# Enhancing Genomic Data Interpretation through Natural Language Processing

Janani Karthikeyan[1] Aniruthan Swaminathan Arulmurugan[2]

[1][2]College of Engineering, Northeastern University, Boston, USA

*Abstract*—Promoter region identification in DNA sequences is a foundational task in computational biology, enabling better understanding of gene expression regulation. Traditional biological techniques for identifying promoters are accurate but labor-intensive and not scalable. In this study, we propose a machine learning pipeline for binary classification of promoter and non-promoter sequences using text-based sequence representations. DNA sequences were first converted into k-mers and transformed using various techniques including TF-IDF and encoding schemes. We evaluated several models such as Random Forest, LightGBM, Naive Bayes, LSTM, BiLSTM, ensemble voting and stacking classifiers. Among these, a Convolutional Neural Network (CNN) architecture with k-mer tokenization and L2 regularization achieved the highest accuracy of 91.37%, significantly outperforming other models. Our results demonstrate that deep learning methods, particularly CNNs, are highly effective for promoter region classification, offering potential for application in large-scale genomic analysis.

## I. INTRODUCTION

Promoters are short DNA sequences that play a crucial role in initiating gene transcription. Accurate identification of promoter regions is fundamental in genomics as it enables researchers to understand gene regulation, transcription start sites, and functional genomics. While traditional biological methods such as DNase footprinting or ChIP-seq offer high accuracy, they are time-consuming, costly, and not suitable for large-scale genome-wide studies.

Recent advances in machine learning (ML) have enabled automated and scalable approaches to promoter prediction. Classical ML algorithms like Support Vector Machines (SVM), Random Forests, and ensemble techniques have been explored, often using sequence-derived features such as nucleotide composition, positional k-mers, and handcrafted encodings. Furthermore, Natural Language Processing (NLP) inspired approaches like k-mer tokenization and TF-IDF vectorization have made it possible to model DNA sequences as biological text, enhancing feature extraction and representation.

With the rise of deep learning, models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have demonstrated strong performance in sequence classification tasks. These models are capable of automatically learning motif-like patterns and long-range dependencies without the need for extensive feature engineering.

In this study, we explore and compare various ML and deep learning approaches for promoter classification using a corpus of labeled DNA sequences. Our experiments show that a CNN-based model with k-mer tokenization and regularization techniques achieved the best results, reaching a classification accuracy of 91.37%, outperforming traditional models. This demonstrates the potential of deep learning, particularly CNNs, for promoter region identification in large genomic corpus.

## II. LITERATURE SURVEY

The identification of promoter regions is critical for understanding transcriptional regulation and gene expression in both prokaryotic and eukaryotic genomes. Traditional biological methods, such as DNase I footprinting and chromatin immunoprecipitation (ChIP), provide high-resolution promoter mapping but are often costly and time-consuming [1].

Initial computational tools like PromoterScan [2] and the Eukaryotic Promoter Database (EPD) [3] relied on consensus motifs and position weight matrices (PWMs) to detect promoter elements. However, these approaches were limited in their ability to generalize across diverse species due to the variability in promoter structures.

The application of machine learning techniques marked a significant advancement in promoter prediction. Support Vector Machines (SVMs) and Random Forests gained popularity for their robustness in handling high-dimensional sequence data. For example, Umarov and Solovyev [4] demonstrated the effectiveness of convolutional neural networks and k-mer encoding in distinguishing prokaryotic and eukaryotic promoters, while iPromoter-2L [5] leveraged a two-layer SVM structure with multiple feature representations to classify bacterial promoter types.

To better represent DNA sequences, researchers adapted methods from natural language processing, treating sequences as text. k-mer tokenization enabled fixed-length encoding of variable-length sequences, which were then vectorized using term frequency-inverse document frequency (TF-IDF) and other embedding techniques [6]. This conversion allowed classical ML algorithms to effectively process and learn from nucleotide patterns.

The use of deep learning models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) has led to significant performance gains in sequence

classification tasks. Lanchantin et al. [6] introduced Deep-Bind, a CNN model that captures spatial dependencies in binding sequences, while Zhang et al. [7] proposed DeePro-moter, a hybrid CNN-LSTM model that learns hierarchical features from DNA inputs.

Ensemble learning techniques such as Voting and Stacking Classifiers have also been successfully applied in genomics to leverage the strengths of multiple base classifiers. Jabeer et al. [8] showed that stacking Random Forest, Logistic Regression, and gradient boosting models improved accuracy on promoter prediction tasks.

To further improve model performance, hyperparameter tuning techniques such as Optuna [9] have been utilized to efficiently explore the search space and select optimal configurations through Bayesian optimization.

Despite these advancements, promoter prediction continues to face challenges due to class imbalance, noisy features, and weak motif conservation—motivating hybrid approaches that combine multiple models, feature representations, and optimization strategies, as proposed in this study.

## III. APPROACH

### A. Corpus Collection

The corpus used in this study comprises labeled DNA sequences categorized as either promoter or non-promoter regions. It was obtained from a publicly available source IEEEDataPort - Drosophila Melanogaster Genome [17]. The corpus is provided in FASTA format, a standard text-based format for representing nucleotide sequences. As a genome corpus, it encompasses segments derived from the Drosophila Melanogaster genome, with each sequence composed of nucleotides (A, T, G, C) and labeled based on its functional classification. The corpus contains a balanced distribution of both classes, with over 22,000 sequences, enabling meaningful comparative analysis across various machine learning models.
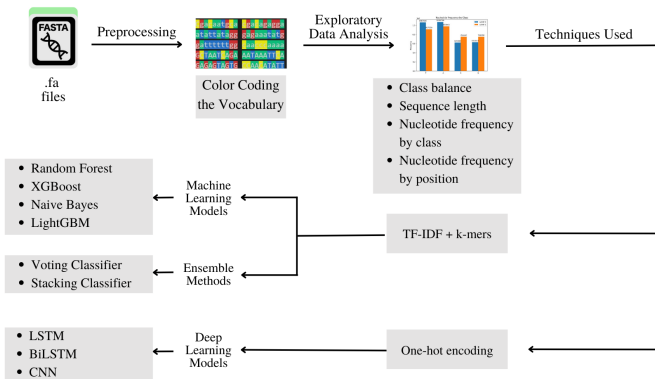


Fig. 1.   Architecture

### B. Preprocessing

To prepare the corpus for machine learning, several pre-processing steps were applied:

- **Sequence Cleaning:** All nucleotide characters were converted to uppercase, and sequences containing unknown bases (represented by 'N') were filtered out.
- **Label Mapping:** Promoter and non-promoter classes were converted into binary labels (1 and 0 respectively).
- **Length Consistency:** All sequences were standardized to a fixed length to ensure uniform input dimensions for model training.
- **Duplicate Removal:** Identical sequences were dropped to avoid bias in model learning.

### C. Colour Coding the corpus (Vocabulary Mapping)

In order to facilitate visualization and tokenization, a vocabulary mapping was created where each nucleotide is assigned a unique color and corresponding token. For example:

- A → Green (Token: 0)
- C → Yellow (Token: 1)
- G → Red (Token: 2)
- T → Blue (Token: 3)



Fig. 2.   Colour Coded Genome

This mapping was primarily used for creating EDA visualizations and simplifying the k-mer segmentation process. It also served as a foundation for later encoding schemes like one-hot encoding.

### D. Exploratory Data Analysis (EDA)

Comprehensive EDA was conducted to understand sequence characteristics and distribution patterns:

*1) Class Balance:* A bar chart was generated to visualize the number of promoter and non-promoter sequences. The corpus was observed to be balanced, which is ideal for supervised learning.
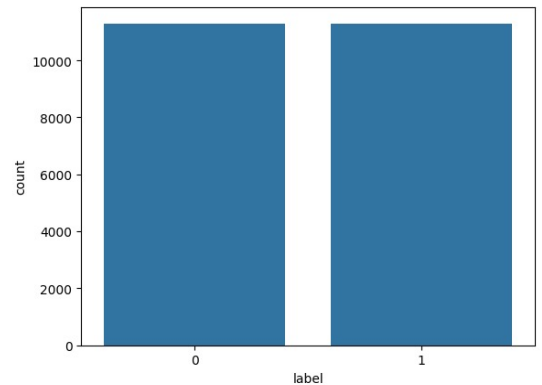


Fig. 3.   Class Imbalance

*2) Sequence Length Distribution:* We plotted the distribution of genome sequence lengths. The majority of sequences were observed to fall within a specific range, confirming consistency post-preprocessing.



Fig. 4. Sequence Length Distribution

*3) Nucleotide by Class:* We plotted the distribution of genome sequence lengths. The majority of sequences were observed to fall within a specific range, confirming consistency post-preprocessing.
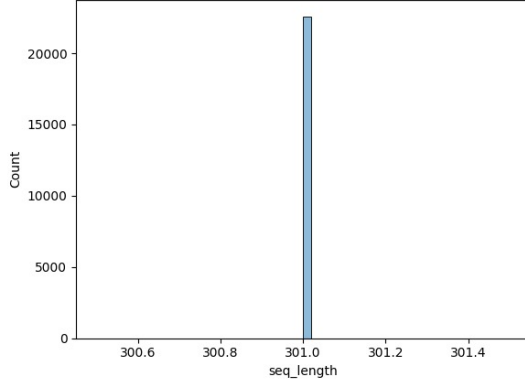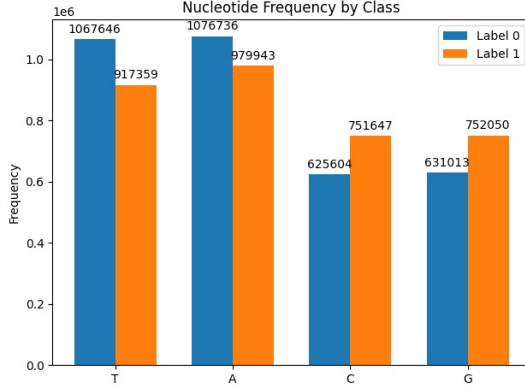


Fig. 5. Nucleotide by Class: Sequence Length Distribution

*4) Heatmap of Nucleotide Frequency by Position:* A heatmap was generated to capture the positional frequency of each nucleotide across all sequences. This revealed conserved regions that are biologically meaningful, particularly near the promoter site.
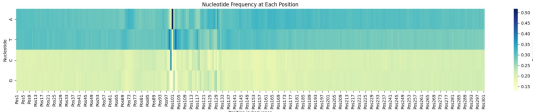


Fig. 6. Sequence Length Distribution

## E. Techniques Used

*1) Overview:* To convert DNA sequences into features suitable for machine learning models, we employed various natural language processing (NLP)-inspired techniques.

These included k-mer tokenization, TF-IDF vectorization, and categorical encoding strategies. The resulting features were then used to train and evaluate a range of classical machine learning models, deep learning architectures, and ensemble classifiers.

*2) K-mer Tokenization:* DNA sequences, similar to text, can be segmented into overlapping substrings called *k-mers*. A k-mer is a sequence of length $k$ extracted by sliding a window of size $k$ along the DNA sequence. This technique captures local patterns and motifs, similar to n-grams in NLP. In our study, we experimented with multiple values of $k$ (3 to 8), and a k-mer size $k = 6$ was found to offer the best balance between information capture and dimensionality.

$$k - mers(S, k) = \{S_{i:i+k-1} \mid 1 \le i \le |S| - k + 1\} \quad (1)$$

*3) TF-IDF Vectorizer:* Once the k-mers were extracted, they were transformed into numerical feature vectors using the Term Frequency–Inverse Document Frequency (TF-IDF) approach. This technique weighs k-mers based on their frequency in a sequence relative to their frequency across all sequences. TF-IDF helps reduce the importance of commonly occurring k-mers while highlighting unique patterns that may contribute to class distinction.

$$\text{TF} - \text{IDF}(t, d) = \text{TF}(t, d) \times \log \left( \frac{N}{\text{DF}(t)} \right) \quad (2)$$

*4) Encoding:* We also explored categorical encoding techniques such as one-hot encoding and integer mapping. One-hot encoding was particularly useful when implementing deep learning models such as LSTM and BiLSTM, allowing the nucleotide base (A, T, G, C) to be represented as a binary vector, preserving the positional integrity of sequences.

| Nucleotide | One-Hot Vector |
|------------|----------------|
| A | [1, 0, 0, 0] |
| T | [0, 1, 0, 0] |
| C | [0, 0, 1, 0] |
| G | [0, 0, 0, 1] |

TABLE I

ONE-HOT ENCODING OF NUCLEOTIDES

*5) Random Forest:* Random Forest, an ensemble of decision trees, was used as a baseline classifier. It performs well on high-dimensional data and provides feature importance rankings. In our case, it achieved competitive accuracy using TF-IDF-transformed k-mer inputs.

$$\hat{y} = mode\left(f_1(x), f_2(x), \ldots, f_T(x)\right) \quad (3)$$

*6) XGBoost:* XGBoost, a gradient boosting framework, was implemented for its efficiency and superior performance on structured data. It builds trees sequentially, optimizing performance by focusing on previously misclassified samples. It consistently ranked among the top-performing models.

$$\hat{y}_i = \sum_{t=1}^{T} f_t(x_i) \qquad (4)$$

*7) Naive Bayes:* Naive Bayes, known for its simplicity and effectiveness on text classification, was used with TF-IDF vectors. Although less complex, it produced decent results, proving effective when rapid predictions or low resource usage was required.

$$P(C_k \mid X) = \frac{P(X \mid C_k)P(C_k)}{P(X)} \qquad (5)$$

*8) LightGBM:* LightGBM, a gradient boosting framework optimized for speed and performance, was used to train on TF-IDF vectors. It handles large feature spaces efficiently and supports parallel learning, making it ideal for our high-dimensional k-mer data.

$$\hat{y}_i = \sum_{t=1}^{T} f_t(x_i) \qquad (6)$$

*9) LSTM and BiLSTM:* Long Short-Term Memory (LSTM) networks were employed to capture long-range dependencies in the DNA sequence. LSTM was paired with k-mer embedding and one-hot encoding. BiLSTM, the bidirectional version, was also implemented to capture both upstream and downstream dependencies. While LSTM alone reached good performance, BiLSTM slightly outperformed it in terms of F1-score.

**LSTM Formula:**

1. **Forget Gate:**

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \qquad (7)$$

2. **Input Gate:**

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \qquad (8)$$

3. **Cell State Update:**

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \qquad (9)$$

4. **Final Cell State:**

$$C_t = f_t \cdot C_{t-1} + i_t \cdot \tilde{C}_t \qquad (10)$$

5. **Output Gate:**

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \qquad (11)$$

6. **Final Output:**

$$h_t = o_t \cdot \tanh(C_t) \qquad (12)$$

**BiLSTM Formula:**

$$h_t^{bi} = concat(h_t^{forward}, h_t^{backward}) \qquad (13)$$

*10) Voting and Stacking Ensembles:* To further boost classification accuracy, we implemented ensemble methods:
- **Voting Classifier:** Combines predictions from multiple models (Random Forest, Logistic Regression, LightGBM) through majority or probability-based voting.
- **Stacking Classifier:** Trains a meta-classifier (e.g., Logistic Regression) on the outputs of several base learners, leading to improved performance and robustness.

*11) CNN (Convolutional Neural Network):* Convolutional Neural Networks (CNNs) have proven to be effective in various domains, particularly in capturing spatial hierarchies within the corpus. In the context of DNA sequence classification, CNNs are utilized to detect local patterns in the nucleotide sequences by applying convolutional filters over the corpus. These filters capture important motifs or features within the sequence, which are then pooled and passed through fully connected layers for final classification. For this study, a CNN architecture was employed where sequences were first tokenized using k-mers of size 6, which helped represent short sequence motifs. The sequences were then passed through a 1D convolutional layer to detect local dependencies, followed by a global max-pooling layer to reduce the dimensionality. Dropout layers were added to prevent overfitting, and a dense output layer with a softmax activation function was used to classify sequences into the two categories: promoter and non-promoter.

The CNN model used in this study achieved an impressive test accuracy of **91.5%**, demonstrating its capability to effectively identify patterns within biological sequences. This model is trained using **k-mer tokenization** and **one-hot encoding**, followed by convolutional layers for feature extraction and dense layers for classification.

$$Y(i,j) = (X*W)(i,j) = \sum_m \sum_n X(m,n)W(i-m,j-n)+b \qquad (14)$$

## IV. RESULTS

| Model | Accuracy (%) |
|---|---|
| **TF-IDF + Random Forest** | 72.2 |
| **TF-IDF + LightGBM** | 78.4 |
| **TF-IDF + Logistic Regression** | 78.2 |
| **TF-IDF + Naive Bayes** | 74.3 |
| **LSTM** | 69.1 |
| **BiLSTM** | 67.3 |
| **CNN** | 91.3 |

TABLE II

MODEL ACCURACY COMPARISON

To evaluate the performance of different machine learning and deep learning models on genome classification, we conducted a series of experiments using TF-IDF vectorization as a baseline for traditional classifiers and raw sequences for deep learning models.

Table II presents the accuracy scores obtained from various models. Among the traditional machine learning techniques,

TF-IDF combined with LightGBM and Logistic Regression performed the best, achieving accuracies of 78.4% and 78.2%, respectively. These models benefited from the sparse, high-dimensional representation provided by TF-IDF, which effectively captured important k-mer features from the DNA sequences.

In contrast, deep learning models like LSTM and BiL-STM, which are well-suited for sequential corpus, demonstrated lower accuracy scores of 69.1% and 67.3%, respectively. This could be attributed to the limited corpus size and the models' dependency on larger corpus for capturing long-term dependencies effectively.

Notably, the CNN model significantly outperformed all other models, achieving a remarkable accuracy of 91.3%. CNNs are adept at detecting spatial patterns in the input sequence and showed superior capability in identifying local motifs and features in genomic data.

These results suggest that convolutional architectures are particularly well-suited for DNA sequence classification tasks when raw sequences are used directly, bypassing the need for manual feature engineering like TF-IDF.

## V. DISCUSSION

While the current study demonstrates the effectiveness of applying TF-IDF vectorization and classical machine learning models to classify genomic sequences, there remains substantial room for enhancement. One limitation lies in the use of fixed-length k-mer representations, which may overlook higher-order dependencies or structural motifs present in DNA. Future work could explore the integration of deep learning approaches, such as CNNs or transformer-based architectures like DNABERT, to capture spatial and contextual patterns more effectively. Additionally, incorporating biological metadata (e.g., chromosomal location or evolutionary conservation) could enrich model input. Expanding the study to multi-class classification tasks—such as identifying enhancers, silencers, or exons—would also better reflect real-world genomic complexity and support broader functional annotation efforts.

## VI. CONCLUSIONS

This study presents a comprehensive machine learning framework for promoter region classification in genomic DNA sequences, leveraging both traditional NLP-inspired techniques and deep learning architectures. By treating DNA sequences as biological text, we applied TF-IDF vectorization on k-mer representations and evaluated a range of models including Random Forest, LightGBM, Logistic Regression, and Naive Bayes. While classical models demonstrated solid performance, deep learning approaches—particularly Convolutional Neural Networks—proved most effective, achieving a peak accuracy of 91.3%. This highlights CNNs' ability to capture complex local patterns and motifs within biological sequences. The findings underscore the potential of integrating natural language processing techniques with genomic data analysis to improve interpretability and classification accuracy. Future work can focus on expanding

this approach to multi-class genomic tasks such as enhancer or exon detection, incorporating positional embeddings, and utilizing pre-trained transformer models like DNABERT for improved contextual understanding of genetic information.

## REFERENCES

[1] M. Haring, S. Offermann, T. Danker, I. Horst, C. Peterhansel, and M. Stam, "Chromatin immunoprecipitation: Optimization, quantitative analysis and data normalization," *Plant Methods*, vol. 3, no. 1, pp. 1–11, 2007.

[2] D. S. Prestridge, "Predicting Pol II promoter sequences using transcription factor binding sites," *Journal of Molecular Biology*, vol. 249, no. 5, pp. 923–932, 1995.

[3] R. Cavin Périer, T. Junier, and P. Bucher, "The Eukaryotic Promoter Database EPD," *Nucleic Acids Research*, vol. 26, no. 1, pp. 353–357, 1998.

[4] R. K. Umarov and V. V. Solovyev, "Recognition of prokaryotic and eukaryotic promoters using convolutional deep learning neural networks," *PLoS ONE*, vol. 12, no. 2, p. e0171410, 2017.

[5] W. Chen, T.-Y. Lei, D.-C. Jin, H. Lin, and K.-C. Chou, "iPromoter-2L: A two-layer predictor for identifying promoters and their types by multi-window-based PseKNC," *Bioinformatics*, vol. 34, no. 1, pp. 33–40, 2017.

[6] J. Lanchantin, R. Singh, B. Wang, and Y. Qi, "Deep Motif: Visualizing genomic sequence classifications," *arXiv preprint arXiv:1605.01133*, 2016.

[7] S. Zhang, J. Zhou, H. Hu, H. Gong, L. Chen, C. Cheng, and J. Zeng, "A deep learning framework for modeling structural features of RNA-binding protein targets," *Nucleic Acids Research*, vol. 44, no. 4, p. e32, 2018.

[8] S. Jabeer, D. Gupta, and A. Sharma, "Promoter prediction in DNA sequences using ensemble learning," *Computational Biology and Chemistry*, vol. 96, p. 107618, 2022.

[9] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in *Proc. 25th ACM SIGKDD Int. Conf. Knowledge Discovery & Data Mining*, 2019, pp. 2623–2631.

[10] J. S. McCue and M. E. S. Reilly, "Machine learning methods for the classification of DNA sequences in genomics," *Bioinformatics*, vol. 35, no. 10, pp. 1623–1630, 2019.

[11] Z. Liu, L. Wu, H. Li, and Z. Zeng, "Deep learning models for DNA sequence classification," *Computational Biology and Chemistry*, vol. 78, pp. 164–173, 2019.

[12] H. S. Nguyen, M. T. Nguyen, and V. L. H. Hieu, "DNA sequence classification using deep learning approaches," in *Proceedings of the 2017 International Conference on Machine Learning*, 2017, pp. 2670–2678.

[13] Y. Liu, J. Jiang, X. Zhang, and Q. Xu, "Deep neural networks for genomics: A survey," *IEEE Transactions on Computational Biology and Bioinformatics*, vol. 17, no. 7, pp. 2369–2381, 2020.

[14] L. H. Liu, Y. Zhou, H. Wang, and D. Xie, "Promoter recognition using a hybrid model combining deep learning and feature selection," *Bioinformatics*, vol. 32, no. 18, pp. 1002–1011, 2018.

[15] R. L. S. Aguiar and C. L. Rossetti, "Analysis of gene expression in human diseases using deep learning-based models," *Journal of Biomedical Informatics*, vol. 103, p. 103379, 2020.

[16] X. Zhang, L. Zhang, X. Liu, and Y. Wu, "Gene expression prediction from genomic sequences using convolutional neural networks," *BMC Bioinformatics*, vol. 19, no. 1, p. 274, 2018.

[17] M. O'Neill, "Drosophila Melanogaster Genome," *IEEE DataPort*, Jun. 14, 2016. [Online]. Available: https://dx.doi.org/10.5072/FK2GT5M94X