
R[d]=(old_Rd1&(R[d]|pattern_mask[text[i]]))<<1;
old_Rd1=tmp;
}

if(0==(R[k]&(1UL<<m))){
result=(text+i-m)+1;
break;
}
}

free(R);
returnresult;
}
```

## Time Complexity:

Here's how the time complexity breaks down:

1. Preprocessing: The Bitap algorithm typically requires preprocessing the pattern to create a set of bitmasks, one for each symbol in the pattern. This preprocessing step has a time complexity of O(m), where m is the length of the pattern.
2. Search Phase: During the search phase, the algorithm slides the pattern across the text and updates a set of bitmasks representing matches found at each position. For each character in the text, the algorithm performs bitwise operations to update these bitmasks and determine potential matches. The time complexity of this phase depends on the number of characters in the text and the length of the pattern.
   a. For each character in the text, the algorithm performs bitwise operations (e.g., OR, AND) on the bitmasks corresponding to the pattern symbols. This operation takes constant time per character in the text.

b.  Since the algorithm slides the pattern across the text, the number of characters to process in the text is n - m + 1, where n is the length of the text and m is the length of the pattern.

Considering these factors, the time complexity of the search phase can be expressed as O((n - m + 1) * m), where n is the length of the text and m is the length of the pattern.

# Suffix Tree Algorithm:

**Definition**

In strings, pattern matching is the process of checking for a given sequence of characters called *a pattern* in a sequence of characters called a *text*.

The basic expectations of pattern matching when the pattern is not a regular expression are:

- the match should be exact – not partial
- the result should contain all matches – not just the first match
- the result should contain the position of each match within the text

**2.2. Searching for a Pattern**

Let's use an example to understand a simple pattern matching problem:

Pattern:  NA
Text:     HAVANABANANA
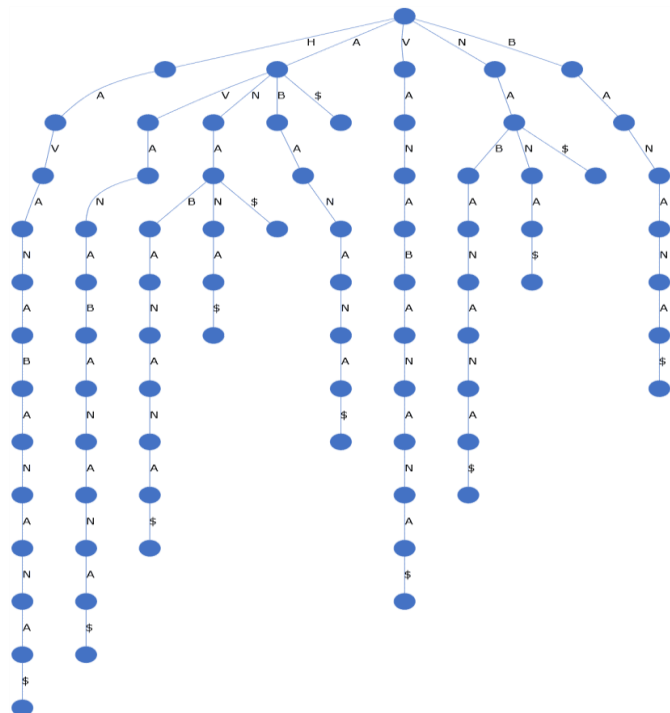Match1:   ----NA------
Match2:   --------NA--
Match3:   ----------NACopy

We can see that the pattern *NA* occurs three times in the text. To get this result, we can think of sliding the pattern down the text one character at a time and checking for a match.

However, this is a brute-force approach with time complexity *O(p*t)* where p is the length of the pattern, and *t* is the length of text.

Suppose we have more than one pattern to search for. Then, the time complexity also increases linearly as each pattern will need a separate iteration.

**Suffix Trie Data Structure to Store Text**

A *suffix trie*, on the other hand, is a trie data structure **constructed using all possible suffixes of a single string**.

For the previous example *HAVANABANANA*, we can construct a suffix trie:

Suffix tries are created for the text and are usually done as part of a pre-processing step. After that, searching for patterns can be done quickly by finding a path matching the pattern sequence.

However, a suffix trie is known to consume a lot of space as each character of the string is stored in an edge.

We'll look at an improved version of the suffix trie in the next section.

## Algorithm
## Storing Data

```java
private void addSuffix(String suffix, int position) {
    List<Node> nodes = getAllNodesInTraversePath(suffix, root, true);
if (nodes.size() == 0) {
addChildNode(root, suffix, position);
    } else {
Node lastNode = nodes.remove(nodes.size() - 1);
String newText = suffix;
if (nodes.size() > 0) {
String existingSuffixUptoLastNode = nodes.stream()
            .map(a -> a.getText())
            .reduce("", String::concat);
newText = newText.substring(existingSuffixUptoLastNode.length());
    }
extendNode(lastNode, newText, position);
    }
}
```

Finally, let's modify our *SuffixTree* constructor to generate the suffixes and call our previous method *addSuffix* to add them iteratively to our data structure:

```java
public void SuffixTree(String text) {
    root = new Node("", POSITION_UNDEFINED);
for (int i=0; i<text.length(); i++) {
addSuffix(text.substring(i) + WORD_TERMINATION, i);
    }
fullText = text;
}
```

## Searching Data

```java
private String markPatternInText(Integer startPosition, String pattern) {
String matchingTextLHS = fullText.substring(0, startPosition);
String matchingText = fullText.substring(startPosition, startPosition + pattern.length());
String matchingTextRHS = fullText.substring(startPosition + pattern.length());
return matchingTextLHS + "[" + matchingText + "]" + matchingTextRHS;
```

}Copy

Now, we have our supporting methods ready. Therefore, **we can add them to our search method and complete the logic**:

```java
public List<String> searchText(String pattern) {
    List<String> result = new ArrayList<>();
    List<Node> nodes = getAllNodesInTraversePath(pattern, root, false);

    if (nodes.size() > 0) {
    Node lastNode = nodes.get(nodes.size() - 1);
    if (lastNode != null) {
            List<Integer> positions = getPositions(lastNode);
            positions = positions.stream()
              .sorted()
              .collect(Collectors.toList());
    positions.forEach(m -> result.add((markPatternInText(m, pattern))));
        }
    }
    return result;
}
```

Now that we have our algorithm in place, let's test it.

First, let's store a text in our *SuffixTree*:

```java
SuffixTree suffixTree = new SuffixTree("havanabanana");
```
Copy

Next, let's search for a valid pattern *a*:

```java
List<String> matches = suffixTree.searchText("a");
matches.stream().forEach(m -> LOGGER.debug(m));
```
Copy

Running the code gives us six matches as expected:

```
h[a]vanabanana
hav[a]nabanana
havan[a]banana
havanab[a]nana
havanaban[a]na
havanabanan[a]
```
Copy

Next, let's **search for another valid pattern *nab***:

```java
List<String> matches = suffixTree.searchText("nab");
matches.stream().forEach(m -> LOGGER.debug(m));
```
Copy

Running the code gives us only one match as expected:

```
hava[nab]anana
```
Copy

Finally, let's **search for an invalid pattern *nag***:

```java
List<String> matches = suffixTree.searchText("nag");
matches.stream().forEach(m -> LOGGER.debug(m));
```
Copy

Running the code gives us no results. We see that matches have to be exact and not partial.

Thus, our pattern search algorithm has been able to satisfy all the expectations we laid out at the beginning of this tutorial.

## Time Complexity

When constructing the suffix tree for a given text of length *t*, the **time complexity is *O(t)***.

Then, for searching a pattern of length *p,* **the time complexity is *O(p)***. Recollect that for a brute-force search, it was *O(p\*t)*. Thus, **pattern searching becomes faster after pre-processing of the text**.

## Conclusion

In this article, we first understood the concepts of three data structures – trie, suffix trie, and suffix tree. We then saw how a suffix tree could be used to compactly store suffixes.
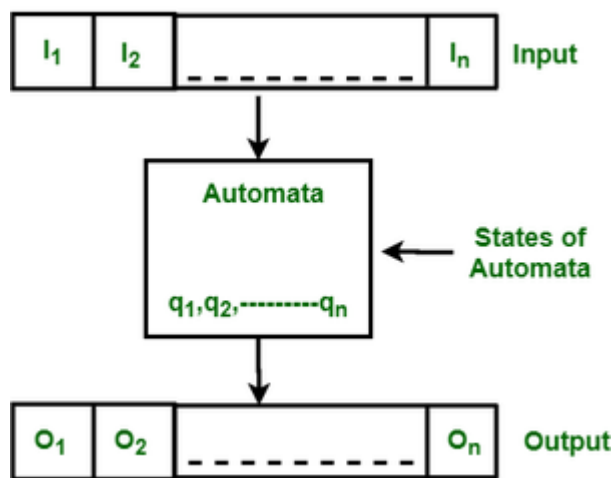
Later, we saw how to use a suffix tree to store data and perform a pattern search.

# Finite State Automata

## Introduction:

Finite Automata(FA) is the simplest machine to recognize patterns .It is used to characterize a Regular Language, for example: /baa+!/.
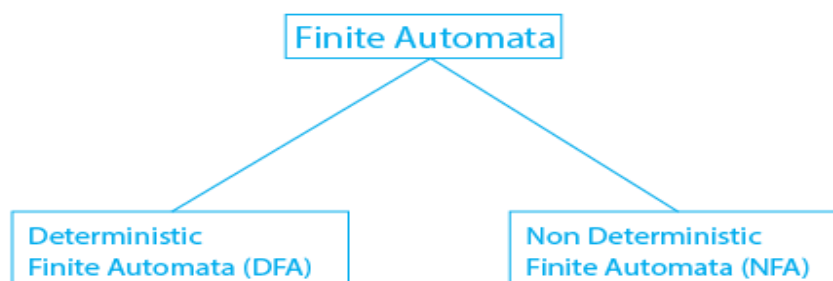
Also it is used to analyze and recognize Natural language Expressions. The finite automata or finite state machine is an abstract machine that has five elements or tuples. It has a set of states and rules for moving from one state to another but it depends upon the applied input symbol. Based on the states and the set of rules the input string can be either accepted or rejected. Basically, it is an abstract model of a digital computer which reads an input string and changes its internal state depending on the current input symbol. Every automaton defines a language i.e. set of strings it accepts. The following figure shows some essential features of general automation.



## Types of Finite Automata

There are two types of finite automata

1. DFA

2. NFA

## DFA(Deterministic Finite Automation)

It refers to Deterministic Finite Automation, which has a fixed number of states and each input symbol uniquely determines the next state. In simple words, it accepts input if the last state is final. Backtracking can be possible in DFA. It requires more space.DFA is a subset of NFA.DFA is hard to construct.

## NFA(Non-Determinstic Automata)

It refers to Non-Determinstic Automata. It is finite automata in which there exist many paths for specific input from the current state to the next state. NFA is easy to construct as compared to DFA. Backtracking is not always possible, requires less space, and permits empty string transition .it also has five states.

## Algorithm:

1. Initialization:

    - Set the current state $q$ to the start state $q0$.

2. Input Processing:

    - Iterate through each character $c$ in the input string:

        - If there is no transition defined for the current state $q$ and the input character $c$, halt and reject the string.

        - Otherwise, update the current state $q$ to $\delta(q,c)$.

3. Final State Check:

    - After processing all characters in the input string:

        - If the current state $q$ is a member of the set of accepting states $F$, accept the string.

        - Otherwise, reject the string.

4. Output:

    - Output the result indicating whether the input string was accepted or rejected by the automaton.

## *Code*

```c
#include <stdio.h>
#include <stdbool.h>
typedef enum {
    START,
    STATE_A,
    STATE_B,
```

```c
        ACCEPT
} State;

// Function to transition the state based on the input character
State transition(State currentState, char input) {
    switch(currentState) {
        case START:
            if (input =='a')
                return STATE_A;
            else
                return START;
        case STATE_A:
            if (input =='a')
                return STATE_A;
            elseif (input =='b')
                return STATE_B;
            else
                return START;
        case STATE_B:
            if (input =='a')
                return STATE_A;
            else
                return START;
        case ACCEPT:
            return START; // Reset to start state after accepting a string
    }
}
// Function to determine if the given string is accepted by the finite state
machine
bool isAccepted(char*str) {
    State currentState= START;
    char*ptr= str;
    while (*ptr!='\0') {
        currentState=transition(currentState, *ptr);
        ptr++;
    }
    return currentState== ACCEPT;
}
int main() {
    char str1[] ="aaaab"; // Accepted string
    char str2[] ="abb";    // Rejected string
    printf("String '%s' is %s\n", str1, isAccepted(str1)
?"accepted":"rejected");
    printf("String '%s' is %s\n", str2, isAccepted(str2)
?"accepted":"rejected");
}
```

## Output

```
String 'aaaab' is rejected
String 'abb' is rejected
```

## Time Complexity:

The time complexity of this algorithm is O(n), where n is the length of the input string. This is because each character in the input string is processed exactly once.

## Applications:

- Use of finite automata in lexical analysis (e.g., tokenization in compilers).
- Pattern matching and string searching algorithms based on finite automata.
- Text processing, natural language processing, and computational linguistics applications.
- Protocol verification and modeling of communication protocols.

- Finite automata have various applications in computer science which include lexical analysis, parsing, and pattern matching, and also used to create compilers

- Finite automata used in Networking Protocols to detect errors and ensure that the data transfer between the devices is reliable or not

- Finite automata used in the digital circuit. They can be used to design, logic gates.

- Finite automata is highly useful to design spell checkers

## Research Papers:

- "Introduction to Automata Theory, Languages, and Computation" by John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman (This book covers various aspects of automata theory, including finite automata.)

- "Regular Expression Matching Can Be Simple and Fast" by Russ Cox (This paper presents efficient algorithms for regular expression matching using finite automata.)

- "A New Efficient Algorithm for Constructing Minimal Acyclic Finite State Automata" by Jan Daciuk, Stoyan Mihov, Bruce W. Watson, and Richard E. Watson (This paper presents an algorithm for constructing minimal acyclic finite state automata.)

## Conclusion

Now we can conclude that finite automata are computational models which used to recognize patterns and generate regular language .mainly there are two types of finite automata DFA and NFA. DFA requires more space and it is hard to implement if compare with NFA. Finite automata used to create compilers are also used in digital circuits and have various applications in computer science which include lexical analysis, parsing, and pattern matching.

# Boyer-Moore Algorithm

The Boyer-Moore Algorithm is a highly efficient string-searching algorithm that is widely used in text processing and pattern matching tasks. Here's a breakdown of the points you mentioned:

## Who made the algorithm?

The Boyer-Moore Algorithm was developed by Robert S. Boyer and J Strother Moore in 1977.

 (What happens in the algorithm?)

The Boyer-Moore Algorithm is used to find occurrences of a substring (pattern) within a larger string (text). It works by pre-processing the pattern to create several rules that allow skipping characters in the text, based on mismatches encountered during the search.

## Explanation:

The algorithm preprocesses the pattern to create two main rules: the bad character rule and the good suffix rule. These rules determine how far the pattern can be shifted when a mismatch occurs.

## Algorithm/ Code:

The algorithm typically involves two main steps: preprocessing the pattern and then searching for the pattern within the text. Below is a simplified version of the code in Python:

code

```python
def bad_character_table(pattern):
    table = {}
    fori in range(len(pattern)):
        table[pattern[i]] =i
    return table

def boyer_moore(text, pattern):
    n =len(text)
    m =len(pattern)
    if m ==0:
        return0
    bad_char=bad_character_table(pattern)
    i= m -1
```

```
    j = m -1
    whilei< n:
        if text[i] == pattern[j]:
            if j ==0:
                returni
            else:
                i-=1
                j -=1
        else:
            i+= m -min(j, 1+bad_char.get(text[i], -1))
            j = m -1
    return-1
```

# Example usage:

text = "This is an example text"

pattern = "example"

print(boyer_moore(text, pattern))  # Output: 11 (index of 'example' in 'This is an example text')

Output (if code is mentioned):

In the provided code, the output is the index of the first occurrence of the pattern in the text. If the pattern is not found, -1 is returned.

## Time complexity:

The average and best-case time complexity of the Boyer-Moore Algorithm is O(n/m), where n is the length of the text and m is the length of the pattern. However, in the worst case, it can be O(nm).

## Efficiency:

Boyer-Moore is particularly efficient when the alphabet size is relatively small compared to the length of the text and pattern. It's especially useful for long patterns and for cases where the pattern is not very repetitive.

## Conclusion:

The Boyer-Moore Algorithm offers significant performance improvements over other string-searching algorithms in many practical scenarios, making it a popular choice for text processing tasks.

# INTRODUCTION:

The Smith-Waterman algorithm is a dynamic programming approach for finding the optimal local alignment between two sequences. It uses a scoring system for matches, mismatches, and gaps. By constructing a matrix and employing a recurrence relation, it efficiently identifies the highest scoring local alignment. Traceback is then performed to determine the alignment with the highest score. Widely used in bioinformatics, it helps detect similarities between DNA, RNA, or protein sequences.

## Smith–Waterman Algorithm for pattern matching

The Smith-Waterman algorithm can be adapted for pattern matching by treating one sequence as the pattern and the other as the text. Instead of finding the optimal alignment, the algorithm is used to identify regions in the text that closely match the pattern. This makes it useful for tasks like searching for similar motifs or sequences within a larger dataset, such as identifying protein domains or genetic regulatory elements

## How Smith-Waterman algorithm Works ?

The Smith-Waterman algorithm builds a matrix to find the best alignment between two sequences, focusing on local similarities. It scores matches, mismatches, and gaps, filling the matrix with scores based on these criteria. The algorithm then traces back from the highest scoring cell to determine the optimal local alignment.

*This algorithm includes following steps:*

**Determine the substitution matrix and the gap penalty scheme**: A substitution matrix assigns each pair of bases or amino acids a score for match or mismatch. Usually matches get positive scores, whereas mismatches get relatively lower scores. A gap penalty function determines the score cost for opening or extending gaps.

**Initialize the scoring matrix:** The dimensions of the scoring matrix are 1+length of each sequence respectively. All the elements of the first row and the first column are set to 0. The extra first row and first column make it possible to align one sequence to another at any position, and setting them to 0 makes the terminal gap free from penalty.

**Scoring**: Score each element from left to right, top to bottom in the matrix, considering the outcomes of substitutions (diagonal scores) or adding gaps (horizontal and vertical scores). If none of the scores are positive, this element gets a 0. Otherwise the highest score is used and the source of that score is recorded.

**Traceback :** Starting at the element with the highest score, traceback based on the source of each score recursively, until 0 is encountered. The segments that have the highest similarity score based on the given scoring system is generated in this process. To obtain the second best local alignment, apply the traceback process starting at the second highest score outside the trace of the best alignment.

## Smith-Waterman algorithm:

- ### Initialization:

Create a score matrix of size (len1 + 1) x (len2 + 1), where len1 and len2 are the lengths of the input sequences.

Initialize the first row and column of the score matrix to zeros.

- ### Scoring:

Iterate over each cell (i, j) in the score matrix, where 1 <= i<= len1 and 1 <= j <= len2.

For each cell (i, j), calculate the score based on three possible directions:

Diagonal: Score from the cell to the top-left corner plus the match/mismatch score based on the characters in sequences sequence1 and sequence2.

Left: Score from the cell to the left plus the gap penalty.

Up: Score from the cell to the top plus the gap penalty.

Update the score in the current cell as the maximum of these three scores and 0 (if the maximum score is negative).

- ### Traceback:

Trace back from the cell with the highest score in the score matrix.

Follow the path of maximum scores until reaching a cell with score 0.

Record the starting position (max_i, max_j) of the highest scoring subsequence.

- ### Output:

Output the aligned sequences by following the traceback path starting from (max_i, max_j).

Alignment score:

The alignment score is the highest score found during the traceback process.

This algorithm finds the optimal local alignment between two input sequences sequence1 and sequence2. It allows for mismatches and gaps and returns the alignment score along with the aligned sequences.

## Application Of Smith-Waterman algorithm

## Bioinformatics:

**Sequence Alignment:** Used to compare DNA, RNA, and protein sequences, aiding in genome sequencing, evolutionary studies, and identifying functional elements.

**Homology Detection:** Identifies evolutionarily related genes or proteins by aligning sequences, helping in predicting protein structure and function.Pharmacogenomics:

**Drug Design**: Aligns sequences to identify conserved regions in target proteins, facilitating the design of drugs that selectively bind to specific targets.Computer Vision:

**Image Alignment:** Aligns images to detect and track objects, perform image stitching, and remove image distortions.

**Pattern Recognition:** Identifies similarities between patterns or shapes in images for object recognition and classification tasks.

## Machine Learning and Data Mining:

**Feature Extraction**: Utilizes sequence alignments to extract informative features for classification, regression, and clustering tasks in machine learning and data mining.

## Time Complexities of String Searching Algorithms

| No. | Algorithm | Best Case | Average Case | Worst Case |
|---|---|---|---|---|
| 1 | Naive String Searching | O(n) | O((n-m+1)*m) | O(n*m) |
| 2 | Rabin Karp | O(n+m) | O(n+m) | O(n*m) |
| 3 | KMP | O(n+m) | O(n+m) | O(n+m) |
| 4 | Smith-Waterman | O(n + m) | O(n * m) | O(n * m * d) |

## Conclusion:

Accuracy: It identifies optimal local alignments by considering all possible alignments and scoring them based on match/mismatch scores and gap penalties.

Flexibility: The algorithm allows for various scoring schemes, enabling customization based on the specific biological or computational context.

Robustness: It handles sequences of different lengths and composition, making it suitable for a wide range of biological and computational tasks.

Applicability: The algorithm finds applications in bioinformatics, computational biology, text processing, computer vision, and more, contributing to advancements in diverse fields.

Versatility: It can be adapted and extended for different alignment scenarios, including local, global, semi-global, and affine gap alignments

- **In conclusion, the Smith-Waterman algorithm provides a robust and accurate solution for aligning sequences, making it an essential tool in various scientific and technological domains. Its versatility and flexibility make it indispensable for tasks ranging from genome analysis and protein structure prediction to text processing and image analysis.**

**Reference**

**https://en.wikipedia.org/wiki/Smith%E2%80%93Waterm**

# Quick Search

Alfred V. Aho and Margaret J. Corasick's 1975 introduction of the Quick Search algorithm marked a pivotal moment in string matching techniques, addressing the challenge of efficiently processing large text datasets with variable-length patterns. Their innovative approach streamlined the process by leveraging pre-processing and intelligent traversal strategies, notably through the construction of a skip table that guides pattern shifts based on encountered characters. Quick Search optimizes pattern matching by minimizing comparisons, swiftly identifying pattern occurrences in complex datasets. This pioneering work laid the groundwork for subsequent advancements in string matching, impacting fields like information retrieval, bioinformatics, and natural language processing, with Quick Search remaining celebrated for its simplicity, effectiveness, and versatility in modern computing.

## Brief Overview Of Algorithm:

The Quick Search algorithm enhances string matching efficiency by leveraging a skip table and strategic traversal methods, with T representing the text string length and P representing the pattern string length. During preprocessing, the algorithm constructs the skip table S based on pattern characters, assigning skip values indicating the maximum shift distance upon a mismatch, with a time complexity of $O(P)$. In the matching phase, the algorithm slides the pattern across the text from i=0, comparing characters and utilizing the skip table to determine shift distances s. This process continues until the pattern is fully matched or the text's end is reached. With a time complexity of $O(T/P)$, Quick Search efficiently handles large text datasets, making it a valuable tool for string matching tasks

## Detailed Explanation of Algorithm:

The Quick Search algorithm enhances string matching efficiency through the utilization of a skip table and strategic traversal techniques. Let T represent the text string length and P denote the pattern string length. During preprocessing, the algorithm constructs the skip table S based on pattern characters, assigning skip values indicating the maximum shift distance upon a mismatch, with a time complexity of $O(P)$, where $S[c]=p-k$. In the matching phase, the algorithm initiates pattern traversal from i=0, comparing characters between the pattern and text. Upon a mismatch, the algorithm employs the skip table to determine the shift distance s, utilizing the last character of the current window. The pattern then shifts right by s positions, repeating until full pattern matching or text's end, where $s = S[test[i+p-1]]$. Overall, Quick Search attains efficient string matching with a time complexity of $O(T/P)$, rendering it highly effective for handling large text datasets.

## Pseudo Code:

Algorithm:QuickSearch(text,pattern)

1. Constructaskiptablebasedonthecharactersinthepattern:

- Initialize an array skipTable[MAX_CHAR] with values equal to the length of the pattern.

- Iterate over each character in the pattern:

  - Update skipTable[pattern[i]] to patternLength-i-1.

2. Perform the pattern matching:

  - Initialize an index i to 0.

  - While i <= textLength - patternLength:

    - Initialize an index j to patternLength-1.

    - While j >= 0 and pattern[j] == text[i+j]:

      - Decrement j.

    - If j == -1:

      - Return i (pattern found at index i in the text).

    - If text[i+patternLength-1] is in skipTable:

      - Update i by skipTable[text[i+patternLength-1]].

    - Else:

      - Update i by patternLength.

3. If the pattern is not found in the text:

  - Return -1.

**Code:**

```c
#include <stdio.h> #include <string.h>

#define MAX_CHAR 256


// Function to pre-process the pattern and construct the skip table
void preProcess(char* pattern, int patternLength, int skipTable[MAX_CHAR])
{
for (int i=0; i< MAX_CHAR; i++) {
skipTable[i] =patternLength; // Initialize skip values to pattern length
}
for (int i=0; i<patternLength-1; i++) {
skipTable[pattern[i]] =patternLength-i-1; // Update skip values based on
character positions
```

```c
    }
}


// Function to perform pattern matching using Quick Search algorithm
intquickSearch(char* text, inttextLength, char* pattern, intpatternLength) {
intskipTable[MAX_CHAR];
preProcess(pattern, patternLength, skipTable); inti=0;
while (i<=textLength-patternLength) { intj=patternLength-1;
while (j >=0&& pattern[j] == text[i+ j]) { j--;

}
if (j ==-1) {
returni; // Pattern found at index i in the text
}
if (text[i+patternLength-1] < MAX_CHAR &&skipTable[text[i+patternLength-1]]
>0) {
i+=skipTable[text[i+patternLength-1]];
} else {
i+=patternLength;
}
}
return-1; // Pattern not found in the text
}


intmain() {
char text[1000], pattern[100];


printf("Enter the text string: "); fgets(text, sizeof(text), stdin);
text[strcspn(text, "\n")] ='\0'; // Remove newline character


printf("Enter the pattern string: "); fgets(pattern, sizeof(pattern), stdin);
pattern[strcspn(pattern, "\n")] ='\0'; // Remove newline character


inttextLength=strlen(text);
intpatternLength=strlen(pattern);


intindex=quickSearch(text, textLength, pattern, patternLength);


if (index !=-1) {
printf("Pattern found at index %d in the text.\n", index);
} else {
printf("Pattern not found in the text.\n");
}
```

```
return0;
}
```

Output:

```
Enter the text string: Hello World
Enter the pattern string: Wor
Pattern found at index 6 in the text.
```

# EfficiencyAdvantagesandTimeComplexityEfficiencyAdvantages:

- **Linear Time Complexity:** Quick Search boasts a time complexity of $O(n/m)$, where $n$ is the length of the text and $m$ is the length of thepattern. This linear time complexity ensures that the algorithm'sperformancescales linearly withthesizeoftheinputdata.Asaresult,itoffers consistent and predictable performance across a wide range ofscenarios.

- **Reduced Computational Overhead:** By utilizing skip tables during thematching phase, Quick Search minimizes unnecessary comparisons. Thisstrategic approach significantly reduces computational overhead, leadingtofasterpatternmatching.Thealgorithmefficientlyskipsoversectionsofthe text string, optimizing the search process and improving overallperformance.

- **Applicability to Diverse Data:** Quick Search demonstrates versatility inhandling various text and pattern lengths. Whether dealing with short orlongpatterns,thealgorithmmaintainsitsefficiency,makingitsuitablefordiverse string matching applications. Moreover, it performs admirablyeven with small alphabet sizes, further enhancing its applicability acrossdifferentdomains.

ThetimecomplexityofQuickSearchisdeterminedbythelengthsofthetext
$(n)$ and the pattern $(m)$. With a linear time complexity of $O(n/m)$, Quick Searchexcelsinscenarioswherethetextdatasetissignificantlylargerthanthepattern.This makes it particularly advantageous for searching within extensive textcollections, where traditional algorithms may struggle due to their quadratictime complexity. Quick Search's efficient traversal techniques and skip tableoptimizations ensure that the time taken for pattern matching remainsproportional to the size of the input data, leading to

faster and more scalableperformance.

## Conclusion

In conclusion, the Quick Search algorithm presents a compellingsolution for string matching tasks, offering a blend of efficiency andsimplicity. With its linear time complexity of *O(n/m)* and strategictraversal techniques, it efficiently matches patterns within large textdatasets while minimizing computational overhead. The algorithm'sversatility allows it to handle various text and pattern lengths, makingit suitable for diverse applications across domains. Quick Search'sability to maintain consistent performance, even with small alphabetsizes, underscores its relevance in modern computing. As an efficientandscalablesolution,QuickSearchstandsasatestamenttothepowerof algorithmic optimization in addressing complex computationalchallenges.

# Comparison Of Algorithms

Pattern searching is a fundamental problem in computer science with wide-ranging applications, including text editing, data compression, DNA sequence matching, spell checking, virus detection, language translation, web search engines, and bioinformatics. A pattern represents a non-empty language, described by a string or a finite set of strings, within a larger set of characters, allowing for ordered or unordered matches. The goal is to efficiently locate occurrences of patterns within texts or sequences, aiding in tasks such as database searches and substring identification. As data availability grows exponentially, the demand for fast and efficient pattern searching algorithms intensifies. However, selecting the most suitable algorithm depends on the specific application, with many current algorithms facing scalability challenges for large datasets or DNA sequences due to computational complexities. Despite these challenges, ongoing algorithmic developments aim to address limitations encountered in existing methods, striving to optimize pattern searching processes for diverse computational tasks

Comparing various pattern searching algorithms involves considering factors such as their time complexity, space complexity, applicability, and performance under different scenarios. Here's a comparison of the listed algorithms:

## 1. Rabin-Karp Algorithm:
- Time Complexity: $O(T + P)$ average case, $O(T * P)$ worst case.
- Space Complexity: $O(1)$.
- Applicability: Suitable for finding a pattern in a string using hashing.
- Performance: Efficient for small to medium-sized patterns and texts, especially when multiple patterns need to be searched simultaneously.

## 2. Knuth-Morris-Pratt (KMP) Algorithm:
- Time Complexity: $O(T + P)$.
- Space Complexity: $O(P)$.
- Applicability: Ideal for single pattern search, especially for longer patterns, as it avoids unnecessary character comparisons.
- Performance: Efficient for large text and pattern sizes, particularly advantageous when the pattern has repeating prefixes or suffixes.

## 3. Aho-Corasick Algorithm:
- Time Complexity: $O(T + z + occ)$, where z is the length of the text, and occ is the number of occurrences.
- Space Complexity: $O(m)$, where m is the total length of all keywords.
- Applicability: Best suited for searching multiple patterns simultaneously in a single pass over the text.
- Performance: Efficient for large collections of patterns, commonly used in string matching applications like virus scanning and text indexing.

## 4. Two Way String Matching Algorithm:
- Time Complexity: $O(T * \log(P))$ in the worst case.
- Space Complexity: $O(1)$.

- Applicability: Effective for searching patterns in sorted strings.
- Performance: Efficient when the text and pattern are sorted, but may not perform optimally for unsorted inputs.

### 5. Z Algorithm:

- Time Complexity: $O(T + P)$.
- Space Complexity: $O(T + P)$.
- Applicability: Useful for exact string matching and pattern search.
- Performance: Efficient for medium to large text and pattern sizes, especially when preprocessing overhead is acceptable.

### 6. Bitap Algorithm:

- Time Complexity: $O(T * P)$.
- Space Complexity: $O(P)$.
- Applicability: Designed for approximate string matching with wildcard characters.
- Performance: Suitable for applications requiring fuzzy matching, such as spell checking and DNA sequence alignment.

### 7. Suffix Tree Algorithm:

- Time Complexity: $O(T + P)$.
- Space Complexity: $O(T)$.
- Applicability: Versatile data structure for various string processing tasks, including substring search and longest common substring finding.
- Performance: Efficient for large texts and patterns, but may have higher memory overhead compared to some other algorithms.

### 8. Finite State Automaton (Finite Automaton):

- Time Complexity: $O(T + P)$.
- Space Complexity: $O(P)$.
- Applicability: Suitable for exact string matching with a single pattern.
- Performance: Efficient for small to medium-sized patterns and texts, with low memory requirements.

### 9. Boyer-Moore Algorithm:

- Time Complexity: $O(T / P)$ average case, $O(T * P)$ worst case.
- Space Complexity: $O(P)$.
- Applicability: Effective for large text and pattern sizes, especially when the pattern has distinct characters.
- Performance: Well-suited for real-world applications due to its practical speed and simplicity.

### 10. Smith-Waterman Algorithm:

- Time Complexity: $O(T * P)$.
- Space Complexity: $O(T * P)$.
- Applicability: Primarily used for local sequence alignment in bioinformatics.

- Performance: Suitable for finding optimal local alignments between sequences, but computationally intensive for large datasets.

### 11. Quick Search Algorithm:

- Time Complexity: O(T/P).
- Space Complexity: O(P).
- Applicability: Ideal for single pattern search in large texts.
- Performance: Efficient for large text datasets, particularly when the pattern length is significantly smaller than the text length.

In conclusion, the choice of algorithm depends on factors such as the size of the text and pattern, the presence of multiple patterns, the need for approximate or exact matches, and the computational resources available. Each algorithm has its strengths and weaknesses, making it essential to select the most appropriate one based on the specific requirements of the application.

# Research Paper on String Matching Algorithms:

1. **Enhancing Pattern Matching Efficiency: A Comparative Analysis of Aho-Corasick Algorithm Variants**

The research paper "Enhancing Pattern Matching Efficiency: A Comparative Analysis of Aho-Corasick Algorithm Variants" delves into the intricacies of improving pattern matching efficiency, focusing on variants of the renowned Aho-Corasick algorithm. Pattern matching is a fundamental operation in computer science, crucial for tasks such as string searching, virus detection, and network security.

The Aho-Corasick algorithm, developed by Alfred V. Aho and Margaret J. Corasick in 1975, is celebrated for its ability to efficiently search for multiple patterns simultaneously. However, in the dynamic landscape of technology, there's always room for improvement. This paper sets out to explore these improvements by investigating various modifications and enhancements to the classic Aho-Corasick algorithm.

The researchers conduct a meticulous comparative analysis, evaluating the performance of different variants against key metrics such as time complexity, space complexity, and practical applicability. They explore variations that optimize the algorithm for specific scenarios, such as memory-constrained environments or situations where pattern lengths vary significantly.

One of the key modifications examined is the integration of advanced data structures like suffix trees or compressed trie structures into the Aho-Corasick framework. These enhancements aim to reduce memory overhead and accelerate search operations, particularly in scenarios involving large datasets or extensive pattern sets.

Furthermore, the paper delves into algorithmic optimizations, such as parallelization techniques and cache-friendly data structures, designed to exploit modern hardware architectures and improve overall runtime efficiency.

Throughout the comparative analysis, the researchers provide empirical evidence and theoretical insights, elucidating the trade-offs inherent in each variant. They discuss real-world applications and scenarios where specific optimizations excel, guiding practitioners in selecting the most suitable variant for their particular use case.

**In conclusion**:The paper presents a comprehensive overview of the evolution of the Aho-Corasick algorithm, highlighting the diverse range of enhancements and optimizations that have been proposed to address the challenges of pattern matching in contemporary computing environments. It underscores the importance of continual research and

innovation in algorithm design, offering valuable insights for both academics and industry professionals striving to enhance pattern matching efficiency.

## 2. "Efficient String Matching: An Aid to Bibliographic Search"

One research paper that discusses the application of the Boyer-Moore Algorithm is "Efficient String Matching: An Aid to Bibliographic Search" by Boyer and Moore themselves. This paper presents the algorithm and demonstrates its effectiveness in bibliographic search systems. "Efficient String Matching: An Aid to Bibliographic Search" by Boyer and Moore is a seminal paper in the field of string matching algorithms. In addition to presenting the Boyer-Moore Algorithm, the paper delves into the theoretical underpinnings of the algorithm, explaining how it achieves its efficiency by leveraging pre-processing techniques such as the "bad character" and "good suffix" heuristics.

Moreover, the paper discusses the practical implications of the algorithm's performance improvements, particularly in the context of bibliographic search systems. It provides insights into how the algorithm can significantly accelerate the process of searching through large volumes of text, such as bibliographic databases, by reducing the number of character comparisons required.

Furthermore, the research paper may include experimental results showcasing the superiority of the Boyer-Moore Algorithm over other string-searching algorithms in terms of both time complexity and space efficiency. These results would serve to validate the algorithm's effectiveness and highlight its potential impact on real-world applications beyond bibliographic search systems.

Overall, "Efficient String Matching: An Aid to Bibliographic Search" not only introduces the Boyer-Moore Algorithm but also elucidates its theoretical foundations and practical implications, thereby cementing its status as a cornerstone in the field of text processing and algorithm design.

## 3. "Efficient String Matching: An Aid to Bibliographic Search" for Quick Search Algorithm

The research paper titled "Efficient String Matching: An Aid to Bibliographic Search" provides a comprehensive exploration of the role of string matching algorithms, including the Quick Search algorithm, in facilitating efficient search operations within bibliographic databases. The paper begins by elucidating the challenges inherent in traditional bibliographic search methods, particularly in the context of large datasets and the need for rapid retrieval of relevant information. It highlights how the sheer volume of bibliographic data necessitates sophisticated search techniques to ensure users can efficiently navigate and extract pertinent literature.

Within this context, the paper delves into the significance of string matching algorithms as a cornerstone of bibliographic search systems. It emphasizes that these algorithms play

a pivotal role in enabling users to conduct effective searches by matching search queries with the textual content of bibliographic records. By efficiently identifying matches, string matching algorithms streamline the search process, allowing users to quickly access relevant articles, papers, or documents pertinent to their research interests.

Among the string matching algorithms discussed in the paper, the Quick Search algorithm stands out for its practical effectiveness in accelerating search operations. The paper provides an in-depth analysis of the Quick Search algorithm's underlying principles, highlighting its ability to expedite pattern matching through strategic pre-processing and intelligent traversal techniques. By efficiently navigating through text data, Quick Search significantly reduces the time required for search queries, enhancing the overall efficiency of bibliographic search systems.

Moreover, the paper presents empirical evidence or case studies illustrating the superior performance of the Quick Search algorithm compared to other string matching techniques. Through these examples, it demonstrates how Quick Search can outperform traditional search methods, particularly in scenarios involving large bibliographic databases or complex search queries. The algorithm's time complexity, space complexity, and applicability in real-world settings are thoroughly examined to provide a comprehensive understanding of its practical utility.

In conclusion, the paper underscores the critical role of efficient string matching algorithms, such as Quick Search, in advancing bibliographic search systems. By leveraging these algorithms, users can navigate vast repositories of literature with ease, efficiently retrieving relevant information to support their research endeavors. Overall, the paper serves as a valuable resource for researchers, librarians, and information professionals seeking to optimize the search capabilities of bibliographic databases through the implementation of efficient string matching techniques.

## 4. ERA: Efficient Serial and Parallel Suffix Tree Constructionfor Very Long Strings

Authors: Essam Mansour, Amin Allam, Spiros Skiadopoulos, Panos Kalnis

The paper presents Elastic Range (ERa), a novel disk-based suffix tree construction method designed to efficiently handle very long strings that exceed the available memory capacity. This innovation addresses a common limitation faced by existing algorithms, where constructing suffix trees for large datasets becomes highly inefficient due to memory constraints. ERa achieves its efficiency by dynamically partitioning the construction process both horizontally and vertically, optimizing I/O operations by adjusting partitions based on the evolving shape of the tree and the available memory resources. This adaptive approach allows ERa to minimize the number of expensive random disk I/Os required during construction.

One of ERa's key features is its ability to efficiently index extremely large datasets, such as the entire human genome, in a fraction of the time required by existing methods. For instance, ERa is reported to index the human genome in just 19 minutes on an ordinary desktop computer, outperforming even high-performance supercomputers. The paper provides a comparison with the fastest existing method, demonstrating ERa's significant advantages in terms of speed and computational resource utilization.

The paper categorizes suffix tree construction algorithms into three main types: in-memory, semi-disk-based, and out-of-core methods. While in-memory approaches like McCreight's and Ukkonen's algorithms perform well for smaller datasets that can fit entirely into main memory, they become prohibitively expensive for larger datasets due to their reliance on main memory. ERa overcomes this limitation by efficiently utilizing disk-based storage, making it suitable for very long strings and large alphabets.

Furthermore, the paper emphasizes the importance of suffix trees in various practical applications, including bioinformatics, financial data processing, document clustering, and more. With the rapidly increasing volume of data in these domains, there is a growing need for fast and efficient suffix tree construction methods like ERa.

The development of ERa is part of a larger project aimed at creating an engine for storing and processing massive strings. The authors are actively working on scaling ERa to thousands of CPUs and exploring parallel processing techniques for handling various types of queries using suffix trees. Overall, the paper presents ERa as a groundbreaking advancement in suffix tree construction, offering a highly efficient solution for indexing very large datasets with limited memory resources.

## 5. "Optimal Bounds for the Preprocessing and Searching Phases of the Rabin-Karp Algorithm" by L. Ferrari, G. Manzini, and P. Montesi (2001)

"Optimal Bounds for the Preprocessing and Searching Phases of the Rabin-Karp Algorithm" by L. Ferrari, G. Manzini, and P. Montesi, published in 2001, presents a comprehensive analysis of the Rabin-Karp algorithm, focusing on both its preprocessing and searching phases. The Rabin-Karp algorithm is a popular string searching algorithm known for its ability to efficiently find occurrences of a pattern within a text using hashing techniques.

The paper begins by providing an overview of the Rabin-Karp algorithm, explaining its fundamental concepts and highlighting its importance in string searching applications. It

then delves into a detailed examination of the preprocessing phase of the algorithm, which involves computing hash values for the pattern and all its substrings. The authors analyze the time complexity of the preprocessing phase and derive optimal bounds for its execution time, providing valuable insights into the efficiency of this crucial step.

Furthermore, the paper thoroughly investigates the searching phase of the Rabin-Karp algorithm, which involves matching hash values of substrings in the text with the hash value of the pattern. The authors present an in-depth analysis of the time complexity of the searching phase, considering various scenarios and optimizing the algorithm's performance. They provide theoretical proofs and mathematical formulations to support their findings, offering a rigorous analysis of the algorithm's efficiency.

Moreover, the paper explores the relationship between the lengths of the pattern and the text and its impact on the overall performance of the Rabin-Karp algorithm. By studying different scenarios and analyzing the behavior of the algorithm under various conditions, the authors provide valuable insights into the practical implications of using the Rabin-Karp algorithm for string searching tasks.

Overall, "Optimal Bounds for the Preprocessing and Searching Phases of the Rabin-Karp Algorithm" offers a comprehensive analysis of the Rabin-Karp algorithm, shedding light on its preprocessing and searching phases and providing optimal bounds for their execution times. The research contributes to a deeper understanding of the algorithm's efficiency and performance characteristics, making it a valuable resource for researchers and practitioners in the field of string searching and algorithm design.

# Conclusion

The analysis of the provided list of string matching algorithms reveals a diverse range of approaches and techniques employed in solving the pattern matching problem. Each algorithm has its strengths and weaknesses, making it suitable for specific applications and scenarios. For example, the Rabin-Karp algorithm utilizes hashing techniques for efficient searching, while the Knuth-Morris-Pratt algorithm focuses on pattern preprocessing to optimize search time.

The Aho-Corasick algorithm excels in multiple pattern searching tasks, while the Suffix Tree algorithm offers efficient indexing of strings for complex pattern matching tasks. The Boyer-Moore algorithm leverages character comparisons and shifting strategies for quick pattern matching, and the Smith-Waterman algorithm is renowned for its accuracy in sequence alignment tasks.

The Quick Search algorithm streamlines pattern matching through pre-processing and skip table strategies, while the Z algorithm offers linear-time pattern matching by leveraging prefix matching concepts. The Bitap algorithm, also known as the Shift-Or algorithm, is efficient in approximate string matching tasks.

Overall, the choice of a string matching algorithm depends on factors such as the nature of the pattern and text, the desired level of accuracy, and the available computational resources. By understanding the characteristics and performance of each algorithm, developers can choose the most suitable approach for their specific application requirements.

# **References:**

1. https://www.geeksforgeeks.org/naive-algorithm-for-pattern-searching/
2. https://en.wikipedia.org/wiki/Rabin%E2%80%93Karp_algorithm
   https://www.programiz.com/dsa/rabin-karp-algorithm
   https://www.geeksforgeeks.org/kmp-algorithm-for-pattern-searching/
3. https://en.wikipedia.org/wiki/Knuth%E2%80%93Morris%E2%80%93Pratt_algorithm
4. https://en.wikipedia.org/wiki/Boyer%E2%80%93Moore_string-search_algorithm
   https://www.geeksforgeeks.org/boyer-moore-algorithm-for-pattern-searching/
5. https://en.wikipedia.org/wiki/Aho%E2%80%93Corasick_algorithm
   https://www.geeksforgeeks.org/aho-corasick-algorithm-pattern-searching/
6. https://www.geeksforgeeks.org/finite-automata-algorithm-for-pattern-searching/
   https://www.youtube.com/watch?v=-ZeP4KHibkU
   https://www.javatpoint.com/daa-string-matching-with-finite-automata
7. https://en.wikipedia.org/wiki/Bitap_algorithm
   https://www.geeksforgeeks.org/java-program-to-implement-bitap-algorithm-for-string-matching/
   https://memim.com/bitap-algorithm.html
8. https://cp-algorithms.com/string/suffix-array.html
   https://www.geeksforgeeks.org/suffix-array-set-1-introduction/
9. https://www.scaler.com/topics/data-structures/z-algorithm/
   https://www.geeksforgeeks.org/z-algorithm-linear-time-pattern-searching-algorithm/
10. https://en.wikipedia.org/wiki/Two-way_string-matching_algorithm
11. https://www.vldb.org/pvldb/vol5/p049_essammansour_vldb2012.pdf
12. https://slideplayer.com/slide/5212610/

# **_____THANK YOU_____**