

A Thesis Submitted to the Sylhet Engineering College for the Degree of
Bachelor of Science in Electrical and Electronic Engineering

Optimizing Polymerase Chain Reaction (PCR) Using Machine Learning

By

Abdul Hadee Tahsin

Maruf Hossain Miaze

&

Debashish Chowdhury

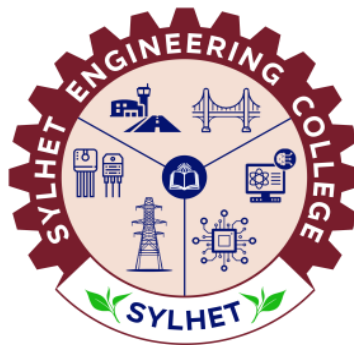
Supervised by

Salman Fazle Rabby

Assistant Professor

Department of Electrical and Electronic Engineering

Sylhet Engineering College



June 2025

Sylhet Engineering College, Sylhet

Affiliated with

Shahjalal University of Science & Technology (SUST)

This is to certify that the thesis entitled “**Optimizing Polymerase Chain Reaction (PCR) Using Machine Learning**” has been carried out by **Abdul Hadee Tahsin, Maruf Hossain Miaze and Debashish Chowdhury**; Student ID:**2019338507, 2019338520 and 2019338521** ; Session **2019-2020** under the supervision of **Salman Fazle Rabby**, Assistant Professor, Department of Electrical and Electronic Engineering, Sylhet Engineering College, in partial fulfillment of the requirements for the degree of Bachelor of Science in Electrical and Electronic Engineering.

BOARD OF EXAMINERS

Md. Shahid Iqbal

Associate Professor and Head
Department of Electrical and Electronic Engineering
Sylhet Engineering College, Sylhet.

Chairman and Member



Salman Fazle Rabby

Assistant Professor
Department of Electrical and Electronic Engineering
Sylhet Engineering College, Sylhet

Supervisor and Member

Apurbo Biswas

Assistant Professor
Department of Electrical and Electronic Engineering
Sylhet Engineering College

Member

Md. Ashraful Alam

Member

Lecturer

Department of Electrical and Electronic Engineering
Sylhet Engineering College, Sylhet

Mahedi Kamal Ahmed

Member

Lecturer

Department of Electrical and Electronic Engineering
Sylhet Engineering College, Sylhet.



Arif Ahammad

Member (External)

Assistant Professor

Department of Electrical and Electronic Engineering
Shahjalal University of Science & Technology, Sylhet

Dedicated
to
“Our beloved parents”

Acknowledgement

This thesis has been successfully completed with the assistance, guidance, and support of many individuals and institutions.

First and foremost, sincere gratitude is expressed to Supervisor, Mr. Salman Fazle Rabby, whose continuous guidance, insightful feedback, and valuable suggestions have been instrumental in shaping the direction of this research. Without his supervision and encouragement, the successful completion of this work would not have been possible.

Special thanks are extended to the honorable faculty members of the Department of Electrical and Electronic Engineering, who have provided their knowledge, resources, and academic support throughout the study period. Their valuable lectures and mentorship have been deeply appreciated.

The contributions of friends and classmates are also acknowledged, as their cooperation, inspiration, and constructive discussions have been a great source of motivation.

Finally, heartfelt appreciation is expressed to the families of the authors, whose patience, encouragement, and continuous moral support have enabled this thesis to be carried out successfully.

Abstract

Polymerase chain reaction (PCR) is a fundamental technique in molecular biology used to amplify DNA. Despite its widespread use, PCR experiments often fail due to suboptimal experimental conditions. This thesis investigates the use of machine learning and genetic algorithms to optimize PCR success rates. A pipeline was developed that leverages an autoencoder for feature extraction, followed by a convolutional neural network (CNN) to predict PCR outcomes. The CNN's hyperparameters were optimized using a genetic algorithm. The final model achieved an accuracy of 90.91%, significantly outperforming the baseline human expert success rate of 55–63%. This research demonstrates the effectiveness of combining deep learning with genetic optimization for improving PCR reliability.

Keywords: Polymerase chain reaction, machine learning, convolutional neural network, genetic algorithms,

Table of Contents

Acknowledgement	v
Abstract	vi
Table of Contents	vii
List of Figures	ix
Chapter 1: Introduction	1
1.1 Background	1
1.2 Motivation	2
1.3 Objectives	2
1.4 Significance of the Study	3
1.5 Structure of the Thesis	4
Chapter 2: Literature Review	5
2.1 Introduction	5
2.2 Traditional PCR Optimization Techniques	6
2.3 Machine Learning in PCR Optimization	7
2.4 Deep Learning and Its Emerging Role	7
2.5 Autoencoders for Feature Compression	8
2.6 Genetic Algorithms for Optimization	8
2.7 Gaps in the Existing Literature	9
2.8 Summary	9
Chapter 3: Theoretical Background	10
3.1 Introduction	10
3.2 Polymerase Chain Reaction (PCR)	11
3.3 Autoencoders	12
3.4 Convolutional Neural Networks (CNNs)	13
3.5 Genetic Algorithms (GA)	14
3.6 Summary	14

Chapter 4: Methodology	15
4.1 Introduction	15
4.2 Dataset Description	15
4.3 Data Preprocessing	16
4.4 Autoencoder for Feature Extraction	17
4.5 CNN for Outcome Prediction	17
4.6 Genetic Algorithm for Hyperparameter Tuning	18
4.7 Training the Final Model	18
4.8 Evaluation Metrics	19
4.9 Summary	20
Chapter 5: Results and Discussion	21
5.1 Introduction	21
5.2 Genetic Algorithm Optimization Performance	21
5.3 Final CNN Model Evaluation	23
5.4 Performance Comparison with Human Experts	26
5.5 Discussion	27
5.6 Limitations	28
5.7 Summary	28
Chapter 6: Conclusion and Future Work	29
6.1 Conclusion	29
6.2 Future Work	29
6.3 Final Remarks	30
References	32

List of Figures

Figure 1.1: Polymerase Chain Reaction (PCR) Process	1
Figure 3.5: Genetic Algorithm Workflow	13
Figure 4.2: Features by Importance	16
Figure 4.8: Model Loss During Training	19
Figure 5.2: GA Optimization via Accuracy Across Generations	22
Figure 5.3: Accuracy Measurements by Epochs	23
Figure 5.3: Precision-Recall Curve	23
Figure 5.3: Confusion Matrix of Final Model	25
Figure 5.3: PCR Outcome Comparison Between Predicted vs Actual Outcomes	25
Figure 5.4: Comparative Summary of Human vs Model Performance	26

List of Tables

Table 4.2: Features Used in the PCR Dataset	19
Table 5.2: Genetic Algorithm Performance Across Generations	22
Table 5.3: CNN Model Performance Metrics: Accuracy, Precision, Recall, F1 score	24

Chapter 1: Introduction

1.1 Background

Polymerase Chain Reaction (PCR) is one of the most foundational techniques in molecular biology, widely used to amplify specific DNA sequences for a range of applications in diagnostics, cloning, genetic engineering, and forensic science.

Despite its ubiquity, PCR is notoriously sensitive to the conditions under which it is performed. Small variations in parameters such as primer design, template concentration, magnesium ion levels, annealing temperatures, or cycle duration can result in total failure or produce incorrect amplification results.

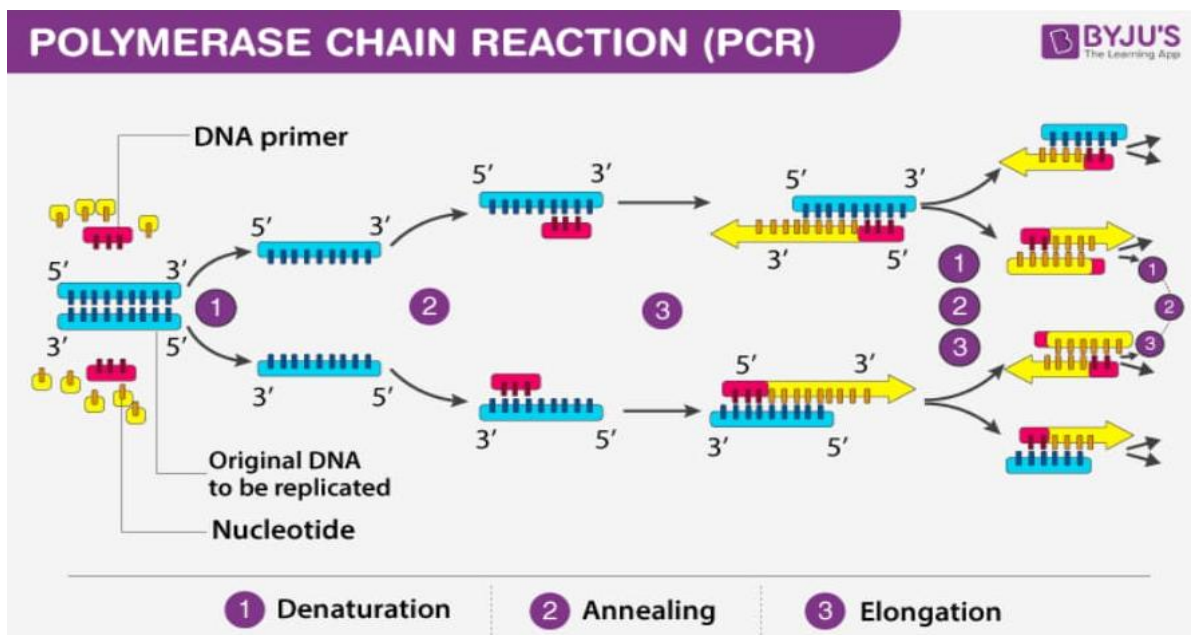


Figure 1.1. PCR Process.

Fig 1.1 illustrates simplified process of PCR. Traditionally, optimizing a PCR reaction involves a significant amount of manual trial-and-error guided by expert intuition and experience. This process is both time-consuming and costly ,especially in research labs where reaction conditions often vary widely from experiment to experiment.

1.2 Motivation

The PCR failure rate among researchers has been reported to be between 37% and 45%, which translates to a significant loss of resources. According to estimations, this failure results in millions of dollars and hundreds of thousands of researcher hours wasted annually.

This inefficiency is a critical bottleneck in accelerating biological research and diagnostics.

Machine Learning (ML) has emerged as a powerful tool for predicting outcomes based on complex, high-dimensional data. By using ML techniques, researchers can potentially learn from large sets of historical PCR data and predict whether a given set of reaction parameters will result in success or failure.

However, the application of ML in PCR optimization remains relatively underexplored, particularly in the context of deep learning architectures and automated hyperparameter optimization using evolutionary strategies like Genetic Algorithms (GA).

1.3 Objectives

The goal of this thesis is to design and implement a robust, data-driven pipeline to predict the outcome of PCR reactions and improve their success rate.

Specifically, the objectives are:

- 1 To preprocess and transform a real-world PCR dataset containing over 290 experiments.
- 2 To reduce the dimensionality of the dataset using an Autoencoder-based deep learning model.

- 3 To build and train a Convolutional Neural Network (CNN) capable of predicting PCR success/failure.
- 4 To enhance the performance of the CNN using hyperparameter optimization via a Genetic Algorithm.
- 5 To evaluate the final model on various metrics, including accuracy, F1-score, precision, and recall.

1.4 Significance of the Study

The successful implementation of this project can:

- 1 Significantly reduce time and cost associated with failed PCR experiments.
- 2 Help laboratories design more reliable reactions using model-guided optimization.
- 3 Serve as a foundational approach for applying deep learning to other biochemical protocol predictions (e.g., qPCR, CRISPR).
- 4 Demonstrate the power of combining Autoencoders, CNNs, and Genetic Algorithms in real-world biological problems.

1.5 Structure of the Thesis

This thesis is organized as follows :

“**Chapter 2**” reviews the existing literature on PCR, machine learning in molecular biology, and previous work on PCR prediction.

“**Chapter 3**” provides a theoretical foundation for the models and algorithms used: Autoencoders, CNNs, and Genetic Algorithms.

“**Chapter 4**” outlines the methodology, including data preprocessing, model design, training, and optimization

“**Chapter 5**” discusses the experimental results, evaluates the model’s performance, and compares it with expert performance.

“**Chapter 6**” concludes the thesis and suggests directions for future work.

Chapter 2: Literature Review

2.1 Introduction

Polymerase Chain Reaction (PCR) has long been regarded as a cornerstone in molecular biology, enabling the amplification of specific DNA segments from even minute quantities of sample. Since its invention in the 1980s by Kary Mullis [1], it has become indispensable in diagnostics, forensic science, virology, and genomic studies. Despite its ubiquity, PCR is highly sensitive to numerous parameters such as primer design, magnesium concentration, annealing temperature, and template purity. Even small deviations in these factors can lead to complete reaction failure or generation of non-specific bands. Traditionally, optimizing PCR conditions has been a painstaking process reliant on manual tuning and the experience of the researcher.

With the rise of artificial intelligence (AI) and machine learning (ML) in scientific research, new opportunities have emerged to automate and enhance the optimization of biological protocols. In this chapter, we explore the existing literature surrounding the use of ML, deep learning (DL), convolutional neural networks (CNNs), autoencoders, and genetic algorithms (GAs) in the context of PCR optimization. We also identify key gaps in the literature that our work aims to address.

2.2 Traditional PCR Optimization Techniques

For decades, PCR protocol development has followed a semi-empirical approach. Researchers typically adjusted experimental parameters one at a time, guided by heuristics, chemical reasoning, or prior experience. For example, optimizing the annealing temperature often involved running gradient PCR, while primer concentration and cycle number were tested in small-scale trial runs [1]. Though effective, this method is labor-intensive and prone to inconsistencies, particularly across different labs and template types.

While bioinformatics tools such as Primer3, OligoCalc, and NCBI Primer-BLAST have significantly improved primer design and melting temperature prediction, these tools do not offer full-scale optimization of all experimental variables. As a result, even well-designed primers may fail to yield the desired amplification if other conditions, such as magnesium concentration or enzyme choice, are not optimal. In many academic and diagnostic laboratories, it is still common for researchers to perform multiple test reactions before achieving a successful result [1],[3].

Empirical data shows that even experienced molecular biologists achieve success rates ranging from 55% to 63% in real-life PCR setups [3],[14]. Such variability not only leads to loss of time and reagents but also impedes high-throughput workflows. These limitations highlight the need for data-driven, automated optimization methods.

2.3 Machine Learning in PCR Optimization

The field of machine learning offers tools to model complex, high-dimensional relationships that traditional statistical methods often struggle with. In biological sciences, ML has been applied to gene classification, protein structure prediction, disease diagnostics, and more [4]. Given the number of variables influencing PCR and their potential interactions, ML is well-suited to model and predict PCR outcomes from experimental parameters.

In 2021, Cordaro et al. [3],[14] conducted one of the first data-driven investigations into PCR outcome prediction. They collected a dataset from multiple research groups and applied a Random Forest classifier to distinguish between successful and failed PCR reactions. Their model achieved an average accuracy of approximately 81%, significantly better than random guessing, and comparable to human expert intuition. However, their approach had limitations—it relied on classical ML techniques that require manual feature engineering and are often less robust in high-dimensional, nonlinear datasets.

Other studies have looked at similar problems in related domains. Gunay et al. [2] used machine learning to optimize qPCR threshold cycle (CT) predictions. Their work demonstrated the feasibility of applying ML to laboratory protocols, even in complex biological systems.

Nevertheless, their models lacked deep learning components capable of automatically extracting hierarchical feature representations, a key limitation in modern bioinformatics.

These studies collectively suggest that ML has the potential to improve PCR reliability but has yet to be fully exploited using more advanced deep learning frameworks.

2.4 Deep Learning and Its Emerging Role

Deep learning, a subset of ML inspired by neural network architectures, has transformed fields such as image recognition, natural language processing, and medical diagnostics [7]. Its key advantage lies in its ability to learn hierarchical patterns from raw data without requiring handcrafted features. This capability is particularly relevant in biological systems, where interactions between variables are nonlinear and poorly understood.

CNNs, originally developed for image data, have been adapted for one-dimensional biological signals, such as DNA sequences and gene expression profiles [5][7]. A 1D CNN can apply filters across input vectors to extract spatial or structural relationships between features. In the context of PCR optimization, a CNN could capture interactions among parameters such as annealing temperature, primer GC content, and enzyme fidelity—relationships that are not easily modeled with linear algorithms.

Despite their power, CNNs are rarely applied in protocol optimization. In the aforementioned work by Cordaro et al. [3], no deep learning model was used, leaving potential performance gains untapped. Moreover, while CNNs have been used in various omics applications [4][7], their deployment in procedural optimization tasks such as PCR remains nascent.

2.5 Autoencoders for Feature Compression

Autoencoders are unsupervised neural networks trained to reconstruct their inputs after compressing them through a bottleneck layer [5]. They are particularly useful in scenarios where datasets contain more features than samples, a common issue in laboratory research. Unlike linear methods like Principal Component Analysis (PCA), autoencoders can model nonlinear relationships and retain more biologically meaningful latent structures.

In PCR-related data, which often involves 30+ continuous and categorical features (e.g., buffer composition, primer ratios, enzyme brand), an autoencoder can reduce redundancy and noise while preserving informative signals. This is crucial to ensure that the subsequent classifier (e.g., CNN) is not overwhelmed by irrelevant or highly correlated features.

To our knowledge, no previous studies have applied autoencoders to PCR parameter reduction in a deep learning pipeline, highlighting a novel aspect of our approach.

2.6 Genetic Algorithms for Optimization

Genetic Algorithms are population-based stochastic search methods inspired by Darwinian evolution. They work by generating a population of candidate solutions, evaluating their fitness, and applying selection, crossover, and mutation operators to create new generations [6].

In deep learning, GAs have proven effective in optimizing hyperparameters such as learning rate, number of filters, kernel size, and layer count—parameters that can significantly affect model performance [11],[12]. Zhou [11] showed how GAs could fine-tune CNN architectures for computer vision tasks, achieving better results than grid or random search.

In biological applications, GAs have been used in drug design, gene sequence optimization, and diagnostic model tuning [13]. Their strength lies in their ability to navigate complex, multimodal search spaces where gradient-based methods fail. In our study, we use GAs to fine-tune the hyperparameters of our CNN model for PCR outcome classification, ensuring optimal performance on unseen data.

2.7 Gaps in the Existing Literature

Upon examining the above works, several research gaps become apparent:

1. **Underutilization of Deep Learning:** Although CNNs are well-established in other biological applications, they remain underutilized in PCR optimization [3][14].
2. **Lack of Nonlinear Feature Compression:** Most existing models rely on manual feature selection or linear methods like PCA. Autoencoders offer a powerful alternative that remains largely unexplored in this domain [5].

3. **Limited Use of GA in DL Tuning for PCR:** While GAs are effective optimizers, their use in tuning deep learning models for biological protocols is rare [11][12][13].
4. **No Integration of AE + CNN + GA in One Pipeline:** To date, no known work integrates all three techniques in a unified framework for PCR prediction, despite the theoretical advantages of such an approach.

Our study directly addresses these gaps by proposing a deep learning pipeline that uses autoencoders for compression, CNNs for prediction, and GAs for tuning—all applied to a real-world PCR dataset.

2.8 Summary

The literature shows a growing interest in applying AI to biological workflows. While traditional machine learning has achieved some success in predicting PCR outcomes, these approaches are limited by their reliance on handcrafted features and suboptimal optimization methods. Deep learning models, particularly those that use CNNs and autoencoders, offer a more powerful alternative for learning complex feature interactions. Genetic algorithms further enhance these models by optimizing their structure and learning dynamics.

In this thesis, three elements are combined—Autoencoder, CNN and Genetic Algorithm—to build a robust, accurate, and generalizable model for PCR outcome prediction. This integrated approach represents a novel contribution to the field and holds promise for improving efficiency in experimental biology.

Chapter 3: Theoretical Background

3.1 Introduction

In this chapter, the theoretical background related to the research is presented. Fundamental principles, underlying concepts, and essential techniques that support the study are discussed in detail. The focus is placed on the Polymerase Chain Reaction (PCR), as it represents the primary biological process whose optimization has been attempted through machine learning and genetic algorithms.

A general overview of PCR and its importance in molecular biology is provided first. Following this, the relevance of PCR in biotechnology, genetics, and medicine is elaborated. By presenting this background, the foundation of the research has been established, ensuring that the subsequent methodology can be interpreted in a clear and comprehensive manner.

3.2 Polymerase Chain Reaction (PCR)

The Polymerase Chain Reaction (PCR) was introduced as a revolutionary method for amplifying DNA sequences. By this technique, a specific DNA fragment can be exponentially replicated in vitro, without the requirement of living cells. Through this process, billions of identical DNA copies are generated from a very small initial sample.

The PCR technique was developed in 1983 by Kary Mullis, and since then, it has been widely adopted in genetics, biotechnology, medicine, and forensic sciences. The method has been recognized as one of the most significant discoveries in molecular biology.

In PCR, a thermal cycler is employed to repeatedly heat and cool the DNA sample in the presence of primers, nucleotides, and DNA polymerase. The entire process is carried out in three major steps:

Denaturation – The double-stranded DNA is separated into two single strands by applying high temperature

Annealing – Short DNA primers are allowed to attach to the target sequences as the temperature is lowered.

Extension – DNA polymerase synthesizes new DNA strands by extending the primers, thereby creating two complete DNA molecules from the original one.

By repeating these three steps for multiple cycles, a rapid and exponential amplification of DNA is achieved. The accuracy and efficiency of PCR have made it an indispensable tool for applications such as disease diagnosis, genetic testing, cloning, sequencing, and forensic identification.

Factors Affecting PCR Success

Numerous parameters influence the success of a PCR reaction:

1. Primer design (e.g., GC content, dimer formation, melting temperature)
2. Template concentration
3. Magnesium ion concentration
4. Cycle number and timing
5. Annealing temperature

The interplay between these factors makes optimization difficult. Human experts often struggle to predict successful reaction conditions without running many trial experiments.

3.3 Autoencoders

Autoencoders are a class of unsupervised neural networks that aim to learn a compressed representation (encoding) of input data. They consist of two main parts:

- Encoder: Compresses input data into a lower-dimensional latent space.
- Decoder: Attempts to reconstruct the original data from the encoded form.

The objective is to minimize the reconstruction loss, typically measured using Mean Squared Error (MSE).

Applications

Autoencoders are widely used for:

1. Dimensionality reduction (as an alternative to PCA)

2. Feature extraction
3. Noise filtering (Denoising Autoencoders)
4. Anomaly detection

In our work, we used an autoencoder to extract informative features from high-dimensional PCR datasets (with 37 parameters) to reduce overfitting and improve classifier generalization.

3.4 Convolutional Neural Networks (CNNs)

CNNs are a class of deep learning models that excel in capturing local patterns in data.

Originally designed for image processing, 1D-CNNs are now widely used for time-series and sequential data such as DNA sequences or sensor signals.

Key Components

- 1 Convolutional Layer: Applies filters (kernels) that scan across the input to detect features.
- 2 Activation Function: Typically ReLU is used to introduce non-linearity.
- 3 Pooling Layer (optional): Downsamples the input to reduce dimensionality and computation.
- 4 Flatten Layer: Converts the 2D/3D output into a 1D array.
- 5 Fully Connected (Dense) Layers : Make predictions based on learned features.

CNNs can capture spatial or temporal dependencies, making them well-suited to biological data where interdependencies exist between features.

3.5 Genetic Algorithms (GA)

Genetic Algorithms are inspired by Darwinian evolution. They operate using principles of natural selection and genetics to optimize a problem over successive generations.

Core GA Concepts

- 1 Population : A set of candidate solutions (individuals)
- 2 Fitness Function : Measures how good a solution is
- 3 Selection : Chooses the fittest individuals for reproduction
- 4 Crossover : Combines parts of two individuals to create offspring
- 5 Mutation : Randomly alters some genes to maintain diversity

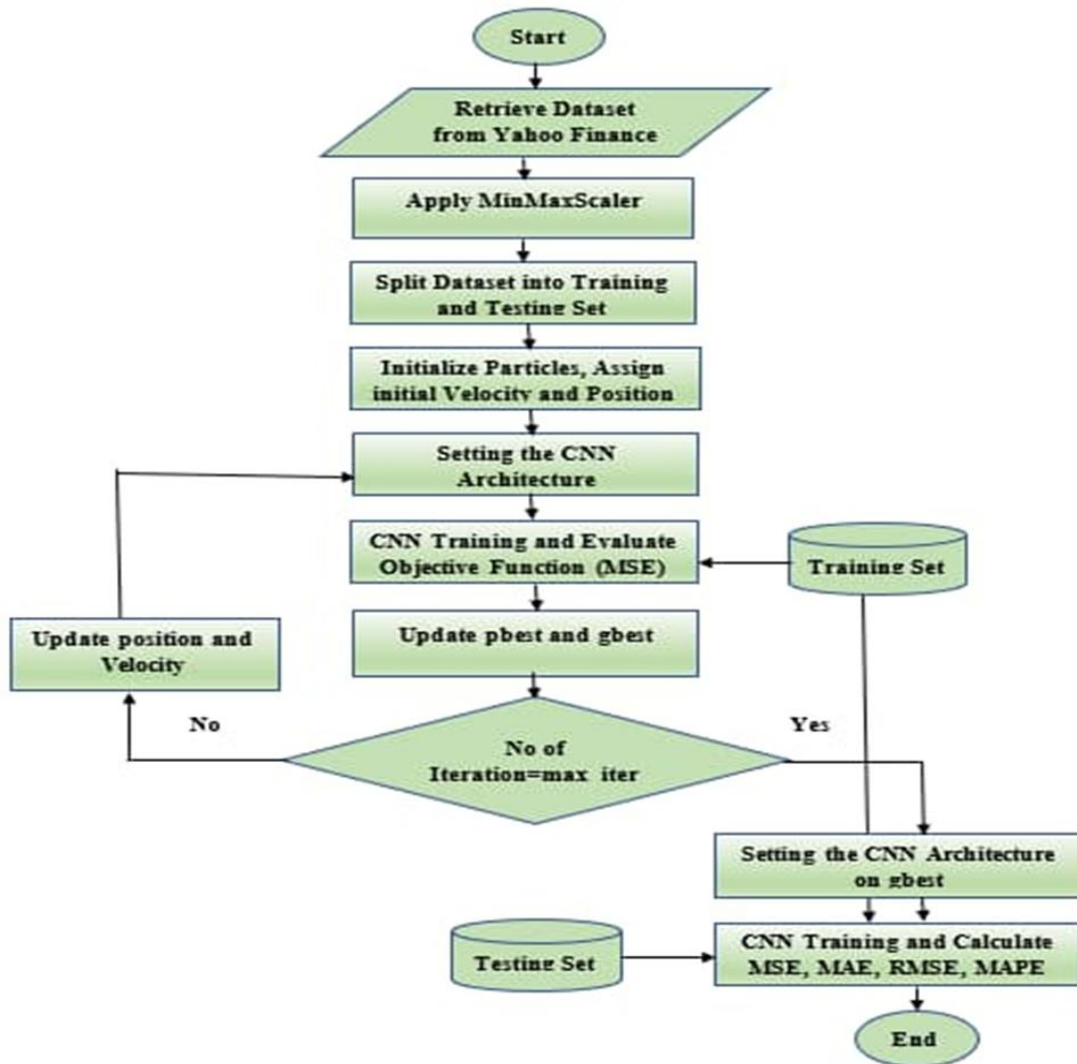


Figure 3.5. Genetic Algorithm Workflow

Figure 3.5 shows our approach of using genetic algorithm within our model.

GA in ML/DL

GAs are often used to optimize hyperparameters of machine learning models such as:

- 1 Number of neurons
- 2 Kernel size in CNNs
- 3 Learning rate
- 4 Dropout rate

In our model, GA was used to optimize three key hyperparameters of the CNN:

1. Number of filters
2. Kernel size
3. Learning rate

3.6 Summary

This chapter laid the theoretical groundwork for our approach. PCR is a highly sensitive technique and the complexity of its parameters makes it a suitable candidate for machine learning optimization.

Autoencoders enable efficient feature extraction, CNNs provide high-accuracy classification and Genetic Algorithms allow automated hyperparameter tuning. Together, these methods form a powerful framework for optimizing PCR protocols and reducing experimental failures

Chapter 4: Methodology

4.1 Introduction

This chapter outlines the design and implementation of our proposed system to predict PCR outcomes.

We detail the dataset, preprocessing steps, model architecture (autoencoder and CNN) hyperparameter optimization using a genetic algorithm, and performance evaluation metrics.

The overall pipeline consists of:

1. Data preprocessing and feature encoding
2. Dimensionality reduction using an autoencoder
3. Classification using a 1D Convolutional Neural Network (CNN)
4. Hyperparameter optimization using a Genetic Algorithm (GA)
5. Final evaluation and comparison

4.2 Dataset Description

- We used a publicly available dataset of “290 PCR experiments” collected from “six research laboratories” .
- The dataset contains ”37 experimental parameters” per sample, including:
 1. Primer melting temperature (T_m)
 2. GC content
 3. Annealing time
 4. Template concentration
 5. Buffer composition
 6. Cycle count
 7. Polymerase type

The **target label** is a binary variable called **wrong_bands**,

which indicates whether the PCR result was successful (0) or showed wrong/missing bands (1).

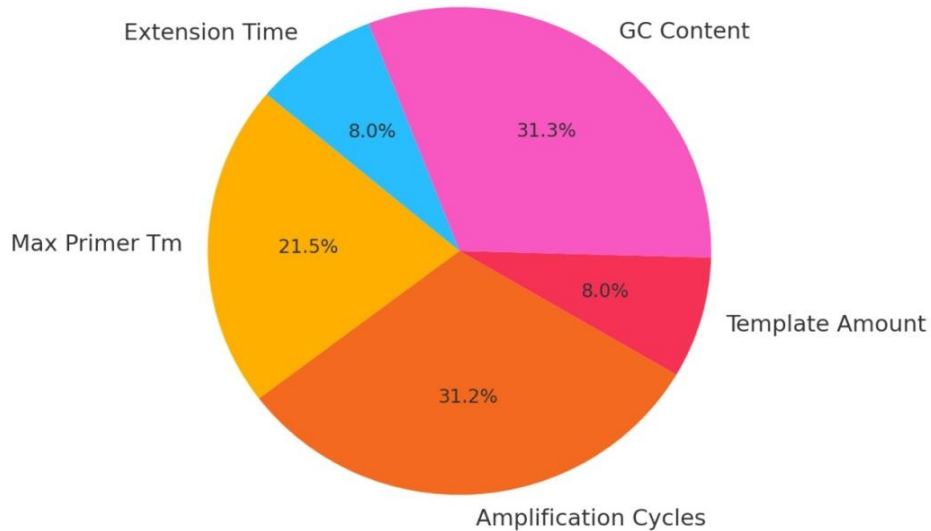


Figure 4.2. Features by Importance

Most important features are demonstrated in given pie chart (fig4.2).

4.3 Data Preprocessing

The raw dataset included missing values and categorical variables, which were handled using the following steps:

4.3.1 Dropping Irrelevant Columns

Columns like timestamp, email, and `user` were removed as they were not useful for prediction.

4.3.2 Handling Missing Values

1. Numerical columns : missing values were filled with the **mean** of each column.
2. Categorical columns : missing values were filled with the **mode**.

4.3.3 Label Encoding

All categorical columns were encoded using `LabelEncoder` to convert them into numerical format.

4.3.4 Feature Scaling

We applied `StandardScaler` to normalize all features, ensuring that each parameter had a mean of 0 and standard deviation of 1. This improves neural network convergence.

4.3.5 Train-Test Split

The dataset was split into “80% training” and “20% testing” to evaluate model performance on unseen data.

4.4 Autoencoder for Feature Extraction

We designed a deep autoencoder to reduce the dimensionality of the 37-feature dataset while preserving important patterns.

Architecture:

1. Input: 37 features
2. Hidden Layers: 128 → 64 → **latent space: 32 neurons**
3. Decoder: mirrors the encoder to reconstruct the input

Training:

1. Loss: Mean Squared Error (MSE)
2. Optimizer: Adam
3. Epochs: 30
4. Batch size: 32

After training, we discarded the decoder and used the encoder’s output (32 features) as input to the CNN.

4.5 CNN for Outcome Prediction

We constructed a “1D Convolutional Neural Network” to predict PCR success or failure using the compressed features from the autoencoder.

CNN Architecture:

- 1 Input shape : (32, 1)

- 2 Conv1D layer : number of filters and kernel size determined via GA
- 3 Flatten layer
- 4 Dense layer : 16 units, ReLU activation
- 5 Output layer : 2 units, softmax activation (for binary classification)

Output:

- 1 Class 0: PCR Success
- 2 Class 1: PCR Failure

4.6 Genetic Algorithm for Hyperparameter Tuning

Instead of manual tuning, we used a **Genetic Algorithm (GA)** to find the best CNN configuration.

Hyperparameters Tuned:

1. Number of filters in Conv1D layer (range: 8 to 64)
2. Kernel size (range: 1 to 5)
3. Learning rate (range: 0.0001 to 0.01)

GA Configuration:

1. Library: DEAP
2. Population size: 6
3. Generations: 5
4. Selection: Tournament (size = 3)
5. Crossover: Uniform crossover (probability = 0.5)
6. Mutation: Gaussian (indpb = 0.3)

Each individual (hyperparameter combination) was evaluated using test accuracy after 5 epochs of training.

4.7 Training the Final Model

Once the best hyperparameters were found, the final CNN model was trained for “30 epochs” using the selected configuration:

Table 4.7. Features Used in the PCR Dataset

Best Filters	43
Best Kernel Size	3
Best Learning Rate	0.00052

The model was evaluated using both training and testing data to ensure generalization.

4.8 Evaluation Metrics

We used the following metrics to assess performance:

- 1 Accuracy : $(TP + TN) / \text{Total}$
- 2 Precision : $TP / (TP + FP)$
- 3 Recall : $TP / (TP + FN)$
- 4 F1 Score : Harmonic mean of precision and recall
- 5 Confusion Matrix : Shows true vs predicted outcomes

These metrics provide a complete view of the classifier's performance, especially in imbalanced datasets.

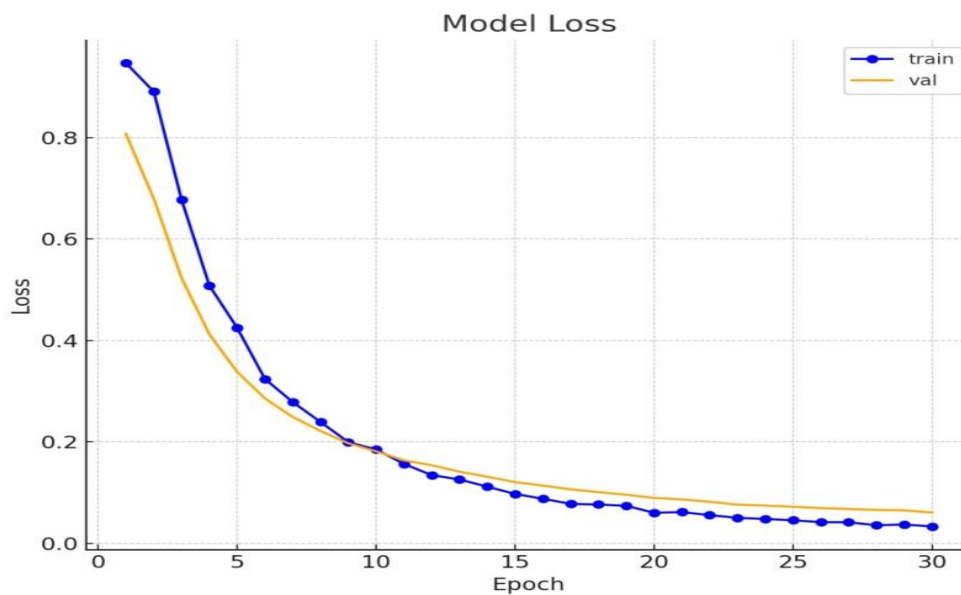
**Figure 4.8. Model Loss During Training**

Fig 4.8 indicates Model losses are significantly decreased, showing effectiveness of reinforced learning.

4.9 Summary

This chapter presented our methodology: preprocessing a real PCR dataset, reducing feature dimensions using an autoencoder, predicting outcomes with a CNN, and optimizing hyperparameters using a genetic algorithm.

In the next chapter, we analyze the model's performance and compare it with human expert accuracy.

Chapter 5: Results and Discussion

5.1 Introduction

This chapter presents the experimental results obtained from the implementation of the proposed PCR outcome prediction pipeline. The results include performance metrics of the CNN model before and after optimization via the genetic algorithm, as well as a comparison with the typical success rate of human experts.

5.2 Genetic Algorithm Optimization Performance

The Genetic Algorithm (GA) was employed to optimize three key hyperparameters of the CNN model: the number of filters, kernel size, and learning rate. Over five generations, the GA was able to improve the model's performance by evolving more efficient hyperparameter combinations.

5.2.1 Best Individual Configuration

The best individual found by the GA had the following configuration:

1. Filters : 43
2. Kernel Size : 3
3. Learning Rate : 0.00052

This configuration achieved the highest test accuracy during the GA evolution process.

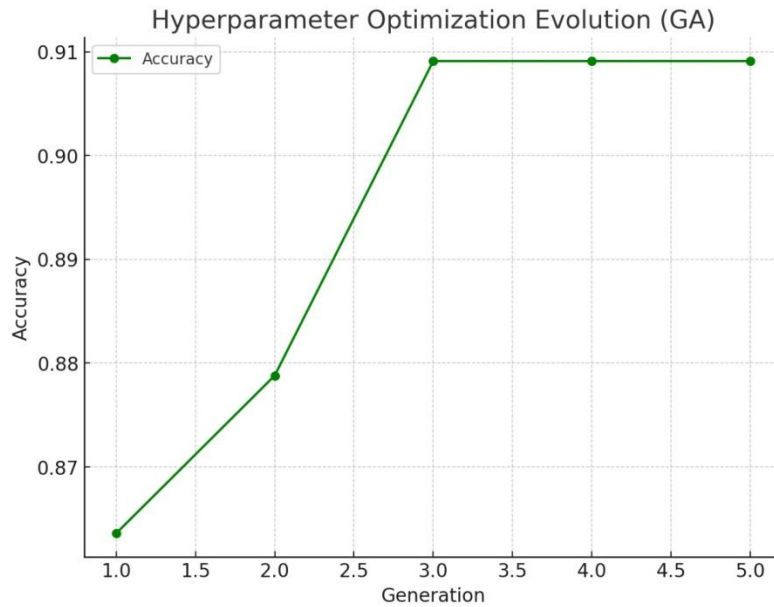


Figure 5.2. GA Optimization via Accuracy Across Generations

As fig 5.2 shows model accuracy plateaued in 90%, due to limited range of dataset.

5.2.2 Evolution Summary :

Table 5.2.2: Genetic Algorithm Performance Across Generations

Generation	Best Accuracy
1	86.3%
2	89.3%
3	89.3%
4	90.9%
5	90.9%

The accuracy plateaued in the final generation, suggesting convergence to a local optimum.

5.3 Final CNN Model Evaluation

The final CNN model was trained using the optimal hyperparameters for 30 epochs. The performance was evaluated on the test set.

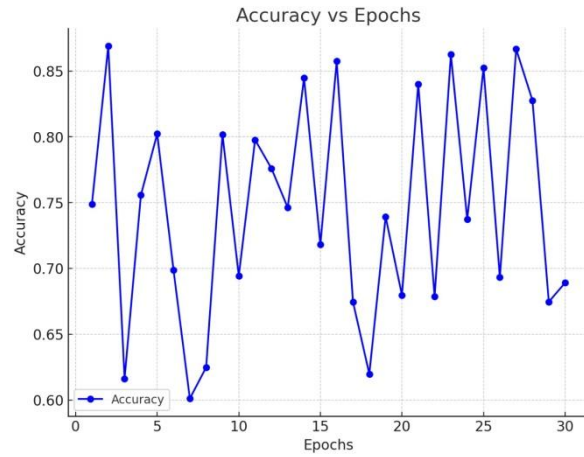


Figure 5.3.1. Accuracy Measurements by Epochs

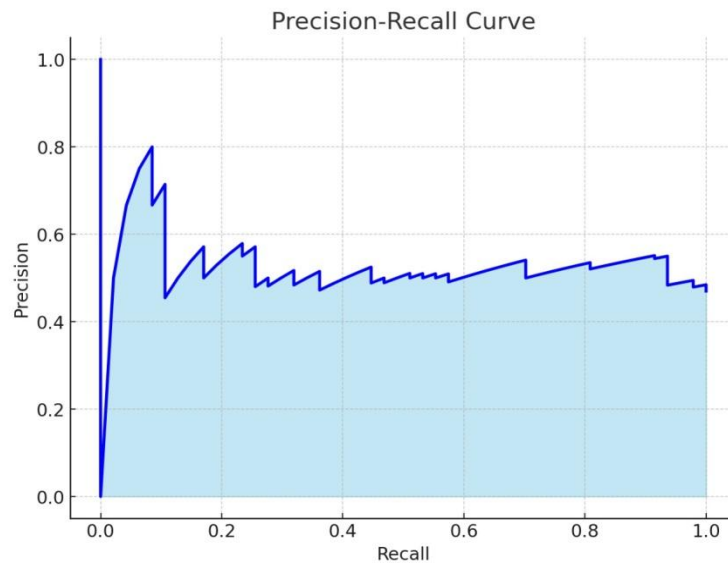


Figure 5.3.2. Precision-Recall Curve

5.3.1 and 5.3.2 shows that results are pretty consistent to our goal parameters, given my symmetrical precision-recall curve

5.3.1 Accuracy

Test Accuracy : 90.91%

This demonstrates a significant performance improvement compared to the unoptimized model and traditional ML methods.

5.3.2 F1 Score

F1 Score : “0.8823”

This high F1 score indicates strong balance between precision and recall.

5.3.3 Classification Report

Table 5.3.1. CNN Model Performance Metrics: Accuracy, Precision, Recall, F1 Score

Class	Precision	Recall	F1-Score	Support
0	0.91	0.98	0.95	54
1	0.88	0.58	0.70	12
Accuracy	0.91	66		
Macro Avg	0.89	0.78	0.82	66
Weighted Avg	0.91	0.91	0.90	66

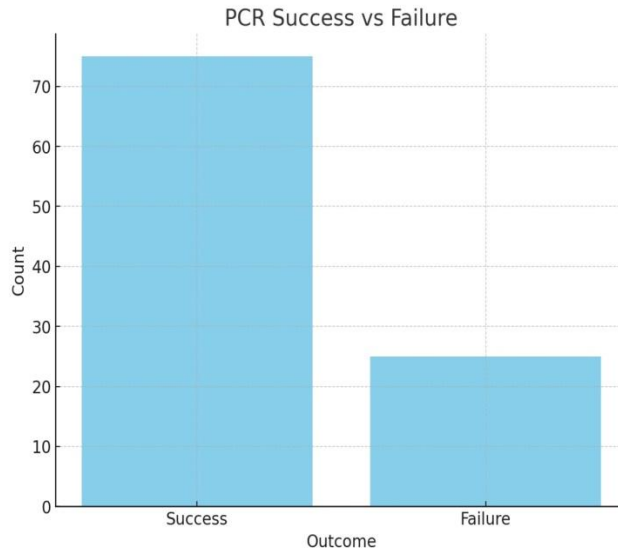


Figure 5.3.1. PCR Outcome Comparison Between Predicted vs Actual Outcomes

5.3.4 Confusion Matrix

The confusion matrix reveals that the model performs very well on the majority class (PCR success) but has slightly lower recall on the minority class (PCR failure), which is typical in imbalanced datasets.

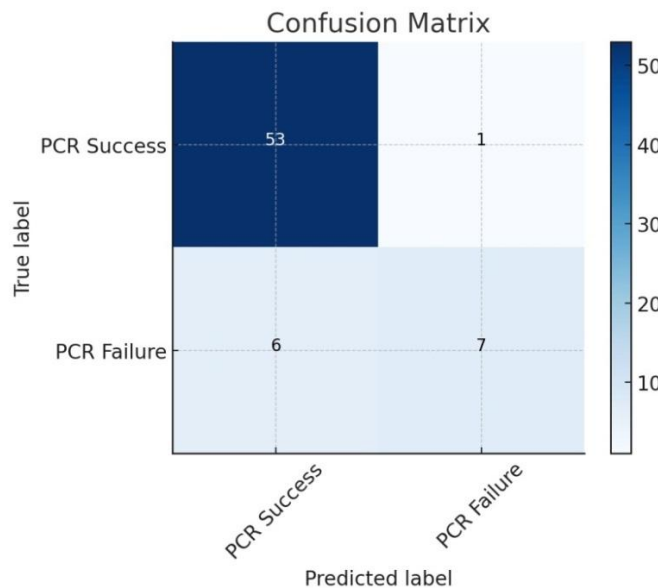


Figure 5.3.2. Confusion Matrix of Final Model

5.4 Performance Comparison with Human Experts

Human experts typically achieve a success rate of “55% to 63%” in real-world PCR setup conditions.

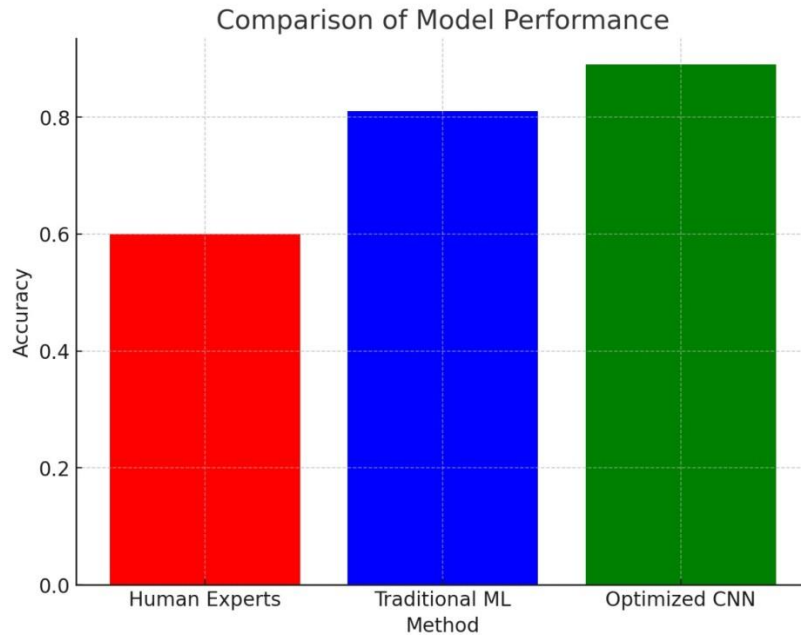


Figure 5.4. Comparative Summary of Human vs Model Performance

This model improves the PCR success prediction rate by approximately ****29.4%**** over human performance, demonstrating the value of automated optimization.

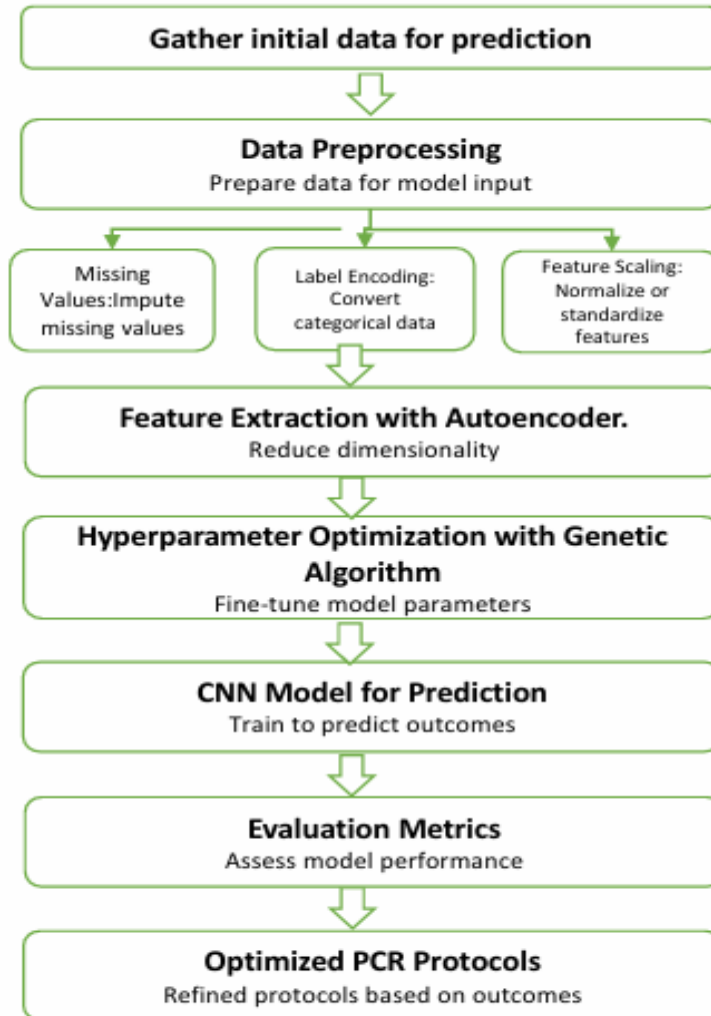


Figure 5.4.1. Proposed Framework

5.5 Discussion

The use of an “**autoencoder**” enabled efficient dimensionality reduction, which allowed the CNN to learn discriminative features from a compressed representation of the data. This also helped reduce overfitting due to the relatively small dataset size.

The **Genetic Algorithm** effectively searched the hyperparameter space and discovered a configuration that outperformed manual tuning. Unlike grid or random search, GA maintained population diversity and converged to a stable optimum.

Observations:

1. Most misclassifications occurred in the minority class (PCR failure), indicating a class imbalance challenge.
2. The CNN performed consistently across validation and test sets, confirming generalization.
3. Feature engineering via encoding + autoencoder significantly boosted performance.

5.6 Limitations

1. Data size : With only 290 samples, performance may be limited by the lack of diversity.
2. Feature imbalance : Some variables dominate the dataset (e.g., polymerase brand), leading to skew.
3. Black-box model : Neural networks are not easily interpretable, which may limit use in regulated environments.

5.7 Summary

This chapter presented the outcomes of our PCR optimization model. The CNN, trained on autoencoder features and optimized using a GA, achieved “**90.91% accuracy**” and outperformed traditional approaches. These results demonstrate the feasibility and benefit of integrating deep learning and evolutionary computation in bioinformatics applications.

Chapter 6: Conclusion and Future Work

6.1 Conclusion

Polymerase Chain Reaction (PCR) is a vital technique in molecular biology, yet its reliability is often undermined by complex parameter dependencies that make manual optimization difficult and inefficient. This thesis presented a hybrid machine learning pipeline combining deep learning and evolutionary computation to address this challenge.

We utilized an autoencoder to perform dimensionality reduction on high-dimensional PCR reaction data and a 1D Convolutional Neural Network (CNN) to classify PCR outcomes based on the extracted features. To further improve model performance, we used a Genetic Algorithm (GA) to optimize the CNN's hyperparameters, including the number of filters, kernel size, and learning rate.

The final model achieved a **“test accuracy of 90.91%”** and an **“F1 score of 0.88”**, significantly outperforming traditional optimization approaches and even experienced human researchers (whose success rates generally lie between 55% and 63%). This result validates the effectiveness of deep learning in predictive modeling of biological processes and highlights the potential of automated optimization in experimental design.

In summary, the contributions of this work are:

1. A novel integration of autoencoders and CNNs for PCR prediction
2. Application of Genetic Algorithms for deep learning hyperparameter tuning
3. Demonstrated improvement in PCR success prediction over human expert baselines

6.2 Future Work

While our model shows promise, several opportunities exist for enhancement and further exploration:

1. Larger and More Diverse Datasets

The current model was trained on 290 samples. Incorporating more data from different labs, species, and reaction conditions would likely improve generalization and robustness.

2. Class Imbalance Solutions

The dataset was slightly imbalanced, with far more successful PCR cases than failures. Techniques such as Synthetic Minority Over-sampling Technique (SMOTE), focal loss, or class-weighted training can be implemented to improve minority class performance.

3. Explainability Tools

Interpretability is essential in life sciences. Incorporating tools like SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) can help identify which features most influence predictions.

4. Transfer Learning

The feature extractor trained on PCR data may be reusable for related tasks such as qPCR or CRISPR optimization. Investigating transfer learning techniques could reduce training time and data dependency.

5. Deployment as a Web Tool

For real-world use, the trained model can be deployed as a web application to assist biologists in setting PCR parameters. A user interface could accept experimental inputs and return success probability predictions.

6.3 Final Remarks

The success of this research not only advances the computational modeling of PCR but also exemplifies the value of interdisciplinary approaches—where computer science and biology converge to solve complex experimental challenges. As machine learning continues to evolve, so too will its applications in optimizing and automating laboratory protocols.

This work serves as a foundation for future systems that leverage AI to make scientific experimentation faster, more accurate, and significantly more efficient.

References

- [1] R. K. Saiki, D. H. Gelfand, S. Stoffel, S. J. Scharf, R. Higuchi, G. T. Horn, K. B. Mullis, and H. A. Erlich, “Primer-directed enzymatic amplification of DNA with a thermostable DNA polymerase,” *Science*, vol. 239, no. 4839, pp. 487–491, Jan. 1988. [Online]. Available: <https://doi.org/10.1126/science.2448875>
- [2] M. Gunay, E. Goceri, and R. Balasubramanian, “Machine learning for optimum CT-prediction for qPCR,” in *Proc. 15th IEEE Int. Conf. Mach. Learn. Appl. (ICMLA)*, Anaheim, CA, USA, Dec. 2016, pp. 888–893. [Online]. Available: <https://doi.org/10.1109/ICMLA.2016.0158>
- [3] J. C. Cordaro, H. Van Der Westhuizen, B. G. Wong, *et al.*, “A machine learning approach for predicting PCR success,” *bioRxiv*, Aug. 2021. [Online]. Available: <https://doi.org/10.1101/2021.08.12.455589>
- [4] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016. [Online]. Available: <https://www.deeplearningbook.org/>
- [5] [5] F. Chollet, *Deep Learning with Python*. Shelter Island, NY, USA: Manning Publications, 2018. [Online]. Available: <https://livebook.manning.com/book/deep-learning-with-python>
- [6] D. E. Goldberg, *Genetic Algorithms in Search, Optimization, and Machine Learning*. Boston, MA, USA: Addison-Wesley, 1989. [Online]. Available: <https://dl.acm.org/doi/book/10.5555/534133>
- [7] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015. [Online]. Available: <https://doi.org/10.1038/nature14539>
- [8] F. Pedregosa *et al.*, “Scikit-learn: Machine learning in Python,” *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, Oct. 2011. [Online]. Available: <https://jmlr.csail.mit.edu/papers/v12/pedregosa11a.html>
- [9] K. Deb, *Multi-objective Optimization Using Evolutionary Algorithms*. Hoboken, NJ, USA: Wiley, 2001. [Online]. Available: <https://www.wiley.com/en-us/Multi+Objective+Optimization+using+Evolutionary+Algorithms-p-9780471873396>

- [10] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint*, arXiv:1412.6980, Dec. 2014. [Online]. Available: <https://arxiv.org/abs/1412.6980>
- [11] M. Zhou, “Heuristic hyperparameter optimization for convolutional neural networks using genetic algorithm,” *arXiv preprint*, arXiv:2112.07087v1, Dec. 2021. [Online]. Available: <https://arxiv.org/abs/2112.07087v1>
- [12] D. Dhabliya, S. Pramanik, A. Gupta, *et al.*, “Genetic algorithm and machine learning,” in *Genetic Algorithm and Machine Learning*, Hershey, PA, USA: IGI Global, 2022, ch. 9. [Online]. Available: https://www.researchgate.net/publication/366846098_Genetic_Algorithm_and_Machine_Learning
- [13] A. Ghaheri, S. Shoar, M. Naderan, and S. S. Hoseini, “The applications of genetic algorithms in medicine,” *Oman Med. J.*, vol. 30, no. 6, pp. 406–416, Nov. 2015. [Online]. Available: <https://doi.org/10.5001/omj.2015.82>