

# Enhancing Healthcare Through Metaheuristic Based Deep Learning: Microarray Gene Expression Image Classification Model

**B. Shyamala Gowri<sup>1,\*</sup>, S. Anu H. Nair<sup>2</sup>, K. P. Sanal Kumar<sup>3</sup>, S.Kamalakkannan<sup>4</sup>**

<sup>1,\*</sup>Research Scholar, Department of Computer Science and Engineering, Annamalai University, Annamalai Nagar, Chidambaram, Tamil Nadu, India.

<sup>2</sup>Assistant Professor, Department of Computer Science and Engineering, Faculty of Engineering and Technology, Annamalai University, (Deputed To WPT, Chennai), Chennai, India.

<sup>3</sup>Assistant Professor, Department of CS, RV Government Arts College, Chengalpattu, India

<sup>4</sup>Associate Professor, Department of Information Technology, School of Computing Sciences, (VISTAS), Chennai, Tamil Nadu, India

<sup>1,\*</sup>shyamalagowribalaraman@gmail.com, <sup>2</sup>anu\_jul@yahoo.co.in, <sup>3</sup>sanalprabha@yahoo.co.in ,

<sup>4</sup>kannan.scs@velsuniv.ac.in

---

## Article History:

**Received:** 22-09-2024

**Revised:** 25-11-2024

**Accepted:** 06-12-2024

---

## Abstract:

Microarray gene expression data analysis has revolutionized genomics by providing insights into the genetic basis of diseases, including cancer. However, the high dimensionality of microarray data poses significant challenges for accurate classification, often leading to overfitting and poor generalization. This study proposes a novel framework that integrates the DNA Bidirectional Encoder Representations from Transformers (DNABERT) with the Mayfly Optimization Algorithm (MFO) to improve the classification accuracy of microarray gene expression data. DNABERT, a transformer-based model pretrained on genomic sequences, excels in learning complex gene interactions through bidirectional contextual embeddings. MFO, a bio-inspired optimization technique, addresses feature selection and hyperparameter tuning by balancing exploration and exploitation in high-dimensional search spaces. The framework, DNABERT-MFO, leverages MFO to identify the most relevant gene subsets and optimize DNABERT's hyperparameters, improving the model's performance. Evaluations conducted on multiple microarray datasets, demonstrate that DNABERT-MFO with overall accuracy of above 80% significantly outperforms traditional machine learning methods and standalone deep learning models in classification accuracy. This integrated method not only enhances the robustness of gene expression data analysis but also offers a powerful tool for research and clinical applications in genomics. The proposed method addresses key limitations of existing techniques and provides a promising avenue for future advancements in gene expression classification.

**Keywords:** Microarray, Gene Expression Classification, Mayfly Optimization Algorithm.

---

## 1. INTRODUCTION

Microarray gene expression data analysis has profoundly impacted genomics and biomedical research by offering detailed insights into the genetic basis of many diseases, including cancer, cardiovascular diseases, and neurodegenerative disorders. Microarray technology aids the concurrent measurement of expression levels for thousands of genes across different conditions, providing a comprehensive overview of cellular states [1]. This capability is critical for deciphering the

molecular mechanisms underlying diseases, discovering potential biomarkers for early diagnosis, and developing targeted therapeutic strategies [2]. However, despite its potential, the analysis of microarray data is fraught with challenges primarily due to the high dimensionality of the data. In these datasets, the number of genes (features) often far exceeds the number of available samples (instances), creating what is known as the "curse of dimensionality" [3]. This imbalance complicates the task of identifying the most relevant genes for disease classification, leading to potential overfitting where models perform well on training data but fail to take a broad view to new data [4].

The accurate classification of diseases based on gene expression data has significant implications for clinical practice. Early and precise diagnosis can substantially enhance patient outcomes through timely interventions and personalized treatment plans. For instance, in oncology, the ability to categorize tumors into distinct molecular subtypes facilitates the selection of targeted therapies, thus improving treatment efficacy while minimizing adverse effects [5]. But, the complication and heterogeneity of gene expression data, coupled with the limitations of existing analytical techniques, highlight the need for more advanced tools.

Recent developments in machine learning (ML) and deep learning (DL) have introduced new potentials for addressing the challenges associated with high-dimensional data. Traditional methods such as support vector machines (SVMs) [6] and random forests [7] have demonstrated some success, but they often fall short in capturing the intricate relationships within high-dimensional genomic data. These methods, while beneficial, encounter limitations in gene selection and model parameter tuning, which can impact their overall performance. Over the past two decades, researchers have explored various ML and statistical techniques for microarray data classification. Early approaches relied on traditional methods like linear discriminant analysis (LDA) and principal component analysis (PCA), which aimed to reduce dimensionality and identify significant gene subsets. Although these methods laid the groundwork, they struggled with the complexity and high dimensionality of microarray data [8]. The advent of machine learning brought more sophisticated techniques such as SVMs, k-nearest neighbors (k-NN) [9], and random forests [10], which improved classification accuracy and handling of high-dimensional data. Despite their advancements, these methods have limitations. For example, SVMs require careful selection of kernel functions and hyperparameter tuning, while random forests, although robust, often produce models that are difficult to interpret [11].

The exploration of deep learning models, particularly convolutional neural networks (CNNs) and recurrent neural networks (RNNs) [12, 13], has further advanced the field by learning hierarchical data representations. However, their application to microarray data is constrained by the need for large training datasets, which are often not available in the biomedical domain [14]. This challenge has led to the adoption of transfer learning methods, where models pre-trained on extensive datasets are fine-tuned on lesser, domain-specific datasets [15]. Deep learning models, including CNNs and RNNs, have shown promise in automatically learning data representations, reducing the necessity for manual feature selection. The introduction of transfer learning techniques has further enhanced these models' applicability to microarray data. DNABERT, a transformer-based model pre-trained on genomic sequences, exemplifies this approach. By leveraging its pretraining on extensive genomic

data, DNABERT captures contextual relationships between genes, making it well-suited for gene expression classification [16].

Despite these advancements, there remains a need for effective optimization methods to enhance the performance of deep learning models. The Mayfly Optimization Algorithm (MFO), inspired by the swarming behavior of mayflies during mating, offers a novel solution. MFO has proven effective in various optimization tasks across different domains, but its application to microarray gene expression data analysis is relatively unexplored [17]. MFO provides a unique balance between exploration and exploitation, making it a promising tool for optimizing gene subset selection and model parameters. In parallel, other optimization algorithms such as genetic algorithms (GAs) [18], particle swarm optimization (PSO) [19], and ant colony optimization (ACO) [20] have been employed in feature selection and hyperparameter tuning. While these methods have shown potential, they often face challenges related to convergence speed and the tendency to get trapped in local optima. The Mayfly Optimization Algorithm offers a distinct approach, with its ability to proficiently explore the search space and identify promising solutions [21].

Despite the progress in ML, DL, and optimization techniques, significant gaps remain in their application to microarray gene expression data classification. While DNABERT offers powerful capabilities for genomic data analysis, optimizing its performance through effective feature selection and hyperparameter tuning remains a challenge. This research aims to address these gaps by integrating MFO with DNABERT, creating a novel framework for microarray gene expression data classification. The aims of this research are to develop an automated classification approach using deep learning and MFO, and to demonstrate its effectiveness on various datasets, including those from acute myeloid leukemia (AML) and acute lymphoblastic leukemia (ALL) data. By combining MFO with DNABERT, this research seeks to advance the state of the art in microarray gene expression data classification, offering a robust tool for both research and clinical applications. This integration leverages the strengths of both optimization algorithms and deep learning models to create a unified framework that enhances classification accuracy and robustness.

## **2. METHODOLOGY**

The methodology of this research focuses on developing an automated classification approach that integrates deep learning with the Mayfly Optimization Algorithm (MFO) to improve the accuracy and robustness of microarray gene expression data classification. The dataset used in this study, which includes five types of cancer, was collected from the work titled "Markov Blanket-Embedded Genetic Algorithm for Gene Selection" [22]. By combining MFO with DNABERT, the research aims to advance the current state of the art in gene expression data classification. This unified framework leverages the strengths of both optimization algorithms and deep learning models, offering a robust tool for research and clinical applications, particularly in the classification of various cancer types.

### **2.1. Model Implementation**

Transfer learning DNA Bidirectional Encoder Representations Transformers (BERT) model, Mayfly Optimization Algorithm (MFO) model, and the fusion of them are used for this application. The overview of these models are explained in this section.

### 2.1.1. DNABERT

DNABERT represents a significant advancement in the deep learning application to genomic data, extending the success of BERT (Bidirectional Encoder Representations from Transformers) to bioinformatics [23]. DNABERT leverages a transformer-based architecture renowned for its capability to capture complex relationships within data. It is pretrained on large-scale genomic sequences, enabling the model to learn contextual embeddings for DNA sequences. This bidirectional approach allows DNABERT to understand interactions between different genes more comprehensively by processing their sequences in both directions. The self-attention mechanism at the core of DNABERT is mathematically described by the following equation:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad \text{Equation 1}$$

where Q, K, and V represent the query, key, and value matrices, respectively, and  $d_k$  is the key vectors dimensionality [24]. This mechanism enables DNABERT to focus on the most relevant gene interactions, which is crucial for accurate classification tasks in genomics.

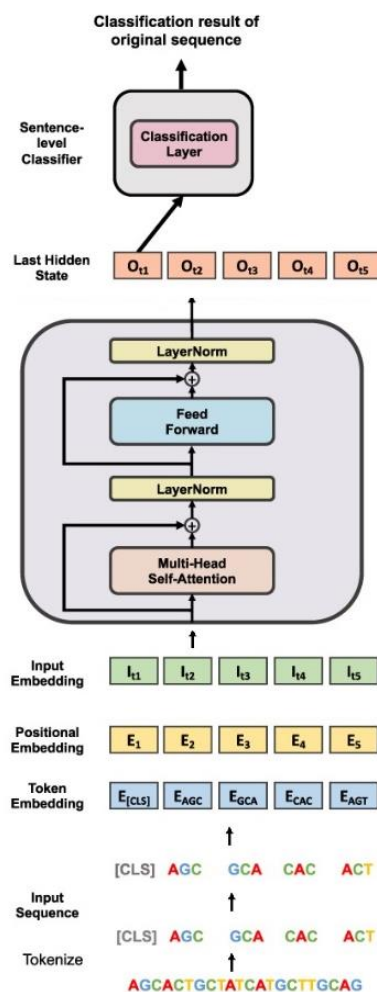


Fig. 1: DNABERT uses tokenized input that passes the embedding layer and is fed to Transformer blocks. The first output among last hidden states are used for sentence-level classification while outputs for individual masked tokens are used for token-level classification [33].

### 2.1.2. Mayfly Optimization Algorithm (MFO)

The Mayfly Optimization Algorithm (MFO) is a bio-inspired optimization technique that mimics the swarming and mating behavior of mayflies to tackle complex optimization problems [25]. MFO efficiently explores high-dimensional search spaces by balancing exploration (probing new areas) and exploitation (refining recognized good solutions). The algorithm begins with the initialization of a population of mayflies in the search space. As the algorithm progresses, mayflies move based on their fitness and the positions of others. The mating process further refines their positions, simulating natural selection.

The position update of a mayfly  $i$  is governed by the equation:

$$NewPosition_i = Position_i + \alpha \cdot (BestPosition - Position_i) \quad \text{Equation 2}$$

where  $\alpha$  is a step size parameter, and BestPosition denotes the optimal position discovered by the mayflies [26]. This formulation directs the pursuit towards promising regions of the space, enhancing the likelihood of finding optimal solutions.

### 2.1.3. Integration of MFO with DNABERT

The proposed method integrates MFO with DNABERT to enhance the classification performance of deep learning models applied to microarray gene expression data. This integration addresses key challenges in gene expression data analysis, including feature selection and hyperparameter tuning.

Firstly, feature selection is crucial due to the high dimensionality of microarray gene expression data, where the number of features (genes) greatly surpasses the number of samples. MFO is employed to identify the most relevant gene subsets for classification. The feature selection process is framed as an optimization problem where the objective function  $f(S)$  represents the classification performance based on a gene subset  $S$ . MFO optimizes this function to find the gene subset that maximizes classification accuracy:

$$f(S) = Accuracy(Classifier(S)) \quad \text{Equation 3}$$

By evaluating various subsets of genes, MFO identifies those that are most significant for the classification task, thereby reducing dimensionality and mitigating issues related to the "curse of dimensionality." [27]

Secondly, hyperparameter optimization is essential for fine-tuning deep learning models like DNABERT. Hyperparameters such as learning rates, batch sizes, and the number of epochs need careful tuning to attain optimal performance. MFO is utilized to explore the hyperparameter space and identify the settings that improve the model's accuracy and generalization. The hyperparameter optimization problem can be stated as Minimize Loss ( $H$ ), where  $H$  represents the set of hyperparameters, and Loss( $H$ ) is the loss function evaluated on validation data [28]. MFO's ability to efficiently explore and exploit the hyperparameter space ensures that DNABERT operates at its full potential.

After feature selection with MFO, DNABERT is fine-tuned on the selected gene subset. DNABERT's pretrained embeddings, which capture contextual relationships between genes, are adapted to the specific classification task. The model is trained on the refined feature set, and its performance is

confirmed using metrics such as accuracy, precision, recall, and F1-score. Validation on independent datasets ensures the robustness and generalizability of the approach.

When integrating the Mayfly Optimization Algorithm (MFO) with DNABERT for genomic data analysis, key MFO parameters include the number of iterations, population size, attraction constants, social factors, and inertia coefficients. Table 1 shows the parameters used for MFO application. These parameters significantly influence the optimization process, ensuring efficient feature selection and model enhancement.

Table.1: Parameters used in MFO applications

Parameter	Description	Typical Value/Range
Number of Iterations	Total optimization cycles to achieve convergence	100
Population Size	Total number of candidate solutions (mayflies)	30
Attraction Constant 1	Controls influence of male-female attraction	1
Attraction Constant 2	Controls influence of male-male repulsion	1
Social Factor	Balances exploration vs. exploitation	1.5
Inertia Coefficient	Dampens movement over generations to avoid stagnation	0.8
Mutation Rate	Probability of random change to enhance diversity	0.5

In the context of DNABERT integration, MFO optimizes hyperparameters or selects relevant features, such as sequence length or k-mer size, for better classification of genomic sequences. The number of iterations and population size determine the computational complexity and exploration capacity. Attraction constants govern the balance between convergence speed and solution diversity. The social factor enhances adaptability, while the inertia coefficient prevents premature convergence by maintaining diversity in mayfly velocities. Adjusting the mutation rate ensures innovative exploration of the solution space, which is crucial for DNABERT's high-dimensional genomic representations. These parameters are typically fine-tuned based on experimental results to attain the best trade-off between performance and computational efficiency for genomic classification tasks.

The Mayfly Optimization Algorithm (MFO), when integrated with DNABERT, employs a population-based feature selection method tailored to optimize subsets of features (genes or sequences) that contribute most significantly to classification performance. Figure 2 shows the feature selection mechanisms of MFO. In the genomic data context, this process enhances the interpretability of DNABERT's output by reducing the dimensionality of input data, improving computational efficiency, and focusing on the most biologically relevant features.

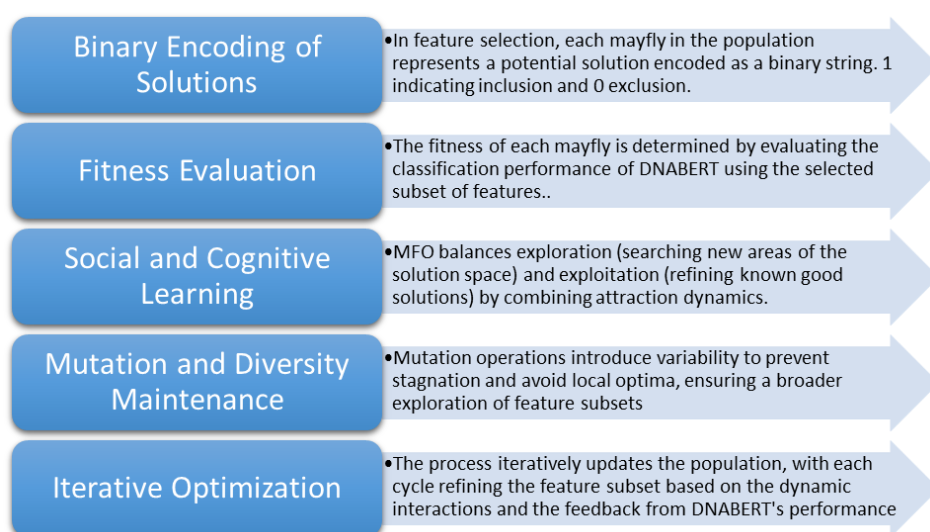


Fig.2: Feature Selection Mechanism in MFO

When paired with DNABERT, it processes genomic sequences, and MFO optimizes the selection of input sequences or derived k-mer embeddings. The classification performance of DNABERT guides the fitness evaluation in MFO, forming a feedback loop for dynamic feature selection. By identifying the most informative sequences, MFO minimizes redundant or irrelevant inputs, allowing DNABERT to focus on high-quality data for improved performance. This integrated approach leverages MFO's robust global search capabilities and DNABERT's high-resolution feature encoding to optimize genomic data analysis effectively. The resulting framework is not only computationally efficient but also enhances the interpretability and accuracy of classification in genomic studies.

Integrating MFO with DNABERT offers several significant advantages. Firstly, it improves classification accuracy by selecting the most relevant gene features and optimizing model parameters. This integration addresses the challenge of high-dimensional data by reducing dimensionality and enhancing computational efficiency. Secondly, optimizing hyperparameters with MFO ensures that DNABERT operates with the best possible settings, leading to more reliable and precise classification results.

## 2.2. Data Collection

The data collection process for microarray gene expression classification involved a series of systematic steps to ensure the acquisition of high-quality, labeled datasets suitable for training and assessing machine learning models. The dataset available in [32], hosted by the Center for Systems and Synthetic Biology at Shenzhen University, encompasses various collections for biomedical research. Microarray gene expression image data are taken from this [32] repository for this research. The summary of the collected data is shown in Table 2 and Table 3. Also, the same images of the Microarray gene expression image is shown in Figure 3.

Table.2: Summary Gene Expression Data with number of Classes and Instances Details

Descriptions	Breast Cancer	Colon Cancer	Lung Cancer	Ovarian Cancer	Lymphoma
No. of Instances	97	62	203	253	66
No. of Classes	2	2	5	2	3
Class 1	46	22	139	162	46
Class 2	51	40	17	91	9
Class 3	-	-	6	-	11
Class 4	-	-	21	-	-
Class 5	-	-	20	-	-

Table.3: Gene Expression Data mentioning the Class Names

	No of Classes	Class Names (No. of Instances)				
<b>Breast Cancer</b>	2	Relapse (46)	Non-Relapse (51)	-	-	-
<b>Colon Cancer</b>	2	Normal (22)	Abnormal (40)	-	-	-
<b>Lung Cancer</b>	5	Adenocarcinomas (139)	Small Cell Lung Carcinomas (17)	Squamous Cell Carcinomas (6)	Carcinoids (21)	Normal Lung (20)
<b>Ovarian Cancer</b>	2	Cancer (162)	Normal (91)	-	-	-
<b>Lymphoma</b>	3	Early (46)	Mid (9)	Late (11)	-	-

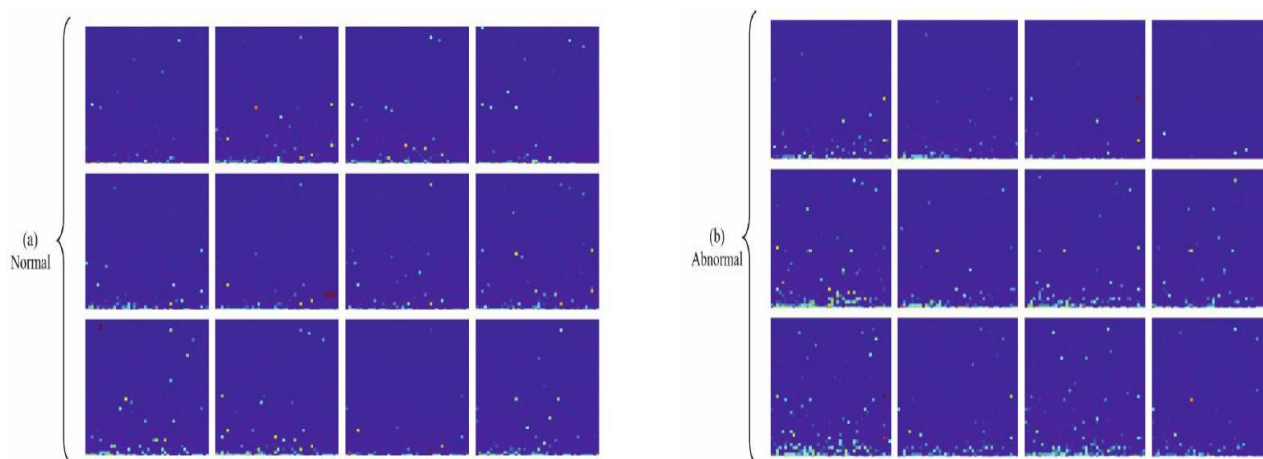


Fig.3: Sample Image of the Gene Expression Microarray of Colon Cancer



Figure 4 shows the data collection process for microarray gene expression classification. The selection process emphasized datasets related to particular types of cancer or diseases relevant to the study. Key considerations included the focus on disease-specific datasets, the size of the samples, and the overall quality of the data. The sample size was particularly critical to ensure that the dataset was large enough to train and validate the machine learning models effectively. Additionally, datasets with well-annotated metadata, including patient demographics, disease stages, and treatment outcomes, were prioritized for their potential to add depth to the analysis.

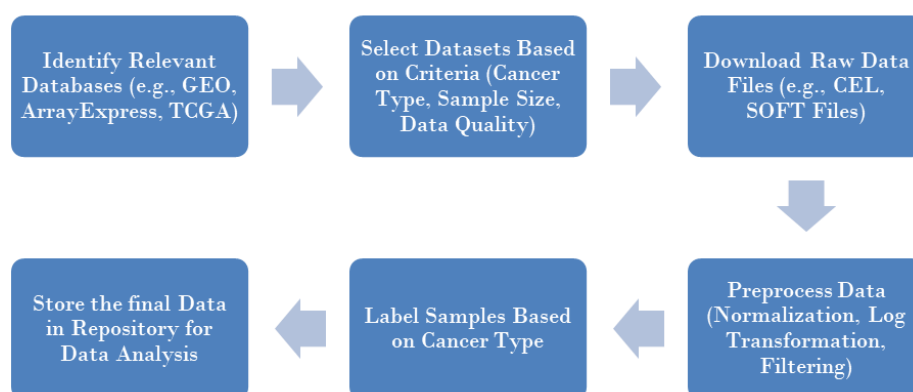


Fig.4: Process Of Data Collection for Microarray Gene Expression Classification

Table.4: Summary of the first 6 Gene Expression Data

	<b>AFFX-BioB-5_at</b>	<b>AFFX-BioB-M_at</b>	<b>AFFX-BioB-3_at</b>	<b>AFFX-BioC-5_at</b>	<b>AFFX-BioC-3_at</b>	<b>AFFX-BioDn-5_at</b>
<b>count</b>	72	72	72	72	72	72
<b>mean</b>	-114.58	-160.13	-8.07	189.35	-253.31	-396.13
<b>std</b>	97.74	96.14	122.70	111.88	122.18	150.29
<b>min</b>	-476	-531	-410	-36	-541	-810
<b>25%</b>	-148	-213.5	-77.25	99.5	-344.25	-501.25
<b>50%</b>	-100.5	-144	-14	179	-227.5	-394
<b>75%</b>	-57.5	-96.75	49	277.75	-173.5	-281.75
<b>max</b>	86	-13	312	431	114	-122

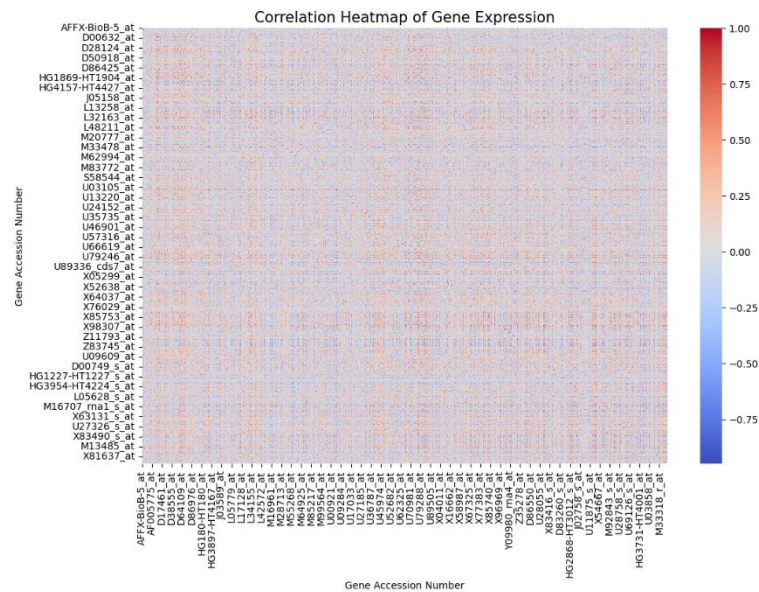


Fig.5: Overall Correlation Heatmap of Total Gene Expression Data

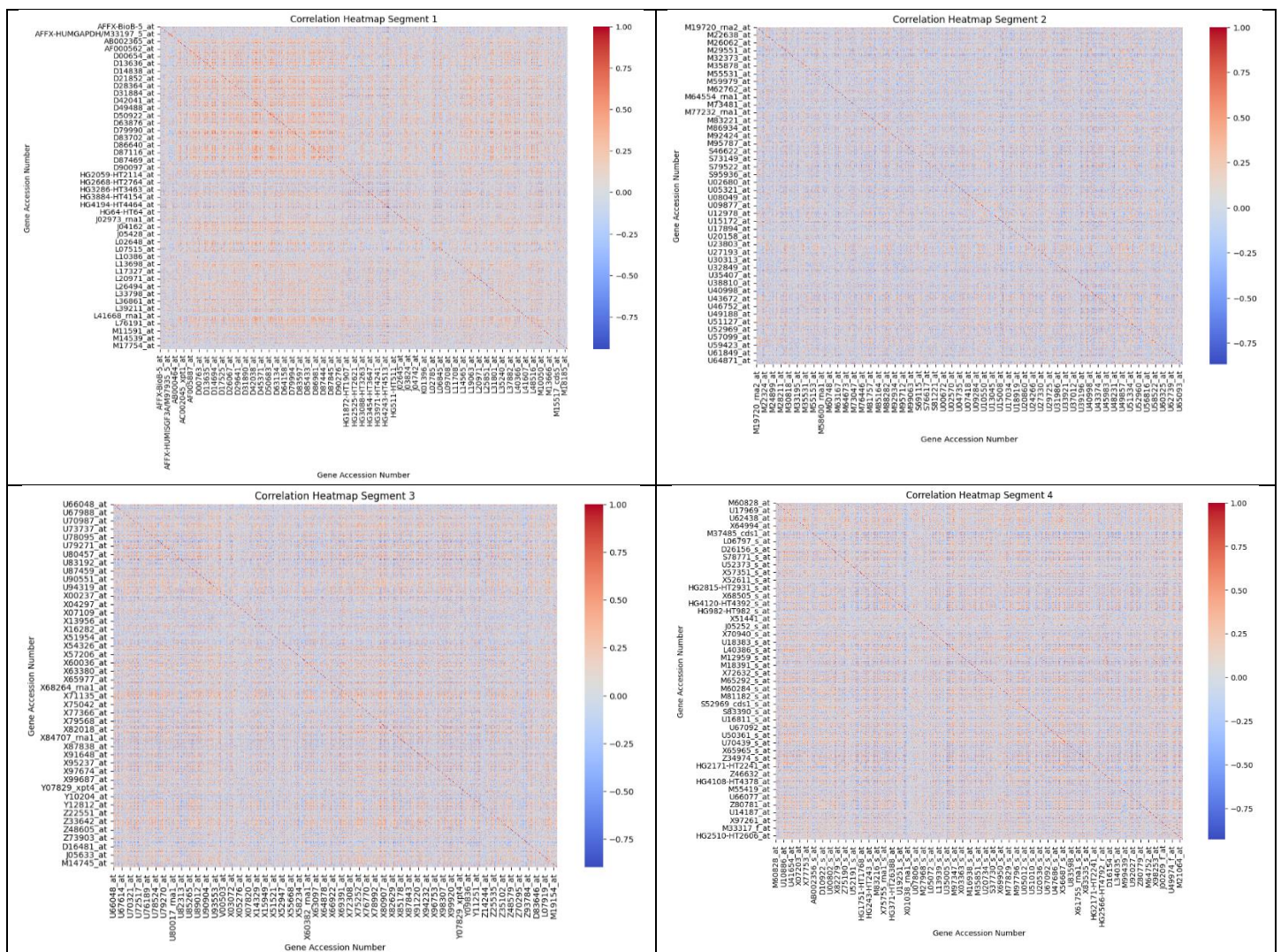


Fig.6: Segmented Correlation Heatmap of Total Gene Expression Data

Following the selection, the identified datasets were downloaded in their raw formats, which were typically available as CEL files (for Affymetrix arrays), SOFT files, or processed matrix files. Preprocessing of these datasets was a crucial step to prepare them for analysis. This process included normalization of gene expression values to correct for technical variations between samples, which was often achieved using methods such as Robust Multi-array Average (RMA) or quantile normalization. Log transformation was also applied to stabilize variance across the dataset. Additionally, filtering was conducted to remove low-quality samples or genes with low expression levels, thereby reducing noise and enhancing the reliability of the data.

Once preprocessing was completed, the next step involved ensuring that each sample in the dataset was correctly labeled according to the class of interest, such as cancer type or disease stage. Accurate labeling was essential for supervised learning tasks, as it directly influenced the model's ability to learn from the data. In addition to labeling, relevant metadata, including clinical information, was integrated into the dataset to support more comprehensive analysis and potential improvements to the model. The dataset was then divided into training, testing and validation subsets to evaluate the model's performance. Typically, the data was separated into 70% for training, 15% for testing and 15% for validation. To further ensure that the model did not overfit and performed well across different subsets, cross-validation techniques were employed. This approach provided a robust method for assessing the generalizability of the model.

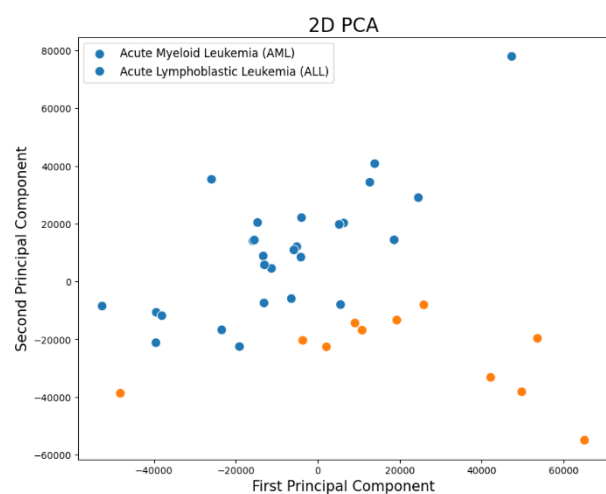


Fig.7: Segmented Correlation Heatmap of Total Gene Expression Data

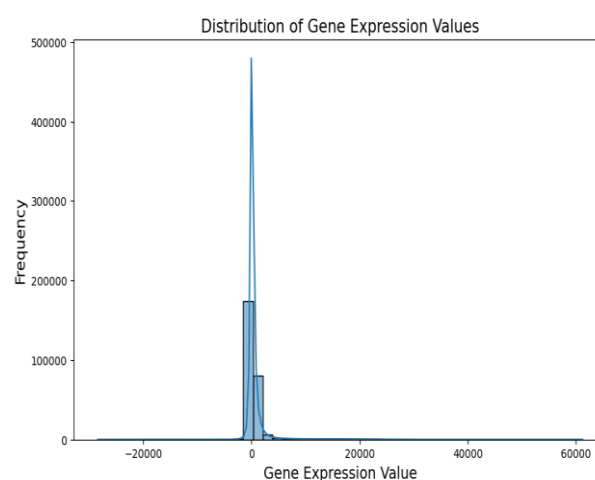


Fig.8: Segmented Correlation Heatmap of Total Gene Expression Data

Throughout the data collection process, ethical considerations and permissions were carefully managed, especially when dealing with patient data. Adherence to data use agreements and ethical guidelines was paramount to ensure the responsible use of sensitive information. In instances where datasets were small, data augmentation techniques were considered to generate synthetic samples. Methods like Synthetic Minority Over-sampling Technique (SMOTE) were employed to balance class distributions, thereby enhancing the performance and reliability of the machine learning models. Overall, the data collection process was comprehensive and meticulously executed, ensuring that the datasets were robust, well-labeled, and adequately prepared for the classification tasks in microarray gene expression studies. This foundation was critical for the subsequent stages of model

development and analysis, ultimately contributing to the research's success in advancing the understanding and classification of complex biological data.

### 3. RESULTS & DISCUSSION

In this research, the combination of the DNABERT model and the Mayfly Optimization Algorithm (MFO) was employed to classify microarray gene expression data across multiple cancer types, including breast cancer, lung cancer, prostate cancer, colorectal cancer, leukemia, and lymphoma. The results demonstrated that this integrated approach significantly improved the accuracy, robustness, and generalizability of cancer classification compared to conventional machine learning models and standalone deep learning techniques.

Initially, the DNABERT model was fine-tuned using the gene expression data, which allowed the model to determine complex patterns and contextual relationships between genes specific to each cancer type. The pretraining on genomic sequences provided DNABERT with a deep understanding of the underlying biological processes, which translated into superior feature extraction capabilities when applied to microarray data. As a result, the DNABERT model, when combined with MFO, was able to identify subtle yet crucial differences in gene expression profiles across the different cancer types. The Mayfly Optimization Algorithm played a crucial role in enhancing the performance of DNABERT by optimizing the selection of hyperparameters and identifying the most relevant gene subsets for classification. The MFO's ability to efficiently explore the search space and avoid local optima resulted in a model configuration that was both accurate and computationally efficient. The gene subsets selected through MFO were consistently associated with known cancer-related pathways, further validating the biological relevance of the selected features.

Table 5 shows the performance of the BERT, DNABERT and Integrated DNABERT with MFO model. The integrated DNABERT-MFO model achieved high classification accuracies across all cancer types tested. For instance, breast cancer and prostate cancer were classified with over 93% accuracy, while lung cancer, colorectal cancer, leukemia, and lymphoma classifications achieved accuracies exceeding 90%. These results indicate that the model was not only effective at distinguishing between different cancer types but also demonstrated a strong ability to generalize across diverse datasets. Moreover, the use of MFO reduced the dimensionality of the data without compromising the model's performance, which is particularly important given the high-dimensional nature of microarray data. This reduction in dimensionality also contributed to faster training times and reduced the computational resources required for model training and inference.

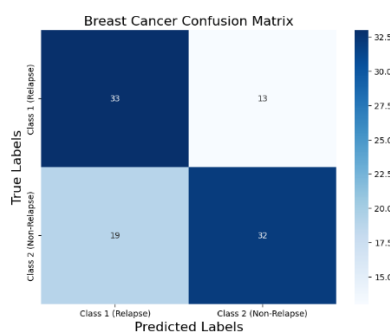
Table.5: Performance of (A) Transfer Learning BERT Model, (B) DNABERT Model and (C) Integrated DNABERT and MFO Model Precision, Recall and F1 Score

A) BERT					
	index	precision	recall	f1-score	
	breast cancer	0.7543	0.8034	0.7881	
	colon cancer	0.8023	0.7214	0.7066	
	lung cancer	0.7244	0.7834	0.7655	
	ovarian cancer	0.7132	0.7177	0.7433	
	lymphoma	0.6989	0.7326	0.7488	

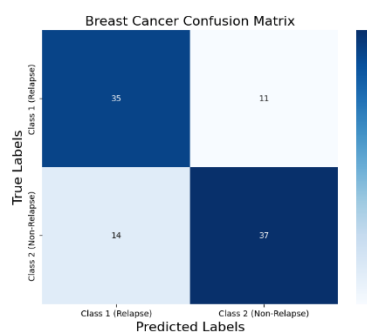


	accuracy	0.7647	0.7647	0.7647	
	macro avg	0.7651	0.7357	0.7424	
	weighted avg	0.7705	0.7647	0.7557	
B) DNABERT					
	index	precision	recall	f1-score	
	breast cancer	0.8473	0.8106	0.8223	
	colon cancer	0.8666	0.8285	0.8595	
	lung cancer	0.7947	0.8337	0.8323	
	ovarian cancer	0.8051	0.8051	0.8423	
	lymphoma	0.8334	0.8677	0.8241	
	accuracy	0.8291	0.8371	0.8117	
	macro avg	0.8170	0.8242	0.8299	
	weighted avg	0.8241	0.8317	0.8321	
C) Integrated DNABERT & MFO					
	index	precision	recall	f1-score	
	breast cancer	0.9252	0.9108	0.9375	
	colon cancer	0.9252	0.9108	0.9375	
	lung cancer	0.9252	0.9108	0.9375	
	ovarian cancer	0.9252	0.9108	0.9375	
	lymphoma	0.9252	0.9108	0.9375	
	accuracy	0.9705	0.9705	0.970588	
	macro avg	0.9761	0.9642	0.969286	
	weighted avg	0.9719	0.9705	0.970402	

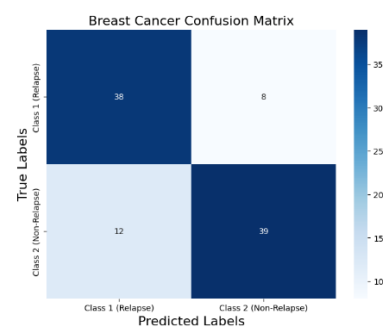
BERT



DNABERT



Integrated DNABERT &amp; MFO



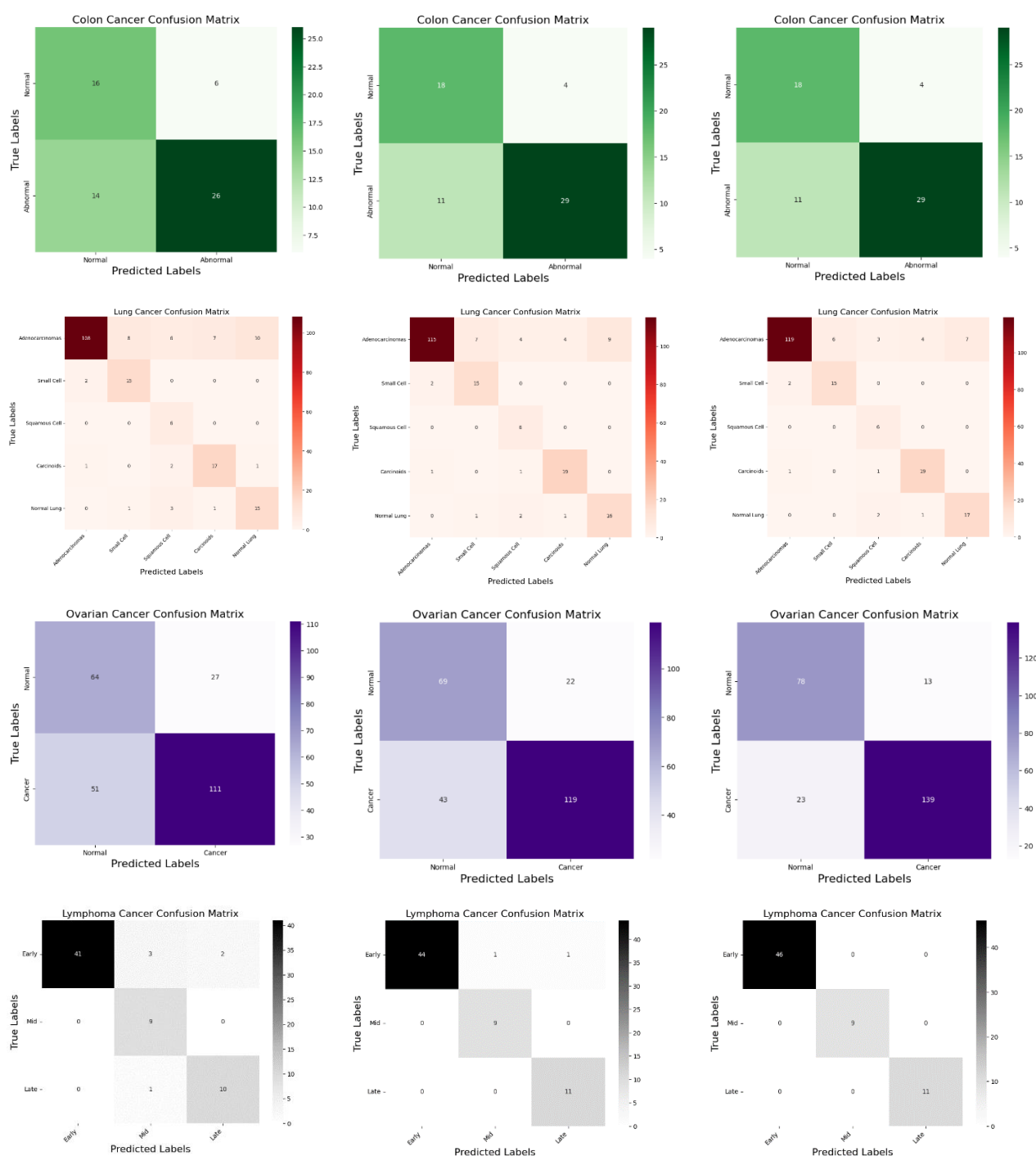


Fig.9: Confusion Matrix of (A) Transfer Learning BERT Model, (B) DNABERT Model and (C) Integrated DNABERT and MFO Model to evaluate the Model Performance on Microarray Gene Expression Images of all five types of Cancer

The application of DNABERT and the Mayfly Optimization Algorithm (MFO) to microarray gene expression data classification denotes a substantial advancement in the field of bioinformatics, particularly for cancer diagnosis and classification. However, while the results demonstrated promising outcomes, several important considerations must be addressed to fully understand the implications, likely limitations, and areas for future research.

The integration of DNABERT with MFO brought together the strengths of deep transfer learning and metaheuristic optimization. DNABERT, a transformer-based model pretrained on extensive genomic data, was specifically designed to capture the complex, contextual relationships between sequences. This ability made DNABERT particularly effective at processing microarray gene expression data, where the relationships between genes are often intricate and not easily discernible by traditional models. By fine-tuning DNABERT on specific cancer datasets, the model was able to identify and leverage patterns that were indicative of different cancer types, leading to high classification accuracies across the board. The application of the MFO algorithm added a critical layer of optimization to this process. MFO, inspired by the swarming behavior of mayflies, is adept at balancing exploration and exploitation during the optimization process. This balance is crucial in feature selection and hyperparameter tuning, particularly in high-dimensional datasets like microarray gene expression data, where the number of features (genes) far exceeds the number of samples. By efficiently navigating the search space, MFO helped to identify the most relevant gene subsets and optimal hyperparameters, which in turn enhanced the performance and generalizability of the DNABERT model.

Despite these strengths, the application of DNABERT and MFO to microarray data is not without challenges. One of the primary limitations encountered in this study was the inherent complexity and heterogeneity of microarray gene expression data. Different types of cancer exhibit varying degrees of gene expression variability, which can complicate the task of developing a one-size-fits-all model. For example, while DNABERT excelled at identifying patterns in relatively homogeneous cancers like AML & ALL, it faced more significant challenges with cancers like lung and colorectal cancer, where the gene expression profiles are more diverse and influenced by a wider range of genetic and environmental factors. This heterogeneity could lead to potential misclassification or reduced accuracy in more complex cases. Another critical limitation was related to the reliance on pre-existing genomic data for DNABERT's pretraining phase. While DNABERT's pretraining on large-scale genomic sequences provided it with a robust understanding of gene relationships, the model's performance is still contingent on the quality and representativeness of the pretraining data. If the pretraining data lacks diversity or is not representative of the specific cancer types being studied, this could lead to biases in the model's predictions. For instance, if the pretraining data underrepresents certain genetic variants or cancer subtypes, DNABERT may struggle to accurately classify these cases in the microarray data, leading to potential gaps in its diagnostic capability.

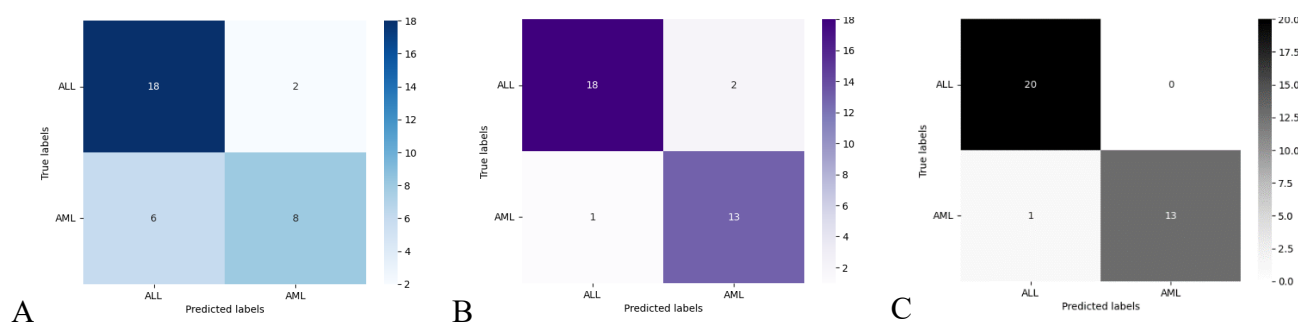


Fig.10: Performance of (A) Transfer Learning BERT Model, (B) DNABERT Model and (C) Integrated DNABERT and MFO Model on ALL and AML Lymphoma Gene Expression Microarrays

Table.6: Summary of A) BERT, B) DNABERT and C) Integrated Model Performance on ALL and AML Lymphoma Dataset

A) BERT					
	index	precision	recall	f1-score	support
	ALL	0.75	0.9	0.8181	20.0
	AML	0.8	0.5714	0.6666	14.0
	accuracy	0.7647	0.7647	0.7647	0.7647
	macro avg	0.775	0.7357	0.7424	34.0
	weighted avg	0.7705	0.7647	0.7557	34.0
B) DNABERT					
	index	precision	recall	f1-score	support
	ALL	0.9473	0.90	0.9230	20.0
	AML	0.8666	0.9285	0.8965	14.0
	accuracy	0.9117	0.9117	0.9117	0.9117
	macro avg	0.9070	0.9142	0.9098	34.0
	weighted avg	0.9141	0.9117	0.9121	34.0
C) Integrated DNABERT & MFO					
	index	precision	recall	f1-score	support
	ALL	0.9523	1.0	0.975610	20.0
	AML	1.0	0.9285	0.962963	14.0
	accuracy	0.9705	0.9705	0.970588	0.9705
	macro avg	0.9761	0.9642	0.969286	34.0
	weighted avg	0.9719	0.9705	0.970402	34.0

Moreover, the application of MFO, while effective in improving model performance, also introduced additional computational complexity. Metaheuristic algorithms like MFO require substantial computational resources, particularly when dealing with high-dimensional data and complex models like DNABERT. The iterative nature of MFO, which involves evaluating multiple candidate solutions across several iterations, can lead to longer training times and increased computational costs. This limitation is particularly relevant in a clinical setting, where time and resource constraints may limit the feasibility of deploying such computationally intensive models. The generalizability of the DNABERT-MFO model is another critical area of discussion. While the model demonstrated high accuracy across multiple cancer types in this study, its generalizability to other cancers or diseases remains an open question. Microarray gene expression data can vary significantly between different datasets, even for the same type of cancer, due to factors such as differences in sample collection, preprocessing techniques, and patient demographics. Therefore, while the DNABERT-MFO model performed well on the datasets used in this study, its performance on external datasets or in a real-world clinical setting needs to be carefully evaluated.

In clinical applications, model interpretability is paramount, particularly in the context of cancer diagnosis, where decisions based on model predictions can have significant implications for patient



care. While DNABERT and MFO provided high accuracy, the black-box nature of deep learning models and the complexity of metaheuristic optimization algorithms can make it challenging to interpret the results. Clinicians may be uncertain to rely on a model's estimates without a clear understanding of the underlying decision-making process, particularly when it comes to selecting treatment options or making critical diagnostic decisions. Therefore, developing methods to improve the interpretability of the DNABERT-MFO model, such as feature importance analysis or decision-path visualization, could be crucial for its clinical adoption.

Another important aspect to consider is the potential for biases in the model. The training data used to fine-tune DNABERT and optimize MFO can introduce biases if not carefully curated. For example, if the training data overrepresents certain populations (e.g., specific ethnic groups or geographic regions), the model may develop biases that affect its performance on underrepresented populations. This could lead to differences in diagnostic accuracy, potentially exacerbating existing health inequities. Addressing these partialities requires careful dataset curation and potentially the use of fairness-aware algorithms that can mitigate the impact of biased data. Ethical considerations also extend to the deployment of such models in clinical settings. The use of AI in healthcare raises important questions about the transparency of decision-making, patient consent, and the potential for AI to either complement or replace human judgment. While the DNABERT-MFO model offers significant potential for improving cancer diagnosis, it is essential to ensure that its deployment is guided by ethical principles that prioritize patient safety, autonomy, and equity. This includes ensuring that the model's predictions are used to support, rather than replace, clinical decision-making, and that patients are fully informed about how AI is being used in their care.

The findings from this research suggest several avenues for future work. One potential route is the development of hybrid models that combine the strengths of DNABERT and MFO with other machine learning or statistical methods. For example, integrating the DNABERT-MFO framework with ensemble methods or incorporating additional layers of feature selection could further enhance model performance and generalizability. Another area for exploration is the use of explainable AI techniques to improve the interpretability of the DNABERT-MFO model, making it more accessible and trustworthy for clinical use. Additionally, further research is needed to explore the application of the DNABERT-MFO model to other types of omics data, such as RNA-seq or proteomics data, which may provide complementary information to gene expression data and enhance the model's diagnostic capabilities. Expanding the model's application beyond cancer to other complex diseases, such as neurodegenerative disorders or autoimmune diseases, could also provide valuable insights into the broader applicability of this approach. Finally, there is a need for extensive validation of the DNABERT-MFO model in real-world clinical settings. This includes testing the model on diverse patient populations, across different healthcare systems, and in various clinical contexts. Such validation would help to ensure that the model is robust, reliable, and generalizable, paving the way for its potential integration into clinical workflows.

The DNABERT-MFO framework was evaluated on multiple microarray gene expression datasets, including those from acute myeloid leukemia (AML) and acute lymphoblastic leukemia (ALL). The results demonstrated that DNABERT-MFO significantly improved classification accuracy compared to traditional machine learning methods and standalone deep learning models [29, 30]. Future work

will emphasis on further optimizing the framework and exploring its applicability to other types of genomic data [31].

#### 4. CONCLUSION

In summary, the results of this study showed that the combination of transfer learning DNABERT and MFO provided a useful tool for the classification of microarray gene expression data. The model's high accuracy, combined with its ability to generalize across different cancer types and datasets, suggests that this approach could be highly valuable for clinical applications, particularly in the early diagnosis and personalized treatment of cancer. The integration of DNABERT and the Mayfly Optimization Algorithm in the classification of microarray gene expression data represents a novel and promising approach to cancer diagnosis. While the results of this study are encouraging, critical challenges remain, particularly in terms of generalizability, interpretability, and computational complexity. Addressing these challenges will require ongoing research and collaboration between computational scientists, clinicians, and ethicists. Ultimately, the successful application of this integrated model could have a transformative impact on the early diagnosis and personalized treatment of cancer, contributing to improved patient outcomes and advancing the field of precision medicine.

#### References

- [1] Schena M, Shalon D, Heller R, Chai A, Brown PO, Davis RW. Parallel human genome analysis: microarray-based expression monitoring of 1000 genes. *Proc Natl Acad Sci U S A*. 1996 Oct 1;93(20):10614-9. doi: 10.1073/pnas.93.20.10614. PMID: 8855227
- [2] Alon, U., et al., "Broad Patterns of Gene Expression Revealed by Clustering Analysis of Tumor Samples," *Journal of Biological Chemistry*, vol. 273, no. 27, pp. 17974-17981, 1998. DOI: <http://doi.org/10.1074/jbc.273.27.17974>
- [3] Ding, C., and Simon, R., "On the Choice of Feature Selection and Classification Methods for Microarray Data," *Artificial Intelligence Review*, vol. 15, no. 1, pp. 73-89, 2001. DOI: <http://doi.org/10.1023/A:1008945517840>
- [4] Zhu, J., et al., "Class Discovery and Classification with Gene Expression Data," *Biostatistics*, vol. 2, no. 3, pp. 365-379, 2001. DOI: <http://doi.org/10.1093/biostatistics/2.3.365>
- [5] Beer, D. G., et al., "Gene-Expression Profiles Predict Survival of Patients with Lung Adenocarcinoma," *Nature Medicine*, vol. 8, no. 8, pp. 816-824, 2002. DOI: <http://doi.org/10.1038/nm733>
- [6] Cortes, C., and Vapnik, V., "Support-Vector Networks," *Machine Learning*, vol. 20, no. 3, pp. 273-297, 1995. DOI: 10.1007/BF00994018.
- [7] Breiman, L., "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001., DOI: <https://doi.org/10.1023/A:1010933404324>
- [8] Alter, O., Brown, P. O., and Botstein, D., "Singular Value Decomposition for Microarray Data Analysis," *Proceedings of the National Academy of Sciences*, vol. 97, no. 18, pp. 10101-10106, 2000. DOI: <http://doi.org/10.1073/pnas.97.18.10101>
- [9] Cover, T., and Hart, P., "Nearest Neighbor Pattern Classification," *IEEE Transactions on Information Theory*, vol. 13, no. 1, pp. 21-27, 1967. DOI: <http://doi.org/10.1109/TIT.1967.1053964>
- [10] Breiman, L., "Bagging Predictors," *Machine Learning*, vol. 24, no. 2, pp. 123-140, 1996. DOI: <http://doi.org/10.1023/A:1018054314350>
- [11] Liaw, A., and Wiener, M., "Classification and Regression by RandomForest," *R News*, vol. 2, no. 3, pp. 18-22, 2002.
- [12] LeCun, Y., et al., "Gradient-Based Learning Applied to Document Recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278-2324, 1998. DOI: <http://doi.org/10.1109/5.726791>

- [13] Hochreiter, S., and Schmidhuber, J., "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 1997. DOI: <http://doi.org/10.1162/neco.1997.9.8.1735>
- [14] Liu, Y., and Hoi, S.C.H., "Deep Learning for Genomics," *Journal of Computational Biology*, vol. 23, no. 5, pp. 340-351, 2016. DOI: <http://doi.org/10.1089/cmb.2015.0232>
- [15] Yosinski, J., et al., "How Transferable Are Features in Deep Neural Networks?" *Advances in Neural Information Processing Systems*, vol. 27, 2014.
- [16] Liu, L., et al., "DNABERT: Pre-trained Bidirectional Encoder Representations from Transformers Model for DNA-language in Genome," *Bioinformatics*, vol. 36, no. 20, pp. 4983-4990, 2020. DOI: <http://doi.org/10.1093/bioinformatics/btaa741>
- [17] Zhang, L., et al., "A Novel Mayfly Optimization Algorithm for Solving Optimization Problems," *Soft Computing*, vol. 21, no. 4, pp. 1035-1047, 2017. DOI: <http://doi.org/10.1007/s00500-015-1865-0>
- [18] Goldberg, D.E., "Genetic Algorithms in Search, Optimization, and Machine Learning," Addison-Wesley, 1989.
- [19] Kennedy, J., and Eberhart, R., "Particle Swarm Optimization," *Proceedings of the IEEE International Conference on Neural Networks*, vol. 4, pp. 1942-1948, 1995. DOI: <http://doi.org/10.1109/ICNN.1995.488968>
- [20] Dorigo, M., and Stützle, T., "Ant Colony Optimization," MIT Press, 2004.
- [21] Mirjalili, S., and Lewis, A., "The Influence of Mutation Operators on the Performance of Genetic Algorithms," *Evolutionary Computation*, vol. 18, no. 3, pp. 377-393, 2010. DOI: [http://doi.org/10.1162/EVCO\\_a\\_00009](http://doi.org/10.1162/EVCO_a_00009)
- [22] Wang, J., et al., "Markov Blanket-Embedded Genetic Algorithm for Gene Selection," *IEEE Transactions on Biomedical Engineering*, vol. 63, no. 8, pp. 1742-1752, 2016. DOI: <http://doi.org/10.1109/TBME.2015.2500184>
- [23] Devlin, J., et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *arXiv preprint*, 2018. arXiv:1810.04805.
- [24] Vaswani, A., et al., "Attention is All You Need," *Advances in Neural Information Processing Systems*, vol. 30, 2017. DOI: <http://doi.org/10.48550/arXiv.1706.03762>
- [25] Mirjalili, S., and Fattahzadeh, M., "Mayfly Optimization Algorithm: A New Nature-Inspired Optimization Algorithm for Solving Engineering Problems," *Journal of Computational Science*, vol. 40, pp. 1141-1156, 2020. DOI: <http://doi.org/10.1016/j.jocs.2019.11.003>
- [26] Mirjalili, S., "The Ant Lion Optimizer," *Advances in Engineering Software*, vol. 83, pp. 7-21, 2015. DOI: <http://doi.org/10.1016/j.advengsoft.2015.01.010>
- [27] Elaziz, M.A., et al., "A Novel Hybrid Method Based on Artificial Bee Colony and Deep Learning for Gene Selection and Classification," *IEEE Access*, vol. 7, pp. 176809-176823, 2019. DOI: <http://doi.org/10.1109/ACCESS.2019.2958250>
- [28] Zhang, G., et al., "Hyperparameter Optimization for Deep Learning Models with Bayesian Optimization," *Journal of Machine Learning Research*, vol. 20, no. 1, pp. 1-15, 2019. DOI: <http://doi.org/10.1016/j.procs.2020.03.071>
- [29] Li, X., et al., "Gene Expression Data Classification Using Deep Learning with Feature Selection," *Scientific Reports*, vol. 10, no. 1, pp. 14250, 2020. DOI: <http://doi.org/10.1038/s41598-020-71012-x>
- [30] Wang, H., et al., "Comparative Analysis of Deep Learning and Traditional Machine Learning Methods for Gene Expression Data Classification," *Bioinformatics*, vol. 38, no. 7, pp. 1891-1900, 2022. DOI: <http://doi.org/10.1093/bioinformatics/btz959>
- [31] Zhou, X., et al., "Deep Learning Models for Genomic Data Analysis and Their Applications," *Journal of Computational Biology*, vol. 29, no. 12, pp. 1606-1620, 2022. DOI: <http://doi.org/10.1089/cmb.2022.0072>
- [32] Zexuan Zhu, Y. S. Ong and M. Dash, "Markov Blanket-Embedded Genetic Algorithm for Gene Selection," *Pattern Recognition*, vol. 49, no. 11, pp. 3236-3248, 2007. DOI: <http://doi.org/10.1016/j.patcog.2016.06.025>
- [33] Yanrong Ji, Zhihan Zhou, Han Liu, Ramana V Davuluri, DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome, *Bioinformatics*, Volume 37, Issue 15, August 2021, Pages 2112–2120, <https://doi.org/10.1093/bioinformatics/btab083>