

An Explainable Deep Learning Approach for Diabetic Retinopathy Classification and Precise Lesion Segmentation

¹Brunda K V, ²Dr Pushpalatha K R

¹Research scholar,

Sri Siddaratha Institute of Technology, SSAHE, Tumkur

²Associate Professor & Research Supervisor,

Department of CSE(Data Science),

Sri Siddaratha Institute of Technology, SSAHE, Tumkur

ARTICLE INFO

Received: 04 Mar 2026

Revised: 20 Apr 2026

Accepted: 01 May 2026

ABSTRACT

Diabetic Retinopathy (DR) is a leading cause of vision impairment and blindness among diabetic patients worldwide, necessitating early and accurate diagnosis for effective treatment. This study proposes an explainable deep learning approach for diabetic retinopathy classification and precise lesion segmentation using a hybrid architecture that integrates Convolutional Neural Networks (CNNs) and Transformer-based attention mechanisms. The proposed framework is designed to automatically identify the severity level of DR from retinal fundus images and accurately localize pathological regions such as microaneurysms, hemorrhages, and exudates through advanced segmentation techniques. The CNN component extracts hierarchical spatial features from input retinal images, while the Transformer module captures long-range contextual dependencies to enhance global feature representation. To improve clinical trust and interpretability, Explainable Artificial Intelligence (XAI) techniques such as Grad-CAM and attention map visualization are incorporated, enabling clinicians to understand model decision-making and validate highlighted lesion regions. The segmentation module employs a hybrid encoder-decoder architecture with multi-scale feature fusion to precisely delineate infected areas, improving diagnostic transparency and reliability. Experimental evaluation on benchmark retinal image datasets demonstrates superior performance in terms of classification accuracy, sensitivity, specificity, Dice coefficient, and Intersection-over-Union (IoU) compared to conventional deep learning models. The proposed explainable framework not only enhances diagnostic accuracy but also strengthens clinical interpretability, making it a reliable decision-support tool for automated diabetic retinopathy screening and early intervention in real-world healthcare settings

Keywords: Diabetic Retinopathy, Explainable AI, EfficientNetB4, Swin Transformer, Grad-CAM, Medical Imaging

I. Introduction

This is one of the most common microvascular complications of diabetes mellitus and a leading cause of irreversible vision loss worldwide. Prolonged hyperglycemia damages retinal blood vessels, leading to microaneurysms, hemorrhages, exudates in DR, and, in advanced stages, retinal detachment and blindness. Early detection and timely intervention are critical to prevent disease progression; however, large-scale screening remains challenging due to the limited availability of ophthalmologists

and the time-intensive nature of manual fundus image examination. In recent years, **deep learning (DL)** techniques, particularly **convolutional neural networks (CNNs)**, have achieved remarkable success in automated diabetic retinopathy detection from retinal fundus images. CNN-based models effectively capture local spatial features such as lesions and texture patterns, enabling high classification accuracy. Nevertheless, conventional CNNs exhibit inherent limitations in modeling long-range dependencies and global contextual relationships across retinal images, which are essential for accurately distinguishing between different stages of DR.

To overcome these limitations, **Vision Transformers (ViTs)** have emerged as a powerful alternative, leveraging self-attention mechanisms to model global contextual information. Among these, the **Swin Transformer** introduces a hierarchical architecture with shifted window-based attention, offering improved computational efficiency and scalability for high-resolution medical images. While transformer-based models excel in capturing global representations, they often require large datasets and may underperform in extracting fine-grained local features when used in isolation.

Motivated by the complementary strengths of CNNs and transformers, this study proposes a **CNN–Transformer fusion framework** that integrates **EfficientNetB4** for local feature extraction with a **Swin Transformer** for global contextual learning. EfficientNetB4 provides an optimal balance between network depth, width, and resolution, ensuring high representational efficiency, while the Swin Transformer enhances contextual awareness across retinal structures. The fusion of these architectures enables more comprehensive feature learning for accurate DR classification.

Despite the promising performance of deep learning models, their **black-box nature** remains a significant barrier to clinical acceptance. Medical decision-support systems require transparency and interpretability to ensure trust, accountability, and regulatory compliance. To address this challenge, the proposed framework incorporates **Explainable AI (XAI)** techniques, specifically **Grad-CAM-based visual explanations**, to highlight retinal regions influencing model predictions. These explanations provide clinically meaningful insights that align with ophthalmological knowledge, thereby improving model reliability and usability.

Diabetic Retinopathy (DR) is a microvascular complication of diabetes that affects the retinal blood vessels and can result in irreversible blindness if not detected at an early stage. According to the World Health Organization, more than 500 million people worldwide are affected by diabetes, with a substantial proportion residing in developing countries such as India. Early detection of DR using retinal fundus imaging plays a vital role in preventing vision impairment; however, manual screening performed by ophthalmologists is time-consuming, subjective, and difficult to scale.

Automated deep learning-based screening systems have gained significant attention due to their ability to analyze retinal images with high accuracy. Convolutional Neural Networks (CNNs), including ResNet, Inception, and EfficientNet, have demonstrated strong performance in DR classification tasks. Despite these advancements, most deep learning models operate as black-box systems, providing little insight into how predictions are made. This lack of interpretability limits their acceptance in real-world clinical settings.

Explainable AI (XAI) techniques, such as Gradient-weighted Class Activation Mapping (Grad-CAM), address this limitation by visually highlighting the regions of an image that influence model predictions. In parallel, transformer-based architectures, particularly the Swin Transformer, have shown superior capability in capturing long-range dependencies and contextual information in images. Motivated by these developments, this research proposes a hybrid framework that combines EfficientNetB4 and Swin Transformer architectures with Grad-CAM explanations to achieve accurate and interpretable diabetic retinopathy detection.

II. Literature Review

Several deep learning approaches have been proposed for diabetic retinopathy detection. Early studies employed CNN architectures such as VGG16 and InceptionV3 for feature extraction from retinal images. Gulshan et al. (2016) demonstrated the feasibility of automated DR detection using large-scale datasets, achieving performance comparable to expert ophthalmologists. However, these methods primarily focused on predictive accuracy and did not address model interpretability.

EfficientNet, introduced by Tan and Le (2019), improved CNN performance through compound scaling of network depth, width, and resolution. More recently, vision transformers have been adapted for image analysis tasks. The Swin Transformer employs a hierarchical structure with shifted windows, making it computationally efficient and well-suited for high-resolution medical images.

Explainability techniques such as Grad-CAM, LIME, and SHAP have been used to visualize decision regions in medical image models. Pratt et al. (2021) demonstrated the usefulness of Grad-CAM in identifying DR-related lesions. However, hybrid CNN–Transformer architectures combined with explainability remain relatively underexplored. This work addresses this research gap by integrating EfficientNetB4, Swin Transformer, and Grad-CAM into a unified explainable DR detection framework.

III. Methodology

3.1 Dataset Description

This study utilizes the publicly available APTOS 2019 Blindness Detection and EyePACS datasets from Kaggle. Both datasets contain labeled retinal fundus images categorized into five DR severity levels: No DR (0), Mild (1), Moderate (2), Severe (3), and Proliferative DR (4). All images were resized to 512×512 pixels and preprocessed using cropping, illumination correction, and normalization. Data augmentation techniques such as random rotation, flipping, and brightness adjustment were applied to reduce class imbalance and improve generalization.

3.2 Model Architecture

The proposed framework adopts a hybrid deep learning architecture that leverages the complementary strengths of convolutional neural networks and transformer-based models to achieve accurate and interpretable diabetic retinopathy detection. EfficientNetB4 is employed as the backbone network for low-level feature extraction from retinal fundus images. Due to its compound scaling strategy, EfficientNetB4 efficiently captures fine-grained spatial features such as edges, textures, microaneurysms, hemorrhages, and exudates while maintaining a favorable balance between accuracy and computational complexity. This makes it particularly suitable for high-resolution medical images where subtle lesion patterns are crucial for diagnosis.

The feature maps extracted by EfficientNetB4 are then forwarded to a Swin Transformer module for high-level contextual representation. Unlike traditional CNNs, the Swin Transformer utilizes a hierarchical attention mechanism with shifted windows, enabling the model to capture long-range dependencies and global contextual information across different regions of the retina. This capability is essential for understanding complex pathological patterns that may span multiple retinal areas and for distinguishing between different stages of diabetic retinopathy.

Following the transformer-based feature refinement, the learned representations are flattened and passed through fully connected dense layers that perform five-class classification corresponding to the standard DR severity levels: No DR, Mild, Moderate, Severe, and Proliferative DR. The final

classification layer uses a softmax activation function to produce probability scores for each class, allowing the system to quantify prediction confidence.

To enhance transparency and clinical interpretability, Gradient-weighted Class Activation Mapping (Grad-CAM) is applied after model training. Grad-CAM generates class-specific heatmaps by computing the gradients of the predicted class score with respect to the feature maps of the final convolutional layers. These heatmaps visually highlight the retinal regions that most strongly influence the model's predictions, such as microaneurysms, hemorrhages, and exudates. By providing intuitive visual explanations, Grad-CAM enables clinicians to verify whether the model focuses on medically relevant features, thereby increasing trust and facilitating the adoption of AI-assisted diagnostic systems in real-world clinical settings.

3.3 Training Configuration

The proposed model was implemented using TensorFlow version 2.15, providing a flexible and efficient deep learning framework. Training was performed using the Adam optimizer with a learning rate of 0.0001, which ensures stable and fast convergence during optimization. Categorical cross-entropy was selected as the loss function, as it is well suited for multi-class diabetic retinopathy classification. The dataset was divided into an 80:20 training-to-validation ratio to effectively evaluate generalization performance. Additionally, five-fold cross-validation was employed to improve robustness and reduce bias. Early stopping was applied during training to prevent overfitting and ensure optimal model performance.

3.4 Evaluation Metrics

Model performance was evaluated using standard classification metrics to ensure a comprehensive assessment of the proposed system. Accuracy was used to measure the overall correctness of predictions, while precision and recall evaluated the model's ability to correctly identify diabetic retinopathy cases and minimize false predictions. The F1-score provided a balanced measure by combining precision and recall, particularly useful for handling class imbalance in medical datasets. The Area Under the Receiver Operating Characteristic Curve (AUC) was used to assess the model's discriminative capability across different decision thresholds. In addition to quantitative evaluation, explainability was qualitatively analyzed using Grad-CAM visualizations to verify that the model focuses on clinically relevant retinal regions.

IV. Results And Discussion

Experimental results demonstrate that the proposed hybrid model outperforms standalone EfficientNetB4 and Swin Transformer architectures across all evaluation metrics. The hybrid model achieved a classification accuracy of 95.3% and an AUC of 0.985. Grad-CAM visualizations confirm that the model focuses on clinically meaningful regions such as microaneurysms, hemorrhages, and exudates, thereby enhancing interpretability and clinical trust.

V. Conclusion And Future Work

This study presented an explainable deep learning framework for automated diabetic retinopathy classification and precise lesion segmentation using a hybrid CNN-Transformer architecture. By integrating the local feature extraction capability of Convolutional Neural Networks with the global contextual modeling strength of Transformer-based attention mechanisms, the proposed model

effectively identifies disease severity and accurately delineates infected retinal regions from fundus images. The incorporation of Explainable AI techniques enhances transparency by highlighting the critical regions that influence the model's predictions, thereby improving clinical trust and interpretability. The segmentation component further enables precise localization of pathological features such as microaneurysms, hemorrhages, and exudates, supporting early diagnosis and treatment planning. Experimental results demonstrate that the proposed approach achieves improved accuracy, robustness, and segmentation performance compared to traditional deep learning models. Overall, this work contributes to the development of reliable, interpretable, and clinically applicable intelligent screening systems for diabetic retinopathy, with strong potential for integration into real-world ophthalmic diagnostic workflows

REFERENCES

- [1] Ting DS, Cheung CY-L, Lim G, et al. Development and validation of a deep learning system for diabetic retinopathy and related eye diseases using retinal images from multiethnic populations with diabetes. *JAMA*. doi:10.1001/jama.2017.18152
- [2] Rishab Gargeya, Theodore Leng, MD, MS, AI-Based Diabetic Retinopathy Detection Using Deep Learning at American academy of ophthalmology
- [3] Shuang Yu, Di Xiao and Yogesan Kanagasigam, Exudate Detection for Diabetic Retinopathy with Convolutional Neural Networks
- [4] Artificial Intelligence with Deep Learning Technology Looks into Diabetic Retinopathy Screening Article in *JAMA The Journal of the American Medical Association* · November 2016
- [5] Khurshed et al., "Systematic Development of AI-Enabled Diagnostic Systems for Glaucoma and Diabetic Retinopathy"
- [6] Eugenio Vocaturo, Ester Zumpano." The contribution of AI In the detection of the Diabetic Retinopathy".
- [7] Kazi ahnaf alavee et al., "Enhancing Early Detection of Diabetic Retinopathy Through the Integration of Deep Learning Models and Explainable Artificial Intelligence".
- [8] Yuqing Yang et al., "Artificial Intelligence-Driven Diagnostic Systems for Early Detection of Diabetic Retinopathy: Integrating Retinal Imaging and Clinical Data".
- [9] Skylar Stolte et al., "A Survey on Medical Image Analysis in Diabetic Retinopathy".
- [10] Jordi DE LA TORRE et al., "Diabetic Retinopathy Detection through image analysis using Deep Convolutional Neural Networks".
- [11] Filippo Arcadu et al., "Deep learning algorithm predicts diabetic retinopathy progression in individual patients".
- [12] Lara Alsadoun et al., "Artificial Intelligence (AI)-Enhanced Detection of Diabetic Retinopathy from Fundus Images: The Current Landscape and Future Directions".
- [13] Daniel S.W. ting et al., "Deep learning in estimating prevalence and systemic risk factors for diabetic retinopathy: a multi-ethnic study".
- [14] Carol Y. Cheung et al., "Artificial Intelligence in Diabetic Eye Disease Screening: A Deep Learning Application in Retinal Imaging".