

# Report on Genetic Sequence Analysis

Fabian Yesith Aguilar Jiménez  
20231020093

September 15, 2024

## 1 Systemic Analysis

### 1.1 Objective

The objective of this analysis is to study genetic sequences to detect motifs, evaluate the complexity of sequences, and understand their chaotic behavior using entropy measures.

### 1.2 Methodology

The approach involves generating a large dataset of random genetic sequences, analyzing them for specific motifs, and evaluating their complexity using Shannon entropy. The data is filtered based on entropy, and subsequent analysis focuses on motif occurrence and execution time.

### 1.3 Tools and Technologies

- **Programming Language:** Java
- **Libraries:** Standard Java I/O for file handling
- **Metrics:** Shannon entropy for complexity analysis

## 2 Complexity Analysis

### 2.1 Entropy of Genetic Sequences

Shannon entropy provides a measure of randomness or uncertainty in a sequence. Sequences with higher entropy have a more uniform distribution of bases (A, C, G, T), while those with lower entropy have more predictable patterns.

### 2.2 Algorithm Complexity

- **Entropy Calculation:** The entropy calculation involves iterating over each sequence and counting base frequencies, which has a time complexity of  $O(n)$ , where  $n$  is the length of the sequence.
- **Database Generation:** The complexity is  $O(n \times m)$  where  $n$  is the number of sequences and  $m$  is the length. Filtering based on entropy is linear in the number of sequences,  $O(n)$ , as each sequence is processed individually.

## 3 Chaos Analysis

### 3.1 Definition of Chaos in Sequences

In the context of genetic sequences, chaos refers to the unpredictability or randomness of the sequence. High entropy indicates high chaos, while low entropy indicates lower chaos.

### 3.2 Entropy as a Measure of Chaos

By filtering sequences based on Shannon entropy, we aim to retain those with higher entropy, hence higher chaos, and remove those with predictable patterns. This helps in focusing on more variable and potentially more interesting sequences.

### 3.3 Impact of Filtering

Filtering sequences with low entropy reduces the dataset size and increases the average entropy of the remaining sequences. This process enriches the dataset with more chaotic sequences.

## 4 Results

### 4.1 Initial Dataset Statistics

- **Number of Sequences:** 200,000
- **Average Length:** 100 bases
- **Initial Database Size:** (20200000 bytes)

### 4.2 Results of Motif Search

- **Occurrences of Motif:** (1154)
- **Time Taken for Search:** (0.1144229 seconds)

### 4.3 Post-Filtering Dataset Statistics

- **Filtered Database Size:** (18648407 bytes)

### 4.4 Results of Motif Search in Filtered Database

- **Occurrences of Motif:** (1076)
- **Time Taken for Search:** (0.0456035 seconds)

## **5 Discussion of Results**

### **5.1 Analysis of Motif Occurrences**

The number of motif occurrences decreased after filtering, indicating that the removed sequences were less likely to contain the motif. This suggests that the filtering process retained sequences more likely to have the motif of interest.

### **5.2 Impact of Entropy Filtering**

Filtering based on entropy effectively reduced the dataset size and increased the average complexity of the sequences. The decrease in total number of sequences but an increase in average entropy supports the hypothesis that high-entropy sequences are more complex and less predictable.

### **5.3 Time Efficiency**

The time taken to search for motifs decreased in the filtered dataset. This improvement suggests that reducing the dataset size through entropy-based filtering made the search process more efficient.

## **6 Conclusions**

### **6.1 Summary**

This analysis demonstrates the effectiveness of using Shannon entropy to filter genetic sequences. The approach not only improves the focus on more complex sequences but also enhances the efficiency of motif searches.

### **6.2 Future Work**

Future work could involve exploring different entropy thresholds to optimize the balance between dataset size and complexity. Additionally, integrating more sophisticated pattern recognition algorithms could further enhance motif detection in chaotic sequences.

### **6.3 Implications**

Understanding the complexity and chaotic nature of genetic sequences has implications for various fields, including genomics and bioinformatics. Improved methods for filtering and analyzing sequences can lead to better insights into genetic patterns and their functional significance.