

# String Matching Algorithms for DNA Sequence Search

CS 371 – Algorithm Design and Analysis

Denison University

November 2025

## Abstract

This paper discovers the foundation and computational application of string matching algorithms in DNA sequence search. DNA can be represented as a string composed of four nucleotides (A, T, C, G). The central challenge lies in identifying specific nucleotide patterns within massive genomic datasets. The study focuses on the Knuth–Morris–Pratt (KMP) algorithm and examines how its prefix-based preprocessing makes it suited for searching biological sequences. Furthermore, the paper connects core concepts from the Algorithm Design and Analysis course, including the Longest Common Subsequence (LCS), hashing, and divide-and-conquer strategies, to genome analysis. Finally, approximate string matching algorithms, such as Needleman–Wunsch and Smith–Wate, are introduced, demonstrating how they extend exact matching techniques to handle biological mutation and variability.

## 1 Introduction

Evolution in DNA sequencing technologies over the past two decades has revolutionized modern biology and medicine. The cost of sequencing a human genome has dropped from millions of dollars in 2001 to much cheaper nowadays, leading to an unprecedented expansion in genomic data. This exponential growth of biological information has made computational methods necessary for extracting meaningful insights. Computer science, especially the field of algorithm design, provides both the theoretical and practical foundations for processing, comparing, and analyzing DNA sequences at scale. Among these computational techniques, string matching plays an important role as a cornerstone for bioinformatics applications such as genome assembly, gene identification, and mutation detection. By efficiently locating specific nucleotide patterns within vast genomic datasets, string matching algorithms allow discoveries that would be impossible through manual analysis.

### 1.1 Overview of String Matching in bioinformatics

In bioinformatics, one of the most basic computational problems is determining whether a specific DNA sequence, or pattern, occurs within a larger genomic database, or text. For example, given a pattern  $P = \text{ATGCG}$  and a genome  $T = \text{ATGCCGTAGCTGAC}$ , the goal is to find all positions  $i$  such that  $T[i \dots i + |P| - 1] = P$ . This is a direct instance of the string matching problem.

String matching algorithms provide the foundation for pattern recognition tasks in both computer science and computational biology, allowing researchers to search for gene fragments or repeated motifs within large DNA datasets [1]. Formally, string matching is the process of locating all occurrences of a smaller string  $P$  (the pattern) within a larger string  $T$  (the text), both composed of symbols from a finite alphabet  $\Sigma$ . In the case of DNA,  $\Sigma = \{A, T, C, G\}$ .

An approach compares  $P$  with every possible substring of  $T$ , resulting in a time complexity of  $O(nm)$  for a text of length  $n$  and a pattern of length  $m$ . This method becomes computationally impossible when dealing with genome-scale data, motivating the need for more efficient algorithms [3].

To address this challenge, several classical algorithms were introduced, such as the Knuth–Morris–Pratt (KMP), the Boyer-Moore algorithm, and the Rabin-Karp algorithm. These methods, summarized as the three “fundamental string matching paradigms,” are now standard approaches in computational sequence analysis [2].

These algorithms form the foundation for many applications in bioinformatics, text retrieval, and data mining. Within DNA analysis, they allow the efficient processing and alignment of massive biological sequences,

supporting genome mapping, motif discovery, and evolutionary comparison.

## 1.2 Motivation of String Matching

The motivation for developing efficient string matching algorithms arises from the scale and complexity of modern biological data. The size of genomic databases increases each year, creating a demand for algorithms that balance speed and accuracy [1]. A single human genome contains approximately three billion base pairs, and sequencing technologies now produce terabytes of new data daily. Manual inspection of such data is impossible, making algorithmic solutions essential for pattern discovery and analysis.

Efficient string matching is therefore motivated by the need to process increasingly large datasets with both speed and accuracy. As genomic repositories expand and sequencing becomes more accessible, the performance of these algorithms plays a significant role in enabling timely analysis, supporting downstream biological interpretation, and bridging theoretical algorithm design with practical computational biology.

## 2 Background and Related Work

### 2.1 Historical Context and Early Applications to DNA

String matching has been a main topic in theoretical computer science since the early 70s, when the initial algorithms were developed to search text without redundant comparisons. As sequencing technologies matured and biological data became digitized, these techniques were rapidly adopted in computational biology to locate short motifs, annotate genes, and support basic sequence search in early genomic databases. The transition from manual inspection to automated search was a major turning point in bioinformatics. Efficient pattern-matching algorithms made it possible to process large amounts of DNA data quickly, allowing researchers to keep up with rapidly growing sequence databases. [1].

### 2.2 Overview of Classical Exact String Matching Algorithms

Three classical paradigms, which form the base of modern practice, are:

- **Knuth-Morris-Pratt (KMP):** KMP achieves linear time  $\Theta(n + m)$  by preprocessing the pattern  $P$  into a failure (prefix) function  $\pi$  and using it to skip redundant comparisons during the scan of the text  $T$  [5]. The table  $\pi[i]$  stores the length of the longest proper prefix of the entire pattern  $P$  that is also a suffix of the substring  $P[0\dots i]$ , allowing the algorithm to fall back to a valid alignment after a mismatch without re-examining characters.
- **Boyer-Moore (BM):** BM compares the pattern right-to-left and uses the bad-character and good-suffix heuristics to skip portions of the text, giving sublinear average-case performance and linear worst-case time [5]. For DNA sequence matching, BM remains effective as an exact method, though its efficiency depends on the characteristics of the underlying data [3].
- **Rabin-Karp (RK):** RK uses a rolling hash to compare length- $m$  windows of the text against the pattern, giving expected  $\Theta(n + m)$  time with an appropriate hash function and  $\Theta(nm)$  in the worst case due to collisions [5]. Hash-based techniques also play an important role in speeding up DNA search by reducing the number of direct pattern comparisons [2].

However, in the remainder of this paper, we focus on KMP because of its deterministic linear running time, low auxiliary space  $O(m)$ , and suitability for exact DNA pattern search in pipelines where reproducibility and predictable performance are critical [1].

### 2.3 From Exact to Approximate Matching: Relationship to KMP

Although exact matching is useful for identifying identical subsequences, real biological sequences often contain mutations, insertions, and deletions. This creates the need for approximate matching methods, where similarity is evaluated using edit operations [3]. Dynamic programming provides a principled way to model these differences by breaking the alignment problem into smaller subproblems and reusing previously computed values [5]. Classical algorithms such as Needleman–Wunsch for global alignment and Smith–Waterman for local alignment apply this framework to score matches, mismatches, and gaps across entire sequences or their most similar subregions [4]. Moreover, efficient matching remains essential for large DNA databases, where both exact and approximate methods play complementary roles in analyzing biological variation [2].

### 2.4 Prior Studies Applying KMP to DNA Search

Several studies evaluate or apply KMP in genomic contexts. Exact string-matching algorithms have been compared on DNA datasets, showing that KMP maintains predictable linear-time performance across varying pattern lengths [1]. Other work highlights the importance of efficient string matching in managing large DNA databases, where methods such as KMP remain essential for locating exact subsequences within extensive collections of sequences [2]. Research on sequence comparison also positions exact algorithms, including KMP, as standard tools for identifying identical subsequences for the use of more flexible approximate matching techniques [3]. Collectively, these studies support KMP as a principled and reliable choice for exact DNA pattern search.

## 3 Understanding DNA Sequence Search

### 3.1 Representation and Storage of DNA Data

Basically, DNA is a linear molecule composed of four nucleotides as adenine (A), thymine (T), cytosine (C), and guanine (G), that can be represented as a string over the finite alphabet  $\Sigma = \{A, T, C, G\}$ . From a computational perspective, each nucleotide is encoded as a character, allowing algorithms developed for text processing to be applied to biological data. This abstraction forms the foundation of bioinformatics, where biological sequences are stored, indexed, and processed as strings. Besides, large-scale DNA databases contain vast numbers of nucleotide sequences, and efficient string matching algorithms are essential for searching these repositories. Studies emphasize the need for indexing and filtering techniques to accelerate exact and approximate matching in such datasets [1, 2].

### 3.2 Computational Workflow of DNA Sequence Matching

DNA sequence search can be viewed as a computational process that identifies relevant regions within large collections of nucleotide data. Because modern DNA databases contain extremely large numbers of sequences, efficient algorithms are required to search feasible [2]. A typical workflow begins with preparing the sequences for comparison and organizing them so that matching operations can be performed efficiently [1]. The core stage is the pattern-matching step, where exact algorithms such as KMP or Boyer-Moore locate occurrences of a query sequence, and approximate methods use dynamic programming to compare sequences that differ by mutations or structural variation [3, 4]. After matches are identified, they are interpreted to determine their biological significance, such as assessing similarity or identifying potential variations. Because pattern matching dominates the computational cost of this workflow, improvements in exact and approximate matching algorithms have a direct impact on the overall efficiency of DNA sequence search [2].

### 3.3 Applications of DNA Pattern Search

The ability to perform efficient pattern matching in DNA sequences has numerous applications in genomics and biomedical research. Among the most significant are:

- **Disease Detection:** Identifying genetic mutations or markers associated with hereditary diseases depends on locating specific nucleotide patterns within an individual's genome. For instance, pattern matching helps detect single-nucleotide polymorphisms (SNPs) linked to cancer susceptibility genes.
- **Genome Alignment:** Comparing genomes across individuals or species involves aligning long DNA sequences to identify conserved regions and evolutionary relationships. Exact matching algorithms like KMP are often used in the initial scanning phase to locate identical segments, followed by approximate alignment methods for refinement.
- **Motif and Gene Discovery:** Biological motifs-short recurring patterns with functional significance, such as transcription factor binding sites-can be detected by repeatedly matching specific sequences within promoter or coding regions [1, 2].
- **Metagenomics and Sequence Classification:** In environmental or microbiome studies, pattern matching allows researchers to classify unknown DNA fragments by comparing them to reference databases, enabling species identification from mixed samples.

As genomic datasets continue to expand, these applications increasingly depend on efficient algorithms and indexing strategies. The integration of classical algorithms like KMP with bioinformatics-specific optimizations underscores how theoretical computer science continues to drive progress in biological discovery.

### 3.4 Worked Example

To illustrate the role of exact string matching in DNA analysis, consider a clinician searching for a short disease-associated marker

$$P = \text{ACGTGGA}$$

within a sequenced genomic region

$$T = \text{TTACGTGGATCAGG...}$$

An algorithmic pattern-matching method can quickly determine whether the marker appears in the genomic sequence and where it occurs. This type of computation often serves as the first step in larger workflows, such as identifying potential variants, flagging regions for closer examination, or connecting the match to known biological annotations.

This simple example highlights how automated string matching provides the foundation for downstream genomic analysis without requiring manual inspection of large DNA sequences.

## 4 The KMP Algorithm in the Context of DNA Search

### 4.1 Algorithmic Overview and Prefix Function Preprocessing

The Knuth-Morris-Pratt (KMP) algorithm, introduced in 1977, represents one of the earliest linear-time solutions to the exact string matching problem. Its focus lies in its use of a prefix function, which enables the algorithm to avoid redundant comparisons after mismatches. Rather than reexamining previously matched characters, KMP precomputes information about how the pattern  $P$  overlaps with itself.

Formally, for a pattern  $P$  of length  $m$ , the prefix table  $\pi[i]$  stores the length of the longest proper prefix of  $P[0..i]$  that is also a suffix of that same substring. During the search phase, when a mismatch occurs between  $P[j]$  and  $T[i]$ , the algorithm shifts the pattern by  $\pi[j - 1]$  positions instead of restarting at the beginning, thus ensuring that each character of  $T$  is compared at most once.

The preprocessing step builds the  $\pi$  table in  $O(m)$  time using the following recurrence:

$$\pi[i] = \begin{cases} k + 1, & \text{if } P[i] = P[k] \\ \pi[k - 1], & \text{otherwise, where } k = \pi[i - 1] \end{cases}$$

and the search phase proceeds in  $O(n)$  time, yielding a total runtime  $O(n + m)$  for a text of length  $n$  and a pattern of length  $m$ . The algorithm's deterministic behavior and fixed memory usage  $O(m)$  make it highly predictable qualities that are essential in biomedical computation, where reproducibility and accuracy are critical [5, 1].

**Prefix Function Pseudocode:** The KMP preprocessing phase is summarized as a procedure that computes the prefix values of the pattern by detecting how many characters match between the current suffix and the earliest prefix of the pattern[1]. Their formulation captures the core logic of the classical KMP table: extending the current prefix length when characters match and falling back to shorter prefixes when they do not. A streamlined version consistent with both [1] and the classical KMP approach is shown below:

---

**Algorithm 1** PREFIX-FUNCTION( $P$ )

---

```

1:  $m \leftarrow \text{length}(P)$ 
2:  $\pi[0] \leftarrow 0, k \leftarrow 0$ 
3: for  $i = 1$  to  $m - 1$  do
4:   while  $k > 0$  and  $P[i] \neq P[k]$  do
5:      $k \leftarrow \pi[k - 1]$ 
6:   end while
7:   if  $P[i] = P[k]$  then
8:      $k \leftarrow k + 1$ 
9:   end if
10:   $\pi[i] \leftarrow k$ 
11: end for
12: return  $\pi$ 

```

---

The array  $\pi$  records, for each position  $i$ , how much of the pattern has been matched so far. When a mismatch occurs, the algorithm does not restart from the beginning; instead, it uses the information stored in  $\pi$  to jump to the longest border of the matched prefix. This mechanism is what reduces the complexity from  $O(nm)$  to  $O(n + m)$ , since only previously verified comparisons are reused[1]. Similarly, the prefix array is the core structure that enables deterministic shifts not only for regular strings but also for indeterminate DNA strings[6].

Using the precomputed prefix table, the search scans the text exactly once. At each position, if the current characters match, both indices advance. If they differ, the pattern index falls back according to  $\pi[j - 1]$ , ensuring no repeated comparison:

---

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

---

```
1:  $n \leftarrow \text{length}(T)$ ,  $m \leftarrow \text{length}(P)$ 
2:  $\pi \leftarrow \text{PREFIX-FUNCTION}(P)$ 
3:  $j \leftarrow 0$ 
4: for  $i = 0$  to  $n - 1$  do
5:   while  $j > 0$  and  $T[i] \neq P[j]$  do
6:      $j \leftarrow \pi[j - 1]$ 
7:   end while
8:   if  $T[i] = P[j]$  then
9:      $j \leftarrow j + 1$ 
10:  end if
11:  if  $j = m$  then
12:    report match at position  $i - m + 1$ 
13:     $j \leftarrow \pi[j - 1]$ 
14:  end if
15: end for
```

---

**Correctness:** The correctness of KMP follows from two standard invariants, supported by the description of prefix computation [1] and the formal treatment of shift logic in [6]:

1. **Prefix Function Invariant:** After processing position  $i$ , the value  $\pi[i]$  equals the length of the longest proper prefix  $P$  that is also a suffix of  $P[0\dots i]$ . This is maintained by extending matches  $P[i] = P[k]$  and falling back to  $\pi[k - 1]$  when mismatches occur.
2. **Search Invariant:** At the beginning of each iteration with the text index  $i$ , the variable  $j$  equals the length of the longest prefix of  $P$  a matching suffix of  $T[0\dots i - 1]$ .  
If characters match, the invariant extends; if they do not, the fallback rule ensures that only feasible prefix lengths are retained. When  $j = m$ , a full match has been found and is correctly reported.

Because the fallback always jumps to a prefix consistent with the previously matched part of the text, KMP never misses a valid occurrence and never reports a false one.

**Running Time:** Both papers report the same asymptotic behavior for the regular (non-indeterminate) KMP algorithm. Each character of the text is processed at most once, and every fallback step moves to a strictly smaller prefix. Therefore, the total number of comparisons is linear in the size of the input:

$$O(n + m).$$

It confirms that KMP achieves the best worst-case performance among classical exact DNA-matching algorithms[1]. Moreover, even in the generalized indeterminate-KMP variant, the classical KMP core maintains its linear-time performance whenever both text and pattern are regular strings, which is precisely the setting of exact DNA matching [6].

## 4.2 Applications of KMP in Genomic Pattern Identification

The simplicity and linear-time performance of KMP make it an ideal candidate for detecting exact gene patterns within genomic data. In bioinformatics pipelines, the genome sequence serves as the text  $T$ , often consisting of billions of bases, and short target sequences - genes, promoters, or motifs - serve as patterns  $P$ . KMP efficiently scans  $T$  to locate all occurrences  $P$  without re-reading any portion of the genome.

For example, in motif detection tasks, a pattern such as  $P = \text{TATAAA}$ , the canonical TATA box sequence, can be matched across an entire chromosome to identify promoter regions associated with gene transcription. Similarly, in variant analysis, KMP can detect short, disease-associated DNA segments with complete precision, serving as a preprocessing step before approximate or alignment-based analyses.

Unlike heuristic algorithms such as Boyer–Moore, which depend on alphabet size for efficiency, KMP’s performance is invariant to the size of  $\Sigma = \{A, T, C, G\}$ . This makes it particularly well-suited for genomic data, where the alphabet is small and repetitive motifs are frequent. Many genomic tools, therefore, incorporate KMP or its derivatives as a fundamental exact-matching component in their search or indexing modules.

### 4.3 Optimizations for Large-Scale DNA Data

Although KMP already operates in linear time, further optimizations are required to handle genome-scale data efficiently. The primary challenges include limited memory bandwidth, the cost of reading massive sequence files, and the need to distribute computation across processors or nodes.

Several strategies have been proposed to address these issues:

- **Parallelization and Data Partitioning:** The genome  $T$  can be divided into independent segments processed in parallel by separate threads or nodes. The boundary overlap of  $|P| - 1$  characters ensures that matches spanning partitions are not missed.
- **Cache-Optimized Implementations:** Since KMP performs a single linear pass, implementations benefit from storing both  $T$  and  $P$  in contiguous memory to reduce cache misses. CPU vectorization (SIMD) can also accelerate comparisons in large DNA scans by evaluating multiple characters simultaneously.
- **Efficient Filtering Before Matching:** Because DNA databases can be extremely large, indexing and filtering techniques are often used to narrow the search space before applying exact string matching. Studies emphasize that organizing sequences and using efficient lookup strategies are essential for reducing unnecessary comparisons in large DNA datasets [2, 1].

These optimizations demonstrate how the theoretical efficiency of KMP translates into practical scalability. As sequencing technologies continue to advance, such adaptations will be crucial for integrating classical algorithms like KMP into high-performance computing environments for genomic analysis.

## 5 Course Related and Extension Concepts

The course emphasizes dynamic programming, hashing, and divide-and-conquer as general algorithmic strategies. In DNA analysis, these concepts naturally appear in both exact matching and sequence alignment, illustrating how classical algorithm design supports modern computational biology [1, 2].

### Dynamic Programming and Longest Common Subsequence

The concept of the Longest Common Subsequence (LCS) problem provides a theoretical foundation for sequence comparison and approximate matching in bioinformatics. The LCS between two sequences  $X$  and  $Y$  is the longest sequence of characters that appears in both strings in the same relative order, though not necessarily contiguously. Formally, given  $X = x_1x_2\dots x_m$  and  $Y = y_1y_2\dots y_n$ , the classical dynamic-

programming solution computes the LCS using the recurrence

$$L[i][j] = \begin{cases} 0, & \text{if } i = 0 \text{ or } j = 0, \\ L[i - 1][j - 1] + 1, & \text{if } x_i = y_j, \\ \max(L[i - 1][j], L[i][j - 1]), & \text{otherwise.} \end{cases}$$

This formulation yields a time complexity of  $O(mn)$  and a space complexity of  $O(mn)$ , with optimized variants reducing memory usage to  $O(\min(m, n))$ . In DNA sequence search, LCS offers a principled way to formalize similarity when sequences differ by substitutions, insertions, or deletions. As a result, LCS serves as a conceptual bridge between theoretical string comparison and practical approximate matching, providing a mathematical basis for quantifying similarity between biological sequences [3, 1].

## Hashing and Divide-and-Conquer for Scalability

Large genomic datasets require strategies that reduce search time beyond pure algorithmic improvements. Hashing enables constant-time lookup of short  $k$ -mers, filtering the genome for candidate regions likely to contain the pattern [2]. Divide-and-conquer techniques further improve scalability by partitioning the genome into overlapping windows, allowing each segment to be processed independently and often in parallel. It highlights that combining hashing with linear-time algorithms, such as KMP, enhances throughput by minimizing unnecessary comparisons and reducing I/O overhead[1]. These ideas closely parallel course topics in hashing-based indexing and decomposing large inputs into tractable subproblems [5].

## Approximate Matching and Its Relation to KMP

Approximate string matching generalizes exact matching by allowing mismatches, insertions, and deletions. Needleman–Wunsch (global alignment) and Smith–Waterman (local alignment) formalize these operations using dynamic-programming matrices that compute similarity across entire sequences or their most similar subregions [4]. Global alignment compares two sequences from beginning to end, penalizing gaps and mismatches along the full length, whereas local alignment focuses only on the highest-scoring matching region within the sequences. Edit distance further quantifies mismatch tolerance by measuring the minimum number of edits needed to transform one string into another. Although KMP performs exact matching, its prefix-function structure still informs approximate methods by illustrating how patterns overlap with themselves, a principle closely related to recurrence design in dynamic programming [5]. In practical genomic analysis, exact matching algorithms such as KMP are often used to filter or locate promising regions of a large genome before more intensive approximate matching methods evaluate their full similarity. Thus, approximate matching builds on the insights of exact pattern matching while relaxing its constraints to handle biological variability.

## Discussion and Future Directions

Modern genomic workflows increasingly integrate exact and approximate methods to balance speed and biological accuracy. Fast exact matching algorithms such as KMP or Boyer–Moore are often used to quickly identify regions of the genome that are likely to be relevant, after which more computationally intensive dynamic-programming alignment methods refine the similarity assessment. Studies of DNA matching systems demonstrate that combining hashing-based filtering, exact matching, and alignment yields strong performance on large-scale datasets [1, 2]. Future directions include improving the handling of repetitive or low-complexity regions, developing more expressive prefix-based heuristics, and incorporating GPU and distributed computing frameworks to accelerate both exact matching and alignment [6]. As genomic databases continue to expand, scalable adaptations of classical algorithms such as KMP will remain essential for efficient and accurate DNA analysis.

## 6 Conclusion

The study of string matching algorithms shows how algorithmic design directly supports biological discovery. Through careful application of computational theory, DNA sequence search has evolved from a purely experimental task into a highly automated, data-driven process. Classical algorithms such as KMP provide the theoretical foundation for efficient exact matching, while dynamic programming techniques like Needleman–Wunsch and Smith–Waterman extend this foundation to handle biological variability.

The principles learned in Algorithm Design and Analysis find direct expression in the design of modern bioinformatics tools. These algorithms are not theoretical constructs; they underpin essential processes such as genome alignment, motif detection, and variant analysis.

As sequencing technologies continue to generate unprecedented volumes of data, the fusion of algorithmic rigor with biological insight will remain central to future innovation. The connection between computer science and biology, embodied in string matching, illustrates how abstract algorithmic thinking can yield concrete impacts in understanding life at the molecular level.

## References

- [1] P. M. Rahate and M. B. Chandak. Comparative study of string matching algorithms for DNA datasets. *International Journal of Computer Sciences and Engineering*, 6(5):1067–1074, 2018.
- [2] Y. Chen. String Matching in DNA Databases. *Open Access Biostatistics & Bioinformatics*, 1(4), 2018.
- [3] I. Alsmadi and M. Nuser. String Matching Evaluation Methods for DNA Comparison. *International Journal of Advanced Science and Technology*, 47:13–31, 2012.
- [4] F. N. Muhamad, R. B. Ahmad, S. M. Asi, and M. N. Murad. Performance analysis of the Needleman–Wunsch algorithm (global) and Smith–Waterman algorithm (local) in reducing search space and time for DNA sequence alignment. *Journal of Physics: Conference Series*, 1019(1):012085, 2018. doi:10.1088/1742-6596/1019/1/012085.
- [5] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*, 4th edition. MIT Press, 2022.
- [6] N. Mhaskar and W. F. Smyth. Simple KMP pattern-matching on indeterminate strings. In *Proceedings of the Prague Stringology Conference*, pages 125–133, 2020.