# Motif Detection in Artificial DNA Sequences

Juan Carlos Quintero Rubiano
*Facultad de Ingeniería*
*Universidad Distrital Francisco José de Caldas*
Bogotá, Colombia
jcquinteror@udistrital.edu.co

*Abstract*—This paper explores motif detection in artificial DNA sequences. Various datasets with different base probabilities and sequence lengths were generated. We applied Shannon entropy to filter sequences, enhancing dataset diversity. The study involved evaluating motif detection performance before and after entropy filtering. The results are summarized in tables showing the characteristics of the dataset, motif sizes, occurrences, and detection times.

*Index Terms*—Motif detection, DNA sequences, Shannon entropy, dataset generation, bioinformatics

## I. INTRODUCTION

Detection of motifs in DNA sequences is a fundamental task in bioinformatics, aiding in understanding genetic patterns and functional elements within biological systems. This study focuses on detecting motifs in artificially generated RNA sequences. The primary objective is to explore how the varying base probabilities and sequence lengths impact motif detection performance.

We generated multiple datasets with different configurations to assess the motif detection capabilities. Additionally, Shannon entropy was employed as a filtering technique to enhance the quality and diversity of sequences. By applying entropy-based filtering, we aimed to reduce redundancy and improve the overall chaos in the datasets, thus refining motif detection.

This report presents a comprehensive analysis of motif detection before and after applying Shannon entropy. The results are organized into tables summarizing the database size, base probabilities, motif sizes, occurrences, and the time required for motif detection. The findings contribute to a better understanding of the effectiveness of entropy-based filtering in motif detection tasks.

## II. SYSTEMIC ANALYSIS

The systemic analysis of the implemented code focuses on understanding the workflow and structure of the motif detection system. The system comprises several key components designed to generate artificial DNA sequences, detect motifs, and evaluate sequence quality using Shannon entropy.

### A. Code Structure and Workflow

The core functionality of the code is encapsulated in the 'BioInformatics' class, which handles sequence generation, motif detection, and entropy-based filtering. The system operates as follows:

- **Sequence Generation:** The code generates DNA sequences based on user-defined probabilities for each base (A, C, G, T) and sequence length. This approach allows for the creation of sequences that mimic realistic DNA data while providing flexibility for users to adjust the parameters according to their needs.
- **Motif Detection:** Motif detection is performed by analyzing generated sequences to identify recurring patterns of specified sizes. The code tracks occurrences of each motif and determines the most frequent patterns. This is achieved through the use of hash maps to store and count motifs.
- **Shannon Entropy Filtering:** Shannon entropy is applied to evaluate the randomness of sequences. Sequences with entropy values below a predefined threshold are considered too repetitive and are discarded. This filtering step aims to enhance the diversity of sequences before motif detection.
- **Performance Metrics:** The code includes mechanisms to measure the time taken for motif detection. This helps in evaluating the efficiency of the motif detection process under different configurations.

### B. Algorithmic Details

The sequence generation algorithm utilizes a random number generator to assign bases according to specified probabilities, ensuring that the sequences exhibit the desired characteristics. Motif detection employs a sliding-window approach to extract and count patterns of varying sizes, leveraging hash maps for efficient counting and retrieval.

The Shannon entropy calculation involves computing the frequency distribution of bases within each sequence and applying the entropy formula:

$$H = -\sum_{i=1}^{n} p_i \log_2(p_i) \tag{1}$$

where $p_i$ represents the probability of base $i$ occurring in the sequence. Sequences with entropy values below the threshold are excluded from the final analysis.

### C. System Performance and Efficiency

The code is designed to efficiently handle varying dataset sizes and configurations. Performance is primarily influenced by factors such as sequence length, motif size, and number of sequences. The filtering step, based on Shannon entropy,

ensures that sequences with entropy values below the user-defined threshold are rewritten to maintain diversity and avoid excessive repetition.

### D. Complexity Analysis

The computational complexity of the system is determined primarily by the sequence generation, motif detection, and entropy filtering steps.

**1. Sequence Generation:** The time complexity of generating a sequence of length $n$ is $O(n)$, as each base in the sequence is selected independently based on user-defined probabilities. For a total of $m$ sequences, the overall complexity is $O(m \cdot n)$.

**2. Motif Detection:** The motif detection process involves searching for predefined patterns within the generated sequences. For each sequence, the algorithm checks for motifs of length $k$. In the worst case, the time complexity of searching for motifs is $O((n - k + 1) \cdot k)$ for each sequence, leading to an overall complexity of $O(m \cdot (n - k + 1) \cdot k)$.

**3. Entropy Filtering:** The entropy filtering step calculates Shannon entropy for each sequence, which has a complexity of $O(n)$ for each sequence, as the algorithm computes the frequency of each base and evaluates the entropy. For $m$ sequences, the overall complexity of entropy filtering is $O(m \cdot n)$.

**Overall Complexity:** Taking all steps into account, the total time complexity of the system is dominated by the motif detection process. Thus, the overall complexity of the system can be approximated as $O(m \cdot (n - k + 1) \cdot k)$. This complexity ensures the system can efficiently handle datasets of varying sizes, though performance is affected by sequence length, motif size, and the number of sequences.

### E. Chaos Analysis

In the context of this system, chaos theory plays a significant role in understanding the behavior of sequence generation, motif detection, and entropy filtering. Several phenomena associated with chaotic systems, such as the butterfly effect, the snowball effect, and the domino effect, can manifest themselves in the code results.

**Butterfly Effect:** The code allows for user-defined probabilities for each base in the DNA sequence. Small changes in these probability values can lead to disproportionately large variations in the final sequences, a hallmark of chaotic systems. A minor alteration in the probability of one base (e.g., changing the occurrence of 'A' from 0.25 to 0.26) could cascade through the sequence generation process, significantly affecting the overall motif detection and entropy filtering steps.

**Snowball Effect:** As motifs are detected across multiple sequences, a small number of errors or incorrect detections can accumulate, creating a snowball effect. If an incorrect motif is identified early in the process, it may propagate through subsequent analyses, leading to a substantial deviation in results. Similarly, sequences with low entropy values may bypass the filtering process and contribute to the formation of patterns that do not align with the expected outcomes, exacerbating errors over time.

**Domino Effect:** In this system, the interdependence between sequence generation, motif detection, and entropy filtering creates a dynamic where changes in one module affect the others. For instance, generating sequences with repetitive motifs (low entropy) will impact both the efficiency of the motif detection algorithm and the performance of the entropy filter. This interconnectedness produces a domino effect, where modifications or errors in one stage can trigger a series of cascading effects across the entire workflow, influencing both performance and accuracy.

**Entropy and Uncertainty:** The integration of Shannon entropy into the system introduces a formal measure of uncertainty and randomness in the sequences. High entropy values correspond to highly variable sequences, while low entropy indicates repetition and predictability. By filtering sequences based on entropy, the system mitigates excessive regularity, but this also introduces an element of chaos, as the threshold set by the user may eliminate useful sequences or retain redundant ones depending on the dataset and threshold value. This balance between randomness and order is key to maintaining the system's robustness and adaptability.

### F. Results

The results of sequence generation, motif detection, and entropy filtering are summarized in two tables: one without entropy filtering and another with the applied entropy threshold. These tables present key metrics such as the number of sequences, base ratios, motif lengths, and the entropy threshold (when applicable). The length interval for each generated sequence is fixed between 5 and 100.

TABLE I
RESULTS WITHOUT ENTROPY FILTERING

| No. of Sequences | A:C:G:T Ratio | Motif & Length | Motif occurrences | Execution Time (ms) |
|---|---|---|---|---|
| 1000 | 0.20, 0.48, 0.67, 1 | TTTTTT (6) | 73 | 18 |
| 2500 | 0.15, 0.56, 0.79, 1 | CCCCC (5) | 1439 | 34 |
| 6000 | 0.25, 0.50, 0.75, 1 | ACTA (4) | 1267 | 46 |
| 7950 | 0.18, 0.36, 0.85, 1 | GGGGGGGGG (9) | 573 | 122 |
| 1020510 | 0.30, 0.50, 0.80, 1 | AAAAGAG (7) | 10585 | 7906 |

TABLE II
RESULTS WITH ENTROPY FILTERING

| No. of Sequences | A:C:G:T Ratio | Motif & Length | instances | Entropy thold. | Exec time (ms) |
|---|---|---|---|---|---|
| 1000 | 0.20, 0.48, 0.67, 1 | TACCGA (6) | 11 | 2 | 8 |
| 2500 | 0.15, 0.56, 0.79, 1 | CCCCC (5) | 1289 | 1.15 | 46 |
| 6000 | 0.25, 0.50, 0.75, 1 | ACGG (4) | 1252 | 1.3 | 46 |
| 7950 | 0.18, 0.36, 0.85, 1 | GGGGGGCGG (9) | 23 | 1.9 | 77 |
| 1020510 | 0.30, 0.50, 0.80, 1 | AGAGGAG (7) | 7051 | 1.95 | 10106 |

### G. Discussion of Results

The results demonstrate how entropy filtering affects the detection of motifs in DNA sequences. Without entropy filtering, sequences with highly skewed base distributions often result in more frequent motifs due to the repetitive nature of the data. For instance, in the case of the 1020510 sequences with a high base ratio skew, the motif detection yielded a large number of

occurrences (10,585) and a significantly higher execution time (7,906 ms).

In contrast, applying entropy filtering reduces the number of sequences by removing those with low entropy, effectively discarding sequences where the base distribution is not uniform. This leads to fewer, but more diverse, motifs. For example, with the same set of 1020510 sequences, applying an entropy limit resulted in 7,051 instances of motifs, showing a more varied set and a longer execution time (10,106 ms), reflecting the increased complexity of processing.

Entropy is maximized at 2 due to the four possible bases, and variations in base distribution affect motif variability. Sequences with a more uniform base distribution tend to produce more varied motifs. This is evident from comparing results with different base ratios, where more balanced ratios (e.g., 0.25, 0.50, 0.75, 1) yield a broader range of motifs compared to highly skewed ratios.

Overall, the results highlight the impact of entropy filtering on motif detection, emphasizing the importance of base distribution in determining motif variety and processing efficiency.

### H. Conclusions

The analysis of sequence generation, motif detection, and entropy filtering reveals several key insights:

Impact of Entropy Filtering: Entropy filtering effectively reduces the number of sequences by removing those with low entropy, leading to more diverse and varied motifs. This process helps in focusing on sequences with a more balanced base distribution, enhancing the overall quality and variety of detected motifs.

Base Distribution and Motif Variety: Sequences with a more uniform base distribution produce a greater variety of motifs. In contrast, skewed base ratios often result in repetitive motifs and longer execution times, highlighting the importance of base distribution in motif detection.

Efficiency and Performance: The efficiency of motif detection is significantly affected by sequence length, motif size, and base distribution. Filtering based on entropy introduces additional computational complexity but provides more relevant results by removing sequences with excessive repetition.

Application of Entropy in Bioinformatics: The use of entropy as a filtering criterion proves valuable in refining sequence data and motif detection processes. By applying entropy thresholds, one can ensure that only sequences with meaningful variability are considered, improving the reliability of the results.