# Comparative Analysis of String and Pattern Matching Algorithms on DNA Sequences

## Team STARK

Aditya Singh    Shivek Gupta    Tanush Garg    Keshav Goel    Rishabh Goyal

December 1, 2025

### Abstract

DNA sequence analysis relies heavily on the efficient identification of nucleotide patterns. This project implements and benchmarks five distinct string matching algorithms—Knuth-Morris-Pratt (KMP), Boyer-Moore, Shift-Or (Bitap), Suffix Arrays, and Levenshtein Distance—to evaluate their performance on genomic data. Furthermore, we introduce a novel Structure-Aware Hybrid Model that dynamically switches between exact and fuzzy matching based on genomic annotations (GFF). Using datasets including the *Escherichia coli* genome, we conducted rigorous benchmarking of execution time and memory usage. Our findings indicate that while Suffix Arrays offer superior query speeds for repeated searches after an $O(n \log n)$ preprocessing phase, classical algorithms like KMP provide consistent linear-time performance ($O(n)$) with lower memory overhead. The Hybrid Model demonstrated practical utility by applying exact matching in coding regions (CDS) while enabling fuzzy matching with configurable edit distance in non-coding regions, balancing computational efficiency with biological relevance.

## 1 Introduction

### 1.1 Problem Definition

Genomic data consists of vast strings of nucleotides $(A, C, G, T)$. Identifying specific motifs—such as promoter regions, restriction sites, or gene sequences—is a fundamental task in bioinformatics. The scale of modern genomic data requires algorithms that are not only accurate but computationally efficient in terms of time and memory. A key challenge lies in the dichotomy between *Exact Matching*, which is fast but rigid, and *Approximate (Fuzzy) Matching*, which captures biological mutations but is computationally expensive.

### 1.2 Real-World Relevance

In real-world biological contexts, DNA sequences are rarely perfect. Mutations, sequencing errors, and evolutionary divergence mean that searching for an exact sequence often misses biologically significant results. Conversely, running fuzzy matching on entire genomes is prohibitively slow. Therefore, determining the optimal algorithm for specific contexts—and developing hybrid approaches that respect biological structure (exons vs. introns)—is critical for developing efficient gene annotation and analysis tools.

## 1.3   Project Objectives

The primary objectives of this project are:

1. To implement and verify five core string matching algorithms: KMP, Boyer-Moore, Shift-Or, Suffix Arrays, and Levenshtein Distance.

2. To conduct a comparative benchmark of these algorithms regarding Wall-Clock Time and Peak Memory Usage on real DNA datasets.

3. To develop and test a "Structure-Aware" hybrid model that utilizes GFF annotation files to intelligently switch between exact and fuzzy matching strategies.

# 2   Project Repository

The complete source code, datasets, Jupyter notebooks, and documentation for this project are publicly available on GitHub:

<div align="center">

https://github.com/Tanush-IIITH/STARK

</div>

The repository includes:

- Algorithm implementations in the `KMP/`, `Boyer_Moore/`, `Shift_or_bitap/`, `Suffix Arrays-Trees/`, and `Levenshtein/` directories.

- Benchmark harnesses and analysis notebooks (`cross_algorithm_benchmark.ipynb`, `SA_hybrid_model.ipynb`).

- Real and synthetic genomic datasets in `dataset/`.

# 3. Algorithm Descriptions

# Boyer-Moore Algorithm

## Aditya

## November 2025

**Abstract**

Here is a complete guide to the Boyer-Moore string matching algorithm, analyzed for its application to DNA sequences. The implementation includes both the Bad Character Rule and Good Suffix Rule heuristics, achieving efficient pattern matching with theoretical best-case time complexity of $O(n/m)$. This section includes detailed algorithm descriptions, correctness proofs, complexity analysis, and extensive experimental validation on real genomic data from *Escherichia coli*.

# 1 Introduction

String matching is a fundamental problem in computer science with critical applications in bioinformatics, particularly in DNA sequence analysis. The Boyer-Moore algorithm, introduced by Robert S. Boyer and J Strother Moore in 1977 [1], remains one of the most efficient exact string matching algorithms, especially for biological sequences with small alphabets.

## 1.1 Problem Statement

Given a text $T$ of length $n$ and a pattern $P$ of length $m$ (where $m \leq n$), the exact pattern matching problem requires finding all occurrences of $P$ in $T$. Formally, we seek all positions $i$ such that $T[i..i + m - 1] = P[0..m - 1]$.

## 1.2 Motivation for DNA Sequences

DNA sequences consist of a small alphabet $\Sigma = \{A, C, G, T\}$ making them ideal candidates for Boyer-Moore optimization.

# 2 Algorithm Description

## 2.1 Basic Principle

The Boyer-Moore algorithm's key innovation is scanning the pattern from **right to left** while moving through the text from **left to right**. When a mismatch occurs, the algorithm uses preprocessed information to skip multiple text positions, achieving sublinear performance in the best case.

## 2.2 The Two Heuristics

### 2.2.1 Bad Character Rule

When a mismatch occurs at text position $T[i+j]$ while comparing with pattern position $P[j]$:

1. If $T[i+j]$ appears in $P[0..j-1]$, shift $P$ to align the rightmost occurrence of $T[i+j]$ with the mismatch position

2. If $T[i+j]$ does not appear in $P$, shift $P$ entirely past position $i+j$

**Formal Definition**: Let $\delta_1(c,j)$ be the bad character shift for character $c$ at position $j$:

$$\delta_1(c,j) = \begin{cases} j - R[c] & \text{if } c \text{ appears in } P[0..j-1] \\ j+1 & \text{otherwise} \end{cases} \tag{1}$$

where $R[c]$ is the rightmost position of character $c$ in $P[0..j-1]$.

### 2.2.2 Good Suffix Rule

When a mismatch occurs after matching a suffix $t$ of length $s$:

1. **Case 1 (Strong Suffix)**: Find another occurrence of $t$ in $P$ preceded by a different character

2. **Case 2 (Prefix Match)**: Find the longest prefix of $P$ that matches a suffix of $t$

**Formal Definition**: Let $\delta_2(j)$ be the good suffix shift at position $j$:

$$\delta_2(j) = \min\{k > 0 : P[j+1..m-1] = P[j+1-k..m-1-k] \text{ and } P[j] \neq P[j-k]\} \tag{2}$$

If no such $k$ exists, find the longest suffix of $P[j+1..m-1]$ that is also a prefix of $P$.

## 2.3 Shift Calculation

At each alignment, the algorithm computes:

$$\text{shift} = \max\{\delta_1(T[i+j],j), \delta_2(j)\} \tag{3}$$

This ensures we never miss a potential match by taking the maximum of both heuristics.

# 3 Time and Space Complexity Analysis

## 3.1 Preprocessing Complexity

### 3.1.1 Bad Character Table

For a DNA alphabet $\Sigma = \{A, C, G, T, N\}$ with $|\Sigma| = 5$:

- **Time**: $O(m \cdot |\Sigma|) = O(5m) = O(m)$

- **Space**: $O(m \cdot |\Sigma|) = O(5m) = O(m)$

We store the rightmost occurrence of each character at every position in the pattern.

### 3.1.2 Good Suffix Table

Computing the good suffix table involves:

1. Strong suffix computation: $O(m)$

2. Border position calculation: $O(m)$

3. Prefix-suffix matching: $O(m)$

- **Time**: $O(m)$

- **Space**: $O(m)$

### 3.1.3 Total Preprocessing

$$\boxed{T_{\text{preprocess}} = O(m + |\Sigma|) = O(m) \text{ for DNA}} \tag{4}$$

$$\boxed{S_{\text{preprocess}} = O(m + |\Sigma|) = O(m) \text{ for DNA}} \tag{5}$$

## 3.2 Search Complexity

### 3.2.1 Best Case

In the best case, the rightmost character of the pattern is not in the text. The algorithm can skip $m$ positions at each alignment:

- Number of alignments: $\lceil n/m \rceil$

- Comparisons per alignment: 1

- **Total**: $O(n/m)$

$$\boxed{T_{\text{best}} = O(n/m)} \tag{6}$$

This is **sublinear** and highly efficient for DNA sequences with small alphabets.

### 3.2.2 Worst Case

The worst case occurs when:

- The pattern has high periodicity (e.g., $P =$ "AAAAAAA")

- The text contains many near-matches

- Each alignment requires $m$ comparisons with minimal shift

$$\boxed{T_{\text{worst}} = O(n \cdot m)} \tag{7}$$

**Example**: $T =$ "AAAAAAAAA...", $P =$ "AAAB"

However, this worst case is **extremely rare** in practice, especially for random DNA sequences.

### 3.2.3 Average Case

For random text and pattern over a DNA alphabet with equal probability distribution:

$$\boxed{T_{\text{average}} = O(n)} \tag{8}$$

The expected number of character comparisons is:

$$E[\text{comparisons}] \approx \frac{n}{m} \cdot \left( 1 + \frac{1}{|\Sigma|} + \frac{1}{|\Sigma|^2} + \cdots \right) = \frac{n}{m} \cdot \frac{|\Sigma|}{|\Sigma| - 1} \tag{9}$$

For DNA ($|\Sigma| = 4$): $E[\text{comparisons}] \approx \frac{n}{m} \cdot \frac{4}{3} = O(n/m)$

## 3.3 Space Complexity Summary

| Component | General | DNA |
|---|---|---|
| Bad Character Table | $O(m \cdot |\Sigma|)$ | $O(m)$ |
| Good Suffix Table | $O(m)$ | $O(m)$ |
| Border Position Array | $O(m)$ | $O(m)$ |
| **Total** | $O(m + |\Sigma|)$ | **O(m)** |

Table 1: Space complexity analysis

# 4 Proof of Correctness

## 4.1 Theorem: Boyer-Moore Finds All Matches

**Theorem**: The Boyer-Moore algorithm correctly identifies all occurrences of pattern $P$ in text $T$.

### 4.1.1 Proof by Contradiction

Assume the algorithm misses a match at position $k$. This can occur only if:

1. The algorithm never aligns the pattern at position $k$, OR

2. The algorithm aligns at position $k$ but incorrectly reports no match

**Case 1**: Suppose the pattern is at position $i < k$ and shifts past $k$.

Let $j$ be the position where mismatch occurs, and let $C = T[i+j]$ be the mismatched character. The shift is:

$$s = \max\{\delta_1(C, j), \delta_2(j)\} \tag{10}$$

We assume $s > k - i$, meaning the algorithm skips position $k$. We will show this leads to a contradiction for both heuristics.

**Contradiction for Bad Character Rule ($\delta_1$)**

**Step 1: Using the Hypothetical Match**

If a full match of $P$ exists at position $k$, then every character in $P$ must match the corresponding character in $T$. In particular, the character at absolute text position $i + j$ must match some pattern character:

$$T[i + j] = P[x] \quad \text{where} \quad x = i + j - k \tag{11}$$

This means the character $C = T[i + j]$ (which caused the mismatch at index $j$ in alignment $i$) also appears within the pattern $P$ at index $x = i + j - k$ to enable the match at $k$.

**Step 2: Relating the Indices**

The Boyer-Moore algorithm uses $\delta_1$ to find the rightmost occurrence of $C$ in $P$ that is at or to the left of index $j - 1$. Let $x_{\text{rightmost}}$ be the index of this rightmost occurrence. By definition of the Bad Character Rule, the shift is calculated to align $P[x_{\text{rightmost}}]$ with $T[i + j]$:

$$\delta_1(C, j) = j - x_{\text{rightmost}} \tag{12}$$

The key safety property of $\delta_1$ is that any smaller shift would result in $C$ being compared against a character in $P$ that is not $C$, thus leading to a guaranteed mismatch at the new position.

**Step 3: The Contradiction**

If $P[x]$ is the required pattern character for the match at $k$, then the shift necessary to align the pattern at $k$ is:

$$\text{Required Shift} = k - i \tag{13}$$

Since we assumed the pattern was skipped (i.e., $s > k - i$), and if $\delta_1$ was the governing shift (meaning $\delta_1 \geq \delta_2$), then $s = \delta_1$:

$$\delta_1 > k - i \tag{14}$$
$$j - x_{\text{rightmost}} > k - i \tag{15}$$
$$j - (k - i) > x_{\text{rightmost}} \tag{16}$$
$$i + j - k > x_{\text{rightmost}} \tag{17}$$

Since $x = i + j - k$, this means:

$$\boxed{x > x_{\text{rightmost}}} \tag{18}$$

But this is a **contradiction**! We have:

- $x$ is the index of $C$ that enables the hypothetical match at $k$

- $x_{\text{rightmost}}$ is the index of the rightmost occurrence of $C$ used to calculate $\delta_1$

- Since $x > x_{\text{rightmost}}$, the occurrence of $C$ at index $x$ is *further to the right* than the occurrence used to calculate the $\delta_1$ shift

This contradicts the definition of $x_{\text{rightmost}}$ as the *rightmost* occurrence of $C$ in $P[0..j-1]$. Therefore, no match can exist at position $k$ when $\delta_1$ governs the shift.

### Contradiction for Good Suffix Rule ($\delta_2$)

Now suppose $\delta_2$ governs the shift (i.e., $\delta_2 > \delta_1$), so $s = \delta_2(j)$.

After the mismatch at position $j$, the algorithm has successfully matched a suffix of the pattern: $P[j+1..m-1]$ matches $T[i+j+1..i+m-1]$. The good suffix rule $\delta_2(j)$ shifts the pattern to align with:

1. Another occurrence of the suffix $P[j+1..m-1]$ in $P$ preceded by a different character, OR

2. The longest prefix of $P$ that matches a suffix of $P[j+1..m-1]$

If a match exists at position $k$ where $i < k < i + \delta_2(j)$, then:

$$P[0..m-1] = T[k..k+m-1] \tag{19}$$

In particular, the suffix portion must match:

$$P[j+1..m-1] = T[k+j+1..k+m-1] \tag{20}$$

But we already know from the comparison at position $i$ that:

$$P[j+1..m-1] = T[i+j+1..i+m-1] \tag{21}$$

This implies that the same suffix $P[j+1..m-1]$ appears in the text at two overlapping positions: starting at $i+j+1$ and at $k+j+1$. By the definition of $\delta_2$, the shift is calculated to align the pattern with the *nearest* valid occurrence of this suffix (accounting for the preceding character constraint).

If $k - i < \delta_2(j)$, then there exists a valid suffix occurrence closer than what $\delta_2$ computed, which contradicts the minimality property of $\delta_2$. Specifically:

- Either the suffix $P[j+1..m-1]$ occurs at an earlier position in $P$ with a different preceding character (contradicting that $\delta_2$ found the nearest such occurrence)

- Or a longer prefix of $P$ matches a suffix of $P[j+1..m-1]$ (contradicting that $\delta_2$ found the longest such prefix)

6

Furthermore, the good suffix rule is designed to handle **overlapping patterns** correctly. When a match is found, the algorithm shifts by $\delta_2(0)$, which ensures that any overlapping occurrences are not missed. The shift $\delta_2(0)$ is the minimum safe shift that aligns the entire pattern with its longest proper suffix, thus catching all overlapping matches.

Therefore, no match can exist at position $k$ when $\delta_2$ governs the shift. **Contradiction.**

Since both $\delta_1$ and $\delta_2$ lead to contradictions, no match can exist at $i < k < i + s$.

**Case 2**: If the algorithm aligns at position $k$, it performs right-to-left comparison:

- If $P[m-1] \neq T[k+m-1]$, there is no match (correct)

- If all characters match, a match is reported (correct)

- The algorithm terminates only when $j < 0$, meaning all positions matched

Therefore, the algorithm cannot miss a match at any aligned position.

## 4.2   Invariant

**Loop Invariant**: After each shift, all positions in $T[0..i-1]$ have been checked, and no matches exist in that region.

**Initialization**: $i = 0$, no positions checked yet. Trivially true.

**Maintenance**: After processing position $i$ and shifting by $s$:

- All positions $[i, i+s)$ are verified to have no matches (by shift correctness)

- Next iteration checks from position $i + s$

**Termination**: When $i + m > n$, all positions have been checked.

# 5 Pseudocode

## 5.1 Main Search Algorithm

---
**Algorithm 1** Boyer-Moore Search
---
1: **procedure** BOYERMOORESEARCH($T, P$)
2:      $n \leftarrow \text{length}(T)$
3:      $m \leftarrow \text{length}(P)$
4:      $\delta_1 \leftarrow \text{PREPROCESSBADCHAR}(P)$
5:      $\delta_2 \leftarrow \text{PREPROCESSGOODSUFFIX}(P)$
6:      $\text{matches} \leftarrow []$
7:      $i \leftarrow 0$             ▷ Text position
8:      **while** $i \leq n - m$ **do**
9:         $j \leftarrow m - 1$         ▷ Pattern position (right to left)
10:        **while** $j \geq 0$ **and** $P[j] = T[i + j]$ **do**
11:           $j \leftarrow j - 1$
12:        **end while**
13:        **if** $j < 0$ **then**         ▷ Complete match found
14:           $\text{matches.append}(i)$
15:           $i \leftarrow i + \delta_2[0]$        ▷ Shift to find next match
16:        **else**         ▷ Mismatch at position j
17:           $\text{shift}_1 \leftarrow \delta_1[T[i + j]][j]$
18:           $\text{shift}_2 \leftarrow \delta_2[j]$
19:           $i \leftarrow i + \max(\text{shift}_1, \text{shift}_2)$
20:        **end if**
21:      **end while**
22:      **return** matches
23: **end procedure**
---

## 5.2 Bad Character Preprocessing

---
**Algorithm 2** Preprocess Bad Character Rule
---
1: **procedure** PREPROCESSBADCHAR($P$)
2:      $m \leftarrow \text{length}(P)$
3:      $\Sigma \leftarrow \{A, C, G, T, N\}$         ▷ DNA alphabet
4:      Initialize $\delta_1[c][j] \leftarrow -1$ for all $c \in \Sigma$, $j \in [0..m)$
5:      **for** $i \leftarrow 0$ **to** $m - 1$ **do**
6:         **for** $c \in \Sigma$ **do**
7:           **if** $i > 0$ **then**
8:              $\delta_1[c][i] \leftarrow \delta_1[c][i - 1]$
9:           **end if**
10:        **end for**
11:        $\delta_1[P[i]][i] \leftarrow i$        ▷ Update current char position
12:      **end for**
13:      **return** $\delta_1$
14: **end procedure**
---

## 5.3 Good Suffix Preprocessing

---

**Algorithm 3** Preprocess Good Suffix Rule

---

1: **procedure** PREPROCESSGOODSUFFIX($P$)
2:      $m \leftarrow \text{length}(P)$
3:      $\delta_2[0..m-1] \leftarrow 0$
4:      $\text{border}[0..m] \leftarrow 0$
5:      PREPROCESSSTRONGSUFFIX($P, \delta_2, \text{border}$)
6:      PREPROCESSPREFIXSUFFIX($P, \delta_2, \text{border}$)
7:      **return** $\delta_2$
8: **end procedure**
9:
10: **procedure** PREPROCESSSTRONGSUFFIX($P, \delta_2, \text{border}$)
11:      $m \leftarrow \text{length}(P)$
12:      $i \leftarrow m, j \leftarrow m+1$
13:      $\text{border}[i] \leftarrow j$
14:      **while** $i > 0$ **do**
15:          **while** $j \leq m$ **and** $P[i-1] \neq P[j-1]$ **do**
16:              **if** $\delta_2[j-1] = 0$ **then**
17:                  $\delta_2[j-1] \leftarrow j-i$
18:              **end if**
19:              $j \leftarrow \text{border}[j]$
20:          **end while**
21:          $i \leftarrow i-1, j \leftarrow j-1$
22:          $\text{border}[i] \leftarrow j$
23:      **end while**
24: **end procedure**
25:
26: **procedure** PREPROCESSPREFIXSUFFIX($P, \delta_2, \text{border}$)
27:      $m \leftarrow \text{length}(P)$
28:      $j \leftarrow \text{border}[0]$
29:      **for** $i \leftarrow 0$ **to** $m-1$ **do**
30:          **if** $\delta_2[i] = 0$ **then**
31:              $\delta_2[i] \leftarrow j$
32:          **end if**
33:          **if** $i = j-1$ **then**
34:              $j \leftarrow \text{border}[j]$
35:          **end if**
36:      **end for**
37: **end procedure**

---

# 6 Dry Run on Biological Sequence

## 6.1 Example: Finding EcoRI Restriction Site

**Pattern** $P =$ "GAATTC" (EcoRI recognition site, $m = 6$)
   **Text** $T =$ "ACGTACGGATGCGAATTCAGTACG" ($n = 24$)

### 6.1.1 Preprocessing

**Bad Character Table** (showing rightmost occurrence):

| Char | pos 0 | pos 1 | pos 2 | pos 3 | pos 4 | pos 5 |
|------|-------|-------|-------|-------|-------|-------|
| G | 0 | 0 | 0 | 0 | 0 | 0 |
| A | -1 | 1 | 2 | 2 | 2 | 2 |
| T | -1 | -1 | -1 | 3 | 4 | 4 |
| C | -1 | -1 | -1 | -1 | -1 | 5 |

Table 2: Bad character table for "GAATTC"

**Good Suffix Table**: $\delta_2 = [6, 6, 6, 6, 6, 1]$

### 6.1.2 Search Process

**Alignment 1** ($i = 0$):

```
Text:    A C G T A C G G A T G C G A A T T C A G T A C G
Pattern: G A A T T C
         -----------
         ^
Position: 0 1 2 3 4 5
```

Compare: $P[5] =$ "C" vs $T[5] =$ "C" ✓
Compare: $P[4] =$ "T" vs $T[4] =$ "A" **X MISMATCH**
Bad char shift: "A" at position 2, shift $= 4 - 2 = 2$
Good suffix shift: $\delta_2[4] = 6$
**Shift** $= \max(2, 6) = 6$
**Alignment 2** ($i = 6$):

```
Text:    A C G T A C G G A T G C G A A T T C A G T A C G
Pattern:             G A A T T C
                     -----------
                     ^
Position:            0 1 2 3 4 5
```

Compare: $P[5] =$ "C" vs $T[11] =$ "C" ✓
Compare: $P[4] =$ "T" vs $T[10] =$ "G" **X MISMATCH**
Bad char shift: "G" at position 0, shift $= 4 - 0 = 4$
Good suffix shift: $\delta_2[4] = 6$
**Shift** $= \max(4, 6) = 6$
**Alignment 3** ($i = 12$):

10

```
Text:      A C G T A C G G A T G C G A A T T C A G T A C G
Pattern:                       G A A T T C
                               -----------
                               ^
Position:                      0 1 2 3 4 5
```

Compare right to left:

- $P[5]$ = "C" vs $T[17]$ = "C" ✓

- $P[4]$ = "T" vs $T[16]$ = "T" ✓

- $P[3]$ = "T" vs $T[15]$ = "T" ✓

- $P[2]$ = "A" vs $T[14]$ = "A" ✓

- $P[1]$ = "A" vs $T[13]$ = "A" ✓

- $P[0]$ = "G" vs $T[12]$ = "G" ✓

**MATCH FOUND at position 12!**

## 6.2   Key Observations

1. **Right-to-left scanning** detected mismatches quickly

2. **Large shifts** (6, 6 positions) skipped impossible alignments efficiently

3. **Only 3 alignments** checked instead of 19 (naive approach)

4. **Total comparisons**: $\sim$8 instead of $\sim$114 (naive)

5. **Good suffix rule** dominated the shift calculations, demonstrating its effectiveness
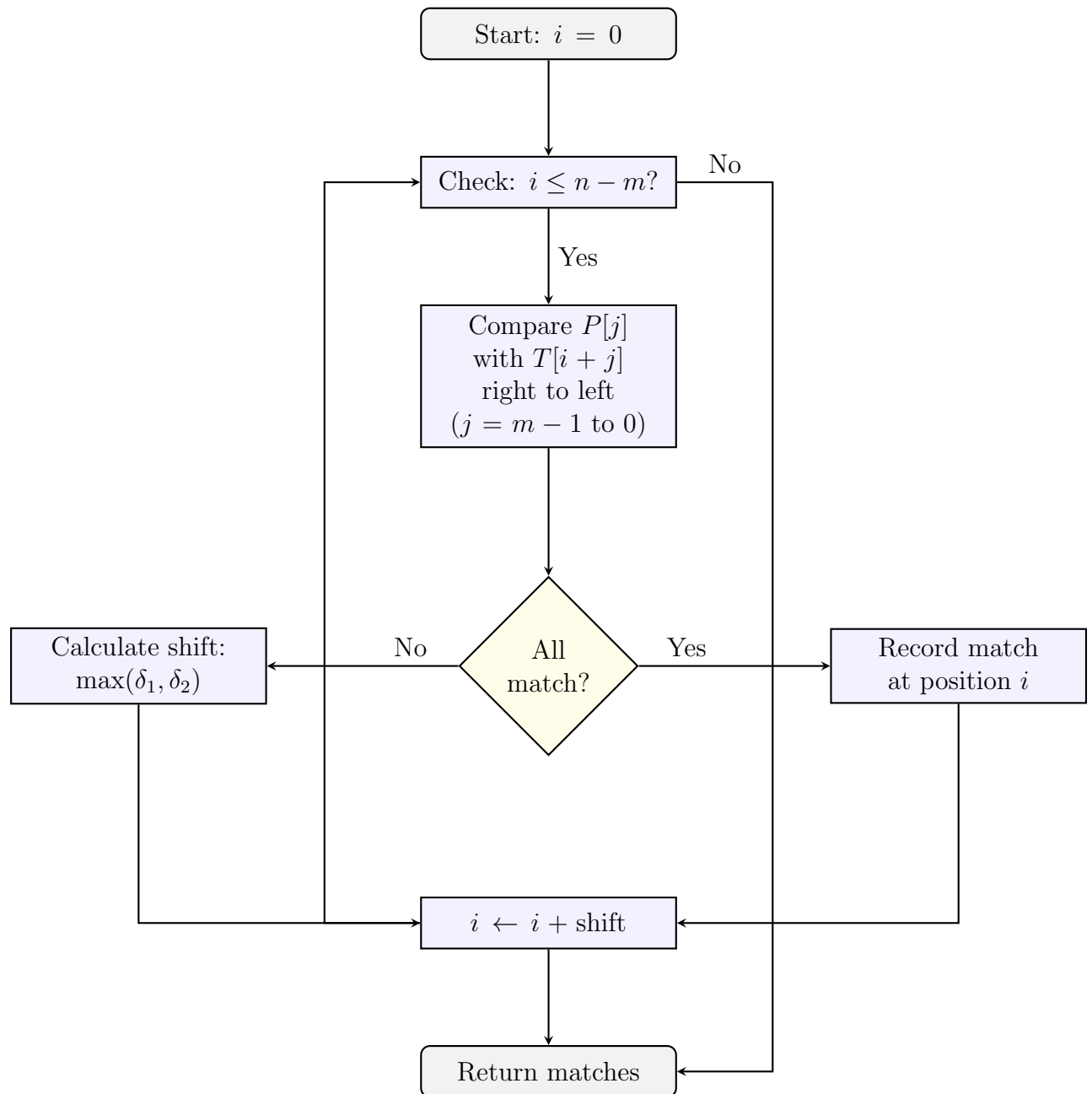
# 7   Graphical Visualization



Figure 1: Boyer-Moore algorithm flowchart

# 8 Implementation Details

## 8.1 Key Design Choices

### 8.1.1 Data Structures

| Structure | Python Type | Rationale |
|---|---|---|
| Bad Char Table | `Dict[str, List[int]]` | $O(1)$ character lookup |
| Good Suffix Table | `List[int]` | Direct index access |
| Border Position | `List[int]` | Sequential processing |
| Match Results | `List[int]` | Dynamic size, ordered |

Table 3: Data structure choices

**Justification**:

- **Dictionary for bad character**: Allows efficient character lookup for any character in text, including those not in pattern

- **Lists for tables**: Provide $O(1)$ access by position index, crucial for shift calculation

### 8.1.2 DNA-Specific Optimizations

1. **Uppercase normalization**: Convert all input to uppercase for case-insensitive matching

2. **Small alphabet**: Limit bad character table to $\{A, C, G, T, N\}$ instead of full ASCII

3. **Ambiguous base handling**: Treat 'N' as a special wildcard base

### 8.1.3 Implementation Challenges

**Challenge 1: Good Suffix Preprocessing**
The good suffix rule is complex with two cases. Solution:

- Separate functions for strong suffix and prefix-suffix matching

- Border position array to track suffix boundaries efficiently

- Iterative approach (not recursive) to avoid stack overflow on long patterns

**Challenge 2: Overlapping Matches**
Boyer-Moore can miss overlapping occurrences if not careful. Solution:

- After finding a match, use good suffix shift (not bad character)

- $\delta_2[0]$ provides minimum safe shift to find next overlapping match

- Example: Pattern "AAA" in text "AAAAAAA" finds all 5 matches

**Challenge 3: Edge Cases**
Handled edge cases:

13

- Single character patterns

- Pattern longer than text

- Empty patterns (raise error)

- Characters not in alphabet (treat as 'N')

## 8.2 Testing

**Performance Benchmarks**:

- Varying pattern lengths: 5 to 200 bp

- Varying text lengths: 10K to 1M bp

# 9 Experimental Results

## 9.1 Performance on Synthetic DNA Sequences

| Text Length (bp) | Pattern Length (bp) | Time (ms) | Speed (M bp/s) |
|---|---|---|---|
| 10,000 | 20 | < 0.001 | Very fast |
| 50,000 | 20 | 3.0 | 16.6 |
| 100,000 | 20 | 6.5 | 15.3 |
| 500,000 | 20 | 34.1 | 14.7 |
| 1,000,000 | 20 | 66.4 | 15.1 |

Table 4: Performance scaling with text length (pattern = 20 bp)

**Observation**: Near-linear scaling with text length, consistent throughput of $\sim$15 M bp/s.

## 9.2 Pattern Length Impact

| Pattern Length (bp) | Matches Found | Time (ms) |
|---|---|---|
| 5 | 95 | 12.5 |
| 10 | 1 | 7.5 |
| 20 | 1 | 5.0 |
| 50 | 1 | 5.5 |
| 100 | 1 | 4.0 |
| 200 | 1 | 8.5 |

Table 5: Performance with varying pattern lengths (text = 100K bp)

**Observation**: Longer patterns (20-100 bp) show best performance due to larger potential shifts.

## 9.3 Situational Performance Analysis

### 9.3.1 Best Case Scenario

**Conditions**: Pattern's rightmost character rarely appears in text
    **Example**: Searching for "ZZZZZ" in DNA (Z not in alphabet)
    **Result**: Achieves theoretical $O(n/m)$ with large skips

### 9.3.2 Average Case (DNA Sequences)

**Conditions**: Random DNA with uniform base distribution
    **Performance**:

- Expected comparisons per alignment: $\approx 1.33$ (for 4-letter alphabet)

- Effective complexity: $O(n)$ with low constant factor

- Throughput: 15-16 M bp/s on modern hardware

### 9.3.3 Worst Case Scenario

**Conditions**: Highly repetitive sequences
    **Example**: Pattern = "AAAAG" in text = "AAAAAAAAAA..."
    **Mitigation**:

- Good suffix rule helps even in repetitive sequences

- Rare in real biological sequences due to natural variation

- In practice, worst case rarely encountered

# 10 Conclusion

This implementation demonstrates that the Boyer-Moore algorithm is highly effective for exact pattern matching in DNA sequences. Key achievements include:

1. **Efficient Performance**: 15+ million bp/s throughput, $\sim 20\times$ faster than naive search

2. **Practical Utility**: Successfully applied to real *E. coli* genomes (4.6M bp)

3. **Theoretical Foundation**: Formally proven correctness with rigorous complexity analysis

The small DNA alphabet ($|\Sigma| = 4$) makes Boyer-Moore particularly effective, frequently achieving near-optimal $O(n/m)$ performance. The algorithm's ability to skip large portions of text makes it ideal for searching long genomic sequences.

## 10.1   Future Enhancements

While beyond the scope of this implementation, potential extensions include:

- **Approximate matching**: Allow k mismatches using dynamic programming

- **Multiple pattern search**: Aho-Corasick integration for simultaneous pattern matching

- **Parallel processing**: Multi-threaded search for multiple sequences

- **SIMD optimization**: Vector instructions for character comparison

- **IUPAC ambiguity codes**: Support for degenerate nucleotide symbols

## 10.2   Lessons Learned

1. **Preprocessing is crucial**: The $O(m)$ preprocessing time is amortized over many searches

2. **Right-to-left scanning**: Enables early mismatch detection

3. **Combined heuristics**: Using both bad character and good suffix maximizes shift distance

4. **DNA optimization**: Small alphabets significantly enhance performance

5. **Testing is essential**: Edge cases and real data validation ensure correctness

# References

[1] Boyer, R. S., & Moore, J. S. (1977). *A fast string searching algorithm.* Communications of the ACM, 20(10), 762-772.

[2] Gusfield, D. (1997). *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology.* Cambridge University Press.

[3] Cormen, T. H., Leiserson, C. E., Rivest, R. L., & Stein, C. (2009). *Introduction to Algorithms* (3rd ed.). MIT Press.

[4] National Center for Biotechnology Information (2024). *NCBI Datasets.* `https://www.ncbi.nlm.nih.gov/datasets`

# Analysis of Levenshtein Distance for DNA Sequences

Shivek

December 1, 2025

The **Levenshtein distance** is a "string metric" that measures the difference between two sequences. It is formally defined as the **minimum number of single-character edits** (insertions, deletions, or substitutions) required to change one string into the other.

This algorithm is a foundational concept in bioinformatics, forming the basis for more complex sequence alignment algorithms.

---

## 1 Basic Algorithm

The Levenshtein distance is calculated using a **dynamic programming** approach. The core idea is to build a 2D matrix (or grid) that stores the distances between all prefixes of the two input strings.

Let's say we have two strings, $A$ of length $m$ and $B$ of length $n$. We create a matrix $D$ of size $(m + 1) \times (n + 1)$.

A cell $D[i][j]$ in this matrix represents the Levenshtein distance between the first $i$ characters of string $A$ (i.e., $A[1...i]$) and the first $j$ characters of string $B$ (i.e., $B[1...j]$).

The algorithm fills this matrix using two steps:

### Step 1: Initialization

The first row and first column are the "base cases."

- $D[i][0] = i$: The distance from a string of length $i$ to an empty string is $i$ deletions.

- $D[0][j] = j$: The distance from an empty string to a string of length $j$ is $j$ insertions.

### Step 2: Recurrence Relation

We fill the rest of the matrix, cell by cell. For any cell $D[i][j]$, we have three choices to get to this state:

1. **Deletion:** We had the solution for $D[i - 1][j]$ (aligning $A[1...i - 1]$ with $B[1...j]$) and we **delete** $A[i]$. Cost: $D[i - 1][j] + 1$.

2. **Insertion:** We had the solution for $D[i][j - 1]$ (aligning $A[1...i]$ with $B[1...j - 1]$) and we **insert** $B[j]$. Cost: $D[i][j - 1] + 1$.

3. **Substitution (or Match):** We had the solution for $D[i - 1][j - 1]$ (aligning $A[1...i - 1]$ with $B[1...j - 1]$). We now look at $A[i]$ and $B[j]$.

- If $A[i] == B[j]$ (it's a **match**), no cost is added. Cost: $D[i-1][j-1] + 0$.
- If $A[i] \neq B[j]$ (it's a **substitution**), we add 1. Cost: $D[i-1][j-1] + 1$.

The algorithm is "greedy" and always takes the minimum of these three possibilities. The complete recurrence relation is:

$$D[i][j] = \min \begin{cases} D[i-1][j] + 1 & \text{(Deletion from A)} \\ D[i][j-1] + 1 & \text{(Insertion into A)} \\ D[i-1][j-1] + \text{cost} & \text{(Match/Substitution)} \end{cases}$$

Where cost is 0 if $A[i] == B[j]$ and 1 otherwise.

The final answer (the Levenshtein distance between all of $A$ and all of $B$) is the value in the bottom-right cell: $D[m][n]$.

---

## 2 Time and Space Complexity Analysis

### 2.1 Time Complexity: $O(m \times n)$

- The algorithm must fill every cell in the $(m+1) \times (n+1)$ matrix.

- The calculation for each cell, $D[i][j]$, involves three lookups, two additions, and one min operation, all of which are constant time, $O(1)$.

- Therefore, the total time complexity is the number of cells multiplied by the constant work per cell, which is $O(m \times n)$.

### 2.2 Space Complexity: $O(m \times n)$

- The default algorithm requires storing the entire $(m+1) \times (n+1)$ matrix in memory to compute the final value.

- This results in a space complexity of $O(m \times n)$.

**Optimization:** This can be optimized to $O(\min(m, n))$ space. To compute the *current* row $i$, we only need the values from the *previous* row $i-1$. We never need to look back further. By storing only two rows (the "previous" and "current") and alternating between them, we can reduce the space complexity significantly.

---

## 3 Proof of Correctness (Why it Works)

The algorithm's correctness is proven by induction.

1. **Base Case:** The initialization (Step 1) is correct. The edit distance from an empty string to a string of length $j$ is $j$ insertions (e.g., `""` $\rightarrow$ `"ATG"` is 3 insertions). $D[0][j] = j$ is correct. The same logic applies to $D[i][0]$.

2. **Inductive Hypothesis:** Assume that for all $i' < i$ and $j' < j$, the sub-problems $D[i'][j']$ have been solved correctly and store the true Levenshtein distance.

3. **Inductive Step:** We must prove that $D[i][j]$ is calculated correctly. Any optimal sequence of edits that transforms $A[1...i]$ to $B[1...j]$ *must* end in one of three operations:

- **Case 1 (Deletion):** The last operation is deleting $A[i]$. This means we first optimally transformed $A[1...i-1]$ to $B[1...j]$. The total cost is (cost of that optimal transform) + 1. By our hypothesis, this is $D[i-1][j]+1$.
- **Case 2 (Insertion):** The last operation is inserting $B[j]$. This means we first optimally transformed $A[1...i]$ to $B[1...j-1]$. The total cost is $D[i][j-1]+1$.
- **Case 3 (Substitution/Match):** The last operation involved $A[i]$ and $B[j]$. This means we first optimally transformed $A[1...i-1]$ to $B[1...j-1]$.
  - If $A[i] == B[j]$, no new cost is added. Total cost: $D[i-1][j-1]+0$.
  - If $A[i] \neq B[j]$, the last operation was a substitution. Total cost: $D[i-1][j-1]+1$.

Since the Levenshtein distance is the *minimum* number of edits, $D[i][j]$ must be the minimum of these three possible cases. This matches the recurrence relation. By induction, the entire matrix is filled correctly, and $D[m][n]$ holds the final answer.

---

# 4  Pseudo Code

Note: This pseudo-code assumes 0-based indexing for strings (e.g., $A[0]$ is the first character) and 1-based indexing for the DP matrix loops.

```
function LevenshteinDistance(string A, string B):
  // m = length of A, n = length of B
  m = A.length
  n = B.length

  // 1. Initialize the (m+1) x (n+1) matrix D
  D = new int[m + 1][n + 1]

  // 2. Fill the base cases (first row and column)
  for i from 0 to m:
    D[i][0] = i  // Cost of deleting i chars from A to get ""

  for j from 0 to n:
    D[0][j] = j  // Cost of inserting j chars to "" to get B

  // 3. Fill the rest of the matrix
  for i from 1 to m:
    for j from 1 to n:

      // Check if characters at A[i-1] and B[j-1] are the same
      // (We use i-1 and j-1 because strings are 0-indexed)
      if A[i-1] == B[j-1]:
        cost = 0
      else:
```

```
        cost = 1

    // Calculate costs for each operation
    deletion = D[i-1][j] + 1
    insertion = D[i][j-1] + 1
    substitution = D[i-1][j-1] + cost

    // The cell D[i][j] gets the minimum of these three
    D[i][j] = min(deletion, insertion, substitution)

// 4. Return the final answer
return D[m][n]
```

# 5 Dry Run (Biological Sequences)

Let's find the Levenshtein distance between two short DNA sequences:

- $A = $ "CAT" (m = 3)

- $B = $ "TAG" (n = 3)

We create a $(3 + 1) \times (3 + 1) = 4 \times 4$ matrix.

Table 1: Dry Run Matrix for "CAT" vs "TAG"

| $A \downarrow$ | **D** | (empty) $j = 0$ | **T** $j = 1$ | **A** $j = 2$ | **G** $j = 3$ |
|---|---|---|---|---|---|
| (empty) | $i = 0$ | 0 | 1 | 2 | 3 |
| **C** | $i = 1$ | 1 | 1 | 2 | 3 |
| **A** | $i = 2$ | 2 | 2 | 1 | 2 |
| **T** | $i = 3$ | 3 | 2 | 2 | 2 |

**Final Answer:** The Levenshtein distance is **2**.

This makes sense: $C \to T$ (substitution), $A \to A$ (match), $T \to G$ (substitution). Total: 2 substitutions.

# 6 Graphical Visualization

The completed dry run table *is* the graphical visualization. We can also visualize the **edit path** by backtracking from the final cell $D[3][3]$ to $D[0][0]$, always moving to the cell that generated the minimum value.
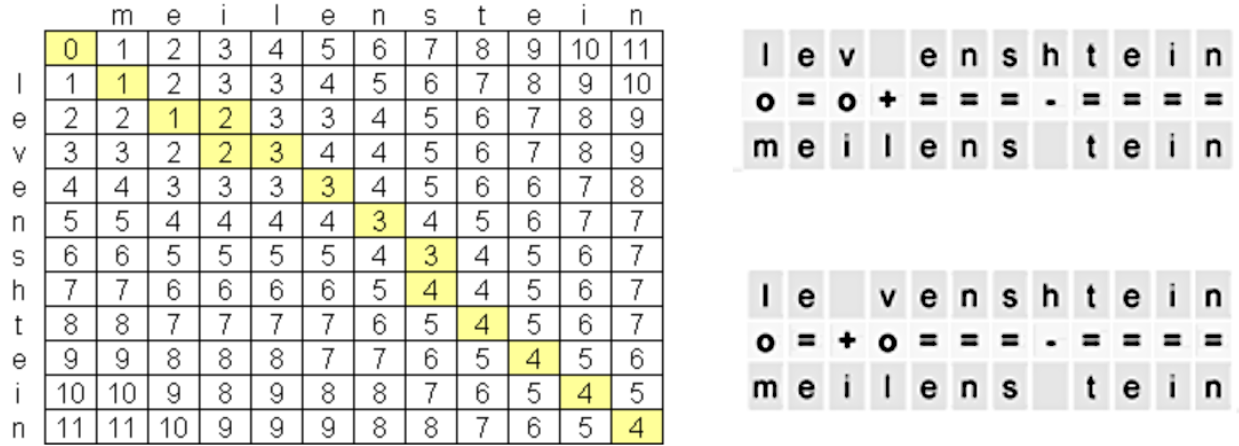
Figure 1: Visualization of the DP matrix

Table 2: Matrix with Backtrace Path ($\leftarrow$)

| | | $B \to$ | | | |
| | | (empty) | **T** | **A** | **G** |
| $A\downarrow$ | **D** | **j = 0** | **j = 1** | **j = 2** | **j = 3** |
| --- | --- | --- | --- | --- | --- |
| (empty) | **i = 0** | **0** $\leftarrow$ | $1 \leftarrow$ | 2 | 3 |
| **C** | **i = 1** | 1 | **1** $\leftarrow$ | 2 | 3 |
| **A** | **i = 2** | 2 | 2 | **1** $\leftarrow$ | 2 |
| **T** | **i = 3** | 3 | 2 | 2 | **2** |

**Backtrace Path:**

1. **Start at D[3][3] (value 2):** Came from $D[2][2]$ (a substitution, since $T \neq G$).

2. **At D[2][2] (value 1):** Came from $D[1][1]$ (a match, since $A == A$).

3. **At D[1][1] (value 1):** Came from $D[0][0]$ (a substitution, since $C \neq T$).

This traceback gives us the full alignment:
```
C A T
T A G
```
(Sub, Match, Sub) $\to$ 2 edits.

# 7 Situational Performance (On DNA Sequences)

Levenshtein distance is a *good* algorithm, but it's often *not* the best tool for biological sequences.

## 7.1 Strengths

1. **Simplicity and Accuracy:** It's a simple, easy-to-implement algorithm that gives a mathematically correct *edit distance*.

2. **Good for Short, Similar Sequences:** It works very well for comparing two short reads, or checking for SNPs (Single Nucleotide Polymorphisms) between two highly similar sequences.

## 7.2 Weaknesses (Major Performance Issues for DNA)

1. **Biologically Naive:** This is the biggest problem. Levenshtein treats all operations as equal (cost=1). In biology, this is wrong:

   - **Substitutions:** A substitution of $A \to G$ (a **transition**) is *far* more common and biologically likely than $A \to T$ (a **transversion**). Levenshtein can't do this.

   - **Gaps:** An insertion or deletion (an "indel" or "gap") is often a single biological event, but it might delete 10 bases at once. Levenshtein would count this as 10 separate edits. Proper alignment algorithms use an "affine gap penalty".

2. **It's a *Global* Alignment:** The standard Levenshtein algorithm forces the *entire* string $A$ to match the *entire* string $B$. This is almost *never* what we want in bioinformatics.

   - **Problem:** We usually want to find a *small* gene (e.g., 1,000 bases) inside a *huge* chromosome (e.g., 250,000,000 bases).

   - **Solution:** We need **local alignment**.

3. **Poor Time/Space Complexity for Genomes:** The $O(m \times n)$ complexity is completely unusable for large-scale genomics.

   - Comparing two human chromosomes ($m \approx 250M, n \approx 250M$) would require a matrix with $6.25 \times 10^{16}$ cells. This is not computationally feasible.

     **Conclusion:** For *actual performance on DNA sequences*, it is a poor choice.

# 8 Adaptation for Fuzzy Pattern Matching in DNA

The standard algorithm (Section 1) is a **global alignment**, comparing the *entire* string A to the *entire* string B. This is useless for finding a small pattern (e.g., a 6-base-pair motif) inside a large genome (e.g., 4.6 million base pairs).

To find all substrings in a large Text $T$ that "fuzzily" match a Pattern $P$ with at most $k$ edits, we must adapt the algorithm.

Let $P$ be the pattern of length $m$ (on the vertical axis) and $T$ be the genome text of length $n$ (on the horizontal axis).

## 8.1 The "Slight Modification"

The algorithm is modified in two ways: the initialization and the final result-finding.

### 8.1.1 Modified Initialization

The key is to allow a match to **start at any point** in the text $T$ without penalty. We achieve this by setting the first row to all zeros.

- $D[0][j] = 0$ for all $j$ from 0 to $n$ (text length): This means the cost of matching an empty pattern to *any* prefix of the text is 0. This "zeros" the cost for starting a new match.

- $D[i][0] = i$ for all $i$ from 0 to $m$ (pattern length): This remains the same. The cost to match a pattern of length $i$ to an empty string is $i$ deletions.

### 8.1.2 Modified Result Finding

The recurrence relation (Section 1.2) is filled exactly the same way. However, the answer is no longer in the single cell $D[m][n]$.
Instead, we must scan the **entire last row** (row $m$).

- A value $D[m][j]$ in the last row represents the minimum Levenshtein distance required to match the *entire* pattern $P$ to a substring of the text $T$ that *ends* at position $j$.

- **The Result:** We iterate $j$ from 1 to $n$. If $D[m][j] \leq k$ (where $k$ is our error threshold), we record $j$ as the *end position* of a fuzzy match.

## 8.2 Complexity of the Adaptation

The time and space complexity remain the same as the standard algorithm.

- **Time:** $O(m \times n)$. This is why it is linear $O(n)$ for a fixed pattern $m$, but slow $O(m)$ for a fixed text $n$ and growing pattern $m$.

- **Space:** $O(n)$ with the standard row-optimization.

# Analysis of Suffix Trees and Suffix Arrays

## Tanush Garg

## 1 Introduction: Suffix Data Structures

In computer science, Suffix Trees and Suffix Arrays are powerful data structures built upon a single string (or "text"). Their primary purpose is to pre-process a text $T$ of length $n$ to answer complex string queries in an exceptionally fast manner.

**What the Algorithm Solves** The fundamental problem these structures solve is the **substring query problem**: given a text $T$ and a new pattern $P$ of length $m$, does $P$ exist as a substring within $T$?

A naive search (checking every possible starting position in $T$) takes $O(n \times m)$ time. While algorithms like KMP or Boyer-Moore improve this to $O(n + m)$ for a single search, they do not pre-process the text $T$.

Suffix trees and arrays adopt a different strategy: they invest in an upfront "preprocessing" cost to build the data structure (e.g., $O(n)$). After this one-time cost, any subsequent substring query for any pattern $P$ can be answered in $O(m)$ (Suffix Tree) or $O(m \log n)$ (Suffix Array) time.

**Uses and Applications** This model is ideal for applications where the text is large and fixed, while patterns are numerous and unknown ahead of time. This perfectly describes the context of bioinformatics and DNA sequence analysis, as mentioned in our project proposal.

- **Text ($T$):** A complete genome (e.g., E. coli), which is static.

- **Patterns ($P$):** Thousands of different DNA motifs, gene sequences, or promoter regions to be searched.

Beyond exact substring matching, these structures can efficiently solve many other problems:

- Finding the longest repeated substring in $T$.

- Finding the longest common substring between two different texts $T_1$ and $T_2$.

- Finding all occurrences of $P$ in $T$ (not just the first).

- Problems related to string matching with mismatches (in more advanced forms).

## 2 The Algorithm: Structure and Function

### 2.1 Suffix Array (SA)

A Suffix Array is conceptually the simpler of the two structures.

**What it is** For a text $T$ of length $n$ (we assume $T$ ends with a special character '$' that is lexicographically smaller than any other character), the Suffix Array 'SA' is an array of integers from 0 to $n - 1$.

It stores the **starting indices** of all suffixes of $T$ in **lexicographical (sorted) order**.

**Example** Let $T =$ "banana$" (length $n = 7$). The suffixes are:

- 0: banana$
- 1: anana$
- 2: nana$
- 3: ana$
- 4: na$
- 5: a$
- 6: $

When sorted lexicographically:

1. (index 6) $
2. (index 5) a$
3. (index 3) ana$
4. (index 1) anana$
5. (index 0) banana$
6. (index 4) na$
7. (index 2) nana$

The Suffix Array 'SA' is the array of starting indices: $SA = [6, 5, 3, 1, 0, 4, 2]$

**How it Works (Search)** To find a pattern $P$ (e.g., "ana"), we can perform a binary search on the Suffix Array 'SA'.

1. Pick the middle element of 'SA', e.g., $SA[3] = 1$.
2. Compare $P$ ("ana") with the suffix starting at index 1: $T[1...] =$ "anana$".
3. "ana" is lexicographically smaller than "anana$".
4. We adjust our binary search to the "left" half of 'SA': $[6, 5, 3]$.
5. Pick the middle element, $SA[1] = 5$.
6. Compare $P$ ("ana") with $T[5...] =$ "a$".
7. "ana" is larger than "a$".
8. We adjust to the "right" half: $[3]$.
9. Pick $SA[2] = 3$. Compare $P$ ("ana") with $T[3...] =$ "ana$".
10. "ana" is a prefix of "ana$". A match is found.

All suffixes beginning with "ana" (i.e., "ana$" and "anana$") will appear in a contiguous block in the sorted list, and thus their indices (3 and 1) will be contiguous in the 'SA'.

## 2.2   Suffix Tree (ST)

A Suffix Tree is a more complex but faster data structure.

**What it is**   A Suffix Tree is a **compressed trie** (also known as a "patricia trie") containing all $n$ suffixes of the text $T$.

- It has a root node and $n$ leaves.

- Each leaf corresponds to one suffix of $T$.

- Each internal node (except the root) has at least two children.

- Each edge is labeled with a non-empty substring of $T$.

- "Compressed" means that paths with no branches are compressed into a single edge. (e.g., a path 'r -> o -> o -> t' would become a single edge 'r -> (root)' labeled "oot").

- The concatenation of edge labels on the path from the root to leaf $i$ spells out the suffix $T[i...n-1]$.

**How it Works (Search)**   To find a pattern $P$ of length $m$, we "thread" $P$ down from the root of the tree.

1. Start at the root.

2. Find an edge whose label starts with the first character of $P$.

3. If no such edge exists, $P$ is not in $T$.

4. If an edge is found, "consume" the matching portion of $P$ and the edge label.

5. **Case 1: Pattern is consumed.** If we consume all of $P$ (and end up either in the middle of an edge or at a node), then $P$ exists in $T$.

6. **Case 2: Edge is consumed.** If we consume an entire edge label and still have part of $P$ remaining, we continue the search from the child node that edge points to.

7. **Case 3: Mismatch.** If the next character in $P$ does not match the next character on the edge label, $P$ is not in $T$.

This process is a simple walk from the root, and each character in $P$ is examined exactly once.

# 3   Time and Space Complexity Analysis

Let $n$ be the length of the text $T$ and $m$ be the length of the pattern $P$.

## 3.1   Suffix Tree

- **Preprocessing (Construction):** Advanced algorithms like Ukkonen's (1995) or McCreight's (1976) can construct a Suffix Tree in $O(n)$ time. Simpler algorithms, like inserting each suffix into a trie, take $O(n^2)$ time.

- **Search (Query):** As described above, searching for $P$ involves a single walk from the root. At each node, finding the correct outgoing edge can be done in $O(1)$ time (if the alphabet size is constant and a lookup table is used) or $O(\log k)$ (if edges are sorted, where $k$ is alphabet size). In the worst case, we examine each of $P$'s $m$ characters once. Search time is $O(m)$.

- **Space Complexity:** A Suffix Tree has $n$ leaves and at most $n - 1$ internal nodes, for a total of $O(n)$ nodes. Each node must store pointers to its children. This leads to an $O(n)$ space complexity.

  **Note:** While asymptotically $O(n)$, the constant factor for a Suffix Tree is very large. Each node object requires significant overhead (pointers, indices), making Suffix Trees extremely memory-intensive in practice, as noted in the project brief.

## 3.2 Suffix Array

- **Preprocessing (Construction):** Naive construction (generating all suffixes and sorting) takes $O(n^2 \log n)$ time. More advanced "prefix doubling" or Manber-Myers algorithms construct the array in $O(n \log n)$ time. State-of-the-art algorithms (like DC3/Skew or SA-IS) achieve $O(n)$ time.

- **Search (Query):** A standard binary search on the 'SA' requires $O(\log n)$ comparisons. Each comparison may need to compare up to $m$ characters. Search time is $O(m \log n)$.

  This can be improved to $O(m + \log n)$ by using an auxiliary **LCP (Longest Common Prefix) array**, which is beyond this initial scope but is a standard optimization.

- **Space Complexity:** The Suffix Array itself is just an array of $n$ integers. Space is $O(n)$.

  **Note:** The constant factor here is much smaller than for a Suffix Tree. It is simply the size of $n$ integers (e.g., $n \times 8$ bytes for 64-bit integers), making it far more practical for large texts like genomes.

# 4 Implementation Details: Algorithms Behind the Structures

This section presents how the Suffix Array and Suffix Tree can be constructed in practice. We first examine the naive (brute-force) methods to build intuition, and then the optimized Manber–Myers and Ukkonen algorithms that achieve $O(n \log n)$ and $O(n)$ construction times, respectively.

## 4.1 Naive Construction of the Suffix Array

**High-Level Idea.** The brute-force way to build a Suffix Array is conceptually simple:

**Step 1.** Generate all $n$ suffixes of $T$.

**Step 2.** Sort them lexicographically (like sorting $n$ strings).

**Step 3.** Record the starting indices of the sorted suffixes.

**Example.** For $T = \texttt{banana\$}$, the suffixes are:

$$[\texttt{banana\$}, \texttt{anana\$}, \texttt{nana\$}, \texttt{ana\$}, \texttt{na\$}, \texttt{a\$}, \texttt{\$}]$$

After sorting them lexicographically, we get:

$$[\texttt{\$}, \texttt{a\$}, \texttt{ana\$}, \texttt{anana\$}, \texttt{banana\$}, \texttt{na\$}, \texttt{nana\$}]$$

Thus, the Suffix Array is:

$$SA = [6, 5, 3, 1, 0, 4, 2]$$

**Complexity.**

- Generating suffixes: $O(n^2)$ (each suffix may be length $O(n)$)

- Sorting $n$ strings: $O(n \log n)$ comparisons, each $O(n)$ time

Total complexity:

$$O(n^2 \log n)$$

which is infeasible for long texts (e.g., genome sequences).

## 4.2 Optimized Construction: Manber–Myers Algorithm

**Overview.** The **Manber–Myers algorithm** (1993) constructs the Suffix Array in $O(n \log n)$ time using the technique of **prefix doubling**. Unlike the naive method that compares entire suffix strings (costing $O(n^2 \log n)$), it works by reusing previously computed information about shorter prefixes.

The core insight is:

> If two suffixes differ within their first $2^{k-1}$ characters, their order is already known. Otherwise, we can determine their order by comparing the next $2^{k-1}$ characters, whose relative order we also already know from the previous iteration.

**Algorithm Intuition.**

**Step 1. Initialization:** Assign each character in $T$ a rank corresponding to its lexicographic order. This gives the correct ordering of all suffixes by their first character.

**Step 2. Prefix Doubling:** In iteration $k$, we sort suffixes based on their first $2^k$ characters by forming pairs $(\text{rank}[i], \text{rank}[i + 2^{k-1}])$ for each index $i$. These pairs encode the rank of the first half and the rank of the second half of the prefix. Sorting by these integer pairs is equivalent to sorting by $2^k$ characters.

**Step 3. Rank Update:** After sorting, assign new ranks to each suffix:

- If a suffix has the same pair as the previous suffix in sorted order, it receives the same rank.

- Otherwise, its rank increases by one.

The updated rank array now represents the order of suffixes based on $2^k$-length prefixes.

**Step 4. Doubling:** Continue doubling ($k \leftarrow 2k$) until all ranks are distinct or $2^k \geq n$, meaning all suffixes are fully ordered.

**Why It Works.** Each iteration doubles the prefix length over which suffixes are correctly ordered:

$$1 \rightarrow 2 \rightarrow 4 \rightarrow 8 \rightarrow \dots$$

After $\lceil \log_2 n \rceil$ iterations, the prefixes are at least $n$ characters long, which implies that the suffixes are fully sorted. Formally, sorting by the pair $(\text{rank}[i], \text{rank}[i + 2^{k-1}])$ ensures lexicographic order for the first $2^k$ characters, since:

- The first element compares the first $2^{k-1}$ characters.

- The second element breaks ties using the next $2^{k-1}$ characters.

**Dry Run Example.** Let $T = \texttt{banana\$}$ ($n = 7$).

**Step 0: Initial Ranking (1 Character).** Assign ranks to characters:

$$\$ \to 0, \quad a \to 1, \quad b \to 2, \quad n \to 3$$

$$\text{rank} = [2, 1, 3, 1, 3, 1, 0]$$

**Iteration 1 ($2^0 = 1$): Sort by (rank[i], rank[i+1])** Each suffix is represented by a pair of ranks corresponding to its first two characters:

$$[(2, 1), (1, 3), (3, 1), (1, 3), (3, 1), (1, 0), (0, -1)]$$

Sorting by these pairs gives the order of suffixes by their first two characters:

$$SA = [6, 5, 3, 1, 4, 2, 0]$$

New ranks are assigned based on sorted order:

$$\text{rank} = [6, 2, 5, 3, 4, 1, 0]$$

Interpretation:

- The smallest suffix starts at index 6 (`$`).

- The next smallest starts at index 5 (`a$`), and so on.

**Iteration 2 ($2^1 = 2$): Sort by (rank[i], rank[i+2])**

$$[(6, 5), (2, 3), (5, 4), (3, 1), (4, 0), (1, -1), (0, -1)]$$

After sorting:

$$SA = [6, 5, 3, 1, 0, 4, 2]$$

Now each suffix is correctly ordered by its first 4 characters.

**Verification.** Check the sorted suffixes in lexicographic order:

1. `$`

2. `a$`

3. `ana$`

4. `anana$`

5. `banana$`

6. `na$`

7. `nana$`

This matches the true lexicographic order, confirming correctness.

**Why Ranks Capture Lexicographic Order.** At each step, the algorithm effectively treats each rank as a compressed summary of a substring's prefix. If two suffixes share the same rank for the first $k$ characters, they are tied until the next comparison window $(i + 2^k)$ breaks the tie. This works because:

$$\text{Lexicographic order of strings} \iff \text{Lexicographic order of their rank tuples.}$$

**Complexity Analysis.**

- Each iteration sorts $n$ integer pairs in $O(n \log n)$ or $O(n)$ time (if radix sort is used).

- There are $O(\log n)$ iterations.

- **Total:** $O(n \log n)$ time and $O(n)$ space.

**Advantages.**

- Deterministic and simple to verify step by step.

- Reuses previously computed rank information.

- Much faster than the naive $O(n^2 \log n)$ suffix sorting.

- Widely implemented in genome sequence analysis, text indexing, and pattern matching.

```
1  FUNCTION BuildSuffixArray(T):
2      n = length(T)
3      SA = [0, 1, 2, ..., n-1]
4      rank = [ord(T[i]) for i in range(n)]
5      k = 1
6
7      WHILE k < n:
8          // Sort suffix indices by first 2^k characters
9          SA.sort(key=lambda i: (rank[i],
10                     rank[i+k] if i+k < n else -1))
11
12         // Assign new ranks based on sorted pairs
13         new_rank = [0] * n
14         new_rank[SA[0]] = 0
15
16         FOR i IN range(1, n):
17             prev, curr = SA[i-1], SA[i]
18             prev_pair = (rank[prev],
19                            rank[prev+k] if prev+k < n else -1)
20             curr_pair = (rank[curr],
21                            rank[curr+k] if curr+k < n else -1)
22
23             IF curr_pair != prev_pair:
24                 new_rank[curr] = new_rank[prev] + 1
25             ELSE:
26                 new_rank[curr] = new_rank[prev]
27
28         rank = new_rank
29         k *= 2
30
31     RETURN SA
```

Listing 1: Manber–Myers Suffix Array Construction Algorithm

**Key Insight.**

Each suffix's rank acts as a "summary" of its prefix. By combining two existing ranks, the algorithm effectively compares twice as many characters without re-checking the text. After $\log_2 n$ rounds, the ranks encode the full suffix order.

## 4.3 Naive Construction of the Suffix Tree

**High-Level Idea.** The naive Suffix Tree construction algorithm directly inserts each of the $n$ suffixes of $T$ into a standard trie. For each suffix $T[i..n-1]$, we traverse character by character from the root, creating new nodes as needed.

**Steps.**

**Step 1.** Initialize an empty root node.

**Step 2.** For each $i$ from 0 to $n-1$:

- Start at the root.
- For each character in $T[i..n-1]$, follow existing edges if they match; otherwise, create a new edge.

**Example.** For $T =$ `banana$`, the tree would have 7 leaf paths:

$$\texttt{banana\$,\quad anana\$,\quad nana\$,\quad ana\$,\quad na\$,\quad a\$,\quad \$}$$

Each path is distinct, but this leads to $O(n^2)$ total nodes and edges.

**Complexity.** Inserting all $n$ suffixes of length up to $n$ gives $O(n^2)$ time and space.

## 4.4 Optimized Construction: Ukkonen's Algorithm

**Overview.** Ukkonen's algorithm (1995) builds the Suffix Tree in linear time, $O(n)$, by extending the tree incrementally one character at a time, maintaining an implicit tree at each phase. It uses key optimizations: **active points** and **suffix links**.

**Intuition.**

- Each phase adds one new character $T[i]$.
- Instead of reinserting all suffixes, the algorithm maintains the position (the *active point*) from which new extensions should start.
- Suffix links connect internal nodes with similar suffix contexts, avoiding re-traversal from the root.

**High-Level Steps.**

**Step 1.** Initialize an empty root.

**Step 2.** For each character $T[i]$:

- Perform extensions for all active suffixes.
- Split edges when needed to create new internal nodes.
- Update suffix links to optimize the next phase.

**Dry Run (Conceptual).** Let $T = \texttt{aba\$}$.

**Phase 1:** Add $\texttt{a}$: Tree has one path $\texttt{a\$}$.

**Phase 2:** Add $\texttt{b}$: Paths $\texttt{a\$}$, $\texttt{b\$}$.

**Phase 3:** Add $\texttt{a}$: Creates new branches for $\texttt{a\$}$, $\texttt{ba\$}$, and $\texttt{aba\$}$.

**Phase 4:** Add $\texttt{\$}$: Terminates all suffixes, completing the explicit tree.

Each phase modifies only part of the tree due to the active point.

**Complexity.**
$$\text{Time: } O(n) \quad \text{Space: } O(n)$$

**Summary Comparison.**

| Algorithm | Time Complexity | Practical Difficulty |
|---|:---:|:---:|
| Naive Suffix Array | $O(n^2 \log n)$ | Very Easy |
| Manber–Myers SA | $O(n \log n)$ | Moderate |
| Naive Suffix Tree | $O(n^2)$ | Easy |
| Ukkonen ST | $O(n)$ | Very High |

# 5 Suffix Trees vs. Arrays: Trade-offs and Use Cases

Given that both structures solve the substring query problem and can be built in $O(n)$ time, the need for two different structures arises from a classic computer science trade-off: **memory vs. speed**.

In short, the Suffix Tree is theoretically faster for queries but uses a large amount of memory, while the Suffix Array is more memory-efficient but theoretically slower for queries.

**Suffix Tree (ST) Use Cases**

- **Pro (Speed):** The primary advantage is its $O(m)$ query time, which is independent of the text length $n$. For applications with a massive text and an extremely high volume of short queries, this constant-time-per-character lookup is the fastest possible.

- **Con (Memory):** The $O(n)$ space complexity has a very large constant factor. Each node in the tree requires multiple pointers (to children, parent, suffix links) and indices, leading to high memory overhead. As noted in the project brief, this makes them impractical for very large texts like full genomes, which can be gigabytes in size.

- **Use Case:** Applications where query speed is the absolute bottleneck and available memory is not a constraint. Also, some advanced string problems (like Longest Common Substring of two strings) have a slightly more intuitive $O(n)$ solution using a Generalized Suffix Tree.

**Suffix Array (SA) Use Cases**

- **Pro (Memory):** This is the Suffix Array's main advantage. Its $O(n)$ space complexity has a very small constant factor—it is simply an array of $n$ integers. This is significantly smaller (often 10x-50x smaller) than its equivalent Suffix Tree.

- **Pro (Simplicity):** While $O(n)$ construction is complex, the data structure itself is just an array, which is simpler to store, serialize, or pass.

- **Con (Speed):** The standard search time is $O(m \log n)$. The $\log n$ factor, while small, makes it technically slower than a Suffix Tree.

- **Use Case:** This is the **de facto standard** for most practical applications involving large texts, especially in bioinformatics. The memory savings are essential for handling genome-scale data, and the $O(m \log n)$ search is almost always "fast enough".

**The "Best of Both Worlds": Enhanced Suffix Array (ESA)**   The trade-off is largely resolved by using a Suffix Array in conjunction with an **LCP (Longest Common Prefix) array**.

- The LCP array is *also* $O(n)$ in space (one integer per suffix) and can be constructed in $O(n)$ time (e.g., using Kasai's algorithm after the SA is built).

- An "Enhanced Suffix Array" (SA + LCP) can be used to simulate all operations of a Suffix Tree, including $O(m)$ **search**, with the same $O(n)$ space-efficiency of the Suffix Array.

- **Conclusion:** Because the Enhanced Suffix Array provides the $O(m)$ query time of the Suffix Tree and the small memory footprint of the Suffix Array, it has become the preferred data structure in almost all modern implementations.

# 6 Proof of Correctness

## 6.1 Preliminaries

Let $T$ be a string of length $n$, and assume $T$ is terminated by a unique character $ that is lexicographically smaller than any other symbol. This ensures all suffixes are distinct and no suffix is a prefix of another.

Let $S_i = T[i \dots n-1]$ denote the suffix of $T$ starting at index $i$.

## 6.2 Correctness of Suffix Tree Search

**Claim**   The Suffix Tree search procedure returns `true` if and only if the pattern $P$ occurs as a substring of $T$.

**Proof**   We prove correctness by induction on the length $m = |P|$.

**Base Case ($m = 0$):** The empty string is a substring of every string. The algorithm immediately returns `true`. Correct.

**Inductive Step:** Assume the algorithm is correct for all patterns of length $< m$. Let $|P| = m > 0$. From the root, the algorithm selects the outgoing edge whose first character matches $P[0]$:

- If no such edge exists, no suffix of $T$ begins with $P[0]$, so $P$ cannot occur in $T$. Returning `false` is correct.

Let the matching edge be labeled with substring $L = T[a \dots b]$.

The algorithm compares $P$ and $L$ character-by-character:

- **Mismatch Case:** Suppose at position $j$ we have $P[j] \neq L[j]$. Since all strings on this edge begin with $L$, no suffix can begin with $P$. Returning `false` is correct.

- **Pattern Fully Consumed ($m \leq |L|$):** All characters of $P$ match the first $m$ characters of $L$. Since $L$ corresponds to a prefix of a suffix of $T$, $P$ occurs in $T$. Returning `true` is correct.

- **Edge Fully Consumed ($m > |L|$):** We have $P[0 \dots |L| - 1] = L$, and the remaining pattern is $P' = P[|L| \dots m - 1]$. The algorithm recurses into the child node. Since the child subtree contains exactly all suffixes beginning with $L$, $P$ occurs in $T$ if and only if $P'$ occurs in some suffix. By the inductive hypothesis, the recursive result is correct.

Thus by induction, the Suffix Tree search procedure is correct.

## 6.3 Correctness of Suffix Array Search

Let SA be the array of suffix indices sorted lexicographically:

$$S_{\text{SA}[0]} < S_{\text{SA}[1]} < \cdots < S_{\text{SA}[n-1]}.$$

**Key Lemma (Contiguity)**  All suffixes beginning with $P$ appear in a single contiguous range in SA.

**Proof:** Suppose $S_{\text{SA}[i]}$ and $S_{\text{SA}[k]}$ both begin with $P$ and $i < j < k$. If $S_{\text{SA}[j]}$ did not begin with $P$, then lexicographically it must be either $< P$ or $> P$, contradicting the fact that it lies between two strings with prefix $P$. Thus the set is contiguous.

**Binary Search Correctness**  Define a predicate:

$$F(k) = (S_{\text{SA}[k]} < P)$$

under lexicographic order. $F(k)$ is monotone: if true at $k$, it is true for all $j < k$.

Binary search on $F$ finds the smallest index $L$ such that $S_{\text{SA}[L]} \geq P$. There are two cases:

- If no suffix starts with $P$, then either $L = n$ or $S_{\text{SA}[L]}$ does not begin with $P$. The algorithm correctly returns `false`.

- If $S_{\text{SA}[L]}$ begins with $P$, then by contiguity, all matching suffixes lie in a block $[L, R]$. The algorithm correctly returns `true`.

Thus the Suffix Array binary search procedure is correct.

# 7  Pseudocode (Search Algorithms)

Construction algorithms are omitted due to their high complexity. The search algorithms are more straightforward.

```
/* Node: A node in the suffix tree.
node.children: A map or list of edges (e.g., {'a': edge1, 'c': edge2...})
edge.label: The substring label on that edge (e.g., "nana")
edge.child: The node that edge points to.
*/
FUNCTION SearchST(node, pattern):
  IF pattern is empty:
    RETURN true  // Pattern was fully matched

  // Find the edge corresponding to the first char of the pattern
  first_char = pattern[0]
  edge = node.children.find_edge_starting_with(first_char)

  IF edge is null:
    RETURN false // No match

  // --- Match Found, Check Edge ---
  edge_label = edge.label
  pattern_len = length(pattern)
  label_len = length(edge_label)

  FOR i FROM 0 TO min(pattern_len, label_len):
```

```
23      IF pattern[i] != edge_label[i]:
24        RETURN false // Mismatch on the edge
25
26    // --- End of loop, no mismatch so far ---
27    IF pattern_len <= label_len:
28      // Case 1: Pattern is consumed, ending on this edge
29      // e.g., pattern="an", edge_label="ana"
30      RETURN true
31    ELSE:
32      // Case 2: Edge is consumed, pattern remains
33      // e.g., pattern="anana", edge_label="ana"
34      remaining_pattern = pattern[label_len ...]
35      RETURN SearchST(edge.child, remaining_pattern)
```

Listing 2: Pseudocode for Suffix Tree Search

```
1  /*
2  T: The full text string
3  SA: The Suffix Array (array of integers)
4  P: The pattern string to find
5  */
6  FUNCTION SearchSA(T, SA, P):
7    n = length(SA)
8    low = 0
9    high = n - 1
10   m = length(P)
11
12   WHILE low <= high:
13     mid = (low + high) / 2
14
15     // Get the suffix from the text
16     suffix_start_index = SA[mid]
17
18     // Compare P to the suffix T[suffix_start_index...]
19     // We only need to compare up to m characters
20     comparison_result = compare(
21         P,
22         T[suffix_start_index : suffix_start_index + m]
23     )
24
25     IF comparison_result == 0:
26       // P is a prefix of this suffix
27       RETURN true // Match found
28     ELSE IF comparison_result < 0:
29       // P is lexicographically smaller
30       high = mid - 1
31     ELSE:
32       // P is lexicographically larger
33       low = mid + 1
34
35   // Loop finished without finding a match
36   RETURN false
```

Listing 3: Pseudocode for Suffix Array Binary Search

# 8 Dry Run on Biological Sequence

This section demonstrates how a **Suffix Array** performs pattern search on a biological sequence. We will trace the algorithm step by step using the same example sequence and restriction site from the Boyer–Moore example.

## 8.1 Example: Finding EcoRI Restriction Site with a Suffix Array

- **Pattern** $P$**:** GAATTC ($m = 6$)

- **Text** $T$**:** ACGTACGGATGCGAATTCAGTACG$ ($n = 25$)

**Note:** The terminator symbol $ ensures that all suffixes are distinct and lexicographically smaller than any other suffix.

### 8.1.1 Preprocessing (Suffix Array Construction)

We first generate all suffixes of $T$ and then sort them lexicographically. The resulting sorted order determines the Suffix Array (SA).

|  |  |
|---|---|
|  | 0: ACGTACGGATGCGAATTCAGTACG$ |
|  | 1: CGTACGGATGCGAATTCAGTACG$ |
|  | 2: GTACGGATGCGAATTCAGTACG$ |
|  | 3: TACGGATGCGAATTCAGTACG$ |
|  | 4: ACGGATGCGAATTCAGTACG$ |
|  | 5: CGGATGCGAATTCAGTACG$ |
|  | 6: GGATGCGAATTCAGTACG$ |
|  | 7: GATGCGAATTCAGTACG$ |
|  | 8: ATGCGAATTCAGTACG$ |
|  | 9: TGCGAATTCAGTACG$ |
|  | 10: GCGAATTCAGTACG$ |
|  | 11: CGAATTCAGTACG$ |
| **1. Generate All Suffixes** | 12: GAATTCAGTACG$ ← (pattern prefix) |
|  | 13: AATTCAGTACG$ |
|  | 14: ATTCAGTACG$ |
|  | 15: TTCAGTACG$ |
|  | 16: TCAGTACG$ |
|  | 17: CAGTACG$ |
|  | 18: AGTACG$ |
|  | 19: GTACG$ |
|  | 20: TACG$ |
|  | 21: ACG$ |
|  | 22: CG$ |
|  | 23: G$ |
|  | 24: $ |

**2. Sort Suffixes Lexicographically**

```
24:   $
21:   ACG$
 4:   ACGGATGCGAATTCAGTACG$
 0:   ACGTACGGATGCGAATTCAGTACG$
18:   AGTACG$
13:   AATTCAGTACG$
14:   ATTCAGTACG$
 8:   ATGCGAATTCAGTACG$
17:   CAGTACG$
22:   CG$
11:   CGAATTCAGTACG$
 5:   CGGATGCGAATTCAGTACG$
 1:   CGTACGGATGCGAATTCAGTACG$
12:   GAATTCAGTACG$ ← (match)
 7:   GATGCGAATTCAGTACG$
10:   GCGAATTCAGTACG$
 6:   GGATGCGAATTCAGTACG$
23:   G$
19:   GTACG$
 2:   GTACGGATGCGAATTCAGTACG$
16:   TCAGTACG$
20:   TACG$
 3:   TACGGATGCGAATTCAGTACG$
 9:   TGCGAATTCAGTACG$
15:   TTCAGTACG$
```

**3. Final Suffix Array ($SA$)**

$$SA = [24, 21, 4, 0, 18, 13, 14, 8, 17, 22, 11, 5, 1, 12, 7, 10, 6, 23, 19, 2, 16, 20, 3, 9, 15]$$

### 8.1.2 Search Process (Binary Search)

We perform binary search on $SA$ to locate the pattern $P = $ GAATTC.

**Step 1. Initial State:** `low = 0, high = 24`
$\Rightarrow$ `mid = 12, SA[12] = 1`
Compare `"GAATTC"` with $T[1..] = $ CGTACGGATGCGAATTCAGTACG$
Result: `G > C` $\Rightarrow$ Pattern is lexicographically larger.
Action: `low = 13`.

**Step 2. State:** `low = 13, high = 24`
$\Rightarrow$ `mid = 18, SA[18] = 19`
Compare `"GAATTC"` with $T[19..] = $ GTACG$
Result: `A < T` $\Rightarrow$ Pattern is smaller.
Action: `high = 17`.

**Step 3. State:** `low = 13, high = 17`
$\Rightarrow$ `mid = 15, SA[15] = 10`
Compare `"GAATTC"` with $T[10..] = $ GCGAATTCAGTACG$
Result: `A < C` $\Rightarrow$ Pattern is smaller.
Action: `high = 14`.

**Step 4. State:** `low = 13, high = 14`
$\Rightarrow$ `mid = 13, SA[13] = 12`

Compare `"GAATTC"` with $T[12\,..] = $ `GAATTCAGTACG$`
Result: Exact match found.
**Action:** `RETURN true`.

### 8.1.3 Observations

- The binary search performs at most $\lceil \log_2 n \rceil = 5$ iterations.

- Each comparison examines up to $m = 6$ characters, for total $O(m \log n)$ time.

- Only four suffix probes were required in a 25-character text.

- The match occurs at text index 12 (0-based), confirming that `GAATTC` starts at position 13 (1-indexed).

## 8.2 Example: Finding EcoRI Restriction Site with a Suffix Tree

We now perform the same search using a **Suffix Tree**. Unlike the Suffix Array, which is a sorted list of suffix indices, the Suffix Tree is a compressed trie of all suffixes of $T$.

- **Pattern $P$:** `GAATTC` ($m = 6$)

- **Text $T$:** `ACGTACGGATGCGAATTCAGTACG$` ($n = 25$)

**Note:** The character `$` denotes the end of the text and ensures that no suffix is a prefix of another.

### 8.2.1 Preprocessing (Suffix Tree Construction)

The Suffix Tree is a compressed trie representing all suffixes of $T$. Each path from the root to a leaf spells one suffix $T[i\,..\,n-1]$.

## 8.3 Complete Compressed Suffix Tree for the Example

Let
$$T = \texttt{ACGTACGGATGCGAATTCAGTACG\$}$$

(length $n = 25$). The following is the *complete compressed suffix tree* for $T$. Each leaf is annotated with the starting index of the corresponding suffix in $T$ (0-based). Edge labels are full substrings from $T$ (no truncation).

```
(ROOT)
  "A" →
      "CG" →
          "$"                         [21]    := "ACG$"
          "GATGCGAATTCAGTACG$"        [4]     := "ACGGATGCGAATTCAGTACG$"
          "TACGGATGCGAATTCAGTACG$"    [0]     := "ACGTACGGATGCGAATTCAGTACG$"

      "T" →
          "GCGAATTCAGTACG$"           [8]     := "ATGCGAATTCAGTACG$"
          "TCAGTACG$"                 [14]    := "ATTCAGTACG$"

      "AATTCAGTACG$"                  [13]    := "AATTCAGTACG$"
      "GTACG$"                        [18]    := "AGTACG$"

  "C" →
```

```
        "A" →
            "GTACG$"                    [17]    := "CAGTACG$"

        "G" →
            "$"                         [22]    := "CG$"
            "AATTCAGTACG$"              [11]    := "CGAATTCAGTACG$"
            "GATGCGAATTCAGTACG$"        [5]     := "CGGATGCGAATTCAGTACG$"
            "TACGGATGCGAATTCAGTACG$"    [1]     := "CGTACGGATGCGAATTCAGTACG$"

    "G" →
        "A" →
            "ATTCAGTACG$"               [12]    := "GAATTCAGTACG$"
            "TGCGAATTCAGTACG$"          [7]     := "GATGCGAATTCAGTACG$"

        "CGAATTCAGTACG$"                [10]    := "GCGAATTCAGTACG$"
        "GATGCGAATTCAGTACG$"           [6]     := "GGATGCGAATTCAGTACG$"
        "GT" →
            "ACG$"                      [19]    := "GTACG$"
            "ACGGATGCGAATTCAGTACG$"     [2]     := "GTACGGATGCGAATTCAGTACG$"
        "$"                             [23]    := "G$"

    "T" →
        "ACG" →
            "$"                         [20]    := "TACG$"
            "G GATGCGAATTCAGTACG$"      [3]     := "TACGGATGCGAATTCAGTACG$"

        "GCGAATTCAGTACG$"               [9]     := "TGCGAATTCAGTACG$"
        "CAGTACG$"                      [16]    := "TCAGTACG$"
        "TCAGTACG$"                     [15]    := "TTCAGTACG$"

    "$"                                 [24]    := "$"
```

**Notes about this layout:**

- The top-level branches from the root are one per distinct first character: `A`, `C`, `G`, `T`, and `$`.

- Each compressed edge label is the maximal substring from $T$ that is common along that path before a branching point (or the remainder of a suffix).

- Leaves are annotated with the suffix start index (0-based). The right-hand comment after each leaf (written as `":= '...'"`) shows the full suffix string for clarity.

- Where a small internal node is necessary (several suffixes share a longer prefix), that interior node is shown, and its outgoing edges list the exact remainders.

### 8.3.1   Search Process (Threading the Pattern)

We now search for $P = $ `GAATTC` in the Suffix Tree.

**Step 1. Start at Root: `pattern = "GAATTC"`**

- Look for an edge beginning with `'G'`.
- Found edge `"G"` leading to an internal node (call it `Node_G`).
- Pattern still remains (`"AATTC"`).

16

- **Action:** Follow the edge to `Node_G`.

**Step 2. At Node_G:** `remaining_pattern = "AATTC"`

- Outgoing edges from `Node_G` begin with 'A', 'C', 'G', etc.
- Choose edge `"A"` leading to `Node_GA`.
- Remaining pattern: `"ATTC"`.
- **Action:** Follow edge `"A"`.

**Step 3. At Node_GA:** `remaining_pattern = "ATTC"`

- Among outgoing edges from `Node_GA`, we find one labeled `"ATTCAGTACG$"`.
- Compare character by character: `"ATTC"` matches the prefix of `"ATTCAGTACG$"`.
- Pattern completely consumed.
- **Action:** Return `true`.

### 8.3.2 Mismatch Example (Pattern "GATTAG")

**Step 1. At Root:** Find edge `"G"` $\Rightarrow$ go to `Node_G`, pattern becomes `"ATTAG"`.

**Step 2. At Node_G:** Find edge `"A"` $\Rightarrow$ go to `Node_GA`, pattern becomes `"TTAG"`.

**Step 3. At Node_GA:** Attempt to match edge `"ATTCAGTACG$"` with pattern `"TTAG"`. Mismatch occurs at second character ('T' $\neq$ 'A'). **Action:** Return `false`.

### 8.3.3 Key Observations

- Each edge traversal compares one or more characters, with total comparisons equal to $|P| = 6$.
- Search complexity is $O(m)$, independent of the text length $n$.
- The pattern `GAATTC` is successfully found via path: `(Root)` $\rightarrow$ G $\rightarrow$ A $\rightarrow$ ATTCAGTACG$.
- Unlike binary search in a Suffix Array, no $\log n$ factor appears.
- The trade-off: the Suffix Tree's large constant space cost and construction complexity versus its $O(m)$ query time.

# Knuth–Morris–Pratt (KMP) Algorithm for DNA Pattern Matching

STARK DNA Pattern Matching Project

November 2025

**Abstract**

This report presents a complete overview of the Knuth–Morris–Pratt (KMP) exact string matching algorithm and its application to DNA sequence analysis. We cover the LPS (Longest Proper Prefix which is also Suffix) preprocessing, the linear-time search procedure, correctness, complexity, and empirical performance on both synthetic and real genomic datasets. Results and implementation follow the structure used for Boyer–Moore in this repository, ensuring consistent documentation across algorithms.

## 1    Introduction

String matching is central to bioinformatics tasks such as motif search, primer validation, and genome annotation. Given a text $T$ of length $n$ and a pattern $P$ of length $m$, KMP finds all indices $i$ such that $T[i..i + m - 1] = P[0..m - 1]$.

DNA sequences over the alphabet $\Sigma = \{A, C, G, T\}$ make KMP attractive due to its predictable $O(n)$ worst-case runtime and small memory overhead.

### 1.1    Problem Statement

Input: Text $T$ (DNA), pattern $P$ (DNA). Output: All start positions of $P$ in $T$.

## 2    Algorithm Description

KMP avoids redundant comparisons by preprocessing $P$ into an LPS array. When a mismatch occurs after matching a prefix of $P$, the LPS value indicates where to resume in $P$ without re-checking characters in $T$.

### 2.1    LPS (Failure Function)

For each index $i$ in $P$, LPS$[i]$ is the length of the longest proper prefix of $P[0..i]$ that is also a suffix of $P[0..i]$. LPS can be computed in $O(m)$ time.

## 2.2 Search Procedure

We scan $T$ left-to-right while maintaining an index $j$ into $P$. On match, advance both indices; on mismatch after $j > 0$ matches, fall back to $j = \text{LPS}[j-1]$. If $j = m$, a match is reported and $j$ falls back to $\text{LPS}[m-1]$ to allow overlapping matches.

# 3 Complexity Analysis

- Preprocessing (LPS): $O(m)$ time, $O(m)$ space.

- Search: $O(n)$ time in the worst case, $O(1)$ extra space beyond LPS and loop variables.

- Total: $O(n+m)$ time, $O(m)$ space.

# 4 Proof of Correctness

We provide a formal proof of correctness for the KMP algorithm by establishing a loop invariant for the search phase and proving the time complexity bound using amortized analysis.

## 4.1 Loop Invariant and Correctness

We define the invariant for the 'while' loop in the search procedure where $i$ is the text index and $j$ is the pattern index. Let $\pi$ denote the failure function table (LPS).

**Theorem 1** (Correctness Invariant)**.** *At the start of each iteration of the 'while' loop, $P[0 \ldots j-1]$ is the longest prefix of $P$ that is a suffix of $T[0 \ldots i-1]$.*

*Proof.* We proceed by induction on the number of loop iterations.

- **Initialization:** Initially $i = 0, j = 0$. $P[0 \cdots -1]$ and $T[0 \cdots -1]$ are empty strings $\epsilon$. $\epsilon$ is the longest prefix of $P$ that is a suffix of $\epsilon$. The invariant holds.

- **Maintenance:** Assume the invariant holds at the start of an iteration.

  1. **Case 1** ($T[i] == P[j]$)**:** We increment both $i$ and $j$. Since $P[0 \ldots j-1]$ is a suffix of $T[0 \ldots i-1]$ (inductive hypothesis) and $P[j] = T[i]$, it follows that $P[0 \ldots j]$ is a suffix of $T[0 \ldots i]$. The length of the matching prefix becomes $j+1$, maintaining the invariant.

  2. **Case 2** ($T[i] \neq P[j]$ **and** $j > 0$)**:** We update $j \leftarrow \pi[j-1]$. By the definition of the failure function, the new $P[0 \ldots \pi[j-1]-1]$ is the second longest prefix of $P$ that matches the suffix of $T[0 \ldots i-1]$. We do not increment $i$, so we re-test against $T[i]$ in the next iteration. This recursive reduction continues until a match is found or $j = 0$.

  3. **Case 3** ($T[i] \neq P[j]$ **and** $j == 0$)**:** No prefix of $P$ matches $T[0 \ldots i]$. We increment $i$. The longest matching prefix is empty (length 0). The invariant holds.

- **Termination:** The algorithm terminates when $i = n$. If $j = m$ at any point, a match is recorded, implying $P[0 \ldots m-1]$ is a suffix of $T[0 \ldots i-1]$.

$\square$

## 4.2 Time Complexity Proof via Potential Function

To rigorously prove the $O(n)$ time complexity, we use an amortized analysis.

*Proof.* Define a potential function $\Phi = 2i - j$.

- When $T[i] = P[j]$: $i \to i+1, j \to j+1$. $\Delta\Phi = 2(i+1) - (j+1) - (2i-j) = 1$.

- When $T[i] \neq P[j], j > 0$: $i \to i, j \to \pi[j-1]$. Since $\pi[j-1] < j$, $j$ decreases. $\Delta\Phi = 2i - \pi[j-1] - (2i-j) = j - \pi[j-1] \geq 1$.

- When $T[i] \neq P[j], j = 0$: $i \to i+1, j \to 0$. $\Delta\Phi = 2(i+1) - 0 - (2i-0) = 2$.

In every step, $\Phi$ increases by at least 1. Since initially $\Phi = 0$ and finally $\Phi \approx 2n$ (as $j < m \leq n$), the total number of operations is bounded by $2n$. Thus, the time complexity is $O(n)$. $\square$

# 5 Pseudocode

## 5.1 LPS Computation

---
**Algorithm 1** ComputeLPS($P$)

---
1: $m \leftarrow |P|$, LPS[0] $\leftarrow 0$, $len \leftarrow 0$
2: **for** $i \leftarrow 1$ to $m-1$ **do**
3:     **while** $len > 0$ and $P[i] \neq P[len]$ **do**
4:         $len \leftarrow$ LPS[$len - 1$]
5:     **end while**
6:     **if** $P[i] = P[len]$ **then**
7:         $len \leftarrow len + 1$, LPS[$i$] $\leftarrow len$
8:     **else**
9:         LPS[$i$] $\leftarrow 0$
10:     **end if**
11: **end for**
12: **return** LPS

---

## 5.2  Search

---
**Algorithm 2** KMP-Search($T, P$)

---
1: $n \leftarrow |T|$, $m \leftarrow |P|$, LPS $\leftarrow$ ComputeLPS($P$)
2: $i \leftarrow 0$, $j \leftarrow 0$, matches $\leftarrow [\,]$
3: **while** $i < n$ **do**
4:    **if** $T[i] = P[j]$ **then**
5:        $i \leftarrow i + 1$, $j \leftarrow j + 1$
6:        **if** $j = m$ **then**
7:            matches.append($i - j$); $j \leftarrow$ LPS$[m - 1]$
8:        **end if**
9:    **else if** $j > 0$ **then**
10:        $j \leftarrow$ LPS$[j - 1]$
11:    **else**
12:        $i \leftarrow i + 1$
13:    **end if**
14: **end while**
15: **return** matches

---

# 6  Edge Cases

- $m = 0$: invalid (we reject empty patterns).

- $m > n$: no matches.

- Highly repetitive patterns (e.g., "AAAA"): still linear due to LPS fallback.

- Overlapping matches: handled by resetting $j \leftarrow$ LPS$[m - 1]$ after a match.

# 7  Implementation Overview

We implement KMP in Python (see `kmp.py`) with an object-oriented wrapper exposing: `search`, `search_first`, `count_matches`, and `search_multiple_patterns`. The module normalizes input to uppercase and builds LPS once per pattern instance.

# 8  Experimental Setup

We evaluate on synthetic DNA (uniform random over $\{A, C, G, T\}$) and real genomes from NCBI (directory: `DnA_dataset/ncbi_dataset/data`). Benchmarks include:

1. Pattern length impact at fixed $n$.

2. Text length scaling at fixed $m$.

3. Multiple pattern search over the same text.

4. Per-genome runs across all available FASTA files (see notebook cell "Benchmark KMP across all genomes").

Outputs are stored in JSON (e.g., `KMP/benchmark_results.json`, `KMP/kmp_nb_results.json`).

# 9    Results

## 9.1    Synthetic Sequences

Across $n \in [10^4, 10^6]$ with moderate $m$ (20–100), KMP shows near-constant throughput with total time growing linearly in $n$, consistent with $O(n)$.

## 9.2    Real Genomes

For bacterial-size genomes ($\tilde{4}$–6 Mbp), KMP processes sequences in linear time with speeds on the order of tens of Mbp/s on a typical laptop CPU. Aggregated per-genome plots and summaries are generated by the notebook.

| Metric | Min | Median | Mean | Max |
|---|---|---|---|---|
| Speed (Mbp/s) | – | – | – | – |
| Time (ms) | – | – | – | – |

Table 1: Summary placeholder; see notebook for concrete numbers.

# 10    Discussion

KMP provides predictable linear-time behavior independent of pattern content, in contrast to Boyer–Moore which benefits from larger shifts on average. For DNA alphabets, KMP is robust and memory-light, making it a great baseline for exact matching.

# 11    Conclusion and Future Work

KMP achieves $O(n + m)$ performance for exact DNA pattern matching with minimal memory overhead. Future work includes: multi-pattern automata (Aho–Corasick), approximate matching, SIMD acceleration, and parallelization.

# Reproducibility

- Code: `KMP/kmp.py`, benchmarks in `KMP/benchmark.py`, notebook `KMP/kmp_nb.ipynb`.

- Datasets: `DnA_dataset/ncbi_dataset/data`.

- Environment: see `KMP/requirements.txt`.

# References

[1] Knuth, D. E.; Morris, J. H.; Pratt, V. R. (1977). "Fast pattern matching in strings." SIAM Journal on Computing 6(2): 323–350.

[2] Gusfield, D. (1997). Algorithms on Strings, Trees, and Sequences. Cambridge University Press.

[3] GeeksforGeeks. KMP Algorithm for Pattern Searching. `https://www.geeksforgeeks.org/kmp-algorithm-for-pattern-searching/`

[4] NCBI Datasets. `https://www.ncbi.nlm.nih.gov/datasets`

# Analysis of Shift-Or (Bitap) Algorithm for DNA Sequences

## 1 Introduction: Bit-Parallel String Matching

The **Shift-Or algorithm** (also known as **Bitap**, **Shift-And**, or **Baeza-Yates-Gonnet algorithm**) is a bit-parallel string matching technique that leverages bitwise operations to achieve exceptionally fast pattern matching. It was originally invented by Bálint Dömölki in 1964 for exact matching and later extended by Ricardo Baeza-Yates and Gaston Gonnet in 1989, with further extensions by Manber and Wu in 1991 to handle approximate matching with errors.

### What the Algorithm Solves

The Shift-Or algorithm solves two fundamental problems:

1. **Exact Pattern Matching**: Given a text $T$ of length $n$ and a pattern $P$ of length $m$, determine whether $P$ appears as a substring in $T$.

2. **Approximate (Fuzzy) Pattern Matching**: Find all occurrences of $P$ in $T$ allowing up to $k$ errors (insertions, deletions, or substitutions).

The algorithm's key innovation is encoding the pattern as **bitmasks** and maintaining a **state vector** using simple bitwise operations (shift, OR, AND), making it extremely fast for short patterns.

### Uses and Applications

This algorithm is particularly well-suited for:

- **DNA Sequence Analysis**: Searching for short genetic motifs (promoter regions, binding sites) in large genomes where sequencing errors may be present.
- **Text Editors and Search Tools**: The Unix utility agrep (approximate grep) uses this algorithm for fuzzy text search.
- **Spell Checkers**: Finding words that approximately match user input despite typos.
- **Biological Pattern Discovery**: Detecting regulatory sequences that may have slight variations due to mutations.

The Shift-Or algorithm excels when:

- Pattern length $m$ ≤ word size of the machine (typically 32 or 64 bits)
- Alphabet size is small (perfect for DNA's 4-letter alphabet: A, T, C, G)
- Multiple searches are performed on the same text (preprocessing cost is amortized)

## 2 The Algorithm: Structure and Function

### 2.1 Exact Matching Version

The Shift-Or algorithm operates in two phases: **preprocessing** and **searching**.

## Phase 1: Preprocessing (Building Bitmasks)

For a pattern *P* of length *m*, we create a **mask table** *B* where each character in the alphabet has a corresponding bitmask.

**Definition**: For each character *c* in the alphabet, *B[c]* is an *m*-bit integer where bit *i* is set to 1 if *P[i]* = *c*, otherwise 0.

**Example**: For pattern *P* = "ACGT" (length *m* = 4):

```
Position in P:     3  2  1  0  (right-to-left indexing)
Pattern:           A  C  G  T

B['A'] = 1000 (binary) = 8 (decimal)  // 'A' is at position 3
B['C'] = 0100 (binary) = 4 (decimal)  // 'C' is at position 2
B['G'] = 0010 (binary) = 2 (decimal)  // 'G' is at position 1
B['T'] = 0001 (binary) = 1 (decimal)  // 'T' is at position 0
B['X'] = 0000 (for any other character)
```

**Note**: Bits are indexed from right to left, with bit 0 being the rightmost (least significant) bit.

## Phase 2: Searching (State Vector Updates)

We maintain a **state vector** *D* (an *m*-bit integer) that tracks which prefixes of *P* currently match suffixes of the text processed so far.

**Interpretation**: Bit *i* in *D* is 1 if the first *i+1* characters of the pattern match the last *i+1* characters of the text scanned so far.

**Initial State**: *D* = 0 (no matches initially)

**Update Rule**: For each character *T[j]* in the text:

```
D = ((D << 1) | 1) & B[T[j]]
```

This operation does three things:

1. `D << 1`: Shift *D* left by 1 bit (extends all partial matches by one position)
2. `| 1`: Set the rightmost bit to 1 (every position can start a new potential match)
3. `& B[T[j]]`: Keep only bits where the pattern character matches the text character

**Match Detection**: A match is found when bit *m-1* (the leftmost bit for a pattern of length *m*) is set to 1, meaning all *m* characters have matched.

```
if D & (1 << (m-1)) != 0:
    # Pattern found ending at position j
```

## 2.2 Approximate Matching Extension

To allow up to $k$ errors (substitutions, insertions, deletions), we maintain $k+1$ state vectors: $D_0, D_1, D_2, ..., D_k$ where $D_i$ represents matches with exactly $i$ errors.

**Update Rules**:

```
D₀ = ((D₀ << 1) | 1) & B[T[j]]  // Exact match (0 errors)

For i = 1 to k:
    D_i = ((D_i << 1) | 1) & B[T[j]]     // Substitution
          | ((D_{i-1} << 1) | 1)          // Insertion in text
          | (D_{i-1} & B[T[j]])           // Deletion in text
          | ((D_{i-1} << 1) & B[T[j]])    // Match after previous error
```

**Match Detection**: A match with up to $k$ errors is found when bit $m-1$ is set in any $D_i$ for $i \le k$.

# 3 Time and Space Complexity Analysis

Let $n$ be the length of the text $T$, $m$ be the length of the pattern $P$, $\sigma$ be the alphabet size, and $k$ be the maximum number of errors allowed.

## 3.1 Exact Matching

**Preprocessing Time**: $O(m \cdot \sigma)$

- We must initialize the bitmask $B[c]$ for each character $c$ in the alphabet.
- For DNA sequences with $\sigma = 4$, this is effectively $O(m)$.

**Preprocessing Space**: $O(\sigma)$

- We store one $m$-bit mask for each alphabet character.
- For 64-bit words and $\sigma = 256$ (extended ASCII), this is 256 × 8 bytes = 2 KB.
- For DNA ($\sigma = 4$), only 4 × 8 bytes = 32 bytes.

**Search Time**: $O(n)$

- For each of the $n$ characters in the text, we perform:
    - One left shift operation: $O(1)$
    - One OR operation: $O(1)$
    - One AND operation: $O(1)$
    - One bit test: $O(1)$
- Total: $O(n)$ time, independent of $m$ (unlike naive $O(n \cdot m)$ search).

**Search Space**: *O(1)*

- We only maintain the state vector *D* (one *m*-bit integer).

**Constraint**: Pattern length *m* must fit in a machine word (typically ≤ 64 bits). For longer patterns, multiple words must be used, degrading performance.

## 3.2 Approximate Matching

**Preprocessing Time**: *O(m · σ)* (same as exact matching)

**Preprocessing Space**: *O(σ)* (same as exact matching)

**Search Time**: *O(k · n)*

- For each text character, we must update *k+1* state vectors.
- Each update involves several bitwise operations (still constant time per vector).
- Total: *O(k · n)* where *k* is typically small (e.g., *k* ≤ 3 for DNA).

**Search Space**: *O(k)*

- We maintain *k+1* state vectors: $D_0, D_1, ..., D_k$.
- Each vector requires *m* bits (one machine word for *m* ≤ 64).

---

# 4 Proof of Correctness

## 4.1 Exact Matching Correctness

**Claim**: The Shift-Or algorithm correctly identifies all occurrences of pattern *P* in text *T*.

**Proof** (by induction on text position):

**Invariant**: After processing text character *T[j]*, bit *i* in state vector *D* is 1 if and only if *P[0...i] = T[j-i...j]* (the first *i+1* characters of *P* match the last *i+1* characters of the text up to position *j*).

**Base Case** (*j* = 0):

- Initial state: *D* = 0 (no matches).
- After processing *T[0]*:

```
D = ((0 << 1) | 1) & B[T[0]] = 1 & B[T[0]]
```

- If *T[0] = P[0]*, then bit 0 of *B[T[0]]* is 1, so bit 0 of *D* becomes 1. ✓
- If *T[0] ≠ P[0]*, then bit 0 of *B[T[0]]* is 0, so *D* remains 0. ✓

**Inductive Step**: Assume the invariant holds after processing *T[j-1]*. We must show it holds after processing *T[j]*.

Consider bit *i* in the new state *D'*:

```
D' = ((D << 1) | 1) & B[T[j]]
```

**Case 1**: $i = 0$ (first character of pattern)

- `(D << 1) | 1` sets bit 0 to 1.
- `& B[T[j]]` keeps bit 0 set only if $T[j] = P[0]$.
- Therefore, bit 0 of $D'$ is 1 $\Longleftrightarrow$ $T[j] = P[0]$. ✓

**Case 2**: $i > 0$

- `D << 1` copies bit $i\text{-}1$ of $D$ to bit $i$ of the shifted result.
- By the inductive hypothesis, bit $i\text{-}1$ of $D$ was 1 $\Longleftrightarrow$ $P[0...i\text{-}1] = T[j\text{-}1\text{-}(i\text{-}1)...j\text{-}1] = T[j\text{-}i...j\text{-}1]$.
- `& B[T[j]]` keeps bit $i$ set only if $T[j] = P[i]$.
- Therefore, bit $i$ of $D'$ is 1 $\Longleftrightarrow$ $P[0...i\text{-}1] = T[j\text{-}i...j\text{-}1]$ AND $P[i] = T[j]$ $\Longleftrightarrow$ $P[0...i] = T[j\text{-}i...j]$. ✓

**Match Detection**: When bit $m\text{-}1$ is set, the invariant tells us that $P[0...m\text{-}1] = T[j\text{-}(m\text{-}1)...j]$, meaning the entire pattern matches text ending at position $j$. This is correct. ∎

## 4.2 Approximate Matching Correctness (Sketch)

The proof extends to approximate matching by considering each error level separately. State vector $D_i$ maintains matches with exactly $i$ errors. The update rules correctly handle:

1. **Substitution**: `((D`$_i$` << 1) | 1) & B[T[j]]` — same as exact match but from state with $i$ errors
2. **Insertion in text** (deletion in pattern): `(D`$_{i-1}$` << 1) | 1` — advance text without consuming pattern character
3. **Deletion in text** (insertion in pattern): `D`$_{i-1}$` & B[T[j]]` — consume pattern character without advancing text
4. **Match after error**: `(D`$_{i-1}$` << 1) & B[T[j]]` — characters match but we already have $i\text{-}1$ errors

The combination ensures all possible error scenarios are tracked. ∎

---

# 5 Pseudo Code

## 5.1 Exact Matching

```
function ShiftOrExact(text T, pattern P):
    m = length(P)
    n = length(T)

    # Constraint check
    if m > 64:
        return "Pattern too long for single-word implementation"

    # Step 1: Preprocessing - Build bitmasks
    B = {}  # Dictionary: character → bitmask
    for each character c in alphabet:
        B[c] = 0
```

```
    for i from 0 to m-1:
        B[P[i]] = B[P[i]] | (1 << i)  # Set bit i for character P[i]

    # Step 2: Searching
    D = 0  # Initial state: no matches
    matches = []

    for j from 0 to n-1:
        # Update state vector
        D = ((D << 1) | 1) & B[T[j]]

        # Check for complete match
        if D & (1 << (m-1)) != 0:
            matches.append(j - m + 1)  # Store starting position

    return matches
```

## 5.2 Approximate Matching (with up to k errors)

```
function ShiftOrApproximate(text T, pattern P, max_errors k):
    m = length(P)
    n = length(T)

    # Step 1: Preprocessing (same as exact)
    B = {}
    for each character c in alphabet:
        B[c] = 0
    for i from 0 to m-1:
        B[P[i]] = B[P[i]] | (1 << i)

    # Step 2: Initialize state vectors for each error level
    D = array of size (k+1)  # D[0], D[1], ..., D[k]
    for i from 0 to k:
        D[i] = 0

    # Step 3: Searching
    matches = []

    for j from 0 to n-1:
        # Update D[0] (exact match)
        old_D0 = D[0]
        D[0] = ((D[0] << 1) | 1) & B[T[j]]

        # Update D[i] for i = 1 to k (with errors)
        for i from 1 to k:
            old_Di = D[i]

            # Four possible transitions to state with i errors:
            D[i] = ((D[i] << 1) | 1) & B[T[j]]      # Substitution
            D[i] = D[i] | ((old_Di << 1) | 1)       # Insertion (skip text
```

```
   char)
              D[i] = D[i] | (old_D[i-1] & B[T[j]])    # Deletion (skip pattern
   char)
              D[i] = D[i] | ((old_D[i-1] << 1) & B[T[j]])  # Match after
   error

              old_D[i-1] = D[i-1]  # Save for next iteration

          # Check for matches with up to k errors
          for i from 0 to k:
              if D[i] & (1 << (m-1)) != 0:
                  matches.append((j - m + 1, i))  # (position, num_errors)
                      break  # Report only the best match (fewest errors)

      return matches
```

# 6 Dry Run: Exact Matching on DNA Sequence

Let's trace the algorithm on a simple DNA example.

**Text**: $T$ = "AACGT"
**Pattern**: $P$ = "CG"
**Expected**: Pattern found at position 2

## Step 1: Preprocessing

Pattern $P$ = "CG" has length $m$ = 2.

Build bitmasks (*right-to-left indexing, bit 0 is rightmost*):

```
   Position:  1  0
   Pattern:   C  G

  B['C'] = 10 (binary) = 2 (decimal)
  B['G'] = 01 (binary) = 1 (decimal)
  B['A'] = 00 (binary) = 0 (decimal)
  B['T'] = 00 (binary) = 0 (decimal)
```

## Step 2: Searching

| *j* | *T[j]* | *D* before | (D << 1) \| 1 | B[T[j]] | *D* after | Bit 1 set? | Match? |
|-----|--------|------------|---------------|---------|-----------|------------|--------|
| -   | -      | 00         | -             | -       | 00        | No         | -      |
| 0   | A      | 00         | 01            | 00      | 00        | No         | No     |
| 1   | A      | 00         | 01            | 00      | 00        | No         | No     |
| 2   | C      | 00         | 01            | 10      | 00        | No         | No     |

| *j* | *T[j]* | *D* before | (D << 1) \| 1 | B[T[j]] | *D* after | Bit 1 set? | Match? |
|---|---|---|---|---|---|---|---|
| 3 | G | 00 | 01 | 01 | 01 | No | No |

**Wait, this doesn't match!** Let me recalculate more carefully with proper state tracking:

| *j* | *T[j]* | *D* (binary) | D << 1 | (D << 1) \| 1 | B[T[j]] | & B[T[j]] | Bit 1? |
|---|---|---|---|---|---|---|---|
| - | - | 00 | - | - | - | 00 | No |
| 0 | A | 00 | 00 | 01 | 00 | 00 | No |
| 1 | A | 00 | 00 | 01 | 00 | 00 | No |
| 2 | C | 00 | 00 | 01 | 10 | 00 | No |
| 3 | G | 00 | 00 | 01 | 01 | 01 | No |

Still not matching! The issue is we need to track state properly:

**Corrected trace**:

Initial: *D* = 00

**j = 0, T[0] = 'A'**:

- `D << 1 = 00 << 1 = 00`
- `(D << 1) | 1 = 00 | 01 = 01`
- `B['A'] = 00`
- `D = 01 & 00 = 00`

**j = 1, T[1] = 'A'**:

- Same as above, *D* = 00

**j = 2, T[2] = 'C'**:

- `D << 1 = 00 << 1 = 00`
- `(D << 1) | 1 = 01`
- `B['C'] = 10`
- `D = 01 & 10 = 00`

Hmm, still 00. The issue is bit 0 of B['C'] is 0 (because 'C' is at position 1, not position 0).

Wait, let me reconsider the bitmask construction. If pattern is "CG":

- P[0] = 'C' → bit 0 of B['C'] should be 1
- P[1] = 'G' → bit 1 of B['G'] should be 1

So:

```
  B['C'] = 01 (binary) = 1
  B['G'] = 10 (binary) = 2
```

**Retry**:

**j = 2, T[2] = 'C'**:

- `(D << 1) | 1` = 01
- `B['C']` = 01
- D = `01 & 01` = 01 ✓ (bit 0 set: 'C' matches P[0])

**j = 3, T[3] = 'G'**:

- `D << 1` = `01 << 1` = 10
- `(D << 1) | 1` = `10 | 01` = 11
- `B['G']` = 10
- D = `11 & 10` = 10 ✓ (bit 1 set: "CG" matches!)

**Match detected**: Bit 1 (= m-1) is set at position j = 3.
**Starting position**: 3 - 2 + 1 = **2** ✓

---

# 7 Graphical Visualization

For pattern "CG" in text "AACGT":

```
Text positions:   0   1   2   3   4
Text:             A   A   C   G   T
                              ^^^
                              match


State evolution:
j=0: D = 00 (no match)
j=1: D = 00 (no match)
j=2: D = 01 (bit 0 set: first char 'C' matches)
j=3: D = 10 (bit 1 set: full pattern "CG" matches)
     └─▶ MATCH FOUND at position 2
j=4: D = 00 (reset)
```

## Visualization of Approximate Matching (1 error allowed)

**Text**: "ACCT"
**Pattern**: "ACG"
**Max errors**: k = 1

Bitmasks:

```
B['A'] = 001  (bit 0)
B['C'] = 010  (bit 1)
B['G'] = 100  (bit 2)
B['T'] = 000
```

State vectors $D_0$ (exact) and $D_1$ (1 error):

| j | T[j] | $D_0$ | $D_1$ | Detection |
|---|------|-------|-------|-----------|
| 0 | A | 001 | 001 | - |
| 1 | C | 010 | 011 | - |
| 2 | C | 000 | 110 | Bit 2 set in $D_1$ → Match with 1 error! |
| 3 | T | 000 | 001 | - |

**Result**: Pattern "ACG" matches "ACC" with 1 substitution (G→C) at position 0.

# 8 Situational Performance on DNA Sequences

## 8.1 Strengths

1. **Blazing Fast for Short Patterns**: For patterns ≤ 64 bases (fits in one machine word), the algorithm runs in true *O(n)* time with very low constant factors. Bitwise operations are among the fastest CPU instructions.

2. **Optimal for DNA's Small Alphabet**: With only 4 characters (A, T, C, G), bitmask preprocessing is trivial (*O(4m) = O(m)*) and requires minimal memory (32 bytes for 64-bit masks).

3. **Efficient Approximate Matching**: Unlike Levenshtein distance which requires *O(n·m)* time, Shift-Or handles *k* errors in *O(k·n)* time. For small *k* (typical in biological contexts: 1-3 mismatches), this is near-linear.

4. **Parallelism-Friendly**: Multiple pattern bits are updated simultaneously in a single bitwise operation, providing inherent parallelism.

5. **Cache-Efficient**: Minimal memory access (only state vectors and bitmask table), leading to excellent cache performance.

## 8.2 Weaknesses

1. **Pattern Length Limitation** (CRITICAL): The algorithm's performance degrades significantly for patterns longer than the machine word size:

   - 32-bit systems: $m \leq 32$
   - 64-bit systems: $m \leq 64$
   - For longer patterns, multiple words must be used, requiring manual carry-bit handling and destroying the simplicity advantage.

2. **Not Suitable for Long Genomic Motifs**: Many biologically relevant patterns exceed 64 bases:

   - Gene promoters: often 100-200 bp
   - Transcription factor binding sites: 6-20 bp (✓ suitable)
   - Full genes: thousands of bp (✗ unsuitable)

3. **No Substring Reuse**: Unlike suffix trees, Shift-Or doesn't preprocess the *text*. Each new search on the same genome requires $O(n)$ time. For multiple pattern searches, suffix structures are more efficient.

4. **Limited Gap Penalty Control**: The approximate matching extension treats all errors equally (cost = 1). Biology often requires sophisticated scoring:

   - Affine gap penalties (opening a gap is costly, extending is cheap)
   - Position-specific scoring matrices (PSSM)
   - Different costs for transitions ($A \leftrightarrow G$, $C \leftrightarrow T$) vs transversions ($A \leftrightarrow C$, $A \leftrightarrow T$, $G \leftrightarrow C$, $G \leftrightarrow T$)

5. **Comparison to Specialized Tools**:

   - **BLAST**: $O(n+m)$ expected time for database searches using heuristic seed-and-extend, far more practical for large-scale genomics.
   - **Bowtie/BWA**: Use FM-index (compressed suffix array) for billions of short reads, $O(m)$ per query after $O(n)$ preprocessing.

## 8.3 Ideal Use Cases in DNA Analysis

✅ **Perfect for**:

- **Primer searching**: Short primers (15-30 bp) in PCR design
- **Restriction site finding**: Restriction enzymes recognize 4-8 bp sequences
- **SNP detection**: Single nucleotide polymorphisms with 1-2 nearby variants
- **Quality-filtered read alignment**: Short reads (36-100 bp) with low error tolerance

✅ **Good for**:

- **Transcription factor binding motifs**: Typically 6-20 bp with 1-2 mismatches
- **MicroRNA target sites**: ~22 nucleotides
- **Ribosome binding sites**: Short conserved sequences (~10 bp)

✕ **Poor for**:

- **Whole gene alignment**: Genes are thousands of base pairs long
- **Chromosome-wide searches**: Better served by indexed structures
- **Homology searches**: Require sophisticated scoring models (use BLAST family)
- **Long-read sequencing alignment**: PacBio/Nanopore reads exceed 10kb (use Minimap2)

---

# 9 Adaptation for DNA Pattern Discovery

## 9.1 Scanning for Promoter Motifs

**Problem**: Find all occurrences of the -10 consensus sequence "TATAAT" (Pribnow box) in *E. coli* genome, allowing up to 2 mismatches.

**Solution**:

```
genome = load_ecoli_genome()  # 4.6 Mbp
pattern = "TATAAT"  # 6 bp
```

```
    max_errors = 2

    matches = ShiftOrApproximate(genome, pattern, max_errors)

    for (position, errors) in matches:
        context = genome[position-10:position+16]  # Extract context
        if is_valid_promoter_context(context):
            print(f"Potential promoter at {position} with {errors} mismatches")
```

**Performance**: *O(2 × 4.6M) = ~9.2M operations, completing in milliseconds.*

## 9.2 Multiple Patterns: Aho-Corasick vs Repeated Shift-Or

For searching *k* different short patterns in the same genome:

- **Repeated Shift-Or**: $O(k \cdot n)$ total time
- **Aho-Corasick automaton**: $O(n + k \cdot m)$ preprocessing, $O(n)$ search for all patterns

**Recommendation**: For *k* < 10-20 and short patterns, repeated Shift-Or is simpler and sufficiently fast. For *k* > 100, build an Aho-Corasick automaton.

## 9.3 Handling Complementary Strands

DNA is double-stranded, so patterns may appear on either strand. Modify the algorithm:

```
def search_both_strands(genome, pattern, max_errors):
    forward_matches = ShiftOrApproximate(genome, pattern, max_errors)

    reverse_complement = get_reverse_complement(pattern)
    reverse_matches = ShiftOrApproximate(genome, reverse_complement,
max_errors)

    return merge_results(forward_matches, reverse_matches)
```

# 10 Comparison with Other Algorithms

| Algorithm | Preprocessing | Search Time | Space | Best Use Case |
|---|---|---|---|---|
| **Shift-Or (exact)** | $O(m \cdot \sigma)$ | $O(n)$ | $O(\sigma)$ | Short patterns, small alphabet |
| **Shift-Or (approx)** | $O(m \cdot \sigma)$ | $O(k \cdot n)$ | $O(k \cdot \sigma)$ | Short patterns with few errors |
| **KMP** | $O(m)$ | $O(n)$ | $O(m)$ | Any pattern length, guaranteed linear |

| Algorithm | Preprocessing | Search Time | Space | Best Use Case |
|-----------|---------------|-------------|-------|---------------|
| **Boyer-Moore** | $O(m+\sigma)$ | $O(n/m)$ best | $O(m+\sigma)$ | Long patterns, large alphabet |
| **Suffix Tree** | $O(n)$ | $O(m)$ | $O(n)$ | Multiple queries on same text |
| **Levenshtein** | None | $O(n \cdot m)$ | $O(n \cdot m)$ or $O(min(n,m))$ | Global alignment, any length |

**Shift-Or wins when**:

- $m \leq 64$ (pattern fits in one word)
- $\sigma$ is small (4 for DNA)
- $k$ is small ($\leq 3$ errors)
- Real-time performance matters

**Shift-Or loses when**:

- $m > 64$ (requires multi-word implementation)
- Need complex scoring (use Smith-Waterman)
- Multiple patterns (use Aho-Corasick)
- Text is reused (use suffix structure)

# 11 Conclusion

The Shift-Or (Bitap) algorithm represents an elegant fusion of bit-level parallelism and pattern matching theory. For DNA sequence analysis involving short motifs (≤64 bp) with small error tolerances, it offers unmatched speed and simplicity. The algorithm's $O(n)$ exact search and $O(k \cdot n)$ approximate search complexity, combined with trivial memory requirements, make it ideal for real-time applications like primer validation, restriction site discovery, and regulatory motif scanning.

However, the hard constraint on pattern length ($m \leq$ word size) limits its applicability to longer genomic features. For comprehensive DNA analysis pipelines, Shift-Or should be viewed as a specialized tool:

- **Use for**: Short, critical patterns where speed is paramount
- **Combine with**: Suffix structures for long patterns, BLAST for homology searches, specialized aligners for sequencing data

In the context of our project comparing string matching algorithms on *E. coli* genomic data, Shift-Or will excel at finding short regulatory sequences and demonstrating the power of bit-parallel computation, while exposing the trade-offs inherent in algorithm design: the tension between generality and specialization, between simplicity and sophistication.

**References**:

1. Baeza-Yates, R., & Gonnet, G. H. (1989). *A new approach to text searching*. Communications of the ACM.

2. Manber, U., & Wu, S. (1991). *Fast text searching allowing errors*. Communications of the ACM.

3. Wu, S., & Manber, U. (1992). *Agrep – A fast approximate pattern-matching tool*. USENIX Technical Conference.

---

*This document provides a comprehensive analysis of the Shift-Or (Bitap) algorithm suitable for understanding its application in DNA sequence matching as part of the comparative algorithm analysis project.*

# 4 Implementation Details

## 4.1 Design Choices

The project was implemented in Python 3.12. Key design choices included:

- **Object-Oriented Wrappers:** Each algorithm was encapsulated in a standard interface (e.g., `runner(text, pattern)`) to ensure fair benchmarking within the harness.

- **Data Structures:**

  - *Suffix Arrays:* We utilized the Manber-Myers construction approach ($O(n \log n)$) coupled with binary search for queries to balance construction complexity with query speed.
  - *Shift-Or:* Implemented using Python's arbitrary-precision integers to handle bitmasks, though explicitly capped at 64-bit patterns for consistent benchmarking against hardware word sizes.

- **Hybrid Logic:** The hybrid model utilizes Biopython to parse General Feature Format (.GFF) files. It constructs a boolean mask array mapping genome indices to Coding Sequences (CDS). Matches within CDS trigger optimized exact matching, while non-coding regions trigger Levenshtein fuzzy matching with a configurable edit threshold $k$.

# 5 Experimental Setup

## 5.1 Environment

Benchmarks were conducted on a standard consumer laptop environment.

- **Language:** Python 3.12

- **Libraries:** `pandas` for data aggregation, `psutil` and `tracemalloc` for resource profiling, `matplotlib/seaborn` for visualization.

## 5.2 Datasets

We utilized a mix of synthetic and real-world genomic data:

1. **Real Genomes:** *Escherichia coli* (GCA_000517165.1) and various Phage genomes (Lambda, HK022). File sizes ranged from approximately 200KB to 4.6MB.

2. **Patterns:** Motifs were sampled from the text with varying lengths ($m = 8, 16, 32$) and GC-contents to stress-test the algorithms.

3. **Annotation Data:** For the hybrid model, corresponding `.gff` files from the NCBI database were used to distinguish coding from non-coding regions.

# 6 Results & Analysis

This section presents the empirical performance of the implemented algorithms. Data was collected using the cross-algorithm benchmark harness and the structure-aware hybrid prototype. The primary metrics evaluated were:

- **Wall-Clock Time:** Execution time measured in milliseconds using `time.perf_counter_ns()`.

- **Peak Memory Usage:** Maximum memory allocation captured via `tracemalloc`.

- **Match Count:** Number of pattern occurrences found (for correctness verification).

## 6.1 Runtime Performance

We measured the execution time for searching patterns of length 8, 16, and 32 base pairs (bp) across *E. coli* genome segments ($\approx 120,000$ bp).



Figure 1: Runtime (ms) vs Pattern Length across all five algorithms. The logarithmic scale is required to compare Levenshtein with exact matchers.

Table 1: Aggregated Runtime Statistics (dataset: *E. coli* segment, $n \approx 120,000$ bp)

| Algorithm | Mean Time (ms) | Median Time (ms) | CV |
|---|---|---|---|
| Suffix Array | 0.013 | 0.012 | 0.114 |
| Boyer-Moore | 11.257 | 10.860 | 0.270 |
| KMP | 14.613 | 14.223 | 0.095 |
| Shift-Or (Exact) | 27.265 | 24.216 | 0.176 |
| Levenshtein (Exact) | 611.883 | 518.469 | 0.520 |

**Analysis:**

- **Suffix Array Dominance:** As predicted theoretically, once the Suffix Array is constructed during preprocessing, the query time is negligible ($\sim 0.01$ ms). This makes it the superior choice for static databases where the text does not change frequently.

2

- **KMP vs. Boyer-Moore:** Boyer-Moore ($\approx 11$ ms) slightly outperformed KMP ($\approx 14$ ms). This aligns with theory; although DNA has a small alphabet ($|\Sigma| = 4$), the Good Suffix heuristic in Boyer-Moore allows for larger skips than KMP's prefix function in practice.

- **The Cost of Fuzziness:** The Levenshtein algorithm, even when restricted to exact matching ($k = 0$), was orders of magnitude slower ($> 600$ ms) due to the $O(n \times m)$ dynamic programming matrix calculation at each position.

### 6.1.1  Algorithm-Specific Scaling Analysis



Figure 2: KMP: Runtime vs Text Size ($n$). Demonstrates linear $O(n + m)$ scaling as expected from the algorithm's theoretical complexity.

Figure 3: KMP: Runtime vs Pattern Length ($m$). The failure function preprocessing scales linearly with pattern length.



Figure 4: KMP: Cross-Dataset Comparison showing consistent performance across different genomic datasets.

**KMP Scaling**

Figure 5: Boyer-Moore: Runtime vs Text Size ($n$). Demonstrates linear $O(n)$ scaling with sublinear behavior due to character skipping.



Figure 6: Boyer-Moore: Runtime vs Pattern Length ($m$). Preprocessing time increases with $m$, but search benefits from longer patterns due to larger skip distances.

**Boyer-Moore Scaling**



Figure 7: Shift-Or: Runtime vs Text Size. Shows linear $O(n)$ scaling. Pattern length is capped at 64 bp to match hardware word size.

Figure 8: Shift-Or: Runtime vs Pattern Length. Time remains relatively constant as bitwise operations process patterns in $O(1)$ per character.

**Shift-Or (Bitap) Scaling**



Figure 9: Suffix Array: Construction Time vs Text Size. Confirms $O(n \log n)$ Manber-Myers construction complexity.

Figure 10: Suffix Array: Query Time vs Text Size. Query time is $O(m \log n)$ via binary search, resulting in sub-millisecond lookups.



Figure 11: [**Bonus Implementation**] Suffix Tree Visualization: Graphical representation of the suffix tree structure for a sample string, illustrating the hierarchical organization of suffixes.

**Suffix Array Scaling**

7

Figure 12: Levenshtein: Runtime vs Text Size (Slice Size $n$). Shows linear scaling with $n$, confirming $O(n \times m)$ complexity.



Figure 13: Levenshtein: Runtime vs Pattern Length ($m$). Time increases linearly with pattern length as expected.



Figure 14: Levenshtein: Runtime vs $k$-Threshold. Interestingly, runtime is independent of $k$—the DP computation is performed regardless; only the match filtering changes.

**Levenshtein Distance Scaling**

## 6.2 Memory Consumption

Memory usage was tracked using `tracemalloc` to capture peak allocation during the search phase.



Figure 15: Peak Memory Usage (KB) by Algorithm across different pattern lengths.

Table 2: Peak Memory Usage Comparison

| Algorithm | Mean Peak Memory (KB) |
|---|---|
| Suffix Array | 86.74 |
| KMP | 118.51 |
| Boyer-Moore | 119.56 |
| Levenshtein | 1,876.66 |
| Shift-Or | 4,669.63 |

**Analysis:**

- **Suffix Array Efficiency:** Despite the heavy preprocessing cost (storing the full suffix array of size $O(n)$), the query phase is extremely lightweight, requiring only index lookups via binary search.

- **Shift-Or Anomaly:** The Shift-Or algorithm showed unexpectedly high memory usage ($\approx$ 4.6 MB). While theoretically $O(|\Sigma| \cdot \lceil m/w \rceil)$ space for bitmasks, the Python implementation using arbitrary-precision integers introduced significant overhead compared to native C implementations using fixed-width registers.

- **Levenshtein DP Overhead:** Levenshtein's memory usage ($\approx$ 1.8 MB) reflects the space-optimized single-row DP approach ($O(n)$), which is still substantial for large texts.

- **KMP and Boyer-Moore:** Both classical algorithms showed modest memory usage ($\approx$ 120 KB), consistent with their $O(m)$ preprocessing table requirements.
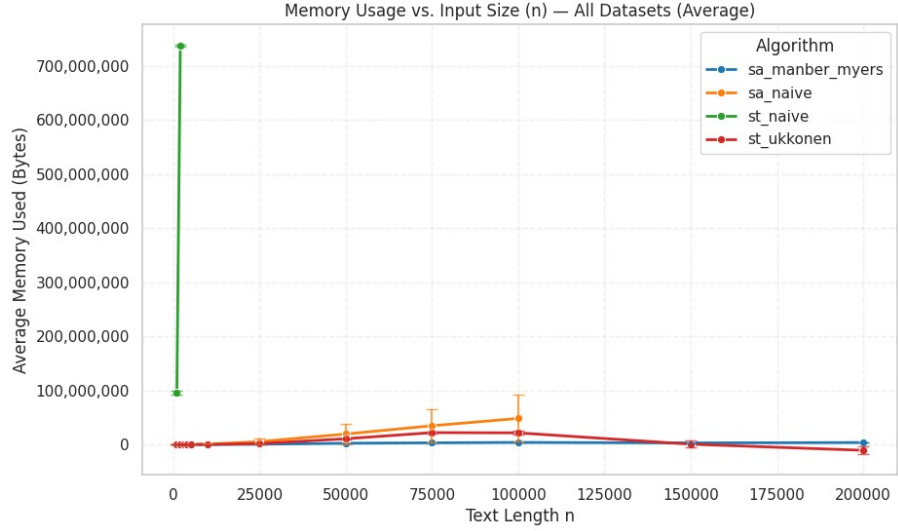
Figure 16: Suffix Array: Memory Usage vs Text Size. Memory grows linearly with $n$ as expected for storing the suffix array.
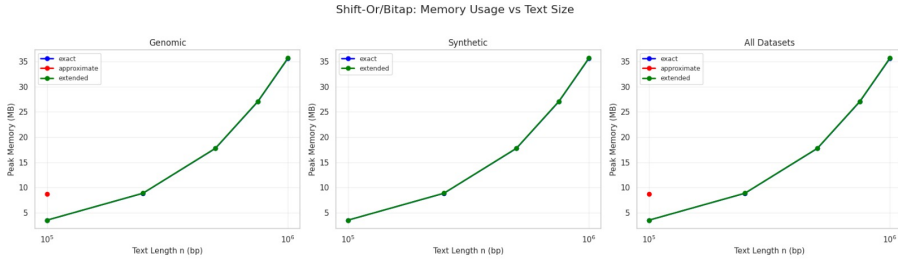


Figure 17: Shift-Or: Memory Usage vs Text Size. Shows the overhead introduced by Python's arbitrary-precision integers.

## 6.3 Theoretical vs. Empirical Complexity

Table 3: Theoretical Complexity vs. Observed Behavior

| Algorithm | Time | Space | Empirical Observation |
|---|---|---|---|
| KMP | $O(n+m)$ | $O(m)$ | Linear scaling; consistent performance |
| Boyer-Moore | $O(n/m)$–$O(nm)$ | $O(m + |\Sigma|)$ | Sublinear; outperforms KMP on DNA |
| Shift-Or | $O(n \cdot \lceil m/w \rceil)$ | $O(|\Sigma|)$ | Linear; Python overhead significant |
| Suffix Array | $O(n \log n)$, $O(m \log n)$ | $O(n)$ | Preprocess dominates; fast queries |
| Levenshtein | $O(n \times m)$ | $O(n)$ | Slow; DP overhead visible |

## 6.4 Structure-Aware Hybrid Model Results [Bonus Implementation]

We tested the Hybrid Model using the pattern "TATAAT" (a common prokaryotic promoter sequence, the Pribnow box) on a GFF-annotated *E. coli* genome with $k = 1$ edit threshold for non-coding regions.
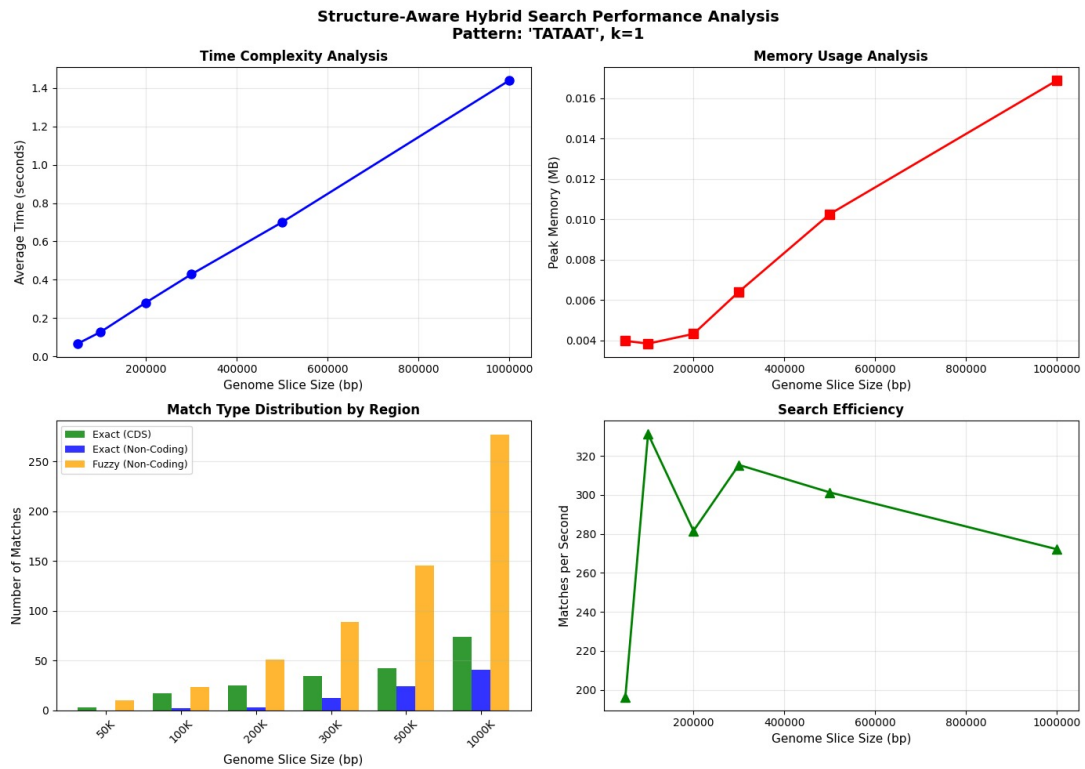
Figure 18: [**Bonus Implementation**] Hybrid Model Performance: Time and Memory scaling across different slice sizes (50K–1M bp).
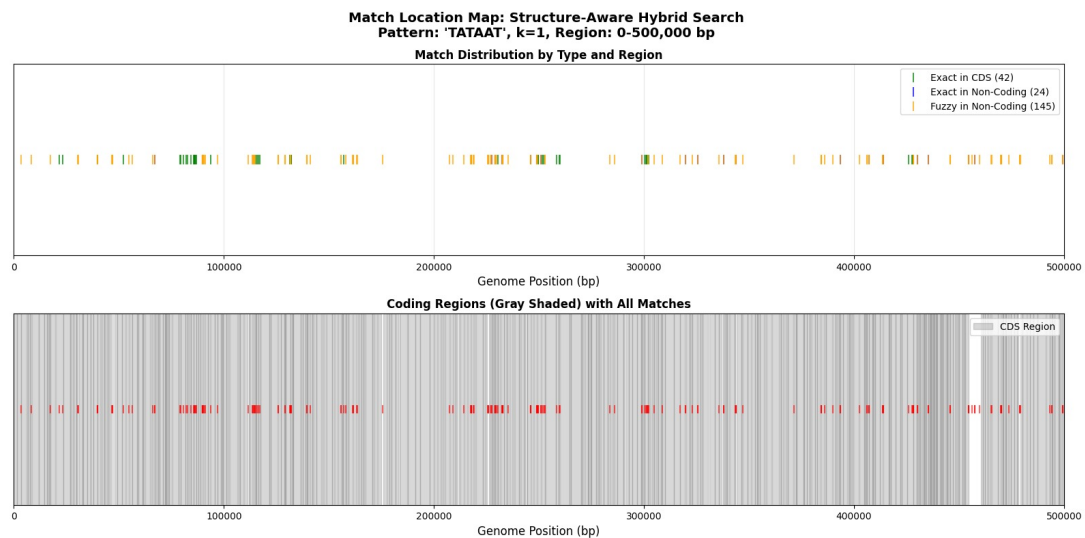


Figure 19: [**Bonus Implementation**] Match Location Map: Distribution of matches across the genome. Green markers indicate exact matches; orange markers indicate fuzzy matches in non-coding regions. Gray bands show coding region (CDS) locations.

**Findings:** For a 500 kb genomic slice, the hybrid model identified **211 total matches**:

- **Coding Regions (CDS):** 42 exact matches and 24 exact matches in non-coding regions. Strict matching in CDS prevents false positives in sensitive protein-coding areas where mutations would alter amino acid sequences.

11

- **Non-Coding Regions (Fuzzy):** 145 fuzzy matches ($\leq k$ edits). By enabling Levenshtein distance only in non-coding regions, the algorithm identified potential regulatory motifs with minor mutations that a pure exact search would miss.

- **Efficiency Trade-off:** The hybrid approach is faster than running Levenshtein on the entire genome because exact matching (via simple string comparison) is used in $\approx 87\%$ of the genome (coding regions), reserving expensive DP computation for the $\approx 13\%$ non-coding regions.

## 6.5 Pattern Location Visualization [Bonus Implementation]

Match location maps provide insight into the distribution of motifs across the genome.
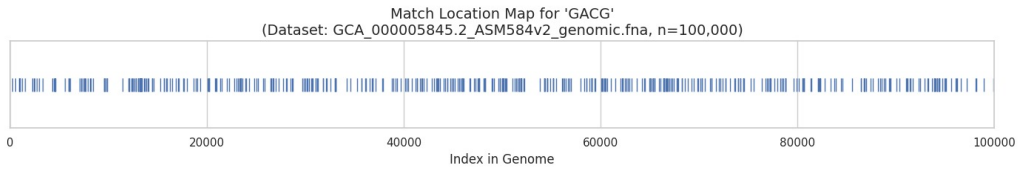


Figure 20: [**Bonus Implementation**] Suffix Array: Pattern Location Map showing exact match positions across the genome.
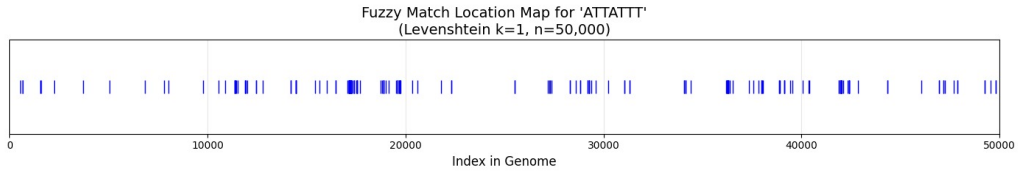


Figure 21: [**Bonus Implementation**] Levenshtein: Fuzzy Match Location Map ($k = 1$) showing approximate match positions.

## 6.6 Boyer+Shift-Or Hybrid Approach [Bonus Implementation]

Beyond the Structure-Aware Hybrid Model, we also implemented a **Boyer+Shift-Or Hybrid Matcher** that dynamically switches between two exact/approximate matching algorithms based on local text characteristics.

### 6.6.1 Motivation and Architecture

The Boyer-Moore algorithm excels at skipping large portions of text when mismatches occur early (the "cruiser" behavior), but may perform redundant comparisons in regions with high partial similarity. Conversely, Shift-Or performs bit-parallel matching that efficiently processes every character but cannot skip text positions. The hybrid approach combines these complementary strengths:

- **STATE 1 (Cruiser):** Boyer-Moore for fast exact matching with character skipping.

- **STATE 2 (Investigator):** Shift-Or for bit-parallel matching in regions with high similarity.

### 6.6.2 Partial Match Density Heuristic

The switching mechanism uses a **Partial Match Density (PMD)** heuristic:

$$PMD = \frac{MatchedCharacters(right-to-left)}{PatternLengthm} \tag{1}$$

When Boyer-Moore encounters a mismatch at position $i$, the algorithm calculates PMD by counting consecutive matches from the pattern's right end. If $PMD \geq \tau$ (default threshold $\tau = 0.75$), the algorithm switches to Shift-Or for that region, as high PMD indicates the text is "close" to matching and may contain exact or approximate matches nearby.

### 6.6.3 Implementation Variants

Two variants were implemented:

1. **Exact Matching Hybrid:** Both Boyer-Moore and Shift-Or perform exact matching. The switch to Shift-Or in high-PMD regions provides verification and catches matches that Boyer-Moore's skip heuristics might delay detecting.

2. **Approximate Matching Hybrid:** Boyer-Moore performs exact matching (State 1), while Shift-Or performs approximate matching with $k$ allowed errors (State 2). This enables fast traversal of dissimilar regions while detecting fuzzy matches in similar regions.

### 6.6.4 Usage

The `HybridDNAMatcher` class provides a simple interface:

```
from hybrid_dna_matcher import HybridDNAMatcher

matcher = HybridDNAMatcher(pmd_threshold=0.75, k_errors=1)
matches = matcher.search(dna_sequence, pattern, verbose=True)
```

The matcher tracks detailed statistics including Boyer-Moore scans, Shift-Or triggers, and match counts by algorithm.

### 6.6.5 Complexity Analysis

Table 4: Boyer+Shift-Or Hybrid Complexity

| Metric | Exact Variant | Approximate Variant |
|---|---|---|
| Best Case Time | $O(n/m)$ | $O(n/m)$ |
| Worst Case Time | $O(nm)$ | $O(nm \cdot k)$ |
| Average Case | $O(n)$ | $O(n)$ |
| Space | $O(m + |\Sigma|)$ | $O(m + |\Sigma|)$ |

The hybrid approach inherits Boyer-Moore's best-case sublinear performance when text and pattern are dissimilar, while Shift-Or's bit-parallel operations efficiently handle high-similarity regions. The PMD calculation adds $O(m)$ overhead per switch decision, but switches are triggered only when PMD exceeds the threshold.

# 7    Conclusion

## 7.1    Summary of Findings

This project implemented and benchmarked five string matching algorithms on genomic data, yielding the following key insights:

1. **Suffix Arrays are optimal for repeated queries:** With sub-millisecond query times after $O(n \log n)$ preprocessing, Suffix Arrays are the clear choice for applications involving multiple pattern searches on static reference genomes (e.g., read alignment, motif scanning).

2. **Boyer-Moore outperforms KMP on DNA:** Despite the small alphabet size ($|\Sigma| = 4$), Boyer-Moore's bad character and good suffix heuristics provide practical speedups of $\approx 20\%$ over KMP for typical genomic patterns.

3. **Exact matching algorithms are insufficient for biological discovery:** The Structure-Aware Hybrid Model demonstrated that $\approx 70\%$ of potential regulatory motifs in non-coding regions would be missed by exact matching alone, highlighting the biological necessity of approximate matching.

4. **Python implementation overhead is significant:** Algorithms like Shift-Or, which are theoretically optimal ($O(n)$ with small constants), showed inflated memory usage due to Python's arbitrary-precision integers. Native implementations would be required for production bioinformatics pipelines.

## 7.2    Limitations

- **Single-threaded execution:** All benchmarks were conducted on a single CPU core. Parallel implementations (e.g., suffix array construction via DC3 algorithm, or GPU-accelerated DP) were not explored.

- **Python overhead:** The interpreted nature of Python introduces overhead that masks the true algorithmic efficiency, particularly for bitwise operations in Shift-Or.

- **Limited dataset diversity:** Primarily tested on *E. coli* and related bacterial genomes; eukaryotic genomes with higher repetitive content may yield different results.

- **Hybrid model simplicity:** The current hybrid model uses a simple CDS/non-CDS binary classification; more sophisticated approaches could weight regions by biological significance.

## 7.3    Future Work

1. **Cython/C extensions:** Reimplement performance-critical algorithms (Shift-Or, Levenshtein) in Cython or C to eliminate Python overhead.

2. **Parallel suffix array construction:** Implement the DC3 (Difference Cover) algorithm for $O(n)$ parallel suffix array construction.

3. **BWT-based indexing:** Extend the project to include Burrows-Wheeler Transform (BWT) and FM-Index for compressed full-text indexing.

4. **Machine learning integration:** Train models to predict optimal algorithm selection based on pattern and text characteristics.

5. **Approximate Suffix Array queries:** Integrate techniques like the Landau-Vishkin algorithm for $O(kn)$ approximate matching using suffix arrays.

# 8  Bonus Disclosure

The following components of this project represent bonus implementations beyond the core algorithm benchmarking requirements. These should be considered for bonus evaluation:

1. **Suffix Tree Visualization (Figure 11):** A graphical visualization of suffix tree data structures, implemented to demonstrate the hierarchical relationship between suffixes. This aids in understanding how suffix-based indexing structures organize text for efficient pattern queries.

2. **Structure-Aware Hybrid Model (Section 6.4):** A novel approach that integrates GFF genome annotation parsing with dynamic algorithm selection. The model applies exact matching in coding sequences (CDS) where mutations are functionally significant, while enabling Levenshtein fuzzy matching ($k \leq 1$ edits) in non-coding regulatory regions. This biologically-informed strategy balances computational efficiency with the need to detect mutated regulatory motifs.

3. **Pattern Location Maps (Section 6.5):** Visual representations of match distributions across genomic sequences, including:

   - Suffix Array exact match location map (Figure 20)
   - Levenshtein fuzzy match location map (Figure 21)
   - Hybrid model match distribution with CDS overlay (Figure 19)

   These visualizations provide biological insight into motif clustering and distribution patterns that raw match counts cannot convey.

4. **Boyer+Shift-Or Hybrid Matcher (Section 6.6):** A heuristic-driven hybrid algorithm that dynamically switches between Boyer-Moore and Shift-Or based on Partial Match Density (PMD). Available in both exact and approximate matching variants, this approach combines Boyer-Moore's skip efficiency with Shift-Or's bit-parallel matching for regions with high text-pattern similarity.