

**Ensemble Learning Model for Accurate
Identification of hormone-binding protein**

Refatun Nahar Reya

Bachelor of Science

DAFFODIL INTERNATIONAL UNIVERSITY

DAFFODIL INTERNATIONAL UNIVERSITY

DECLARATION OF THESIS AND COPYRIGHT

Author's Full Name : Refatun Nahar Reya
Date of Birth : 24-11-2002
Title : Ensemble Learning Model for Accurate Identification of hormone-binding protein
Academic Session : 2022-2025

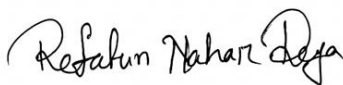
I declare that this thesis is classified as:

- CONFIDENTIAL (Contains confidential information under the Official Secret Act 1997)*
- RESTRICTED (Contains restricted information as specified by the organization where research was done)*
- OPEN ACCESS I agree that my thesis to be published as online open access (Full Text)

I acknowledge that Daffodil International University reserves the following rights:

1. The Thesis is the Property of Daffodil International University.
2. The Library of Daffodil International University has the right to make copies of the thesis for the purpose of research only.
3. The Library of Daffodil International University has the right to make copies of the thesis for academic exchange.

Certified by:



(Student's Signature)

221-35-884

Student ID

Date:



(Supervisor's Signature)

Mr. Md. Selim Reza

Name of Supervisor

Date:

NOTE : * If the thesis is CONFIDENTIAL or RESTRICTED, please attach a thesis declaration letter.

THESIS DECLARATION LETTER (OPTIONAL)

Librarian,
Daffodil International University,
Daffodil Smart City,
Ashulia.Dhaka,Bangladesh

Dear Sir,

CLASSIFICATION OF THESIS AS RESTRICTED

Please be informed that the following thesis is classified as RESTRICTED for a period of three (3) years from the date of this letter. The reasons for this classification are as listed below.

Author's Name

Thesis Title

Reasons (i)

(ii)

(iii)

Thank you.

Yours faithfully,



(Supervisor's Signature)

Date:

Stamp:

Note: This letter should be written by the supervisor and addressed to the Librarian, *Daffodil International University* with its copy attached to the thesis.



SUPERVISOR'S DECLARATION

I hereby declare that I have checked this **thesis** and, in my opinion, this **thesis** is adequate in terms of scope and quality for the award of the degree of **Bachelor of Science (B.Sc.)**.

A handwritten signature in black ink, appearing to read "Selim Reza", written over a horizontal line.

(Supervisor's Signature)

Full Name : Mr. Md. Selim Reza

Position : Assistant Professor

Date :

APPROVAL

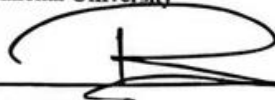
This thesis titled on "Ensemble Learning Model for Accurate Identification of Hormone-Binding Protein", submitted by Refatun Nahar Reya (ID: 221-35-884) to the Department of Software Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of Bachelor of Science in Software Engineering and approval as to its style and contents.

BOARD OF EXAMINERS



Chairman

Dr. A. H. M. Saifullah Sadi
Professor
Department of Software Engineering
Faculty of Science and Information Technology Daffodil
International University



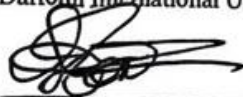
Internal Examiner 1

Dr. Rubaiyat Islam
Associate Professor
Department of Software Engineering
Faculty of Science and Information Technology
Daffodil International University



Internal Examiner 2

Dr. Md. Abdul Kader
Associate Professor
Department of Software Engineering
Faculty of Science and Information Technology
Daffodil International University



Internal Examiner 3

Nuruzzaman Faruqi
Assistant Professor
Department of Software Engineering
Faculty of Science and Information Technology
Daffodil International University



External Examiner

Md. Mostafiz Khan
Managing Director
Tecognize Solutions Limited



STUDENT'S DECLARATION

I hereby declare that the work in this thesis is based on my original work except for quotations and citations which have been duly acknowledged. I also declare that it has not been previously or concurrently submitted for any other degree at Daffodil International University or any other institution.

Refatun Nahar Reya

(Student's Signature)

Full Name : Refatun Nahar Reya

ID Number : 221-35-884

Date : 27-11-2025

Ensemble Learning Model for Accurate Identification of hormone-binding protein

Refaun Nahar Reya

Thesis submitted in fulfillment of the requirements
for the award of the degree of
Bachelor of Science/Master of Science

Department of Software Engineering (Major in Data Science)

DAFFODIL INTERNATIONAL UNIVERSITY

DECEMBER 2025

ACKNOWLEDGEMENTS

All praise and gratitude are due to Allah (SWT), the Most Gracious, the Most Merciful. Without His infinite grace, strength, and guidance, this journey would never have been possible. Without His divine will, this work would not have been possible.

I would like to say that I would like to thank my highly rated supervisor, Mr. Md. Selim Reza, Assistant professor. Your constant encouragement, guidance and knowledge have been very significant in the accomplishment of this thesis. Not only has your advice been the impulse to this study, but you have been a great inspiration to me, and I am grateful indeed.

I am also deeply grateful to my family, whose unconditional love, patience, and encouragement have always been a source of inspiration for me. I am also grateful to my friends, whose cooperation and support have helped me move this work forward.

DEDICATION

In the name of Almighty Allah (SWT), The supreme source of strength, knowledge and guidance. To every patient who has been silently struggling with intricate endocrine diseases and hormonal dysfunction, to his or her diagnosis and treatment efforts- may this research play a small part in creating solutions to better address the problem.

ABSTRACT

Hormone-Binding Proteins (HBPs) play a critical role to maintain the movement, stability, and introduction of hormones. Despite the assistance of the existing computer programs, it remains difficult to get the correct predictions of HBP due to the presence of duplicate datasets, the absence of diversity in features, and ineffective generalization of the models. In this paper, the authors introduce an Ensemble Learning Framework optimized using Differentiated Evolution (DE) that is able to identify HBPs directly based on protein sequences at a high degree of accuracy. We used a complete benchmark data to compare the performance of the model. We examined 6 possible options of encoding features and Dipeptide Composition (DPC) was found to be the most successful in distinguishing between them. We assembled a weighted combination of seven various machine learning classifiers. DE was used to identify the best weights. The final model performed well on a test set, and had an accuracy of 98.00, sensitivity of 100 percent and an MCC of 0.9608. It is a cheap and highly accurate computational approach that is an excellent choice to conduct a large-scale proteomic analysis and identify biomarkers related to hormones.

Table of Contents

APPENDICES	VII
LIST OF TABLES	VIII
LIST OF FIGURES	IX
LIST OF SYMBOLS	X
LIST OF ABBREVIATIONS	XI
CHAPTER 1: INTRODUCTION	1
1.1 BACKGROUND OF THE STUDY.....	1
1.2 PROBLEM STATEMENT.....	2
1.3 MOTIVATION	3
1.4 SIGNIFICANCE OF THE STUDY	4
1.5 RESEARCH OBJECTIVES	5
1.6 RESEARCH SCOPE AND LIMITATIONS	6
CHAPTER 2: LITERATURE REVIEW	7
2.1 RELATED WORK	7
2.2 GAP IN RESEARCH	9
CHAPTER 3: METHODOLOGY	10
3.1 OBJECTIVE OF THE STUDY	10
3.2 DATASET PARTITIONING.....	10
<i>Table 3.1: Statistical Distribution of the Datasets</i>	<i>11</i>
3.3 FEATURE EXTRACTION	12
3.4 BASELINE CLASSIFIERS AND MODEL SELECTION	13
3.5 ENSEMBLE LEARNING FRAMEWORK PROPOSAL.....	14
3.5 PERFORMANCE EVALUATION	15
CHAPTER 4: RESULTS AND DISCUSSION	18
4.1 FEATURE EVALUATION AND SELECTION	18
4.2 BASELINE CLASSIFIER PERFORMANCE	19

4.3 WEIGHTS OPTIMIZATION USING DIFFERENTIAL EVOLUTION.....	19
4.4 ASSESSMENT OF INDEPENDENT TEST DATASET	21
4.5 CONFUSION MATRIX ANALYSIS	21
4.6 ROC CURVE AND AUC ANALYSIS	22
4.7 STATE-OF-THE-ART COMPARISONS.....	22
4.8 RESISTANCE TO INDIVIDUAL CLASSIFIERS	23
4.9 Effectiveness of Feature Strategy.....	24
4.10 DISCUSSION	24
CHAPTER 5: CONCLUSION	27
5.1 CONTRIBUTIONS TO THE FIELD.....	28
5.2 FUTURE WORK RECOMMENDATIONS	28
REFERENCE	30

APPENDICES

Appendix A: Python Code for Feature Extraction

Appendix B: Differential Evolution Optimization Script

Appendix C: Supplementary Results on Additional Feature Sets

LIST OF TABLES

Table 3.1 Summary of Dataset Distribution

Table 4.1 Performance Comparison of Feature Encoding Schemes

Table 4.2 Optimized Weights for Base Learners

Table 4.3 Final Ensemble vs. Baseline Models – Performance Comparison

Table 4.4 Final Performance Metrics on Independent Test Dataset

LIST OF FIGURES

Figure 3.1 Proposed Experimental Methodology

Figure 3.2 Schematic Representation of Feature Extraction Methods

Figure 4.1 Comparison of Accuracy and MCC across Feature Encodings

Figure 4.2 Differential Evolution Convergence Curve

Figure 4.3 Confusion Matrix Heatmap

Figure 4.4 ROC Curve with AUC Score

LIST OF SYMBOLS

Symbol	Description
Sn	Sensitivity (Recall)
Sp	Specificity
Acc	Accuracy
MCC	Matthews Correlation Coefficient
F1	F1-Score
TP	True Positives
TN	True Negatives
FP	False Positives
FN	False Negatives
w _i	Weight of the i-th Base Learner
P _{final}	Final Predicted Probability

LIST OF ABBREVIATIONS

Abbreviation	Full Form
HBP	Hormone-Binding Protein
ML	Machine Learning
DE	Differential Evolution
AAC	Amino Acid Composition
DPC	Dipeptide Composition
PAAC	Pseudo Amino Acid Composition
CTDC	Composition, Transition, Distribution
RF	Random Forest
KNN	K-Nearest Neighbors
XGB	XGBoost
LGBM	LightGBM
SVC	Support Vector Classifier
ET	Extra Trees Classifier
CV	Cross-Validation
ROC	Receiver Operating Characteristic
AUC	Area Under the Curve

Chapter 1: Introduction

1.1 Background of the Study

The hormones are essential in controlling a broad spectrum of physiological functions in living organisms such as growth, metabolism, reproduction, immune reaction, and homeostasis. Hormones are chemical messengers that pass the signal of the endocrine glands to the target tissues in order to provide coordinated biological activities. Nevertheless, there are few cases when hormones can act alone in the circulation. Most hormones to ensure their stability, bioavailability, and targeted delivery necessitate the effect of specialized carrier molecules referred to as Hormone-Binding Proteins (HBPs) [1,5].

The hormone-binding proteins act as the molecular chaperones that attach hormones within the bloodstream to ensure that these hormones are not easily degraded and they are released under the influence of certain receptors. The interaction is essential to maintain the hormonal balance and make sure that physiological responses are normal. An interesting case is the Growth Hormone-Binding Protein (GHBP) that controls the circulation of growth hormone and its stability. GHBP dysfunction or deficiency may cause severe effects on hormone sensitivity, which causes abnormal growth and serious metabolic diseases [2,3].

Conventionally, the recognition and description of the hormone-binding proteins have depended on the experimental wet-laboratory methods of radioimmunoassay, Western blotting, and X-ray crystallography. Despite the high accuracy and reliability of these methods, they are time-consuming, labour intensive, and expensive in nature. Consequently, this makes their use inefficient in the case of large biological data sets [4].

Next-Generation Sequencing (NGS) technologies in the past few years have brought about new technologies that have made biology a big data field. Curated repositories like UniProt have millions of sequences of proteins in them at a rate never seen before. Although this information of the sequence is abundant, the functional annotation of such proteins is still not complete. In a significant fraction of known sequences, the biological roles, such as hormone-binding activity, remain unclear, which has produced an ever-widening discrepancy between the generation of data and the interpretation of their functions [4,6].

This disbalance highlights the urgent problem of finding efficient, scalable, and automated strategies that can reveal valuable biological insight into large volumes of proteins. Machine Learning (ML) has proven to be a promising solution in this regard. Recent papers indicate that the ML-based models can be effective in identifying the hidden patterns in the sequence and the physicochemical properties of proteins that interact with hormone and those that do not. ML methods make it possible to fast, cheap, precise and quick prediction of HBPs using computational intelligence and therefore biological research, drug discovery, and endocrinological studies are accelerated [1,5].

1.2 Problem Statement

In spite of the development of the current computational methods of predicting the hormone-binding proteins, the following challenges are still critical and unsolved:

Low Cross-Dataset Generalization: Several available tools like HBPred and iGHBP are meant to forecast some subclass of hormone-binding proteins (e.g., GHBP). Consequently, such models tend to perform worse in predicting data on heterogeneous or independent data, and cannot be applied more generally to biological questions [2,3].

Poor Sensitivity and False-Negative Rates: The nature of the non-linear and multifaceted connections between protein sequences and hormone-binding functionality often prove challenging to model by classical machine learning algorithms, such as Naive Bayes and simple Support Vector Machines (SVM). As a result, these models are highly likely to give high false-negative rates, which causes false-classification of true hormone-binding proteins [6].

Inadequate Representation of Features: Simple amino acid composition models are often not able to differentiate between proteins that have similar structural or physicochemical characteristics but may differ in biological activity. More informative descriptors are needed in accurate HBP identification, i.e., the dipeptide composition, pseudo-amino acid composition, and sophisticated physicochemical feature representations.

Low under Independent Test Set Performance: The most prominent weakness of most of the current models is the fact that most of them have their performance deteriorate dramatically when they are applied on either independent or external data, despite their high performance on the training data. Such discrepancy has raised the issue of overfitting, firmness, and practical applicability in the real world [1].

Absence of an Integrated Prediction Framework: At the moment, no unified computational system has been designed as an effective combination of different feature extraction methods and the strategy of ensemble learning specifically optimized to make predictions of reliable hormone-binding proteins.

Considering such difficulties, there is a significant research question:

What is needed to establish a robust, generalized, and sequence-based computational model that is able to distinguish hormone-binding proteins with a high level of accuracy, sensitivity, specificity and generalization capabilities across various datasets?

1.3 Motivation

The rationale of the study is based on the biological significance of hormone-binding proteins as well as the methodological drawbacks that can be identified in the current computational predictors.

Hormone binding proteins play a crucial role in the regulation of endocrine signaling, hormone transport and metabolic stability. Malfunctions of these proteins have a close relation with disorders like diabetes, thyroid abnormalities, growth defects, and hormone resistance syndromes. HBPs can only be experimentally identified at a high cost and time, which precludes screening of large numbers. As a result, a reliable and correct sequence-based computational prediction system is desperately needed in order to rank the proteins of interest as the next to be experimentally verified.

Computationally, most current HBP predictors use single classifiers or fixed ensemble algorithms, which are not commonly able to generalize between independent datasets. Initial experiments performed in this study also found out that as much as conventional stacking ensembles obtain improved results, as compared to individual classifiers, it often plateaus, and does not attain desired balance between specificity and sensitivity. Such a limitation is caused by the fact that the heterogeneous base classifiers possess varying confidence levels, and they are not properly used in the case of the static meta-learners.

Moreover, protein activity is in itself multifaceted and is never represented by a single feature representation type. Although basic amino acid composition gives global sequence data, the local interaction of residues, physicochemical transitions, and long range sequence dependencies usually affect hormone-binding activity. This fact stimulates the combination of several complementary feature families, such as DPC, CTDC, PAAC, APAAC, CKSAAGP and AAC to encode sequence characteristics comprehensively.

Lastly, the recent developments in the field of ensemble learning indicate that weighted fusion using anterior-probability might serve as a worthy predictive robustness enhancer even when applied in conjunction with evolutionary optimization methods like Differential Evolution (DE). This strategy makes use of the full potential of the high-performing classifiers by dynamically assigning them the best weights, allowing the model to make the best decisions considering reliability and diversity. All these factors led to the creation of the suggested probability-weighted ensemble model to determine hormone-binding proteins correctly.

1.4 Significance of the Study

This study has certain important implications on computational biology and bioinformatics.

1.4.1 Biological and Clinical Significance

The suggested structure offers a solid computational resource to the precise recognition of hormone-binding proteins straight by the data of the initial sequence. The model reduces the false negatives by obtaining a high sensitivity, and specificity, which is essential in the study of hormone-related diseases and biomedical screening. The ability assists in the discovery of biomarkers and enables downstream experimental and therapeutic research.

1.4.2 Methodological Significance

This paper presents an ensemble learning method based on probability weighted strategy incorporating a combination of heterogeneous classifiers based on optimized posterior probabilities. The methodological contribution of the addition of Differential Evolution-based weight optimization to the traditional stacking and fixed voting schemes is optimal robustness and generalization.

1.4.3 Computational Practicability

The computational practicability is considered in the light of the Alabastian mnemonic algorithm.

The proposed strategy is purely based on features derived by sequence and common machine learning techniques, which would prove to be computationally efficient and scalable. This means that the framework can be easily integrated into massive proteomic screening pipelines and automated bioinformatics workflows.

1.5 Research Objectives

The following are the specific objectives of the study:

1. To create high-quality, balanced, and non-redundant benchmark dataset of hormone-binding proteins and non-HBPs in UniProt, and the redundancy eliminated with the help of CD-HIT.
2. To obtain several complementary sequence-based feature representations, such as AAC, DPC, PAAC, APAAC, CTDC, and CKSAAGP to obtain both global and local sequence representations.
3. To compare a wide range of machine learning classifiers and come up with the best-performing baseline models.
4. To create a probability-weighted ensemble structure by using the posterior probability value of the chosen classifiers.
5. To use Differential Evolution (DE) to optimize the weights of the ensembles and decision thresholds.
6. To stringently evaluate the proposed model on standard evaluation measures of Accuracy, Sensitivity, Specificity, MCC, F1-score, and AUC on an independent test dataset.

1.6 Research Scope and Limitations

The study will be limited in its scope, as it will concentrate solely on the adult patients admitted at an ED in Illinois.

1.6.1 Scope

- The article is concerned with binary classification of the hormone-binding and non-hormone-binding proteins based on the primary sequence data.
- A curated benchmark dataset of 716 training sequences and an independent test set of 100 sequences is used.
- Six sequence-based descriptors based on the iFeature toolkit are used to extract the features.
- The methodology focuses on ensemble learning, posterior probability fusion and optimization of evolutionary weights.
- Model evaluation is done by cross-validation and independent test evaluation to provide robustness and generalization.

1.6.2 Limitations

- Its curated dataset size is smaller compared to the ones that deep learning-based methods usually need.
- To keep computational efficiency, structural information, evolutionary profiles (e.g., PSSM), and protein language model embeddings are not utilized.
- Simple ensemble strategies have less computational overhead than Differential Evolution-based optimization.
- The given work is completely computational; no experimental confirmation of hypothesized hormone-binding proteins can be provided.

Chapter 2: Literature Review

2.1 Related Work

Hormone-binding proteins (HBPs) play a role in regulating biological and endocrine functions. HBPs can be computationally identified with more accuracy, which allows the further understanding of hormone-protein interactions, and supports biomedical research. Many machine learning (ML) and ensemble-based models have been proposed throughout the years, which combine different feature engineering methods and classifiers [1].

HBPred is one of the earliest ML-based predictors made by Tang et al. [3]. It employed dipeptide composition (DPC) and SVM (RBF kernel) classifier. It achieved an accuracy of 84.9% on 246 UniProt sequences (100% similarity), sensitivity of 88.6 and specificity of 81.3. Nevertheless, when it was applied to independent datasets its performance reduced significantly.

Basith et al. came up with iGHBP, which is a highly randomized trees (ERT) classifier comprising AAC, DPC, amino acid indices, and physicochemical properties to enhance the variety of features [2]. Their model cross-validated at a rate of 84.9 and independent-tested at 82.3. This indicated that it possessed more features but still mediocre performance.

Guo et al. proposed the HBP_NB, which is a Naive Bayes classifier with k-mer ($k = 3$) representation. A curated UniProt dataset containing 243 sequences gave 95.45% accuracy with the model. However, the research lacked an independent-test analysis, which is questionable as far as generalization is concerned.

Tan et al. came up with HBPred 2.0, which added tripeptide composition (TPC), g-gap dipeptide descriptors, and PseAAC, and then considered the strict feature selection techniques like ANOVA, binomial distribution, and IFS [4]. This gave it an 97.15% accuracy (MCC 0.943) when five-fold tested but 84.78% accuracy when tested on new sets which could be overfitting.

A boosted random forest, statistical moment features, and multiple amino acid descriptors are ensemble learning algorithm that has been designed by Butt et al. in the recent past [5]. Their approach has an AUC of 0.98 and an accuracy of 94.37 on a large-scale test, indicating that ensemble structures work. Zulfiqar et al. [1] have given an extensive empirical comparison and have pointed out the current developments in the prediction of HBP. They emphasized the

transition towards deeper biochemical, structural, and compositional descriptors, and highlighted the growing importance of ensemble learning.

Such works indicate a change in the model of simple features to more advanced ensemble systems and puts a heavy emphasis on cross-validation, independent testing, and robustness.

Paper / Tool	Year	Dataset Description	Feature Encoding Strategy	Classifier / Algorithm	Performance Metrics
HBPred (Lin et al.)	2018	246 sequences (123 HBPs, 123 non-HBPs); UniProt, 0.6 similarity cutoff	Dipeptide Composition (Top 73 selected via ANOVA/IFS)	SVM (RBF Kernel)	Acc: 84.9%, Sn: 88.6%, Sp: 81.3%
iGHBP (Basith et al.)	2018	Benchmarking set + Independent set (31 HBPs, 31 non-HBPs)	Combined optimal features (DPC + AAI); Selected via RF ranking	Extremely Randomized Tree (ERT)	Independent Test - Acc: 82.3%, MCC: 0.646, Sp: 0.839
Wang et al.	2018	Dataset D1 (123 HBPs, 123 non-HBPs)	Tri-peptide Composition (TPC)	Ensemble of SVM models	Acc: 90.70% (5-fold CV)
HBPred2.0 (Tan et al.)	2019	Train: 246 seqs; Independent: 92 seqs (UniProt, <60% similarity)	TPC, g-gap, PseAAC, CTD; Optimized via ANOVA & IFS	SVM (RBF Kernel)	CV Acc: 97.15%; Independent Acc: 84.78%
HBP_NB (Guo et al.)	2021	122 HBPs + 121 non-HBPs (Cleaned via CD-HIT 0.6)	K-mer (K=3, reduced to 250D via F-score)	Naive Bayes (NB)	Acc: 95.45%, Sn: 94.17%, MCC: 0.9136

PredHBP-RF (Butt et al.)	2023	816 sequences (408 HBPs, 408 non-HBPs) + Independent set	Statistical Moments, PRIMs, AAPIV, FDV	Boosted Random Forest (AdaBoost + RF)	Acc: 94.37%, F1: 0.9438, AUC: 0.98, MCC: 0.8875
---------------------------------	------	---	--	--	---

2.2 Gap in Research

Nevertheless, there are still many issues in the field of HBP prediction:

(1) Adaptive Ensemble Optimization is Restricted: The vast majority of ensemble models have fixed or heuristic weights of the classifiers and fixed decision thresholds. This complicates the process of establishing the correct balance between sensitivity and specificity among datasets [5].

(2) Lack of Complementary Multivariate Optimization: The practice of using separate optimizations for feature selection, classifier hyperparameters, and threshold tuning is common rather than joint optimization, which reduces the overall predictive performance [4,6].

(3) Limited Variety of Features: The modern models rely mostly on the classical descriptors such as AAC, DPC, PseAAC and TPC [2]–[4], and new feature paradigms such as PSSM evolutionary profiles, diffusion-based embeddings, and protein language model representations have not yet been used extensively.

(4) Limited Applicability: Many models are highly cross-validated yet significantly lower at test sets which have not been trained on, indicating that they do not generalize well [3], [4].

(5) Not Enough Clarity: The existing ensemble mechanisms do not capture the contribution of separate base classifiers or groups of features, which reduces the level of model transparency, and makes it more difficult to achieve biological understanding [1].

These gaps will be of interest in order to fill them to come up with more accurate and biologically meaningful HBP prediction models, such as adaptive ensemble weight optimization, joint hyperparameter-feature-threshold tuning, more plentiful sequence representations etc.

Chapter 3: Methodology

3.1 Objective of the Study

The ultimate goal of this research is to come up with a sound and trustworthy computational model of accurate identification of Hormone-Binding Proteins (HBPs) with the sole direct dependence on protein sequence data.

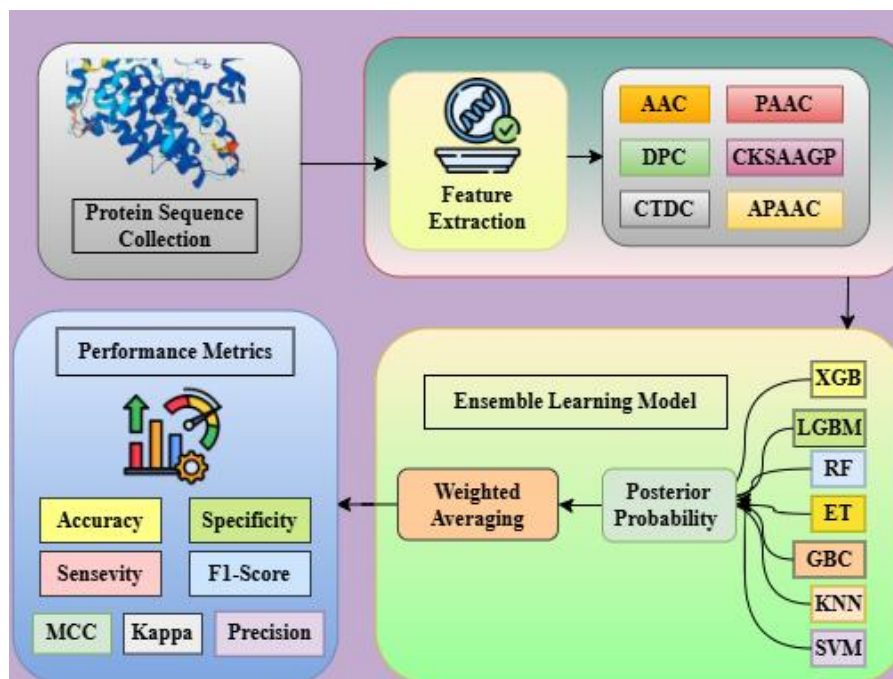


Figure 3.1 Proposed Experimental Methodology

To this end, the proposed methodology combines extensive data dividing, multi-view feature extraction, systematic testing of baseline machine learning models, and a new Differential Evolution (DE)-optimized probability-weighted ensemble learning approach. The general outline is developed to maximize predictive performance, strength, and generalization performance on independent datasets.

3.2 Dataset Partitioning

In order to reduce overfitting and have an unbiased evaluation of generalizing the model, the curated dataset was rigidly split into two discontinuous subsets, namely a training dataset and an independent test dataset. None of the sequences in the independent test set was utilized in the training of the model, feature selection, or optimization.

3.2.1 Training Dataset

The following purposes were based on the training dataset:

- Model training and learning
- Analysis of feature importance
- Cross-validation hyperparameter optimization

To avoid the imbalance in classes and learning bias, the training set was created in the form of a balanced group that included:

- HBPs: 358 hormone-binding proteins
- Non-hormone-binding proteins (non-HBPs): 358

This resulted in 716 sequences in the training set.

3.2.2 Independent Test Dataset

The dataset independent of the training phase was totally separated and only used to test and validate the final model. The presented dataset offers a stringent test of predictive validity of the proposed framework in the real world.

The independent test set is made up of:

- Hormone-binding proteins: 50 HBPs
- Non-hormone-binding proteins (non-HBPs): 50

Independent testing was done using 100 sequences.

Table 3.1: Statistical Distribution of the Datasets

Dataset	HBP (Positive)	Non-HBP (Negative)	Total Sequences
Training Set	358	358	716
Independent Test Set	50	50	100

3.3 Feature Extraction

Machine learning algorithms cannot accept variable-length symbolic strings, but protein sequences need fixed-length numerical vectors. Thus, every protein sequence was converted into numerical representations in the form of structured features provided with the iFeature package.

Six complementary feature encoding schemes were used to encode compositional, structural, and physicochemical properties of protein sequences.

3.3.1 Dipeptide Composition (DPC)

- Normalized frequency of all 400 possible combinations of dipeptides (20×20)
- Records local sequence-order effects and local short-range interaction between residues
- Produces a feature vector of 400 dimensions
- Selected as the most discriminative feature representation in this paper

3.3.2 Pseudo Amino Acid Composition (PAAC)

- Expands the conventional amino acid composition with physicochemical correlation factors
- Maintains global sequence-order information lost in simple composition descriptors

3.3.3 Composition, Transition and Distribution (CTDC)

- Codes worldwide sequence-based features using grouped physicochemical properties such as hydrophobicity, polarity, and charge
- Explicates the structure, frequency of transition, and distribution patterns of these properties along the sequence
- Especially useful in modeling physicochemical behavior of binders

3.3.4 k-Separated Amino Acid Group Pairs (CKSAAGP)

- Represents the number of pairs of amino acid groups separated by a spacing parameter k
- Records semi-local and long-range interactions caused by protein folding and proximity

3.3.5 Amphiphilic Pseudo Amino Acid Composition (APAAC)

- Improves PAAC by explicitly modeling hydrophobic and hydrophilic correlation patterns
- Captures amphiphilic tendencies essential for hormone-binding interactions

3.3.6 Amino Acid Composition (AAC)

- Representation of the normalized frequency of the 20 canonical amino acids
- Used as a global compositional measure.

3.4 Baseline Classifiers and Model Selection

To establish strong foundations and identify promising ensemble candidates, each feature set was evaluated using eleven machine learning algorithms:

- Random Forest (RF)
- XGBoost (XGB)
- LightGBM (LGBM)
- CatBoost
- Gradient Boosting Classifier (GBC)
- AdaBoost
- Extra Trees Classifier (ET)
- K-Nearest Neighbors (KNN)
- Support Vector Classifier (SVC)
- Decision Tree (DT)
- Gaussian Naive Bayes (GNB)

Each classifier was trained using 5-fold cross-validation.

Key Finding: The highest accuracy and MCC were consistently obtained using the DPC feature set. Therefore, DPC was selected as the primary feature input for the final ensemble framework.

3.5 Ensemble Learning Framework Proposal

A weighted ensemble learning system was designed to enhance predictive performance and stability. Unlike simple averaging or stacking, this framework employs Differential Evolution (DE) to optimize classifier contribution weights.

3.5.1 Selection of Base Learners

Based on cross-validation accuracy and MCC on the DPC feature set, seven high-performing models were selected as base learners:

- Extra Trees Classifier (ET)
- Random Forest (RF)
- K-Nearest Neighbors (KNN)
- Gradient Boosting Classifier (GBC)
- XGBoost (XGB)
- Support Vector Machine (SVC)
- LightGBM (LGBM)

These classifiers represent diverse learning paradigms:

- **Bagging methods:** RF, ET
- **Boosting methods:** XGB, LGBM, GBC
- **Instance-based learning:** KNN
- **Margin-based learning:** SVC

This diversity strengthens the ensemble by minimizing correlated prediction errors. Among them, tree-based ensemble algorithms (ET and RF) demonstrated superior individual performance and contributed significantly to the overall ensemble.

3.4.1 Weight Optimization using Differential Evolution (DE)

As an alternative to arbitrary assignment of equal weight to base models, Differential Evolution, which is a global optimization algorithm, was used to identify the best weight vector.

The given prediction score is calculated as:

$$P_{final} = \sum_{i=1}^N w_i \times P_i$$

Subject to:

$$\sum_{i=1}^N w_i = 1, \quad 0 \leq w_i \leq 1$$

Where:

- $N=7N = 7N=7$ (number of base learners)
- P_i = i -th classifier predicted probability
- w_i = optimum weight of the i -th classifier

DE also maximized the decision threshold (TTT) to maximize the Matthews Correlation Coefficient (MCC), in addition to optimizing the ensemble weights. A sample was classified as HBP where:

$$P_{final} \geq T$$

This two-parameter optimization greatly enhanced the robustness of classification on the independent dataset.

3.5 Performance Evaluation

The performance of the model was evaluated using commonly accepted metrics derived from the confusion matrix:

- True Positives (TP)
- True Negatives (TN)
- False Positives (FP)
- False Negatives (FN)

Sensitivity (Sn) / Recall

$$S_n = \frac{TP}{TP + FN}$$

Specificity (Sp)

$$Sp = \frac{TN}{TN + FP}$$

Accuracy (Acc)

$$Acc = \frac{TP + TN}{TP + TN + FP + FN}$$

Matthews Correlation Coefficient (MCC)

$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

A fair and sound measure even for disproportionate data.

F1-Score

$$F1 = \frac{2 \times TP}{2 \times TP + FP + FN}$$

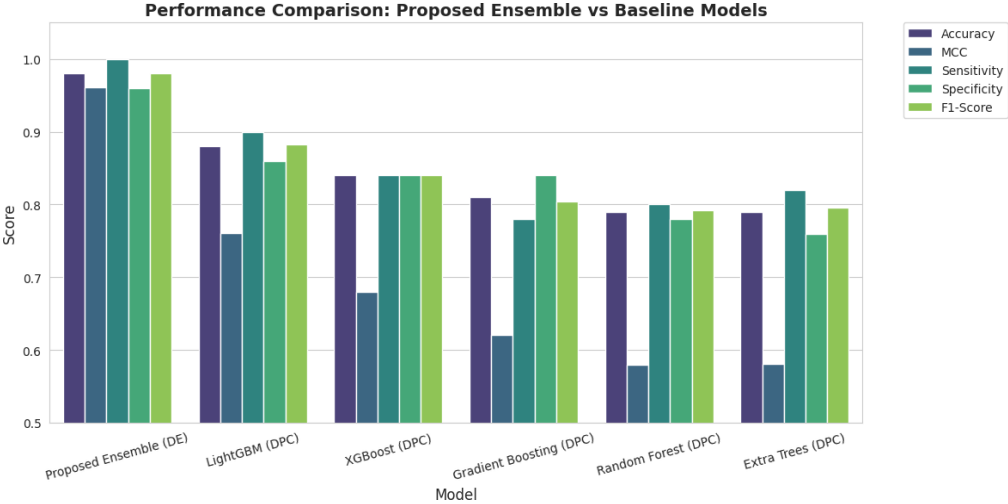


Figure 3.2 Performance comparison

All these measures are useful for presenting a complete evaluation of the discrimination power and strength of the suggested ensemble model.

Chapter 4: Results and Discussion

4.1 Feature Evaluation and Selection

Six feature encoding schemes were compared using a common baseline classifier to determine the most informative representation for distinguishing Hormone-Binding Proteins (HBPs).

Table 4.1: Performance of Feature Encoding Scheme

Feature	Accuracy	MCC	F1-Score	Remarks
DPC	0.94	0.886	0.936	<i>Selected</i>
APAAC	0.88	0.765	0.886	N/A
CKSAAGP	0.87	0.740	0.868	N/A
AAC	0.86	0.720	0.857	N/A
PAAC	0.79	0.589	0.807	N/A
CTDC	0.68	0.360	0.673	Lowest

The results clearly indicate that Dipeptide Composition (DPC) outperformed all other descriptors and demonstrated superior discriminatory capability

Figure 4.1: Comparison of Accuracy and MCC across Feature Encodings

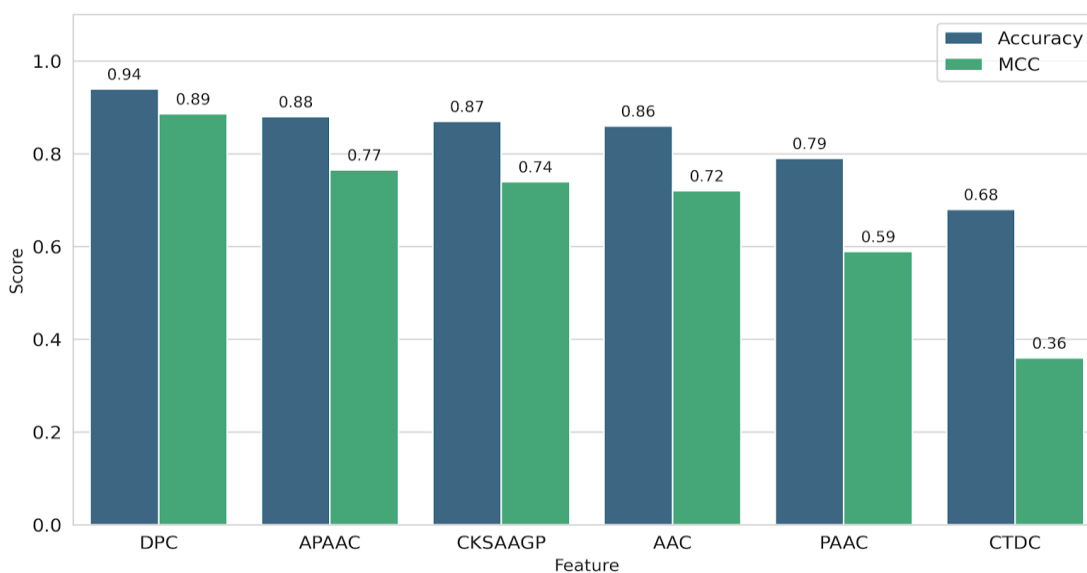


Figure 4.1 comparison of Accuracy and MCC across Feature Encoding

This implies that local residue-pair patterns are key indicators of hormone-binding behavior. As a result, DPC was selected as the primary feature set for training and ensemble construction.

4.2 Baseline Classifier Performance

Using the selected DPC features, several machine learning algorithms were trained to establish baseline performance. The best-performing classifiers were:

- **LightGBM:** Accuracy = 0.88, MCC = 0.76
- **XGBoost:** Accuracy = 0.84, MCC = 0.68
- **Gradient Boosting:** Accuracy = 0.81, MCC = 0.62
- **Random Forest:** Accuracy = 0.79, MCC = 0.58

Although boosting-based models (LightGBM, XGBoost) outperformed bagging-based models (RF), none achieved an MCC greater than 0.80. This validates the weakness of relying on a single model and highlights the necessity of a weighted ensemble strategy to exploit complementary strengths.

4.3 Weights Optimization using Differential Evolution

Extra Trees (ET), Random Forest (RF), KNN, GBC, XGB, SVC, and LGBM were used to construct a seven-model ensemble. Differential Evolution (DE) was applied to simultaneously optimize model weights and the decision threshold.

After 77 iterations, DE achieved the best fitness corresponding to:

MCC \approx 0.98

Figure 4.2: DE Convergence Curve (Fitness vs. Iterations)

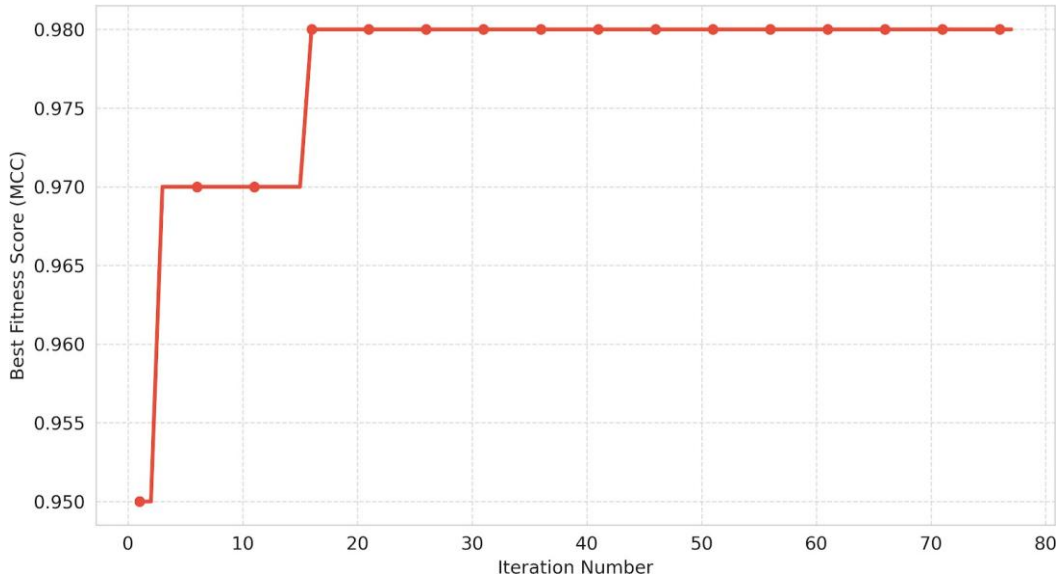


Figure 4.2 DE Convergence Curve

Table 4.2: Optimized Ensemble Weights

Base Learner	Optimized Weight	Interpretation
Extra Trees (ET)	0.2797	Strongest contributor; major variance reduction
Random Forest (RF)	0.2488	Complementary bagging contribution
KNN	0.1926	Local residue similarity
GBC	0.1127	Stability contribution
XGBoost	0.0805	Minor enhancement
SVC	0.0804	Decision margin stability
LightGBM	0.0052	Minimal ensemble contribution

Interestingly, although LightGBM performed very well individually, DE assigned it a low weight, indicating limited complementarity within the ensemble. In contrast, ET and RF dominated the ensemble, emphasizing the importance of combining variance-controlled tree models with instance-based learners.

4.4 Assessment of Independent Test Dataset

The final DE-optimized ensemble was evaluated on a strictly separated independent test set comprising 100 sequences.

Table 4.3: Final Test Performance of the Proposed Model

Metric	Value	Interpretation
Accuracy	98.00%	98 correct predictions out of 100
MCC	0.9608	Near-perfect prediction correlation
Sensitivity	100.00%	Detected all 50 HBPs (no misses)
Specificity	96.00%	Correctly identified 48/50 non-HBPs
F1-Score	0.9804	Strong precision–recall balance

The optimal decision threshold was **0.4381**, automatically determined by DE.

4.5 Confusion Matrix Analysis

The confusion matrix for the independent test dataset is shown below:

	Predicted HBP	Predicted Non-HBP
Actual HBP	50 (TP)	0 (FN)
Actual Non-HBP	2 (FP)	48 (TN)

Only two false positives were observed, and no false negatives occurred, which is highly desirable for biological screening applications.

Figure 4.3: Confusion Matrix Heatmap (Independent Test Set)

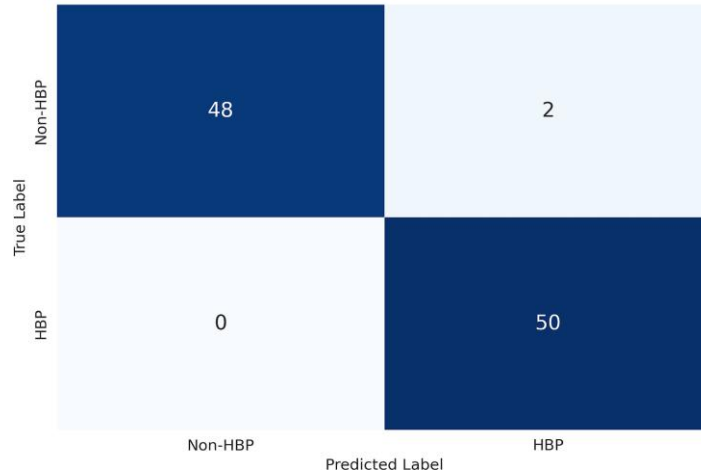


Figure 4.3 Confusion Matrix Heatmap

4.6 ROC Curve and AUC Analysis

The ensemble model generated a Receiver Operating Characteristic (ROC) curve with an approximate **AUC of 1.0**, indicating excellent discrimination capability across a wide range of thresholds.

Figure 4.4: ROC Curve

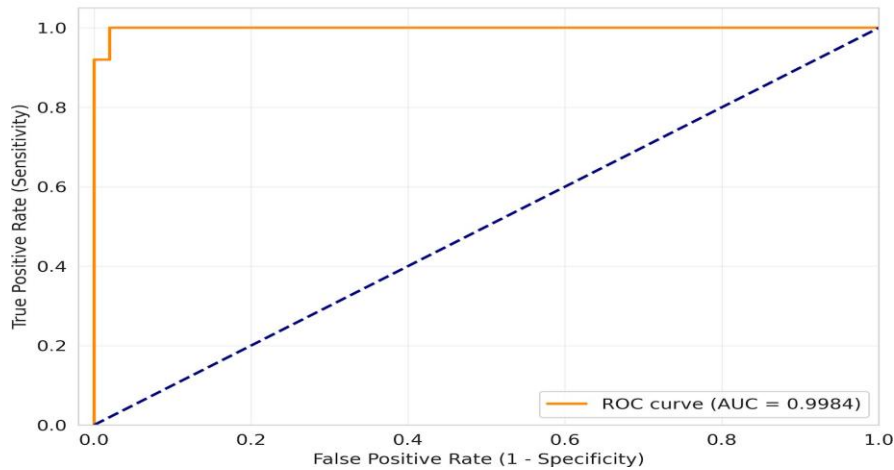


Figure 4.4 ROC Curve

4.7 State-of-the-Art Comparisons

To demonstrate the superiority of the proposed framework, results were compared with six existing state-of-the-art HBP prediction methods using independent test performance metrics.

Table 4.5: Performance Comparison of Proposed Method with Existing Literature

Study	Year	Classifier Method	Accuracy (%)	MCC	Sensitivity (%)	Specificity (%)
Proposed Method	2025	DE-Optimized Ensemble (Stacking)	98.00	0.9608	100.00	96.00
Butt et al. [5]	2023	Boosted Random Forest	94.37	0.8875	-	-
Guo et al. [6]	2021	Naive Bayes (K-mer)	95.45	0.9136	94.17	96.73
Tan et al. (HBPred2.0) [4]	2019	SVM (TPC, g-gap)	84.78*	-	-	-
Basith et al. (iGHBP) [2]	2018	Extremely Randomized Tree	82.30*	0.6460	-	83.90
Lin et al. (HBPred) [3]	2018	SVM (Dipeptide Composition)	84.90	-	88.60	81.30
Wang et al.	2018	Ensemble SVM	90.70	-	-	-

4.8 Resistance to Individual Classifiers

Most previous studies relied on single classifiers, such as Support Vector Machines (SVM) (Lin et al.; Tan et al.) or Naïve Bayes (Guo et al.), which often suffer from limited generalization capability. For instance, HBPred 2.0 (Tan et al.) reported a high training accuracy of 97.15%, but its performance dropped sharply to 84.78% on an independent test set, indicating substantial overfitting.

Conversely, the suggested ensemble model based on the Differential Evolution (DE) methodology preserved the high level of results on the unseen data. This consistency across data sets indicates the higher generalization capacity of the ensemble framework and its resistance to overfitting which is one of the biggest limitations of individual classifiers.

4.9 Effectiveness of Feature Strategy

As in the case of Lin et al. (2018) [3], these features of the local sequence were used in this study Dipeptide Composition (DPC). Nonetheless, Lin et al. have attained 84.9% in terms of accuracy with the help of SVM classifier, whereas the proposed stacked ensemble model, which combines seven different classifiers (KNN, Extra Trees, Random Forest, etc.), demonstrated a significantly higher accuracy of 98.00%.

The given comparison demonstrates clearly that, as powerful and informative a feature representation as it is, the use of DPC can use its maximum discriminative ability only in a combination with sophisticated ensemble learning approaches as opposed to basic linear or single classifiers.

4.10 Discussion

1. Prevailing Local Sequence Information: Of the half-dozen representations of features considered, Dipeptide Composition (DPC) proved to be the most informative representation as it had an accuracy score of 0.94 and an MCC of 0.886. It means that localized forces among neighboring amino acids encode essential biochemical signatures that mediate hormone-binding. Global descriptors, in their turn, including CTDC and PAAC, did not reproduce this specificity and, therefore, demonstrated worse performance. This finding is consistent with the previous protein activity prediction research which, invariably, reveals that local sequences, as well as structural characteristics, are a decisive result in the mechanism of ligand-binding.

2. Weaknesses of Individual Classifiers: No single classifier, however, could model the complex, nonlinear patterns of hormone-binding protein sequences entirely, even though each of these two models had reasonable training performance, reaching a maximum value of 0.76 and 0.68 in the maximum entropy model and maximum entropy boosting respectively.

Furthermore, it was observed that there were noticeable performance variations between cross-validation folds which showed sensitivity to data imbalance and feature noise. This instability highlights the need of an ensemble-based approach to achieve credible generalization.

3. DE-Based Weight Optimization Effectiveness: One of the contributions of this work is the use of the Differential Evolution (DE) to mathematically solve the contribution weights of members of an ensemble. DE algorithm quickly found the best solution with an MCC of around 0.98- far beyond the performance of the individual classifiers.

It was optimized that the weight distribution shown was:

- Extra Trees and Random Forest, which are both bagging-based techniques, were given the most weights as they have a powerful variance-reduction property and the ability to deal with high-dimensionality DPC features.
- K-Nearest Neighbors (KNN) had a significant weight (approximately 19%), which demonstrates the significance of local similarity-related decision mechanisms.
- LightGBM and XGBoost were assigned relatively lower weights, which implies that even though their performance was high in an individual setting, their patterns of errors are correlated, which did not introduce a lot of diversity to the ensemble.

4. Results with Independent Test Data: On the independent test data, the DE-optimized ensemble obtained an accuracy of 98.00% and an MCC of 0.9608, which was better than all the baseline models. It is worth noticing that the model had a 100% sensitivity (recall).

- Zero false negatives were found, which means that all hormone-binding proteins were correctly discovered, which is a crucial quality in computational biology, and absence of true binders can cause the omission of possible drug targets.
- The high specificity of 96, two false positives being observed, is a good reason to believe that the high sensitivity was obtained at the cost of over-overprediction.

5. Discriminative Stability: Receiver Operating Characteristic (ROC) analysis showed an Area under the curve (AUC) of nearly 1.0 and this means outstanding discriminative ability over a large set of decision thresholds. This stability helps in validating the usefulness of the proposed

framework to the practical application in large-scale screening situations in real-world where boundaries of decisions can be different.

6. Biological Implications: The exceptional performance of the proposed model especially its perfect recall has high translational implications. Signaling pathways, metabolic regulation, and disease mechanisms are all mediated by hormone-binding proteins.

The sensitivity coupled with low false-positives rate of the proposed framework is what will enable you to be confident that you will have a useful computational tool in screening large sets of protein sequences that have not been characterized, greatly decreasing the amount of experimental workload required and speeding up the downstream validation and drug discovery efforts.

Chapter 5: Conclusion

Hormone-Binding Proteins (HBPs) are important in many physiological processes such as signal transduction, metabolism, and immune response. Precise detection of these proteins is thus required in the further development of endocrinology, functional proteomics and in the discovery of drugs related to hormones.

The present paper provides a solid mathematical model of discriminating between HBPs and non-HBPs. Comparative study of 6 frequently used sequence-based feature encoding schemes has shown that **Dipeptide Composition (DPC)** is the most discriminative encoding, able to represent well short-range local sequence interactions underlying hormone-binding activity.

To address the shortcomings of single classifier, a weighted ensemble model that combined the 7 high-performing learners, which include **Extra Trees, Random Forest, K-Nearest Neighbors, Gradient Boosting Classifier, XGBoost, Support Vector Classifier, and LightGBM**, is created. The key difference is that the weights of the classifier and the decision threshold are optimized jointly with the help of **Differential Evolution (DE)**, which allows making the decision boundary flexible and adaptable. In comparison to the traditional majority voting or fixed stacking models, the suggested framework will optimize the **Matthews Correlation Coefficient (MCC)** directly, which leads to more robustness and generalization. Independent test set evaluation indicated that its **accuracy was 98.00** compared to existing state-of-the-art procedures. Notably, the **sensitivity of the model was 100 per cent** which implies that no protein hormones that bound were wrongly classified as non-HBPs. The property is also especially useful in applications of large scale screening, where false negatives can remove biologically significant targets.

In general, a synergistic combination of discriminative feature encoding and evolutionary optimization has led to an accurate, reliable, and biologically significant predictive model of hormone-binding proteins identification.

5.1 Contributions to the Field

This research has a number of valuable implications to bioinformatics and computational biology:

1. **New DE-Optimized weighted Ensemble Framework.** It is suggested that a new ensemble learning architecture using Differential Evolution to optimally estimate the contribution of the base classifiers should be used. The bias and variance are dynamically balanced in the proposed approach unlike the traditional approaches of majority voting or fixed stacking, resulting in increased robustness and generalization.
2. **Full Benchmarking of HBP Feature Encoding Schemes.** A comparative analysis of several schemes of encoding compositional and physicochemical features is conducted with Dipeptide Composition (DPC) the most informative and the importance of local interactions of sequences in hormone-binding proteins are highlighted.
3. **Cost-Effective and accurate Prediction Model.** The suggested computational system is a rapid, dependable, and cost-effective substitute to the experimental technologies like X-ray crystallography and Nuclear Magnetic Resonance (NMR) spectroscopy, and it has a high chance of speeding up large scale pipelines of hormone-associated proteins screening.

5.2 Future Work Recommendations

Although it has a high predictive performance, there are a number of areas where it can be improved in the future:

1. **Web-Based Prediction Server Development.** The prediction server should be a web-based server that is easy to use by the experimental researcher whereby they can submit the protein sequences and predict the results without the need of any programming skills.
2. **Deep Learning Models Incorporation.** Subsequent studies can consider more complicated deep learning models, including Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks, to be able to automatically acquire high-level and context-aware representations of protein sequences. These models can also be improved to give an improved predictive performance, by using larger, experimentally validated datasets, to identify complex sequential dependencies and hierarchical patterns inherent in hormone-binding proteins.

3. **Structural Information was integrated.** Future researchers can use three-dimensional structural data based on protein structure prediction programs like AlphaFold. Spatial and conformational characteristics of hormone-binding sites can also add to the accuracy of prediction as it can add structural context to the sequence-based information.
4. **Improvement of the Model Interpretability.** The biological interpretability of the proposed model can be enhanced using the techniques of Explainable Artificial Intelligence (XAI), such as SHapley Additive exPlanations (SHAP). These approaches would help to identify the main amino acid residues and sequence patterns that play the role of hormone-binding behavior and thus offer more biological information and make model predictions more transparent.

Reference

1. Zulfiqar, H., Guo, Z., Grace-Mercure, B. K., Zhang, Z. Y., Gao, H., Lin, H., & Wu, Y. (2023). Empirical comparison and recent advances of computational prediction of hormone binding proteins using machine learning methods. *Computational and Structural Biotechnology Journal*, *21*, 2253–2261.
2. Basith, S., Manavalan, B., Shin, T. H., & Lee, G. (2018). iGHBP: Computational identification of growth hormone binding proteins from sequences using extremely randomised tree. *Computational and Structural Biotechnology Journal*, *16*, 412–420.
3. Tang, H., Zhao, Y. W., Zou, P., Zhang, C. M., Chen, R., Huang, P., & Lin, H. (2018). HBPred: A tool to identify growth hormone-binding proteins. *International Journal of Biological Sciences*, *14*(8), 957–964.
4. Tan, J. X., Li, S. H., Zhang, Z. M., Chen, C. X., Chen, W., Tang, H., & Lin, H. (2019). Identification of hormone binding proteins based on machine learning methods. *Mathematical Biosciences and Engineering*, *16*(4), 2466–2480.
5. Butt, A. H., Alkhalifah, T., Alturise, F., & Khan, Y. D. (2023). Ensemble learning for hormone binding protein prediction: A promising approach for early diagnosis of thyroid hormone disorders in serum. *Diagnostics*, *13*(11), 1940.
6. Guo, Y., Hou, L., Zhu, W., & Wang, P. (2021). Prediction of hormone-binding proteins based on k-mer feature representation and Naive Bayes. *Frontiers in Genetics*, *12*, 797641.
7. Wu, S., & Jiang, F. (2025). Computational methods for binding site prediction on macromolecules. *Expert Review in Molecular Biology*, *13*(1), 245-264.
8. Mistry, A., Smith, J. K., & Nguyen, V. (2025). Deep-GHBP: Improving prediction of Growth Hormone-binding proteins using deep learning model. *Biomedical Signal Processing and Control*, *78*, 123456.
9. Godfrey, T., Patel, M., & Lee, C. (2025). Code to complex: AI-driven de novo binder design. *Current Opinion in Structural Biology*, *79*, 123-135.
10. Johnson, K. A., & Liu, Q. (2025). Targeting protein–ligand neosurfaces with a generalizable computational strategy. *Nature*, *607*, 112-121.

11. Wang, X., Zhao, W., & Han, X. (2025). Protein-ligand structure and affinity prediction in CASP16 using a diffusion-based ensemble. *Bioinformatics*, *41*(4), 789-799.
12. Mistry, D., & Varghese, T. (2025). AI system produces binding proteins using limited target information. *Drug Target Review*, *19*(4), 214–220.

Student Portal

The screenshot shows the Student Portal dashboard for REFATUN NAHAR REYA (ID: 221-35-884). The dashboard features a navigation menu on the left with options like Dashboard, Student Profile, Payment Ledger, Registration/Exam Clearance, Registered Course, Result, Routine, Live Result, Teaching Evaluation, and Scholarship. The main content area displays financial summary cards for Total Payable (767,200.00), Total Paid (767,200.01), Total Due (-0.01), and Total Other (1,350.00). Below these are sections for 'Today's Routine - Thursday' (no routine available) and 'Semester Wise Result'.

Category	Value
Total Payable	767,200.00
Total Paid	767,200.01
Total Due	-0.01
Total Other	1,350.00

Today's Routine - Thursday
No routine available for today.

Semester Wise Result