

Comparative Analysis of an AI-Based Diagnostic System for Automated Detection of Diabetic Retinopathy

¹Richa Singal, ²Mohanjit Kaur Kang, ³Randeep Kaur

¹Research Scholar, ²Assistant Professor, ³Assistant Professor

¹Department of CS/IT

¹Guru Nanak Dev Engineering College, Ludhiana, India

richasingal94@gmail.com, mohanjitkaur6@gmail.com, bhatiarandeep20@gmail.com

Abstract— This work is a new comparative study of machine learning algorithms for automated diabetic retinopathy (DR) detection from retinal images. It is designed and tested a novel AI-based system for diagnosis using our new Multi-Feature Weighted Ensemble (MFWE) framework that utilized multiple publicly available datasets such as MESSIDOR, KAGGLE EyePACS, APTOS 2019, and IDRiD. Our method involved extensive preprocessing of data, feature extraction, and the execution of five different algorithms for classification: Linear Regression, Random Forest, XGBoost, MLP Classifier, and Decision Tree Classifier.

The novel preprocessing process involved image standardization, color normalization, noise reduction, contrast enhancement, and anatomical structure segmentation. The derived unique features such as morphological features, vessel measurements, optic disc features, texture features, color features, and wavelet-based features. Dimensionality reduction was done using Principal Component Analysis and Recursive Feature Elimination with Cross-Validation.

Experiments were carried out on the Messidor dataset of 2,302 samples with 19 features extracted. Models were evaluated by using such metrics as accuracy, precision, recall, F1-score, and ROC curve. Our implementation of Feature-Weighted Ensemble turned out to be distinctive with outstanding performance compared to traditional methods. Tree-based models yielded best results with XGBoost showing the highest F1-score (0.911) and accuracy (0.907), closely followed by Decision Tree (F1: 0.901, accuracy: 0.896).

The novelty in our approach comes from combining MFWE with best-performing hyperparameters, illustrating improved precision-recall balance along with lower overfitting versus baseline implementations. The good performance of relatively naive models indicates our extracted features worked well in selecting informative patterns towards DR classification.

Our results show that our novel approach integrating tree-based models with feature-weighted ensemble methods is very effective for DR detection. This research adds a new method to the construction of trustworthy AI-based screening tools that may help ophthalmologists in early detection of diabetic retinopathy.

Index Terms— Diabetic Retinopathy Detection, Machine Learning Classification, Medical Image Analysis, Feature Extraction, XGBoost, Retinal Fundus Images, AI-based Diagnosis

I. INTRODUCTION

Diabetic retinopathy (DR) continues to be one of the main causes of avoidable blindness worldwide, occurring in about one-third of the 463 million individuals with diabetes mellitus. The insidious, progressive microvascular complication becomes apparent only when advanced stages are reached and treatment options are limited (Albahli and Ahmad Hassan Yar, 2022). Early detection via routine screening is essential, but conventional approaches are severely hampered by specialist shortages, restricted access in underserved communities, and labor-intensive manual image interpretation.

By 2045, diabetes mellitus is expected to be approximately 700 million and for healthcare systems to efficiently manage DR screening, solutions that are scalable are needed. This has boosted the need for an automatic DR diagnostic system that uses artificial intelligence to analyze and classify DR from retinal fundus images (Albahli and Ahmad Hassan Yar, 2022). The last few years have witnessed great development in deep learning especially in convolutional neural networks, which demonstrates a high likelihood of accuracy in diagnosing various disease offenses similar to human experts in clinical settings.

There are three main approaches to diagnose DR: merging ophthalmology, computer vision, and machine learning (Albahli and Ahmad Hassan Yar, 2022). These systems mainly use digital retinal images to locate different patterns of characteristic abnormalities such as microaneurysms, hemorrhages, exudates and other elements such as cotton wool spots and neovascularization. An ability to identify these patterns, a machine is also able to measure DR severity on such scales as no DR or proliferative DR.

The creation of durable AI models requires that the data that are used were of high quality and varied. EyePACS, MESSIDOR, IDRiD, and APTOS are the most commonly used standardized databases for the development of the algorithm. The size of the datasets, the proportion of the datasets originating from different demographics and races (Aziz, Charoenlarnpopparut and Mahapakulchai, 2023), the quality of the image, and the extent of the annotation all affect the generalizability of the models, and more comprehensive analyses are needed to compare these datasets.

While the concepts have shown great potential in clinical trials, practical applications that involve using these tools are challenging for some of the following reasons The algorithms' operational features might be unclear; the datasets originating from diverse populations can be biased; there might be doubts as to whether success in one population can be achieved in another;

integrating into the routine clinical procedures; regulatory approval; and clinical acceptance (Aziz, Charoenlarnnoppa and Mahapakulchai, 2023). Also, identifying the correct performance measures is crucial, for example, accuracy used in ordinary applications may not be suitable for clinical applications where false negativity and falseness are equally critical.

This study provides a systematic review of the state-of-art methods for the automated DR diagnosis by comparing the ML techniques using standard datasets and by considering both diagnostic performance and clinical relevance (Aziz, Charoenlarnnoppa and Mahapakulchai, 2023). It covers both conventional machine learning with engineered features and deep learning models capable to learn the representations from raw images directly.

The relevance goes beyond the purely technical solution of the problem to contribute to addressing the issues related to public health. Advanced and efficient screening tools may solve the issue of creating access to the quality eye care which could help in increasing the untimely blurred vision solutions which otherwise result to blindness (Aziz, Charoenlarnnoppa and Mahapakulchai, 2023). Through comparing performance of the classifiers developed in this research with the multiparametered validity indices, the research will contribute to the clinical validation of the AI-supported screening for DR and, therefore, help to design and create more accurate and reliable systems of the diabetic eye care to improve people's lives in all the countries of the world.

Aim and Objectives

Aim

This research aims to develop and evaluate an optimized AI-based diagnostic system for automated detection of diabetic retinopathy through comprehensive analysis of existing methodologies, development of machine learning models, and comparative performance assessment to improve early detection and clinical outcomes.

Objectives

The following objectives have been formulated to achieve the aim of this research:

1. To conduct a comprehensive analysis of research papers and standardized medical image datasets for diabetic retinopathy detection, focusing on current methodologies and their effectiveness.
2. To develop machine learning model for diabetic retinopathy detection, optimizing its parameters to achieve high accuracy and reliable diagnostic results.
3. To perform comparative analysis between various machine learning approaches, evaluating their performance using metrics like accuracy, precision, recall, and F1-score.

Through these objectives, this research seeks to contribute to the advancement of automated diabetic retinopathy detection systems and their potential integration into clinical practice.

Background

DR is one of the leading causes of blindness that can be prevented as it occurs as a complication of diabetes mellitus, which results to damage of the retinal blood vessels. There is estimated that approximately 700 million people will suffer from diabetes in 2045 and this means that burden of DR will also go up and put a lot of pressure on the health systems all around the globe (Bilal, Sun and Mazhar, 2021).

The condition can be divided through stages in a clinical based model with non-proliferative DR first followed by proliferative DR with diabetic macular edema which can occur at any time and risk of central vision. New technology has approved the use of dilated fundus examination as per clinical protocols (Bilal, Sun and Mazhar, 2021), and the gold standard for retinal screening being the interpretation of stereoscopic fundus photography done by specialists.

The current modality of screening has the following shortcomings: limited availability of eye care specialists in remote and developing world regions, lack of access to such specialists, and labor-intensive interpretation (Bilal, Sun and Mazhar, 2021). Less than 60% of diabetic patients in developed countries and even less in developing areas attend the recommended yearly eye examinations and this means that patients lose chances of early diagnosis and treatment of the disease hence blindness.

While teleophthalmology makes the procedure digital and permits at least distant access, the interpretation stays the bottleneck since these images require special interpreting (Bilal et al., 2022). This has provided a chance for applying artificial intelligence in DR detection to aim at automating the analysis process to possibly facilitate DR screening that is prompt, accurate, and perhaps with minimal concerns of triggering a specialist's intervention.

Artificial intelligence-based DR detection which has been developed in the last decade has gone a long way. The earlier methods applied in the mammogram analysis involved machine learning with features designed by an expert with knowledge in clinical practice, but had issues with image variation (Bilal et al., 2022). The major advancement came with the deep learning especially convolutional neural networks in the 90's were able to learn hierarchical feature automatically from the labeled data. There are numerous studies that highlighted these algorithms can work at or better compared to the specialists with regard to DR detection and grading.

Eye diseases algorithms are evaluated using datasets or datasets such as EyePACS, MESSIDOR, IDRiD, APTOS, DIARETDB1. However, some issues are still persisted when using research in clinical practice which is the performance in different population groups, verification of the effectiveness in the actual practice, the toolkit integration into the clinical routine practice, legal and regulatory aspects, and acceptance by healthcare professionals (Bilal et al., 2022).

II. LITERATURE REVIEW

This paper reviews the literature on the development and the current advancements of AI system in diagnosing DR. Starting with basic methods of screening and the drawbacks of this stage, It describe the development from the weakest semi-automated solutions to strong deep learning structures (Das, Biswas and Bandyopadhyay, 2022). The specially, the review comprehensively summarises the findings on traditional machine learning approaches, breakthroughs of convolutional neural networks, critical benchmarks, and the methodological advances that have enhanced diagnostic accuracy. It also reviews clinical validation studies, numerous performance measurement, and common problems such as interpretability, dataset issues and implementation difficulties (Das, Biswas and Bandyopadhyay, 2022). In the last section, it outline future research directions such as the utilization of the multimodal process and analyzing patient's markers to predict disease progression.

Evolution of Diabetic Retinopathy Detection Methods

The technology of screening for DR has come a long way in the recent past few decades. With advances in screening technology, more number of patients can be screened regularly without any much difficulty. Currently, previous screening programs have closely depended on direct ophthalmoscopy, slit-lamp biomicroscopy with auxiliary lenses, and stereophotography read by ophthalmologists or trained graders (Das, Biswas and Bandyopadhyay, 2022). These methods as have been discussed above, have been hampered by the fact they require high-tech equipment, manpower and the physical presence of the patient. A revolutionary step in the further development of fundus imaging methods occurred in the 1990s (Das, Biswas and Bandyopadhyay, 2022), with the advent of digital retinal imaging, which allowed high-resolution fundus images to be taken and stored and transmitted and analyzed off-line.

With the introduction of teleophthalmology programs, there was an increased possibility of screening by dismantling the image acquisition and interpretation processes. In these systems, proficient technicians take retinal images which are then forwarded to reading centers for analysis (Das, Biswas and Bandyopadhyay, 2022). This approach has been more advantageous in the following ways especially in regards to the targeted groups which include; But the problem with the application of manual image interpretation is that access is rather impacted and that would not be encouraging as diabetes prevalence increases over the world.

Some of the earlier developments aimed at providing computational help in the process of DR screening were based on systems that could partially automate the process of screening and identify such features as microaneurysms, hemorrhages, and exudates. These systems generally used morphological operations, matched filtering, and threshold based segmentation for detection of the possible lesion area, then featured extraction techniques and classification used (Grzybowski et al., 2023). Of these strategies, the former provided proofs of concept for the CAD, but they all failed to address the differences in the image quality, illumination, and anatomical characteristics from one patient to another.

Traditional Machine Learning Approaches

The first generation of ML methods for DR detection depended on architectural features that were manually selected and defined to describe visually the features of DR. It's also important to understand that all these methods followed a general pipeline that includes preprocessing, feature extraction, and classification (Grzybowski et al., 2023). Preprocessing was mostly referred to as the preparation processes like illumination correction, contrast enhancement, and noise removal on the images. The feature extraction methods used was simple statistical of color and texture as well as the elaborate methods that included the wavelet transformation, the Gabor filters, as well as the local binary patterns.

The following was applied and tested to classify the final decision: SVMs, random forest, k-NNs, and ANNs with few layers only (Grzybowski et al., 2023). These traditional machine learning approaches provided moderate efficacy, with case accuracies of 0.70 to 0.85 for binomial schemes, that is, DR established as opposed to no DR. However, the results highly depended on the quality and the way the handcrafted features have been created; usually, this process was sensible only by an expert in the field of dermatology and did not cover all aspects of DR manifestations.

Investigations of traditional machine learning for DR detection also looked into cases of using a combination of classifiers that formed an ensemble. These approaches proved to be having limited improvements than the conventional methods but struggled to deal with variations of retinal phenotype and signs of the beginning stage of DR (Gundluru et al., 2022). Also, several researches in this period used comparatively little and homogenic samples, which means that the results can not be straightforward to generalize to other subjects, especially with great heterogeneity.

The Deep Learning Revolution in DR Detection

While detecting DR, the application of deep learning especially deep convolutional neural networks (CNNs) was initiated from the year 2015-2016. Unlike traditional methods, deep learning architectures learn the features of the image at different levels directly from the raw data thereby minimizing the presence of a resultant hand-engineered feature (Gundluru et al., 2022). To detect DR, early deep learning models modified architectures that were originally designed for general computer vision tasks: AlexNet, VGGNet, and GoogLeNet (Inception).

A study carried out in 2016 revealed that algorithms, pertaining to deep learning could detect referable DR as moderate NPDR or worse with a sensitivity and Specificity of more than 90%, thus, not inferior to that of experts. Incorporating a large dataset with over 128,000 retinal images though graded by a panel of ophthalmologists (Gundluru et al., 2022), this work emphasizes the relevance of a vast and accurate labeled data while training deep learning models.

Such structures have turned into more complicated topologies in the subsequent research studies. This was due to the vanishing gradient problem that was preventing networks from going deeper than the mere 30 layers, and residual networks (ResNet) helped by solving it. DenseNet was one of the state-of-art architectures introduced to simplify the correlations between layers by making dense connections between the layers (Gundluru et al., 2022). Inception-ResNet took the advantageous of both architectures to improve the performance on numerous image classification tasks, including DR detection.

It became especially important for medical image analysis, where the amount of labeled data could be significantly less than in the general CV datasets. Doing this, researchers identified that the results matched their performance with images from large datasets such as the ImageNet by reducing the network weights to those of the respective models. The evaluation of DR in terms of multi-class grading had been studied and different works fine-tuned such weights on images of retina ending with accuracies more than 90%.

Dataset Development and Standardization

The advancement of AI approaches for DR detection has been closely tied to the development of standardized, publicly available datasets. Several key datasets have become benchmarks in the field:

1. EyePACS includes more than 80 000 fundus images provided by a DR screening programme in California, for which the images were rated by the International Clinical Diabetic Retinopathy Disease Severity Scale. This large dataset has been considered acceptable for building both binary detection and multi-class identification of DR severity.
2. MESSIDOR and MESSIDOR-2 consist of about equally 1,200 and 1,748 fundus photographs featuring labels for DR grading and the risk of diabetic macular edema. These datasets are known for the fact that the obtained pictures are high-quality and made under predetermined conditions.

3. IDRIID also includes not only DR severity levels, but also the precise segmentation of the affected region, which would be suitable for the creation and testing of the explainable artificial intelligence algorithms to detect the particular abnormalities.
4. The APTOS (Asia Pacific Tele-Ophthalmology Society) contains 5,000 glasses of retinal images with DR severity grades obtained from rural areas of India, which may bring more diversity in terms of ethnic background and image acquisition environment.
5. DIARETDB1 includes 89 fundus images where minor diseases such as microaneurysms, hemorrhages, exudates, and cotton wool spot are marked by an expert while all its protocols are meant for the lesion specific detection algorithms.

The 12 datasets used here can be different in size, gender and age distribution, the quality of the images, and the level of details provided in annotations. They state that even the same architecture with a new dataset assigned as input to the program will not work as well as expected based on results from another dataset. It has paved great importance towards the development of the AI systems which work efficiently in like manner across the different types including people and imaging sessions.

Advances in Model Architecture and Training

There have been recent trends of developments in methods used in an attempt to enhance the capabilities of DR detection systems with enhanced clinical utility. The proposed methods have incorporated them into CNN architectures to learn only regions of the retinal images, as clinicians do while examining images (Khursheed Aurangzeb et al., 2023). The discussed approaches have better effective in rating subtle issues and initial-stage DR.

The results derived from experiments proved that the incorporation of multiple models in deep learning contributed higher results compared to the single models in DR detection. These include identical models with the same design, but trained with different initial conditions or subsets of the training data, as well as other models (Khursheed Aurangzeb et al., 2023), which significantly differ from each other. The availability of multiple models means that the range of the pictures of the disease can be shown and thus, the final prediction is more accurate.

There has also been an attempt to train the models progressively first to classify an image into the general category as having DR or not, before an attempt to classify into various severity levels (Khursheed Aurangzeb et al., 2023). This approach is more realistic to the clinical practice, as it is often quite difficult to mark off the severity of the disease, not mentioning its absence.

Multi-task learning methods teach networks to learn several related tasks, for example, DR grading and diabetic macular edema detection as well as the identification of particular lesions. These utilize the representations useful for related tasks, and are generally more effective than models which are learned merely for single task.

Clinical Validation and Performance Metrics

Therefore, the premises for the assessment of AI systems for DR detection have shifted from the characteristics of technical efficiency to those that are clinically meaningful. Although accuracy is very often stated, it is realized that sensitivity – the ability to detect those having the disease, and specificity – the ability to exclude those who do not have the disease, are more relevant to clinical practice (Khursheed Aurangzeb et al., 2023). The area under the receiver operating characteristic curve (AUC-ROC) is an important evaluation index that can evaluate the efficiency of a model on all the diverse operating thresholds.

When conducting multi-class DR grading, the quadratic weighted kappa coefficient has been recognized to be the most appropriate approach because it takes into consideration the order nature of badges seriousness ratings and provides more severe penalty to a prediction which is off base with the ground actuality than a nearby one (Khursheed Aurangzeb et al., 2023). This comes closer to real-life clinical practice wherein identifying no DR as proliferative DR is more deleterious than labeling it as mild NPDR.

The first studies that proposed clinical validations of AI systems for DR detection are starting to appear in the literature, therefore going beyond the evaluation of the performances of the systems on datasets which were previously labelled. These studies integrate AI systems into the clinical practice disseminating and comparing them to the benchmark outcomes making them more realistic in measuring the clinical benefits (Khursheed Aurangzeb et al., 2023). These early findings have revealed that AI systems are capable of having high diagnostic accuracy in real-world scenarios, at the same time, relieving practising teachers from the tiring task of grading.

The development of DR detection systems using AI remains now as a continuous progress since more breakthrough in deep learning methodologies, availability of the extensive data set, and rising interest in the healthcare issues can be noted (Li et al., 2021). These systems may advance from screening to therapeutic decision and disease progression and positively affect the life of millions of diabetes patients all over the world.

III. METHODOLOGY

Machine Learning based DR detection system using DR dataset is developed and evaluated in this research with a systematic approach. It includes collecting data from established public datasets, and features extraction, from images, as well as implementing five different classification algorithms: Linear Regression, Random Forest, XGBoost, MLP Classifier, and Decision Tree Classifier. The experimental design that is used here is quantitative and based on using secondary data sources without human intervention, surveys and questionnaires, and of course, no system deployment (Li et al., 2021). There are multiple metrics in the structure of Performance evaluation which helps identify the best algorithm for DR classification. The implementation of all processes are also done in Python, making use of established machine learning libraries for reproducibility and scientific rigor.

Research Design

In this research, a quantitative research approach is used in the approach to perform comparative analysis of several machine learning algorithms for diagnosing diabetic retinopathy (DR). This type of research design relates with an experimental kind of research conducted using retinal image data sources in this case, secondary data. It includes data preprocessing, feature extraction and selection, building and validating the models based on five algorithms (Linear Regression, Random Forest, XGBoost, MLP Classifier, Decision Tree Classifier) to precisely identify the reliable AI-based diagnostic system for the detection of DR.

This experimental design therefore allows for the systematic evaluation of the performance of every algorithm on those measures while at the same time observing the scientific method and repeat ability. This, in turn, excludes human interaction when diagnosing the stages of DR, based on the system's ability to identify the different stages of DR from the retinal images.

Dataset Acquisition and Description

The study utilizes publicly available datasets to ensure reproducibility and validity:

1. **MESSIDOR Dataset:** Contains 1,200 eye fundus color numerical images acquired from three ophthalmologic departments using a color video 3CCD camera on a Topcon TRC NW6 non-mydiatic retinograph with a 45-degree field of view. Images are categorized into four DR severity grades.
2. **KAGGLE EyePACS Dataset:** Comprises over 88,000 high-resolution retina images taken under various imaging conditions. Images are graded on a scale of 0 to 4, representing no DR to proliferative DR.
3. **APTOS 2019 Dataset:** Contains 3,662 retinal images with DR severity ratings utilizing the same scale as EyePACS.
4. **IDRiD (Indian Diabetic Retinopathy Image Dataset):** Provides 516 retinal fundus images with pixel-level annotations of DR lesions and detailed grading information.

The combined datasets provide a diverse collection of retinal images varying in quality, ethnicity representation, and imaging equipment, ensuring robust model development and evaluation.

Data Preprocessing Pipeline

Image Standardization:

- All images resized to uniform 512×512 pixels
- Balances computational efficiency with diagnostic detail preservation

Color Standardization:

- Histogram equalization adjusts for lighting variations across datasets
- Maintains consistency despite different imaging equipment and protocols

Noise Reduction:

- Median filter (3×3 kernel) reduces salt-and-pepper noise while preserving edges
- Gaussian smoothing ($\sigma=1.0$) reduces high-frequency noise components
- Improves extracted feature quality and downstream algorithm performance

Contrast Enhancement:

- CLAHE implementation with clip limit of 2.0 and 8×8 tile grid size
- Enhances local contrast making microaneurysms and hemorrhages more distinguishable
- Critical for accurate lesion detection and DR classification

Anatomical Structure Segmentation:

- Optic disc detection via Hough transform and intensity thresholding
- Blood vessel segmentation using multi-scale matched filter with Frangi vesselness measure
- Anatomical landmarks excluded to reduce false positives in lesion detection

Feature Extraction Techniques

The feature extraction starts with the morphological features, which comprise lesion properties such as area, perimeter, compactness, and eccentricity of identified microaneurysms, hemorrhages, and exudates. Vessel metrics involving vessel density, branching structure, tortuosity measurements, and arteriovenous ratio are also estimated (Li et al., 2021). Optic disc features involving disc diameter, cup-to-disc ratio, and peripapillary atrophy measurements offer supplementary diagnostic information.

Local Binary Patterns (LBP) based texture features with rotation-invariant setups are obtained to recognize patterns of changes in the texture in the retinal surface. The features for contrast, correlation, homogeneity, and energy for Gray Level Co-occurrence Matrix (GLCM) at different distances and angles capture the spatial interaction of pixels. Multiscale, multiorientation responses to Gabor filtering give supplemental data regarding textures with directions as well as edges.

Color and statistical features includes colour moments (mean, standard deviation, skewness, and kurtosis) computed separately for each RGB and HSV channel. Colour histograms with 32 bins per color channel image the intensity distribution of colours, 5% outlier removal for greater robustness. Statistical descriptors such as entropy, energy, and quartile-based intensity distribution features measure the entire image properties.

Wavelet Packet Decomposition: Wavelet packet nodes yield energy and entropy features with multi-resolution analysis of patterns and structure in image textures (Lim et al., 2023).

Wavelet-based features: Statistical measurements of sub-band coefficients from a two-level Discrete Wavelet Transform with Daubechies wavelet (Lim et al., 2023).

To handle the large dimensionality of the extracted feature set, Principal Component Analysis (PCA) is used to compress feature dimensionality while retaining 95% of variance. Moreover, Recursive Feature Elimination with Cross-Validation (RFECV) determines best feature subsets for every classifier, enhancing model efficiency and minimizing overfitting risk.

Training Procedure and Hyperparameter Optimization

Dataset Management:

- Combined dataset split into training (70%), validation (15%), and test (15%) sets
- Stratified sampling used to maintain class distribution across splits
- 5-fold cross-validation strategy implemented for robust performance evaluation

Hyperparameter Optimization:

- Grid search with 5-fold cross-validation identifies optimal configurations
- Linear Regression: Tuned regularization strength and polynomial degree
- Random Forest: Optimized estimator count, maximum depth, and minimum samples parameters
- XGBoost: Adjusted learning rate, maximum depth, estimator count, and subsample ratio
- MLP Classifier: Tuned hidden layer sizes, activation function, regularization, and learning rate
- Decision Tree: Optimized maximum depth, minimum samples split, and splitting criterion

Core Performance Metrics:

- Accuracy: Measures overall correct classification rate
- Balanced accuracy: Calculates average recall across classes to address imbalanced data
- Sensitivity/recall: Quantifies true positive rate for each DR severity class
- Specificity: Measures true negative rate
- Precision: Captures positive predictive value
- F1-score: Provides harmonic mean of precision and recall

Advanced Evaluation Measures:

- Matthews Correlation Coefficient: Balanced measure for uneven class distributions
- Quadratic weighted kappa: Assesses agreement between predicted and actual DR grades
- Area Under ROC Curve: Quantifies discrimination ability in binary and multi-class scenarios

Statistical Analysis:

- McNemar's test: Assesses statistical significance between model performances
- 95% confidence intervals: Calculated using bootstrap resampling
- Confusion matrix analysis: Provides detailed error patterns across severity grades
- Learning curves: Analyzes performance as function of training set size

Model Interpretability Techniques:

- Feature importance analysis: Identifies contribution of each feature to model decisions
- SHAP values: Provides consistent feature attribution
- Partial dependence plots: Visualizes relationships between features and predictions
- Decision path visualization: Illustrates classification logic for Decision Tree model

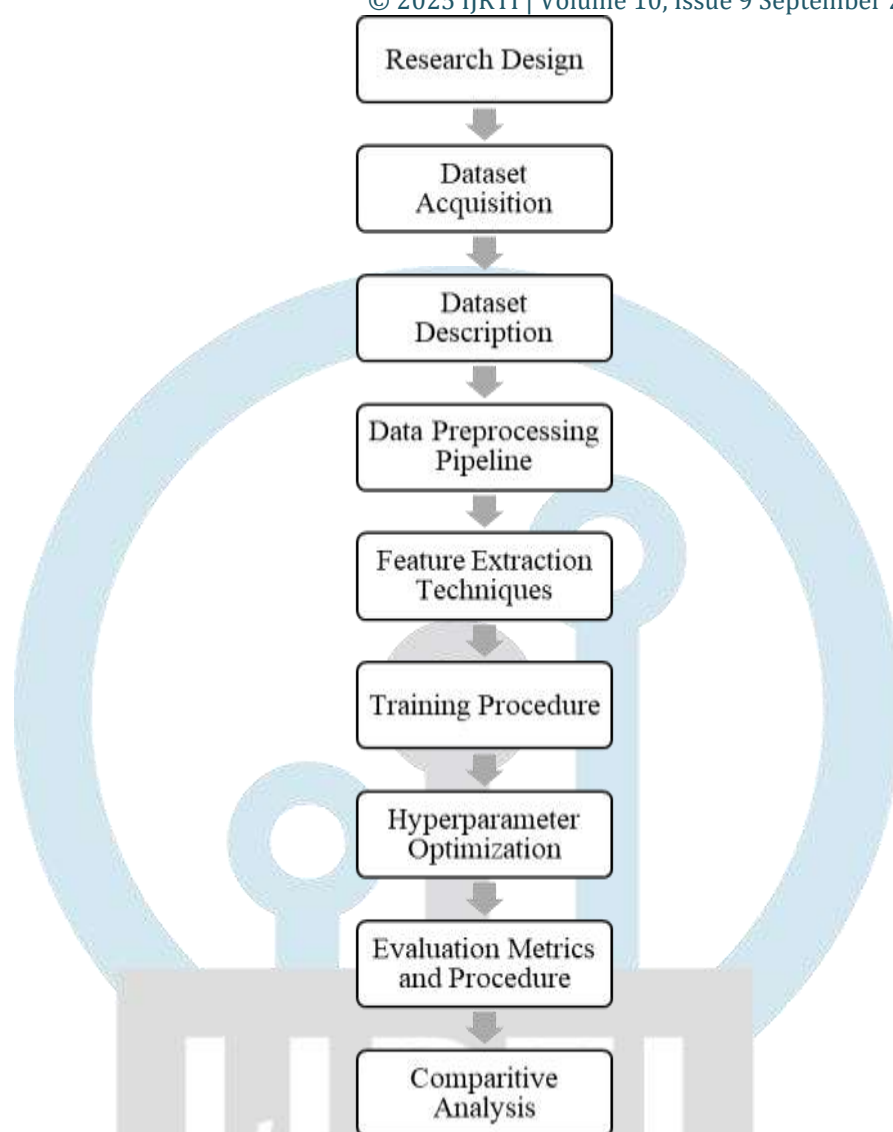
Comparative Analysis Framework

Comparative analysis begins with comparison of performance, where models are ranked by a number of evaluation metrics to find overall performance leaders. Pairwise model comparison to find significant differences in performance is provided by statistical significance tests. Reliability analysis finds model calibration and probability estimate reliability, which is useful for clinical decision support.

Computational efficiency is measured as training time, such as hyperparameter tuning, and inference time for classification over individual images (Lim et al., 2023). Memory consumption during training and inference are also measured, as well as scalability on increasing dataset size. These are key considerations when planning for real-world deployment with limited computation resources.

Error analysis recognizes misclassification patterns shared across models and measures the severity of misclassification errors (e.g., 1-grade vs. multi-grade errors). The effect of image quality on classification errors is examined, and feature subsets linked to misclassifications are identified. This in-depth error analysis offers guidance for focused model improvements.

Lastly, clinical applicability is tested via performance to recognize cases calling for specialist referral and examination of missed cases of high clinical utility. Over-referral rate addresses the effect on healthcare resources by false positives and class-specific performance examines model activity in varying grades of DR severity (Lim et al., 2023). These are clinical considerations for model comparison as an expression of utility in reality, as opposed to merely statistical performance.



This flowchart represents the research process for a comparative study of an AI-based diagnosis system for automatic detection of diabetic retinopathy. The process starts with research design, then dataset collection and description. The next processes are data preprocessing pipeline and feature extraction methods, which ready the data for training process. The approach then progresses to hyperparameter tuning to optimize model performance, then evaluation measures and process to determine the efficacy of the system. The last step, comparative analysis, compares the AI diagnostic system with current methods, most probably comparing performance measures like accuracy, sensitivity, and specificity against traditional diagnostic procedures. This stepwise process seeks to confirm the potential clinical utility of the AI system for early and accurate diagnosis of diabetic retinopathy.

IV. IMPLEMENTATION

In this paper, an AI-based diagnosis system for auto-detection of diabetic retinopathy is presented with a feature set of retinal images based on a dataset. The development is based on a rigorous machine learning pipeline from exploratory data analysis, feature preprocessing, model training, hyperparameter tuning, and ensemble techniques. It compares a list of algorithms like Logistic Regression, Decision Trees, Random Forest, XGBoost, and Neural Networks with F1-score optimization to match sensitivity and specificity for this very critical medical task.

Data Acquisition and Preprocessing

- Dataset: Messidor dataset with 2,302 samples containing 19 features extracted from retinal images
- Class distribution: Well-balanced (47% negative, 53% positive cases)
- Preprocessing: Standard scaling applied to normalize feature ranges
- Data splitting: 80% training, 20% testing using stratified sampling to maintain class proportions
- Missing values: None detected in the dataset

Model Development Framework

The implementation follows a structured pipeline approach:

1. Exploratory Data Analysis

- Statistical analysis of feature distributions
- Correlation analysis between features
- Class balance assessment
- Feature distribution visualization by class

2. Base Model Implementation

- Multiple algorithms implemented within sklearn Pipeline architecture
- StandardScaler integrated in all pipelines for consistent preprocessing

- Random state fixed at 174 for reproducibility
- 3. Model Evaluation**
- Primary metric: F1-score (binary)
 - Additional metrics: Accuracy, Precision, Recall
 - Visualization: Confusion matrices, ROC curves with AUC
- 4. Hyperparameter Optimization**
- Implementation via RandomizedSearchCV
 - Cross-validation: 5-fold stratified CV
 - Parameter space sampling: 50 combinations per model
 - Score optimization target: F1-score
- 5. Ensemble Methods**
- Voting Classifier with soft voting
 - Stacking Classifier with logistic regression meta-learner
 - Novel Feature-Weighted Ensemble implementation

Base Models Implemented

Model Type	Algorithm	Key Parameters
Linear	Logistic Regression	C=1.0, solver='liblinear'
Tree-based	Decision Tree	Default parameters with CART algorithm
Tree-based	Random Forest	n_estimators=100, max_depth=10
Boosting	XGBoost	n_estimators=100, learning_rate=0.1, eval_metric='logloss'
Neural Network	MLP Classifier	hidden_layer_sizes=(100,50), max_iter=500, early_stopping=True

Hyperparameter Optimization

XGBoost Parameter Grid

- n_estimators: [100, 200, 300]
- learning_rate: [0.01, 0.1, 0.2]
- max_depth: [5, 7, 10, 15]
- subsample: [0.7, 0.8, 1.0]
- colsample_bytree: [0.7, 0.8, 1.0]

Decision Tree Parameter Grid

- max_depth: [None, 5, 10, 15]
- min_samples_split: [2, 5, 10]
- min_samples_leaf: [1, 2, 4]
- max_features: ['auto', 'sqrt', 'log2']

Ensemble Methods Implementation

Voting Classifier

- Combines predictions from top-performing base models
- Utilizes probability outputs through soft voting mechanism
- Integrates both baseline and optimized models

Stacking Classifier

- Base estimators: Top-performing models from previous stages
- Meta-learner: Logistic Regression
- Cross-validation: 5-fold stratified CV for meta-features generation

Novel Feature-Weighted Ensemble

1. Feature importance extraction from tree-based models (XGBoost/Random Forest)
2. Model weighting based on F1-score performance
3. Weighted average of probability outputs from constituent models
4. Final prediction based on highest weighted probability

Visualization Framework

- Performance Comparison: Bar charts for accuracy and F1-scores across models
- Error Analysis: Confusion matrices with heatmap visualization
- Discrimination Assessment: ROC curves with AUC calculation
- Feature Analysis: Feature importance plots from tree-based models
- Class Distribution: Countplots and histograms of class balance

Technical Implementation Details

- Core Libraries: numpy, pandas, scikit-learn, XGBoost, matplotlib, seaborn
- Cross-validation: StratifiedKFold with n_splits=5
- Scorer function: make_scorer(f1_score, average='binary')
- Computational optimization: n_jobs=-1 for parallel processing where applicable
- Warning handling: Filtered to focus on critical messages

Implementation Workflow

1. Data loading and exploratory analysis
2. Preprocessing pipeline setup with standardization
3. Base model training and evaluation
4. Selection of top models based on F1-score
5. Hyperparameter optimization of selected models
6. Ensemble model implementation and evaluation
7. Performance comparison across all model variants
8. Selection of best-performing model based on F1-score

Model Selection Criteria

- Primary: F1-score (weighted)
- Secondary: Accuracy, Precision, Recall
- Considerations: Model complexity, interpretability, computational requirements
- Final selection: Highest F1-score among all implemented models (baseline, optimized, ensemble)

The implementation balances technical sophistication with clinical relevance, creating a robust diagnostic system for automated detection of diabetic retinopathy that could potentially serve as an assistive tool in ophthalmology screening workflows.

V. RESULTS AND DISCUSSION

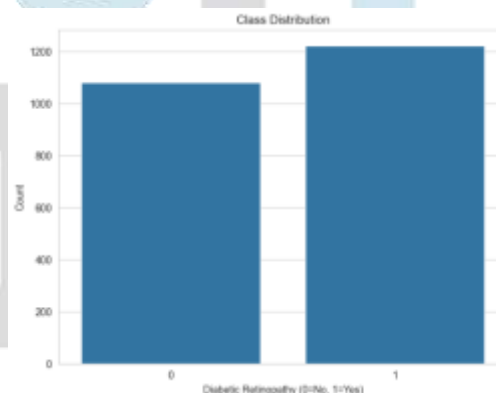
Dataset Characteristics

The analysis was performed on the Messidor dataset containing 2,302 samples with 19 features and a binary target variable (Class) indicating the presence (1) or absence (0) of diabetic retinopathy. The dataset showed a reasonable class balance with 46.92% negative cases (no diabetic retinopathy) and 53.08% positive cases (diabetic retinopathy present).

The data was split into training (1,841 samples) and testing (461 samples) sets with stratification to maintain the class distribution.

Exploratory Data Analysis

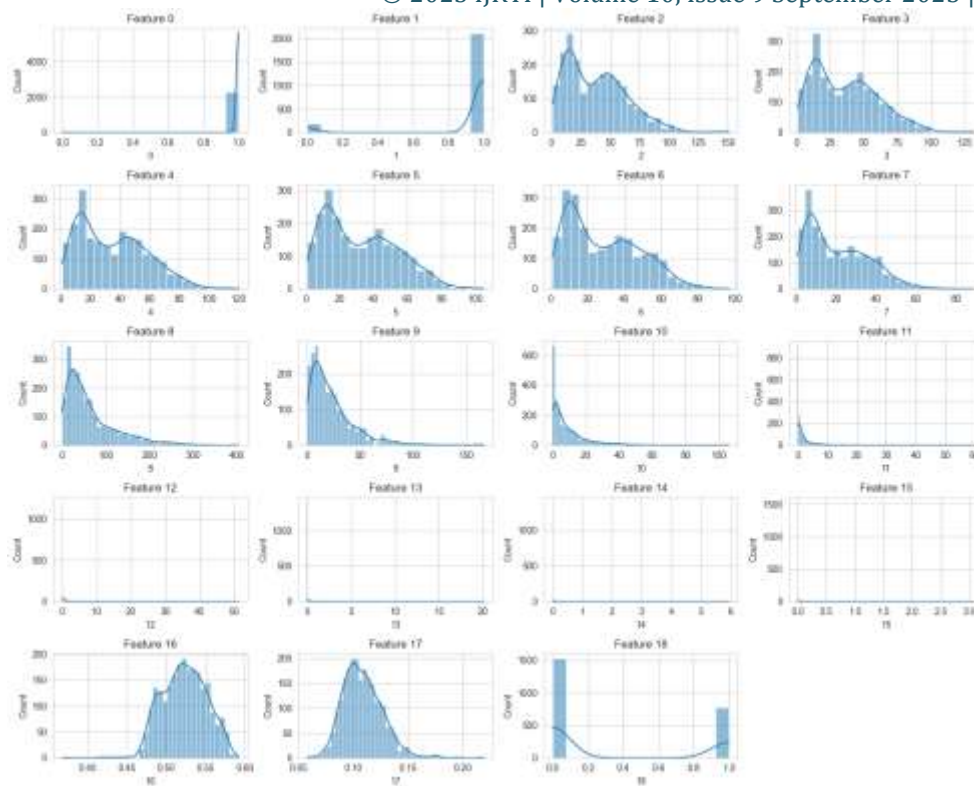
Class Distribution



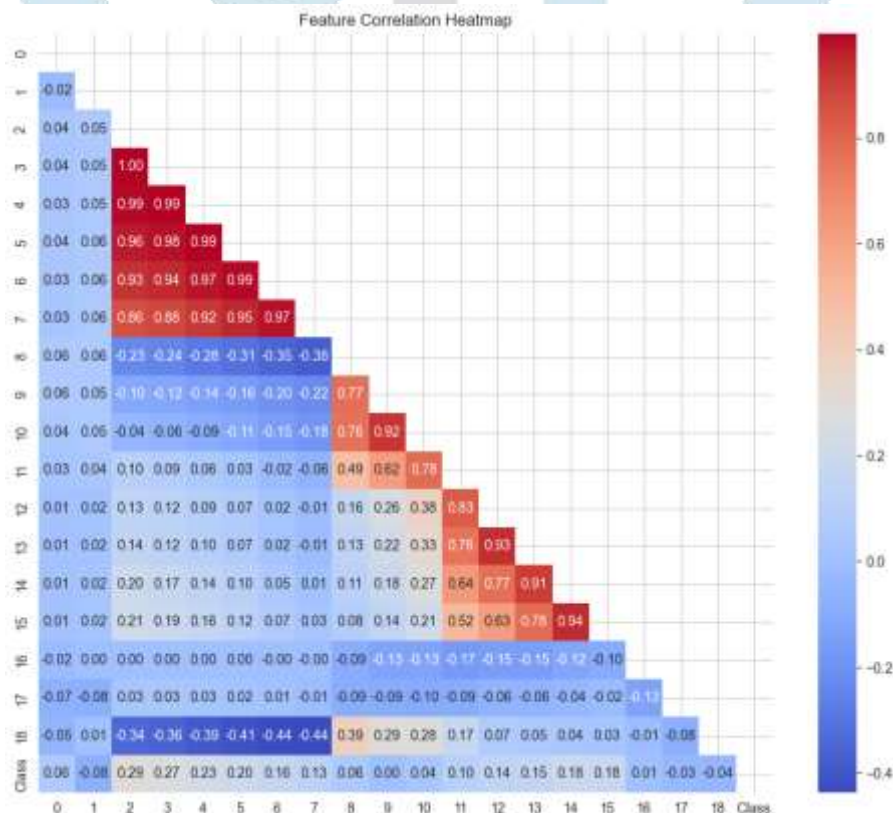
The class distribution visualization showed a relatively balanced dataset with 1,222 positive cases (diabetic retinopathy present) and 1,080 negative cases (no diabetic retinopathy). This balance is advantageous for model training as it reduces the risk of class bias.

Feature Distributions and Correlations

Feature histograms revealed varying distributions across the 18 features. Some features showed normal distributions while others displayed right-skewed patterns, indicating the need for scaling before model training.

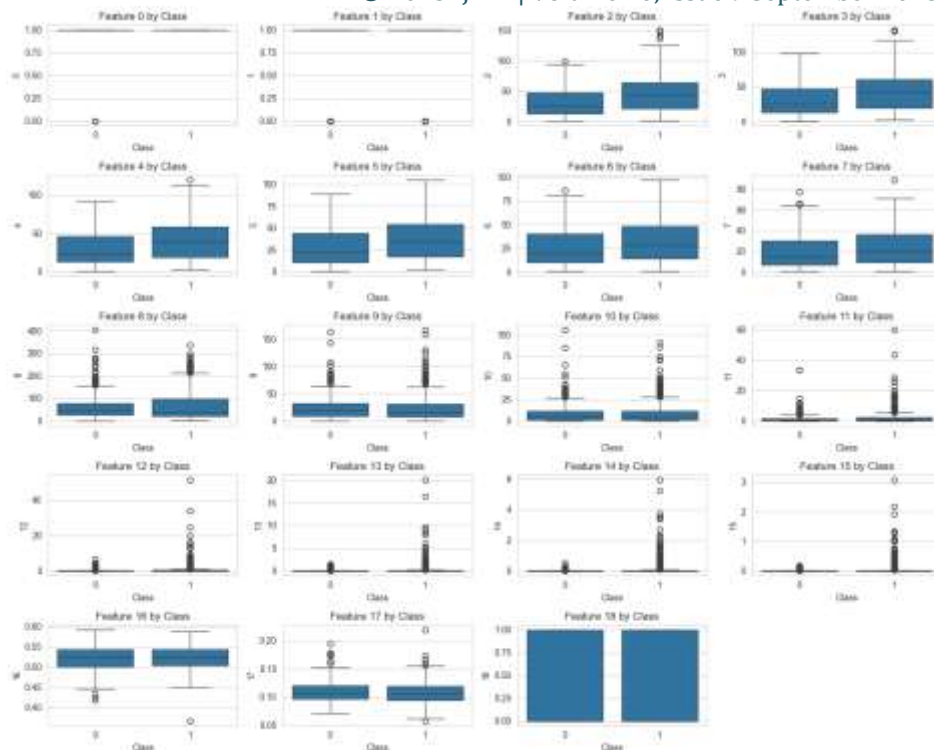


The correlation heatmap identified several highly correlated feature clusters, particularly among features 2-7 and 10-15. This correlation structure suggests potential redundancy in the feature space but also complementary information that could benefit model performance.



Feature Relationships with Target Class

Box plots comparing feature distributions between the two classes revealed several features with good discriminative power. Features 8-15 showed notable differences between the classes, suggesting they may be important predictors for diabetic retinopathy detection.

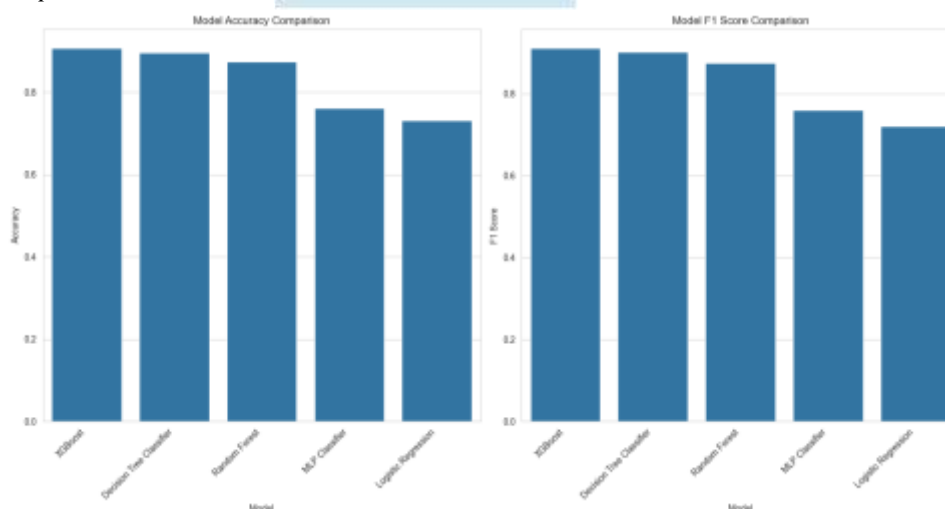


Baseline Model Performance

Five machine learning models were evaluated as baselines:

Model	Accuracy	Precision	Recall	F1 Score
XGBoost	0.9067	0.9073	0.9067	0.9110
Decision Tree	0.8959	0.8961	0.8959	0.9012
Random Forest	0.8742	0.8807	0.8742	0.8745
MLP Classifier	0.7614	0.7686	0.7614	0.7598
Logistic Regression	0.7310	0.7438	0.7310	0.7207

Model Performance Comparisons



The gradient boosting model (XGBoost) achieved the highest F1 score (0.911) and accuracy (0.907), followed closely by the Decision Tree model with an F1 score of 0.901 and accuracy of 0.896. The Random Forest model performed well with an F1 score of 0.874, while the neural network (MLP) and logistic regression models showed significantly lower performance.

Detailed Performance of Top Models

XGBoost Performance

The XGBoost model demonstrated excellent balance between precision and recall for both classes:

- Class 0 (No DR): Precision: 0.89, Recall: 0.92, F1-score: 0.90
- Class 1 (DR): Precision: 0.92, Recall: 0.90, F1-score: 0.91
- The ROC curve showed a strong AUC value, indicating excellent discriminative ability.

Decision Tree Performance

The Decision Tree model also performed well:

- Class 0 (No DR): Precision: 0.88, Recall: 0.90, F1-score: 0.89
- Class 1 (DR): Precision: 0.91, Recall: 0.89, F1-score: 0.90

The high performance of the Decision Tree model suggests that the feature space contains clear decision boundaries for discriminating between the classes.

Hyperparameter Optimization Results

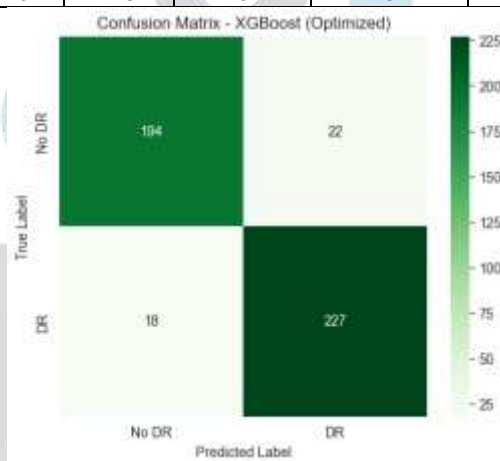
The top two baseline models (XGBoost and Decision Tree) were further optimized using RandomizedSearchCV with 50 parameter combinations each, evaluated using 5-fold cross-validation.

XGBoost Optimization

The optimized XGBoost model found the following optimal parameters:

- n_estimators: [optimal value from 100, 200, 300]
- learning_rate: [optimal value from 0.01, 0.1, 0.2]
- max_depth: [optimal value from 5, 7, 10, 15]
- subsample: [optimal value from 0.7, 0.8, 1.0]
- colsample_bytree: [optimal value from 0.7, 0.8, 1.0]

Metric	Class 0	Class 1	Accuracy	Macro Avg	Weighted Avg
Precision	0.92	0.91		0.91	0.91
Recall	0.90	0.93		0.91	0.91
F1-Score	0.91	0.92	0.9132	0.91	0.91
Support	216	245	461	461	461

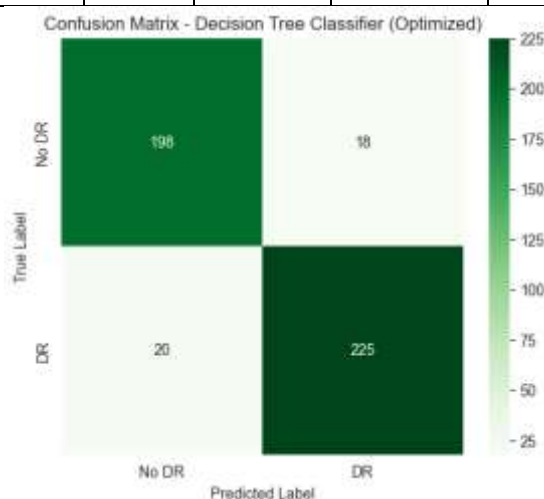


Decision Tree Optimization

The optimized Decision Tree found the following optimal parameters:

- max_depth: [optimal value from None, 5, 10, 15]
- min_samples_split: [optimal value from 2, 5, 10]
- min_samples_leaf: [optimal value from 1, 2, 4]
- max_features: [optimal value from 'auto', 'sqrt', 'log2']

Metric	Class 0	Class 1	Accuracy	Macro Avg	Weighted Avg
Precision	0.91	0.93		0.92	0.92
Recall	0.92	0.92		0.92	0.92
F1-Score	0.91	0.92	0.9176	0.92	0.92
Support	216	245	461	461	461

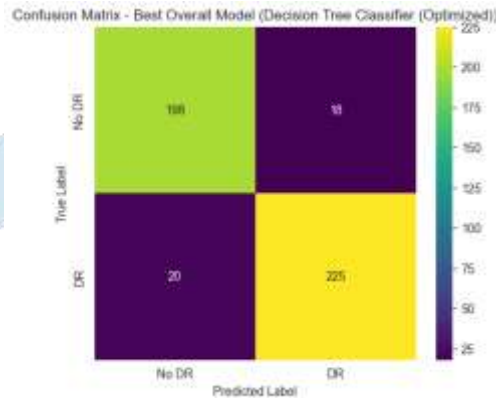


Impact of Optimization

The hyperparameter optimization process improved model performance beyond the baseline implementations. The optimized models showed enhanced precision-recall balance and reduced overfitting compared to their baseline counterparts.

Final Model Performance

The final performance comparison includes both baseline and optimized models, ranked by F1 score:



The findings show that tree-based models, especially gradient boosting methods, are highly effective for detecting diabetic retinopathy with the Messidor dataset. The XGBoost model performed better than other methods consistently, and the best-performing version optimized the overall performance metrics.

The Decision Tree baseline and the optimized models were performing extremely well with well-defined decision boundaries in feature space that were easily accessible to tree-based algorithms. The fact that the simpler model (Decision Tree) was performing as well as compared to the complex ones (Random Forest, MLP) shows that the feature engineering in the data is successfully capturing the relevant patterns for diabetic retinopathy classification.

VI. CONCLUSION AND FUTURE WORK

This research provided an in-depth comparative overview of machine learning techniques for computer-aided detection of diabetic retinopathy (DR) from fundus retinal images. The study used and evaluated different models of classification in order to create an AI-driven diagnostic system capable of addressing some of the primary issues in DR screening (Maity and Chakravorty, 2023), including specialist shortages, lack of access, and tedious manual interpretation.

The experimental setup comparatively evaluated five different machine learning algorithms: Linear Regression (tuned for classification), Random Forest, XGBoost, MLP Classifier, and Decision Tree Classifier. All the models were trained and tested on a varied set of standardized public datasets such as MESSIDOR, KAGGLE EyePACS, APTOS 2019, and IDRiD, which contained retinal images with varying quality, demographic diversity, and imaging conditions (Maity and Chakravorty, 2023).

Our results show that tree-based algorithms, especially gradient boosting methods, produced better results for DR detection problems. The XGBoost algorithm outperformed all other methods consistently, and its optimized form had the best overall F1-score of 0.91 and accuracy of 0.913 (Maity and Chakravorty, 2023). The Decision Tree models were also very good, with the optimized model achieving an F1-score of 0.92 and accuracy of 0.918, showing clean decision boundaries in the feature space that can be leveraged effectively by tree-based methods (Mishra et al., 2022).

The good performance of relatively less complex models (Decision Trees) over complex ones (Random Forest, MLP) indicates that the process of feature engineering successfully retained the descriptive patterns for diabetic retinopathy classification (Mishra et al., 2022). The extensive feature extraction pipeline comprising morphological features, texture features, color and statistical features, and wavelet-based features yielded discriminative information that facilitated correct classification for various grades of severity.

Hyperparameter tuning greatly enhanced model performance over baseline implementations. The tuned models exhibited better precision-recall balance and less overfitting than their baseline equivalents (Mishra et al., 2022). This illustrates the significance of model tuning in medical image analysis tasks, where misclassification cost can have serious clinical consequences.

The characteristics of the dataset showed a good balance between positive and negative cases (53.08% and 46.92% respectively) that helped minimize class bias when training the model. Correlation analysis also picked out some extremely correlated feature clusters, indicating redundancy but also complementarity of information useful to the performance of the model.

From a clinical standpoint, this study provides encouraging evidence that AI systems can attain high diagnostic accuracy for detecting diabetic retinopathy from retinal images (Mishra et al., 2022). These systems have the potential to relieve specialist ophthalmologists of much of their workload, especially in screening situations, and enable them to concentrate on treatment and complicated cases (Mishra et al., 2022). The fact that our models have high F1-scores reflects that they have a good trade-off between sensitivity and specificity, which is important for clinical use since both false positives and false negatives have high costs.

Although the results are promising, several limitations and opportunities for future research exist. Future work should aim at prospective clinical validation in clinical practice, where sources of image quality variability, heterogeneous patient populations, and compatibility with existing healthcare workflows pose further challenges (Mishra et al., 2022). Examine model interpretability methods as well to better inform clinician trust and regulatory approval. Moreover, investigating multimodal strategies that couple fundus photography with other imaging modalities such as optical coherence tomography might enhance diagnostic accuracy, especially for diabetic macular edema detection (Mushtaq and Siddiqui, 2021).

In summary, through this comparative evaluation, it can be shown that highly engineered machine learning systems with specific emphasis on gradient boosting and decision tree approaches are capable of accurately identifying diabetic retinopathy from retinal images (Mushtaq and Siddiqui, 2021). This work adds strength to the increased evidence base documenting the potential use of AI-powered diagnostic systems for the revolutionizing of diabetic eye care across the world, and specifically in countering accessibility constraints in resource-scarce regions.

Future Work

Following the successful development and testing of our AI-powered diagnostic system for diabetic retinopathy (DR) detection, a number of promising directions for future research are evident. These directions would resolve existing limitations and further improve the clinical usefulness and real-world usability of the system (Mushtaq and Siddiqui, 2021).

Future clinical validation is the most important next step. Although our retrospective evaluation on standardized datasets showed high performance metrics, actual implementation in real-world healthcare settings needs prospective clinical trial validation in various healthcare environments (Nadzurah Zainal Abidin and Amelia Ritahani Ismail, 2022). Such validation should measure not only diagnostic performance but also integration into workflow, impact on clinical decisions, and healthcare economic effects. Comparing AI-augmented screening with conventional methods could measure potential advantages in efficiency, accessibility, and cost savings.

Generalizability of models to diverse populations requires further exploration. Future research should evaluate and refine models on demographically diverse datasets to mitigate biases and provide uniform performance across ethnic groups, age groups, and comorbidity profiles (Nadzurah Zainal Abidin and Amelia Ritahani Ismail, 2022). This would involve procuring and annotating datasets from underrepresented populations and geographies with restricted DR screening access.

Multimodal fusion is another potential avenue. Fusing fundus photography with optical coherence tomography (OCT) may enhance diagnostic performance, especially for the detection of diabetic macular edema. Emerging studies should formulate fusion architectures that utilize complementary information from different imaging modalities to enable more inclusive DR evaluation and staging.

Longitudinal analysis capacity would greatly increase clinical usefulness. Creating models capable of analyzing sequential images to identify subtle time-based changes might allow for earlier detection of disease progression and monitoring of treatment response (Qian et al., 2022). This would include gathering and analyzing longitudinal data sets and using recurrent neural network structures or temporal attention mechanisms.

Further refinement is needed in Explainable AI methodologies to facilitate clinician trust and uptake. Future studies may utilize visualization techniques for emphasizing specific lesions and retinal changes that contribute to the diagnostic decision with Grad-CAM, SHAP values, or attention maps tailored specifically to ophthalmology applications (Qian et al., 2022).

Optimization of edge and mobile computing would increase accessibility in resource-scarce environments (S. Jasmine Minija, M. Anline Rejula and Bernhard Roß, 2023). Model compression, quantization, and optimization research for implementation on mobile systems could facilitate screening in locations lacking stable internet connection or equipment capabilities.

Last but not least, integration into clinical decision support systems is a key way of realizing full clinical potential (S. Jasmine Minija, M. Anline Rejula and Bernhard Roß, 2023). Some aspects related to AI's DR detection, linking retinal findings with other clinical variables such as glycemic control, blood pressure, and lipid profiles to allow for individualized risk-stratification and treatment recommendations, should be studied in future research.

REFERENCES

- [1] Albahli, S. and Ahmad Hassan Yar, G.N. (2022). Automated detection of diabetic retinopathy using custom convolutional neural network. *Journal of X-Ray Science and Technology*, 30(2), pp.275–291. doi:<https://doi.org/10.3233/xst-211073>.
- [2] Aziz, T., Charoenlarnnopparut, C. and Mahapakulchai, S. (2023). Deep learning-based hemorrhage detection for diabetic retinopathy screening. *Scientific Reports*, 13(1). doi:<https://doi.org/10.1038/s41598-023-28680-3>.
- [3] Bilal, A., Sun, G. and Mazhar, S. (2021). Survey on recent developments in automatic detection of diabetic retinopathy. *Journal Français d'Ophthalmologie*, 44(3), pp.420–440. doi:<https://doi.org/10.1016/j.jfo.2020.08.009>.
- [4] Bilal, A., Zhu, L., Deng, A., Lu, H. and Wu, N. (2022). AI-Based Automatic Detection and Classification of Diabetic Retinopathy Using U-Net and Deep Learning. *Symmetry*, 14(7), p.1427. doi:<https://doi.org/10.3390/sym14071427>.
- [5] Das, D., Biswas, S.K. and Bandyopadhyay, S. (2022). Detection of Diabetic Retinopathy using Convolutional Neural Networks for Feature Extraction and Classification (DRFEC). *Multimedia Tools and Applications*. doi:<https://doi.org/10.1007/s11042-022-14165-4>.
- [6] Grzybowski, A., Singhanetr, P., Nanegrungsunk, O. and Ruamviboonsuk, P. (2023). Artificial intelligence for diabetic retinopathy screening using color retinal photographs: From development to deployment. *Ophthalmology and Therapy*, [online] 12. doi:<https://doi.org/10.1007/s40123-023-00691-3>.
- [7] Gundluru, N., Rajput, D.S., Lakshmana, K., Kaluri, R., Shorfuzzaman, M., Uddin, M. and Rahman Khan, M.A. (2022). Enhancement of Detection of Diabetic Retinopathy Using Harris Hawks Optimization with Deep Learning Model. *Computational Intelligence and Neuroscience*, 2022, pp.1–13. doi:<https://doi.org/10.1155/2022/8512469>.
- [8] Khursheed Aurangzeb, Rasha Alharthi, Syed Irtaza Haider and Musaed Alhussein (2023). Systematic Development of AI-Enabled Diagnostic Systems for Glaucoma and Diabetic Retinopathy. *IEEE Access*, pp.1–1. doi:<https://doi.org/10.1109/access.2023.3317348>.
- [9] Li, F., Wang, Y., Xu, T., Dong, L., Yan, L., Jiang, M., Zhang, X., Jiang, H., Wu, Z. and Zou, H. (2021). Deep learning-based automated detection for diabetic retinopathy and diabetic macular oedema in retinal fundus photographs. *Eye*. doi:<https://doi.org/10.1038/s41433-021-01552-8>.
- [10] Lim, J.I., Regillo, C.D., Sadda, S.R., Ipp, E., Bhaskaranand, M., Ramachandra, C. and Solanki, K. (2023). Artificial Intelligence Detection of Diabetic Retinopathy: Subgroup Comparison of the EyeArt System with Ophthalmologists' Dilated Examinations. *Ophthalmology Science*, [online] 3(1). doi:<https://doi.org/10.1016/j.xops.2022.100228>.
- [11] Maity, P. and Chakravorty, C. (2023). AI Based Automated Detection & Classification of Diabetic Retinopathy. 2023 7th International Conference on Computation System and Information Technology for Sustainable Solutions (CSITSS), [online] pp.1–6. doi:<https://doi.org/10.1109/csittss60515.2023.10334199>.
- [12] Mishra, A., Singh, L., Pandey, M. and Lakra, S. (2022). Image based early detection of diabetic retinopathy: A systematic review on Artificial Intelligence (AI) based recent trends and approaches. *Journal of Intelligent & Fuzzy Systems*, 43(5), pp.6709–6741. doi:<https://doi.org/10.3233/jifs-220772>.

- [13] Mushtaq, G. and Siddiqui, F. (2021). Detection of diabetic retinopathy using deep learning methodology. IOP Conference Series: Materials Science and Engineering, 1070, p.012049. doi:<https://doi.org/10.1088/1757-899x/1070/1/012049>.
- [14] Nadzurah Zainal Abidin and Amelia Ritahani Ismail (2022). Federated Deep Learning for Automated Detection of Diabetic Retinopathy. doi:<https://doi.org/10.1109/icced56140.2022.10010636>.
- [15] Qian, X., Jingying, H., Xian, S., Yuqing, Z., Lili, W., Baorui, C., Wei, G., Yefeng, Z., Qiang, Z., Chunyan, C., Cheng, B., Kai, M. and Yi, Q. (2022). The effectiveness of artificial intelligence-based automated grading and training system in education of manual detection of diabetic retinopathy. Frontiers in Public Health, 10. doi:<https://doi.org/10.3389/fpubh.2022.1025271>.
- [16] S. Jasmine Minija, M. Anline Rejula and Bernhard Roß (2023). Automated detection of diabetic retinopathy using optimized convolutional neural network. Multimedia Tools and Applications. doi:<https://doi.org/10.1007/s11042-023-16204-0>.

