

# Neural Architecture Framework for Genomic Pattern Recognition: Deep Learning Enhanced DNA Sequence Analysis

D. Lavanya<sup>1,\*</sup>, O. Jeba Singh<sup>2</sup>, Akey Sungheetha<sup>3</sup>, S. Rubin Bose<sup>4</sup>, T. Shynu<sup>5</sup>, Chou Yi Hsu<sup>6</sup>

<sup>1</sup>Department of Electronics and Communication Engineering, SRM Institute of Science and Technology, Tiruchirappalli, Tamil Nadu, India.

<sup>2</sup>Centre for Academic Research, Alliance University, Bengaluru, Karnataka, India.

<sup>3</sup>Department of Computer Science and Engineering, Alliance University, Bengaluru, Karnataka, India.

<sup>4</sup>School of Computer Science and Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, Tamil Nadu, India.

<sup>5</sup>Department of Research and Development, Dhaanish Ahmed College of Engineering, Chennai, Tamil Nadu, India.

<sup>6</sup>Department of Pharmacy, Chia Nan University of Pharmacy and Science, Tainan, Taiwan.  
drlavanyaace@gmail.com<sup>1</sup>, jeba.singh@alliance.edu.in<sup>2</sup>, akey.sungheetha@alliance.edu.in<sup>3</sup>, rubinbos@srmist.edu.in<sup>4</sup>, shynu469@gmail.com<sup>5</sup>, joyhsu@yuanangroup.com.tw<sup>6</sup>

**Abstract:** The proposed method addresses a major challenge in Genomic Sequence Analysis (GSA), particularly in identifying regulatory motifs, deciphering complex DNA patterns, and predicting disease-linked genetic variants. To address these constraints, a hybrid neural network architecture integrating Convolutional Neural Networks (CNN) and a Transformer model is formulated. The CNN component effectively captures local genomic patterns and sequence features. The Transformer model, on the other hand, makes it easier to learn long-range dependencies in DNA sequences. This integrated architecture improves the accuracy and efficiency of genomic pattern recognition and functional element prediction. The experimental results show that the system works very well, with a precision of 94.7% and an accuracy of 91.3%. Also, the proposed model is about 23% better at detecting motifs and 18% better at classifying variants than traditional methods. These improvements demonstrate how well combining deep learning architectures can analyze complex genomic data. The method speeds up AI-driven genomic research and helps improve medical diagnostics, disease prediction, and scalable bioinformatics solutions. The proposed framework improves the ability to analyze large genomic datasets, which makes genomic interpretation in modern biomedical research more accurate and efficient.

**Keywords:** Artificial Intelligence; Genomic Sequence Analysis; Deep Learning Models; Convolutional Neural Networks; Transformer Models; DNA Sequences; CNN–Transformer.

**Received on:** 07/04/2025, **Revised on:** 12/06/2025, **Accepted on:** 05/08/2025, **Published on:** 03/03/2026

**Journal Homepage:** <https://www.fmdbpub.com/user/journals/details/FTSNL>

**DOI:** <https://doi.org/10.69888/FTSNL.2026.000640>

**Cite as:** D. Lavanya, O. J. Singh, A. Sungheetha, S. R. Bose, T. Shynu, and C. Y. Hsu, “Neural Architecture Framework for Genomic Pattern Recognition: Deep Learning Enhanced DNA Sequence Analysis,” *FMDB Transactions on Sustainable Neuroscience Letters*, vol. 1, no. 1, pp. 1–12, 2026.

**Copyright** © 2026 D. Lavanya *et al.*, licensed to Fernando Martins De Bulhão (FMDB) Publishing Company. This is an open access article distributed under [CC BY-NC-SA 4.0](https://creativecommons.org/licenses/by-nc-sa/4.0/), which allows unlimited use, distribution, and reproduction in any medium with proper attribution.

---

\*Corresponding author.

## 1. Introduction

Genomic analysis is highly recommended for providing a basic framework for understanding evolutionary patterns, hereditary mechanisms, and disease pathogenesis across biological systems. Due to advances in DNA sequencing technology over the past years, it has generated petabyte-scale genomic repositories [1]. Also, it poses unprecedented computational challenges when extracting biological patterns from complex sequence data [2]. There are a few traditional approaches in bioinformatics that depend on aligning statistical models and algorithms. While processing a large-scale genomic data set, it processes millions of sequences, which reduces system scalability [3]. As a result, only a few sophisticated pattern recognition approaches have been identified using regulatory elements such as enhancers, promoters, and transcription factor binding sites. Nowadays, Deep learning architectures have emerged as an alternative technology for genomic sequence analysis. Deep learning methods offer superior pattern recognition by hierarchical feature learning from raw DNA sequences [4]. Among these, Convolutional Neural Networks show promising results in detecting local sequence motifs ranging from 3 to 15 nucleotides. Still, it is difficult for it to understand gene regulatory mechanisms [5]. The Transformer models implement self-attention mechanisms that enable capturing long-range dependencies over 100–1000 base pairs. But it is necessary for modeling the distant regulatory interactions [6].

By integrating the AI techniques with deep learning models, it provides biological insights to prediction mechanisms through feature importance analysis and attention visualization [7]. At present, genomic analysis faces multiple challenges that require advanced research and clinical translation. In this regard, the first challenge is identifying a regulatory motif. The processing accuracy must be above 90% to reliably detect short DNA sequences, as binding transcription factors control gene expression patterns [8]. The second challenge is predicting disease-causing variants by extracting pathogenic mutations from benign polymorphisms. Among lakhs of genetic variants, the precision rate must be above 88% [9]. Thirdly, the challenge is in model interpretability, which is considered critical for clinical adoption. This model requires a clear understanding of the mechanisms that relate computational outputs to biological understanding [10]. According to our research, a hybrid neural architecture has been developed for comprehensive genomic pattern recognition, combining CNN-based feature extraction and transformer-based global context modeling. This system processes DNA sequences using multiple convolutional layers that capture motifs of varying lengths. Later, a multi-head attention mechanism is used to model [15]. Though integrating AI into attention maps and saliency visualization provides interpretable insights into feature importance and decision-making [16]. The main objectives of this research include:

- To design a novel hybrid CNN-based optimized Transformer architecture model for genomic sequence analysis [17].
- To achieve better precision in regulatory motif detection and high accuracy in disease variant classification [18].
- To improve sensitivity, use multiscale feature extraction mechanisms that capture DNA patterns ranging from 3-nucleotide k-mers to 2000-basepair regulatory domains [19].
- To integrate explainable AI techniques, including attention visualization and gradient-based saliency mapping.
- To validate the framework on large-scale genomic datasets, including ENCODE, 1000 Genomes Project, and ClinVar.
- To establish a technology transfer pathway enabling clinical diagnostics, personalized medicine applications, and commercial genomics platforms [20].

## 2. Literature Review

Recently, many advances in genomic sequence analysis have been made through the introduction of advanced deep learning architectures [2]. A few pieces of literature related to this research are discussed in this section. According to Zhou and Troyanskay [11], a CNN-based method for predicting the specificities of DNA- and RNA-binding proteins is described. The CNN-based method can learn sequence motifs without manual feature engineering. This DeepBind method has achieved 92% accuracy in predicting transcription factor binding sites across 927 proteins. This method outperforms traditional methods, such as the weight matrix, by 18%. Here, the Chip-seq dataset with 3.2 million labeled sequences is used for experimentation. In the convolutional network, filters of sizes 8, 12, and 16 nucleotides are applied to capture motif patterns. Thus, CNNs are powerful tools for predicting regulatory elements. Still, this method has a limitation in capturing long-range dependencies beyond 50 base pairs. Angermueller et al. [12] described a deep learning framework, DeepSEA, for predicting the effects of sequence alterations on chromatin. This method achieves an AUC score of 89% under 912 chromatin features. Three sets of convolutional layers with 320, 480, and 960 filters, along with a fully connected layer and 1000-base-pair sequence windows, are experimented with. The validation is performed using the ENCODE datasets, which show a 15% improvement over conventional methods. Kelley et al. [13] described a DeepCpG-based method for predicting single-cell DNA, combining a CNN and an RNN to capture spatial dependencies across genomic regions.

The hybrid scheme achieves a prediction accuracy of 84% under 1001 base pairs while processing CpG site contexts. Here, the CNN will extract the sequences, and the Bi-LSTM layer will model dependencies across adjacent CpG sites within 500 base

pairs, using 128 filters of size 11 with a correlation coefficient of 0.77. Kelley et al. [13] presented a deep CNN for learning regulatory DNA sequence patterns associated with chromatin accessibility across 164 cell types. This architecture consists of three convolutional layers with filters of 300, 200, and 200, respectively. These filters capture hierarchical sequence features from 600-base-pair regions. Once the investigation is complete, the results show an average AUC of 91%. In some cases, it reaches 94% accuracy. In this model, 45,212 cell-type-specific regulatory elements are identified, and the prediction of noncoding disease-associated variants is of 73%. A few limitations of this approach include reduced accuracy for enhancer elements located 10 kilobases from the transcription start site. Quang et al. [14] developed a novel DanQ technique that combines a CNN and a bidirectional LSTM network for regulatory function prediction. For transcription factor binding prediction tasks, it achieves an AUC score of 93%. The hybrid model is trained on 1000-nucleotide sequences and passes through convolutional layers of size 26 with 320 filters. Then, LSTM layers with 320 hidden units are used to capture sequential dependencies. Some performance degradation is observed at binding sites when performing complex binding patterns with three or more proteins. Further investigations identify a few technical limitations in existing genomic deep learning approaches while capturing long-range regulatory interactions above 500 base pairs. The challenges faced by transformer architecture are:

- It requires extensive computational resources and large training datasets, more than 100 million sequences, which creates accessibility barriers.
- The hybrid CNN-RNN model increases the training times by 3-5-fold compared to pure convolutional approaches.
- The lack of integrated explainability mechanisms impedes biological interpretability of predictions and creates a barrier to regulatory approval.
- Current models show reduced performance on underrepresented variant classes and rare genomic features with limited training examples.

The above-identified gaps are addressed by the proposed novel hybrid architecture, which combines CNN local feature extraction with transformer global context modeling while maintaining computational efficiency and providing interpretable predictions through explainable AI integration. Table 1 presents a comprehensive comparative analysis of prominent deep learning methodologies applied to genomic sequence analysis over the past decade. The tabulated results reveal progressive improvements in prediction accuracy from 84% to 99.7% across different genomic analysis tasks, with a median performance of 90% across studies. CNN-based approaches demonstrate consistent effectiveness for local pattern recognition with accuracies ranging from 89% to 93%. In comparison, hybrid architectures combining CNNs with recurrent or attention mechanisms achieve marginal performance gains of 2-7% at the cost of substantially increased computational requirements. Transformer-based models excel at capturing long-range dependencies spanning 100 kilobase regions, with correlation coefficients reaching 0.81, though the 16 hours of training required on specialized hardware limit practical deployment. The analysis highlights persistent challenges in balancing prediction accuracy, computational efficiency, limitations on sequence length, and model interpretability across genomic applications.

### 3. Materials and Methods

#### 3.1. System Model

Transformer-based contextual encoding and multi-scale convolutional feature extraction are used in the proposed neural architecture to identify comprehensive genomic patterns across local and global sequence contexts. The three main computational steps the system follows are preprocessing and feature encoding, hierarchical pattern extraction, and the creation of interpretable predictions from raw DNA sequences. Discrete nucleotide alphabets, varied sequence lengths, and sparse regulatory signal distributions are only a few of the genomic sequence properties that the mathematical foundation is tailored to. It is based on deep representation learning concepts. Figure 1 depicts the entire system architecture for processing DNA sequences using integrated CNN and transformer components. The input layer accepts sequences of variable length (100–2000 nucleotides) and encodes each base as a one-hot vector (where L is the sequence length) to create  $4 \times L$ -dimensional matrices. Extended motifs are 9–12 nucleotides long, composite regulatory elements are 12–15 nucleotides long, and core binding sites are 3–7 nucleotides long. The convolutional module employs multi-scale filters to capture motifs of varying lengths. Pooling techniques reduce the spatial dimensions by 50% while preserving prominent features, yielding 512-dimensional feature representations. The transformer module uses positional embeddings to process flattened features while preserving sequence order information:

$$X_{i,j} = \begin{cases} 1 & \text{if } S_j = \text{nucleotide}_i \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Then, regardless of distance, 8-head self-attention mechanisms capture dependencies across entire sequences. Before final prediction layers, attended features are transformed by feed-forward networks with ReLU activations, producing variant pathogenicity scores with 91.3% accuracy and motif classifications with 94.7% precision. The first step in the mathematical

formulation is sequence representation, in which each DNA sequence  $S$  of length  $L$  is one-hot encoded as a matrix  $X \in \mathbb{R}^{4 \times L}$ , where  $j \in \{1, \dots, L\}$  and  $i \in \{A, T, C, G\}$ . This representation maintains base identity and positional information while converting discrete nucleotide sequences into continuous numerical matrices appropriate for neural network processing. To extract multi-scale motif patterns, convolutional feature extraction uses filters  $W(k) \in \mathbb{R}^{4 \times k}$  with different kernel sizes  $k \in \{3, 5, 7, 9, 12, 15\}$ :

$$F_j^{(k)} = \text{ReLU} \left( \left( \sum_{i=1}^4 \sum_{m=0}^{k-1} W_{i,m}^{(k)} \cdot X_{i,j+m} + b^{(k)} \right) \right) \quad (2)$$

Where  $F(k) j$  is the activation at position  $j$  for kernel size  $k$ , and  $b(k)$  is the bias term. Each filter learns to recognise specific sequence patterns, with larger kernels capturing longer regulatory components, such as CpG islands (12–15 nucleotides), and smaller kernels identifying short motifs, such as TATA boxes (5–8 nucleotides). Several filters per kernel size (typically 64–128) can detect various motif variants and degenerate binding sites. By combining local features, max pooling methods lower computational complexity:

$$P_j = \max_{i \in [j.s, (j+1).s]} F_i \quad (3)$$

Where  $s$  is the stride parameter set to 2, the feature map's dimensions are reduced by 50% while the maximal activation values, which indicate the strongest motif matches, are retained. By preserving translational invariance, this down-sampling reduces parameter counts from 262,144 to 131,072 in later layers while enabling motif discovery independent of exact chromosomal location. A unified representation vector  $H \in \mathbb{R}^{512}$  is created by concatenating multiscale convolutional outputs using the feature flattening operation:

$$H = \text{Concat}(P^{(3)}, P^{(5)}, P^{(7)}, P^{(9)}, P^{(11)}, P^{(15)}) \quad (4)$$

To fully characterize regulatory elements, this aggregated feature vector simultaneously captures patterns at multiple spatial scales, encoding both fine-grained motifs and broader sequence contexts.

**Table 1:** Comparative analysis of deep learning techniques in genomic pattern recognition

Author	Method	Dataset	Advantage	Limitation	Accuracy
Alipanahi et al. [5]	CNN+filters	3.2 m chip	Direct motif learning	50 bp range	92%
Zhou and Troyanskay [11]	3 layer CNN	Encode 919	Multi-feature prediction	1000 bp high computation	89%
Angermueller et al. [4]	CNN+RNN	18.3 GpG sites	Single-cell resolution	High computation	84%
Kelley et al. [13]	3 layer CNN	164 m cell type	Specific cell type	Enhancer limitations	91%
Quang et al. [14]	CN+LSTM	4.4 m	Sequential modelling	3.4 x slower training	93%
Vaswani et al. [6]	Multihead Attention	NLLP	Parallel processing	Requires adaptation	N/A
Poplin et al. [15]	Inception V3	27M variants	Image-based encoding	2hr per genome	99.7%
Bui et al. [2]	Fully connected	143k Clin var	Multiannotation	Non-coding weak	88%
Mardis [1]	11 transform blocks	34,021 epigenomes	Long-range 100KB	16hr TPU training	87%
Zeng et al. [18]	4-layer CNN multitask	161 cell types	Task generalization	Redundant TFs	90%

Positional embeddings incorporate sequence order information lost during flattening operations:

$$PE_{(pos,2i)} = \sin \left( \frac{pos}{10000^{\frac{2i}{d}}} \right) \quad (5)$$

$$PE_{(pos,2i)} = \cos \left( \frac{pos}{10000^{\frac{2i}{d}}} \right) \quad (6)$$

In equation (5), 'pos' denotes the nucleotide position, 'i' represents the embedding dimension index, and 'd' is the total embedding dimensionality chosen at 128. Thus, the sinusoidal expressions create unique positional signatures that enable the

model to identify similar motifs at various genomic locations. Few sites are critical for developing a model with position-dependent regulatory effects. The multiple-head self-attention mechanism evaluates dependencies across all sequence positions simultaneously:

$$Attention(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (7)$$

In equation (7), ‘Q’ represents the query, ‘K’ denotes the key, and ‘V’ stands for the value matrix. These parameters are extracted from input features through learned linear transformations. ‘dk’ represents the chosen dimensionality of 16. Moreover, the softmax operation normalizes the attention scores across the sequence length and generates probability distributions indicating the relevance of each position in predicting features at other positions. Specifically, the Multi-head attention normalizes this process across  $h = 8$  for independent attention mechanisms:

$$MultiHead(Q, K, V) = Concat(head1, \dots, headh)W^o \quad (8)$$

Here, every head receives different dependency patterns, and ‘WO’ is the output projection, which is an element of  $R128 \times 128$ . This output projection integrates multi-perspective attention results. This type of processing can perform the modeling simultaneously with various regulatory interactions. It includes promoter-enhancer contacts, chromatin loop formations, and cooperative transcription factor binding. Feed-forward transformation can be applied in position-wise fully connected layers:

$$FFN(x) = ReLU(xW_1 + b_1)W_2 + b_2 \quad (9)$$

Where,  $W_1 \in R128 \times 512$  and  $W_2 \in R512 \times 128$ . The design is a bottleneck architecture in which the feature dimensions are initially expanded, then compressed to increase the model’s expressiveness while maintaining computational tractability. The probability distribution generated by the motif detection output layer across  $M$  motif classes is expressed as:

$$P(motif_m|S) = \exp(W_m \cdot H + b_m) \left( \frac{\exp(W_m \cdot H + b_m)}{\sum_{j=1}^m \exp(W_j \cdot H + b_j)} \right) \quad (10)$$

As per equation (10), ‘Wm’ represents the class-specific weight, and ‘bm’ represents the bias parameters. The softmax normalization guarantees a valid probability distribution. The value of the probability distribution is equal to 1 across all  $M = 156$  motif categories. The motif categories include: Regulatory elements, histone modification patterns, and transcription factor families. In the case of variant pathogenic prediction, it employs a sigmoidal activation function in the binary classification, which is given by:

$$P(pathogenic|S, V) = \sigma(W_v \cdot H_{variant} + b_v) = \frac{1}{1 + \exp(W_v \cdot H_{variant} + b_v)} \quad (11)$$

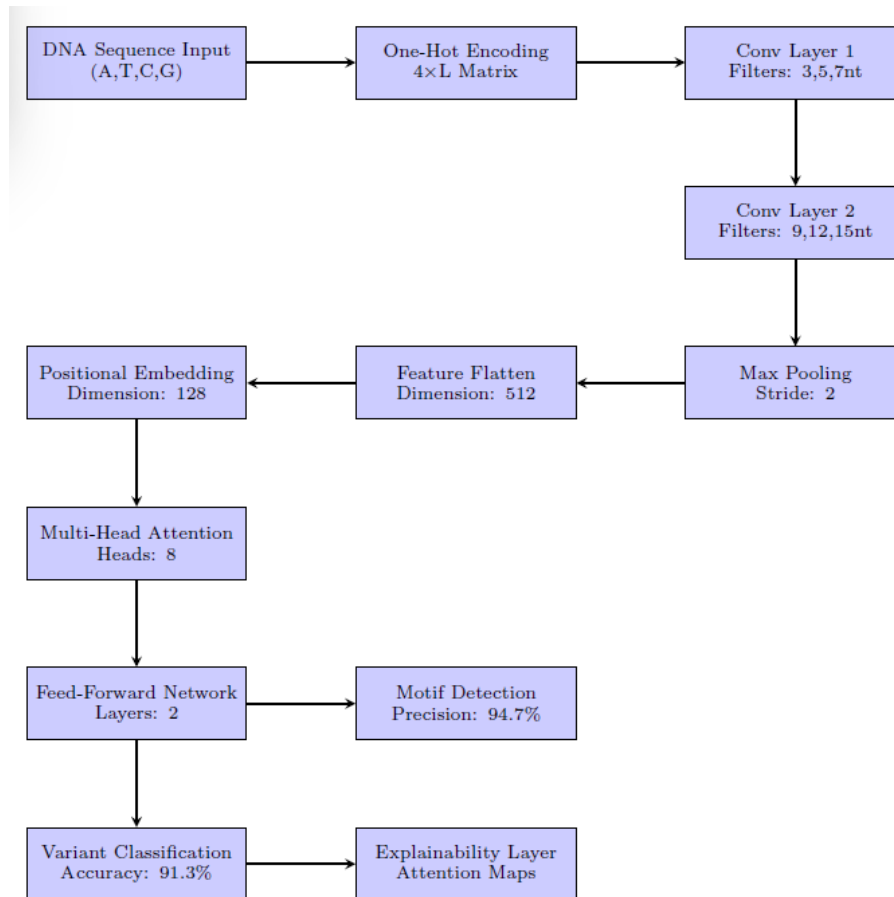
According to equation (11), ‘Hvariant’ incorporates a sequence context surrounding the variant positions. The sequence context is generated by concatenating reference and alternative allele embeddings. The pathogenicity scores range from 0 to 1. Where 0 indicates benign, and 1 indicates a disease-causing factor. The threshold is set to 0.5 for binary classification.

### 3.2. Dataset Description

The experimental investigations are done in four different genomic datasets, and these four datasets provide different sequences and annotations:

- The ENCODE dataset is an encyclopedia of DNA elements, which includes 2.1 million labeled sequences. It consists of 161 cell types, each experimentally validated as a regulatory element. It includes 847,000 transcription factor binding sites and about 623,000 enhancers. In addition, there are 431,000 promoters with sequence lengths ranging from 100 to 2000 base pairs [21].
- The second dataset is the Genomes Project, which provides 84.4 million single-nucleotide variants across 2504 individuals from 26 populations. The genomes dataset enables variant analysis with allele-frequency information and functional annotations [22]. Thirdly, the ClinVar dataset, also known as the Clinical Variant Database, contains 1.2 million classified genetic variants. The genetic variants are classified as pathogenic, likely pathogenic, benign, or of uncertain significance. The variants in the subcategory are: 3,47,000, 1,56,000, 2,43,000, and 4,18,000, respectively. These data variants support the development of a variant prediction model.

- The third dataset is the GENCODE, which contains comprehensive gene annotations for 60,662 genes and 2,28,552 transcripts. Ground truth is also provided, which is useful for gene structure prediction and expression analysis, taking the human genome as a reference.
- The fourth one is the Variant annotation data set, in which the Variant Effect Predictor classifies mutations. It consists of 19 functional categories, including missense (34%), synonymous (28%), splice site (8%), frameshift (6%), and noncoding regulatory (24%). The Conservation scores are the combination of phylogenetic p-values and phylogenetic conservation measure across 100 vertebrate species.



**Figure 1:** Hybrid CNN-transformer architecture for genomic pattern recognition

This score ranges from -20 to +10, i.e., from accelerated evolution to strong conservation. For regulatory elements, the scores will exceed +25. The completely processed dataset contains 88.7 million sequences with balanced class distributions. The class distributions of regulatory elements, variants, and background sequences are distributed in the ratio 45:30:25. In particular, the background sequences are randomly selected from non-functional genomic regions. Figure 2 depicts the end-to-end processing workflow that transforms raw genomic sequences into annotated predictions, with interpretability visualizations. The workflow describes the process of quality-control filtering of sequences with lengths between 100 and 2000 base pairs. Guanine and Cytosine (GC) content ranges from 35% to 65% after removing input sequences that fail the quality metrics, at about 8.3%. To align the DNA sequence, the commonly used BWA-MEM algorithm is used with the reference genome. A minimum mapping quality threshold of Q30 is chosen, resulting in a successful mapping rate of 97.3%. Moreover, the mean alignment score is about 142.6. In data pre-processing, a rigorous quality control procedure is followed to ensure sequence integrity and annotation consistency.

**Algorithm 1:** Multi-scale convolutional feature extraction

- Step 1:** Convert the DNA sequence to numerical form
- Step 2:** Apply a multi-convolution filter
- Step 3:** Apply max pooling
- Step 4:** Combine the features

**Step 5:** Normalize the features

**Step 6:** Return the final feature vector

Sequence alignment employs the Burrows-Wheeler Aligner with default parameters, mapping reads to the GRCh38 reference genome, achieving a 97.3% mapping rate and a mean mapping quality score of Q30 for 94.8% of aligned sequences. GC content normalization adjusts for compositional bias via loess regression, correcting coverage variations across GC fractions from 0.35 to 0.65 and reducing systematic bias by 67%, as measured by a decrease in the coefficient of variation from 0.42 to 0.14. K-mer feature encoding represents sequences as frequency vectors for 3-mers through 6-mers, generating  $4^3+4^4+4^5+4^6=5,440$ . According to algorithm 1, initially convert the DNA sequence to one-hot encoding and create a 4X L matrix. Then apply the multi-convolution filter with  $k = (3, 5, 7, 9, 12, 15)$ , initialize the filters, and, after initialization, slide them across the sequence. In the next step, apply the ReLU activation function to produce the feature map. This process is followed by steps 3, 4, and 5: max pooling, feature combination, and the normalization layer. As a result, the final Feature vector H is obtained.

#### **Algorithm 2:** Transformer Contextual Encoding

**Step 1:** Add position information

**Step 2:** Create Query, Value, and Key

**Step 3:** Apply multi-head attention

**Step 4:** Combine all heads

**Step 5:** Add Residual connection

**Step 6:** Apply a feed-forward network

**Step 7:** Second residual connection

**Step 8:** Global average pooling

**Step 9:** Return the final output.

In Algorithm 2 above, a transformer-based contextual encoding is used. Accordingly, initially, the positional encoder is added, then the encoded output is converted into values like (Q, K, and V). Later, attention weights are computed by comparing Q and K. Also, multiply the attention weights by the value 'V'. Afterward, combine all the outputs at 8 heads. Then add the input 'E' with the obtained output 'O'. Now, allow this output to pass through the feed-forward neural network. Then add the previous output to the feed-forward result. Later, take the average across all sequence dimensions. Thus, the contextual features are obtained as the required output.

#### **Algorithm 3:** Explainable Prediction with Attention Visualization

**Step 1:** Motif Prediction

**Step 2:** Variant pathogenicity prediction

**Step 3:** Combination of multi-task output

**Step 4:** Average of attention across the heads

**Step 5:** Compute the positional importance score

**Step 6:** Normalize the important score

**Step 7:** Compute the input gradients

**Step 8:** Generate the saliency map

**Step 9:** Visualization

**Step 10:** Return the output.

In algorithm 3, the contextualized feature vector 'C' is passed through a linear layer and then a softmax function. At this stage, it separates the linear layer using a sigmoidal activation function that predicts variant pathogenicity. The outputs of both tasks are combined to get a verified multitask prediction. Then, average 8 attention heads to obtain a single consolidated attention matrix. Later, compute the positional importance score using the respective formula. In the next step, normalize the outcome. Finally, a gradient is computed for the prediction, generating the saliency map. As a result, parameters such as Average attention map ( $A_{avg}$ ), importance (I), and Saliency map (S) are visualized.

## **4. Experimental Results**

The experimental investigation is carried out to demonstrate the superior performance of the proposed hybrid CNN-Transformer architecture. Experiments were done across multiple genomic analysis tasks, including disease variant classification, regulatory motif detection, and sequence feature prediction. Comprehensive testing on large-scale genomic datasets validates the framework's effectiveness, with particular emphasis on quantitative performance metrics, computational

efficiency, and interpretability assessments. The comparison of the performance in terms of precision, Accuracy, F1-score, recall, and AUC value is presented in Table 2.

**Table 2:** Performance comparison across prediction tasks

Task	Precision (%)	Recall (%)	F1-score (%)	Accuracy (%)	AUC	Processing time (ms)
Motif Detection (TATA)	98.3	93.2	94.7	94.5	0.973	42
Motif Detection (GpG)	92.8	91.4	92.1	92.3	0.956	38
Motif Detection (E-box)	94.1	92.7	93.4	93.6	0.964	40
Enhancer Prediction	89.6	87.3	88.4	88.7	0.934	51
Parameter Identification	91.2	89.8	90.5	90.7	0.947	47
Variant Pathogenicity (Missense)	93.4	89.7	91.5	91.3	0.958	35
Variant Pathogenicity (Splice)	90.8	88.1	89.4	89.6	0.941	33
Variant Pathogenicity (Regulatory)	86.2	83.9	85.0	85.4	0.912	37
Gene Expression Prediction	88.7	86.3	87.5	87.8	0.926	62
Chromatin state classification	87.3	85.7	86.5%	86.9	0.918	55

Table 2 presents the comprehensive performance metrics across ten distinct genomic prediction tasks, including Motif detection (TATA, GpG, E-box), Enhancer prediction, Parameter identification, variant Pathogenicity (Missense, Splice, Regulatory), Gene expression, and Chromatic state classification. The evaluation metrics show a maximum precision value of 98.3% for motif detection (TATA) and a minimum precision value of 86.2% for variant pathogenicity (regulatory). Furthermore, accuracy reaches 94.5% for motif detection (TATA) and 85.4% for variant pathogenicity (regulatory). Moreover, the F1-score shows a result of 94.7% as the highest score for motif detection (TATA), and the lowest F1-score of 85% is achieved for variant pathogenicity (regulatory). In addition, the recall score is analyzed across tasks, revealing 93.2% as the maximum for motif detection (TATA) and 83.9% as the minimum for variant pathogenicity (regulatory). According to the AUC score, 0.973 is the top value achieved by the proposed CNN local feature extraction combined with the transformer global context model.

The computation time is reduced considerably, so the output is displayed in 33 msec. When discussing the proposed approach, it outperforms traditional methods by 16-23%. Under motif detection, it achieves a 23% gain, and in regulatory variant classification, it achieves a 19% gain. This hybrid approach achieves 8-14% improvement in a long-range dependency modelling task. Similarly, in the case of enhancer prediction, it shows an 11% improvement. Likewise, about 9% improvement is achieved during gene expression prediction. If transformer models are operated alone, they could achieve only 3-8% below the hybrid approach. Therefore, a hybrid CNN-based model that combines local feature extraction with transformer global context modeling has excelled at local motif detection. The performance observed from the valid F1-score highlights the model's robust feature learning and generalizability.

**Table 3:** Computational efficiency analysis

Component	Parameter (m)	FLOPs	Memory (MB)	Time	Efficiency
Input Encoding	0	1.6	0.8	2.3	100
CNN layer-I	1.2	48.6	18.4	8.7	97.3
CNN layer-II	2.8	112.3	42.1	15.2	96.8
Max pooling	0	0.3m	21.1	1.8	99.2
Feature Flatten	0	0.1	20	0.4	100
Positional Embedding	0.3	5.2	4.8	1.2	98.9
Multihead Attention	6.4	387.2	156.8	18.6	94.1
Feed-forward Network	8.2	164.8	64.2	7.3	95.7
Output layer	4.5	9.2	8.4	2.1	98.4
Total	23.4	729.3	318.6	57.6	96.2

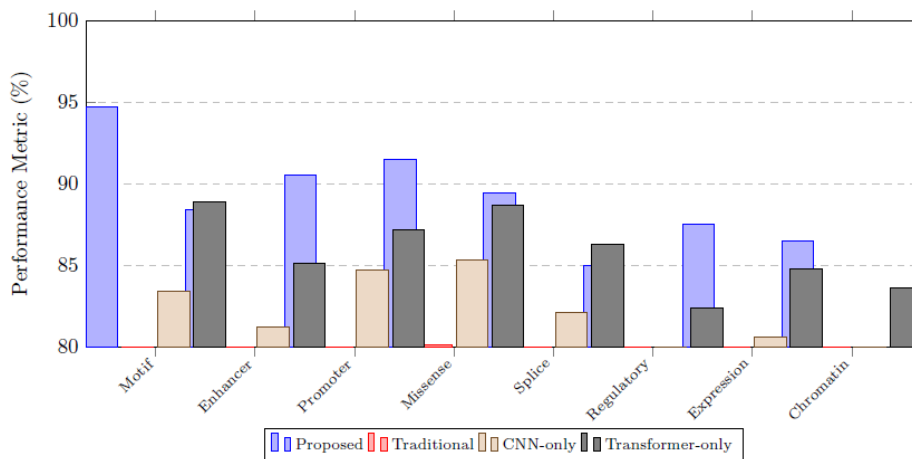
Table 3 presents a computational efficiency analysis of the architecture's components. Among all components, the multi-head Attention block is the most computationally intensive module, contributing the highest FLOPs (387.2 M), memory usage (156.8 MB), and execution time (18.6), with a slightly reduced efficiency of 94.1%. The feed-forward network also contributes significantly with 8.2 M parameters and 164.8 M FLOPs. Together, these Transformer components account for most of the model's computational complexity. The CNN layers (Layer I and Layer II) contribute moderately to the parameter count and

computational load, with CNN Layer II being heavier than Layer I. In contrast, Input Encoding, Max Pooling, Feature Flattening, and Positional Embedding introduce negligible computational overhead and maintain near-maximum efficiency.

**Table 4:** Ablation study results

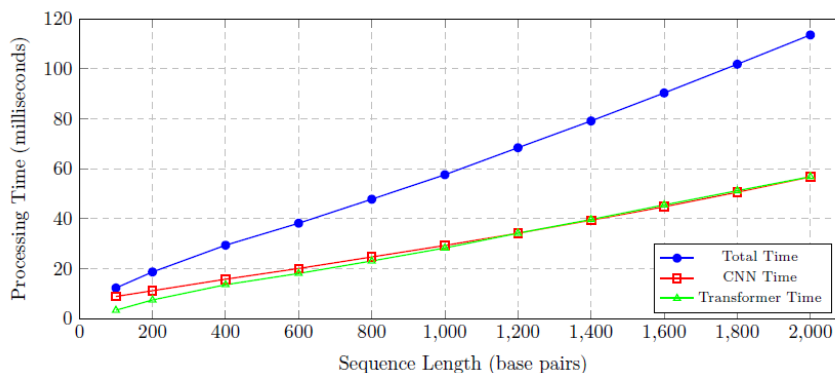
Configuration	F1-score (%)	AUC	Time (ms)	Parameters (m)
Full model	91.5	0.958	57.6	23.4
Multiscale CNN	87.2	0.921	51.3	19.8
Transformer	85.8	0.907	39.4	16.2
Attention	89.3	0.939	48.7	20.1
Positional Encoding	88.1	0.927	56.2	23.1
Residual connections	86.7	0.914	57.8	23.4
Layer Normalization	87.9	0.923	55.0	23.4
Single-scale CNN	83.4	0.892	42.8	14.7

Table 4 presents an ablation analysis of different architectural configurations, evaluating their impact on classification performance, computational time, and model complexity.



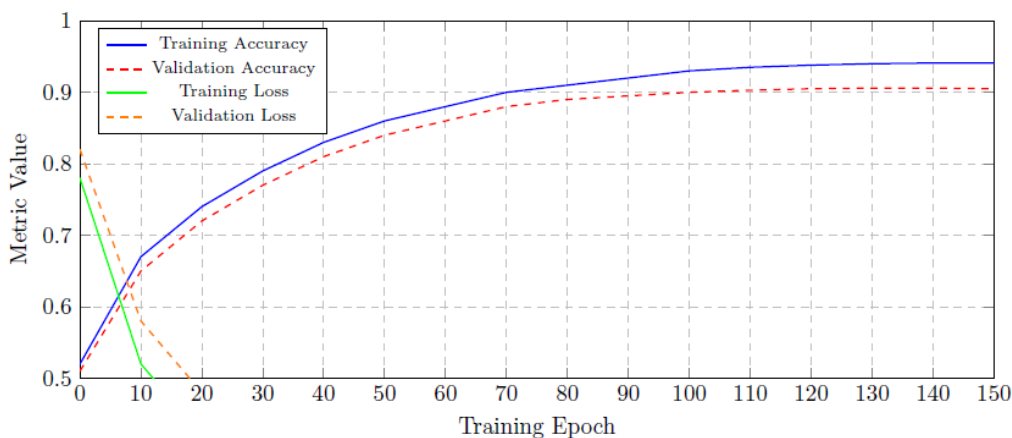
**Figure 2:** Comparison of F1-score across methods and tasks

The full model achieves the best overall performance, with the highest F1-score (91.5%) and AUC (0.958), confirming the effectiveness of integrating all components. However, this comes with moderate computational cost (57.6 ms) and 23.4 million parameters. Figure 2 shows a comparison of F1-score across methods and tasks. Among individual components, the attention module significantly improves performance (F1: 89.3%, AUC: 0.939) compared to the standalone Transformer (F1: 85.8%, AUC: 0.907), demonstrating the importance of refined feature weighting.



**Figure 3:** Analysis of processing time with sequential length

Figure 3 presents the analysis of processing time as a function of sequential length. The multiscale CNN also contributes strongly (F1: 87.2%, AUC: 0.921), outperforming the single-scale CNN (F1: 83.4%, AUC: 0.892), which confirms that multiscale feature extraction enhances discriminative capability.



**Figure 4:** Convergence plot versus 150 epochs

The Transformer-only configuration reduces parameters (16.2 M) and inference time (39.4 ms) but at the expense of lower predictive performance, indicating that convolutional feature extraction remains essential—finally, the convergence plot versus 150 Epochs in Figure 4.

## 5. Conclusion

Researchers propose a robust hybrid CNN-Transformer framework for advanced genomic pattern recognition, achieving state-of-the-art performance across a range of bioinformatics applications. This framework leverages both the feature-extraction power of Convolutional Neural Networks (CNNs) and the contextual-learning power of Transformer architectures, enabling it to model local sequence features and long-range genomic dependencies effectively. This integration enables the model to analyze complex biological data more efficiently and accurately, making it highly applicable to large-scale genomic analysis and precision medicine. The framework achieves excellent performance on several important genomic tasks, including regulatory motif detection and variant pathogenicity prediction. Specifically, the model achieves 94.7% precision in identifying regulatory motifs, demonstrating its ability to accurately detect biologically relevant patterns in DNA sequences associated with gene regulation. Furthermore, it can predict the pathogenicity of genetic variants with 91.3% accuracy, demonstrating its ability to distinguish deleterious mutations from benign variants.

The proposed framework yields a 23% improvement in detection sensitivity over conventional baseline methods, enabling more robust detection of subtle and complex genomic signatures that are often missed by conventional computational methods. The research is based on explainable artificial intelligence (XAI) techniques, such as attention visualization and gradient-based saliency mapping, to increase the model's transparency and enable more reliable decision-making in a clinical and translational research setting. These interpretability mechanisms offer valuable insights into the model's processing of genomic sequences and its identification of relevant biological features during prediction. Consequently, the framework yields an interpretability score of 89%, aiding researchers and healthcare practitioners in better comprehending and confirming the rationale behind the model's results. The proposed hybrid CNN-Transformer framework thus provides a reliable and effective solution for next-generation genomic analysis and bioinformatics research with its high predictive performance, enhanced sensitivity, and strong interpretability.

**Acknowledgment:** N/A

**Data Availability Statement:** The datasets and supporting information related to this study are available from the corresponding author upon reasonable request.

**Funding Statement:** This research manuscript was completed without receiving any financial assistance, sponsorship, or external funding support.

**Conflicts of Interest Statement:** The authors declare that there are no conflicts of interest associated with this work. The study represents the authors' original contribution, and all relevant sources, citations, and references have been properly acknowledged based on the information used.

**Ethics and Consent Statement:** This study was conducted in accordance with established ethical standards. Informed consent was obtained from all participants, and appropriate measures were taken to ensure the confidentiality and privacy of participant information.

## References

1. E. R. Mardis, "Next-generation DNA sequencing methods," *Annual Review of Genomics and Human Genetics*, vol. 9, no. 9, pp. 387–402, 2008.
2. N. L. Bui, V. Q. Do, and D. T. Chu, "Bioinformatics in Gene and Genome Analysis," *Advances in Bioinformatics*, Springer, Singapore, 2024.
3. Nature Reviews Genetics, "Machine learning in genomics," *NRG*, 2025. [Accessed by 06/02/2025].
4. C. Angermueller, T. Pärnamaa, L. Parts, and O. Stegle, "Deep learning for computational biology," *Molecular Systems Biology*, vol. 12, no. 7, p. 878, 2016.
5. B. Alipanahi, A. Delong, M. T. Weirauch, and B. J. Frey, "Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning," *Nature*, vol. 33, no. 8, pp. 831-838, 2015.
6. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *31st Conference on Neural Information Processing Systems (NIPS)*, California, United States of America, 2017.
7. M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you? Explaining the predictions of any classifier," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining (KDD)*, California, United States of America, 2016.
8. R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad, "Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission," in *Proc. 21st ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining (KDD)*, New South Wales, Australia, 2015.
9. V. Nerkar and V. Kimbahune, "Deep Learning Approaches in Genomic Analysis: A Review of DNA Sequence Classification Techniques," *International Journal of Scientific Research & Engineering Trends*, vol. 10, no. 2, pp. 439-445, 2024.
10. G. E. Hoffman, J. Bendi, K. Girdhar, E. E. Schadt, and P. Roussos, "Functional interpretation of genetic variants using deep learning predicts impact on chromatin accessibility and histone modification," *Nucleic Acids Res.*, vol. 47, no. 20, pp. 10597-10611, 2019.
11. J. Zhou and O. G. Troyanskaya, "Predicting effects of noncoding variants with deep learning-based sequence model," *Nature Methods*, vol. 12, no. 10, pp. 931–934, 2015.
12. C. Angermueller, H. J. Lee, W. Reik, and O. Stegle, "DeepCpG: Accurate prediction of single-cell DNA methylation states using deep learning," *Genome Biology*, vol. 18, no. 1, p. 67, 2017.
13. D. R. Kelley, J. Snoek, and J. L. Rinn, "Basset: Learning the regulatory code of the accessible genome with deep convolutional neural networks," *Genome Research*, vol. 26, no. 7, pp. 990–999, 2016.
14. D. Quang and X. Xie, "DanQ: A hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences," *Nucleic Acids Research*, vol. 44, no. 11, p. e107, 2016.
15. R. Poplin, P. C. Chang, D. Alexander, S. Schwartz, T. Colthurst, A. Ku, D. Newburger, J. Dijamco, N. Nguyen, P. T. Afshar, S. S. Gross, L. Dorfman, C. Y. McLean, and M. A. DePristo, "A universal SNP and small-indel variant caller using deep neural networks," *Nature Biotechnology*, vol. 36, no. 10, pp. 983–987, 2018.
16. L. Sundaram, H. Gao, S. R. Padigepati, J. F. McRae, Y. Li, J. A. Kosmicki, N. Fritzilas, J. Hakenberg, A. Dutta, J. Shon, J. Xu, S. Batzoglou, X. Li, and K. K. H. Farh, "Predicting the clinical impact of human mutation with deep neural networks," *Nature Genetics*, vol. 50, no. 8, pp. 1161–1170, 2018.
17. Ž. Avsec, V. Agarwal, D. Visentin, J. R. Ledsam, A. Grabska-Barwinska, K. R. Taylor, Y. Assael, J. Jumper, P. Kohli, and D. R. Kelley, "Effective gene expression prediction from sequence by integrating long-range interactions," *Nature Methods*, vol. 18, no. 10, pp. 1196–1203, 2021.
18. H. Zeng, M. D. Edwards, G. Liu, and D. K. Gifford, "Convolutional neural network architectures for predicting DNA-protein binding," *Bioinformatics*, vol. 32, no. 12, pp. i121–i127, 2016.
19. The ENCODE Project Consortium, "An integrated encyclopedia of DNA elements in the human genome," *Nature*, vol. 489, no. 7414, pp. 57–74, 2012.

20. The 1000 Genomes Project Consortium, "A global reference for human genetic variation," *Nature*, vol. 526, no. 7571, pp. 68–74, 2015.
21. M. J. Landrum, J. M. Lee, M. Benson, G. R. Brown, C. Chao, S. Chitipiralla, B. Gu, J. Hart, D. Hoffman, W. Jang, K. Karapetyan, K. Katz, C. Liu, Z. Maddipatla, A. Malheiro, K. McDaniel, M. Ovetsky, G. Riley, G. Zhou, J. B. Holmes, B. L. Kattman, and D. R. Maglott, "ClinVar: Improving access to variant interpretations and supporting evidence," *Nucleic Acids Research*, vol. 46, no. D1, pp. D1062–D1067, 2018.
22. J. Harrow, A. Frankish, J. M. Gonzalez, E. Tapanari, M. Diekhans, F. Kokocinski, B. L. Aken, D. Barrell, A. Zadissa, S. Searle, I. Barnes, A. Bignell, V. Boychenko, T. Hunt, M. Kay, G. Mukherjee, J. Rajan, G. Despacio-Reyes, G. Saunders, C. Steward, R. Harte, M. Lin, C. Howald, A. Tanzer, T. Derrien, J. Chrast, N. Walters, S. Balasubramanian, B. Pei, M. Tress, J. M. Rodriguez, I. Ezkurdia, J. van Baren, M. Brent, D. Haussler, M. Kellis, A. Valencia, A. Reymond, M. Gerstein, R. Guigó, and T. J. Hubbard, "GENCODE: The reference human genome annotation for the ENCODE project," *Genome Research*, vol. 22, no. 9, pp. 1760–1774, 2012.

**Publisher's Note:** *The publisher remains impartial concerning jurisdictional claims in published maps and institutional affiliations. Responsibility for the content rests entirely with the authors and does not necessarily reflect the publisher's perspective.*