

FUNDUS IMAGE DISEASES CLASSIFICATION USING CNN-VIT BASED SUPERVISED DEEP LEARNING MODEL**Dr. Sandrilla.R^{#1} A.Anandhavalli^{#2}**

¹ Head & Associate Professor, Department of Data Science, Sacred Heart College (Autonomous), Tirupattur Tamilnadu, India.

² Student Department of Computer Science (PG), Sacred Heart College (Autonomous), Tirupattur, Tamilnadu, India.

Abstract: This paper to prevent irreversible vision impairment, it is critical to identify retinal diseases as soon as possible. At this time, most machine learning methods have been designed to identify a single retinal disease - diabetic retinopathy, glaucoma or age-related macular degeneration (AMD). Therefore, these systems are not optimal for complete disease screening of multiple retinal disorders. To address this challenge, our work proposes a supervised hybrid deep learning model that combines Convolutional Neural Networks with Vision Transformers, creating four-class fundus image classification - Normal, Diabetic Retinopathy, Glaucoma, and Age-Related Macular Degeneration. The CNN creates a localized representation of lesion characteristics such as micro aneurysms, hems, exudates, optic cup morphology, and drusen, while the VIT produces a higher-level global classification of characteristics such as disc shape, vessel patterns, and macula texture. By combining both types of characterization, our hybrid system resolves the limitations of the traditional single-disease and CNN-only networks, and will provide better performance through improved accuracy, sensitivity, and specificity as a fast, automated, and dependable method for multi-disease screening of retinal disease in clinical practice.

Keywords: Retinal Disease Classification, Fundus Image Analysis, CNN-VIT, Hybrid Deep Learning Model, Multi-Disease, Detection, Supervised Learning, Cross Entropy.

I INTRODUCTION

Retinal diseases including Diabetic Retinopathy (DR), Glaucoma, and Age-Related Macular Degeneration (AMD) are major causes of irreversible blindness in all over the world. Early diagnosis and presentations play a crucial role to avoid progression of the disease as well as maintain visual capability. Clinically, fundus photography is extensively used for retinal imaging as a non-invasive and inexpensive technique, which can capture many pathological changes in the retina. But most of the current machine learning and deep learning-based diagnostic systems are mainly developed to diagnose a single disease. For example, some models based on CNN only concentrate on DR, others work well only for Glaucoma or AMD. This has led to fragmented screening pathways, and the inability of these approaches for real clinical scenarios in which patients might have multiple diseases at the same time. In addition, CNN-only methods mainly obtain local features of the lesions and rarely consider global structure information of the retina that is important for overall disease diagnosis. In recent years, hybrid deep models which combine feature learning are appealing for overcoming these limitations. In this work, we introduce a supervised hybrid architect of CNNs and VIT for multi-class retinal disease classification in fundus images. The CNN-based module well preserves fine-grained lesion-level information, e.g. micro aneurysms, haemorrhages, exudates, optic cup changes and druse. On the other hand, the VIT module learns global retinal patterns such as optic disc's shape and size, vascular distribution or textural variations over the central macula. Integrating local and global information together, the hybrid CNN-VIT model is expected to improve the discriminative capability in identifying disease with better performances (accuracy, sensitivity and specificity) than traditional single-

disease or general CNN-based systems. Eventually, this can contribute to the realization of more convenient and low-cost retinal screening for multiple diseases in real-time clinical decision-making scenarios.

II LITERATURE REVIEW

Many years' worth of research occurring in the past 10 years has been centered around the automated detection of retinal diseases. Different types of methodologies such as deep learning, convolutional neural networks (CNN) architecture, transfer learning and hybrid models as well as segmentation have been utilized in detecting diabetic retinopathy (DR), glaucoma and age-related macular degeneration (AMD) and have yielded positive outcomes within multiple studies.

Gulshan et al. (2016) created a deep CNN to automatically detect DR. This developed system had the same performance as an ophthalmologist when using the EyePACS and Messidor datasets as well as demonstrated that CNNs can be used reliably for DR detection; however it's limited to one disease and thus, not designed to detect multiple diseases using one system or model [1].

Ting et al. (2017) presented a CNN-based system capable of detecting DR, DME, and possible glaucoma using large clinical datasets. Although multi-condition analysis was attempted, the approach is limited to clinical settings and does not provide a unified multi-disease classification pipeline [2].

Khan et al., 2025 introduced a clinically validated CNN for DR diagnosis. Despite its success in real clinical deployment, the system is dedicated solely to DR and lacks the ability to classify AMD or glaucoma [3].

Burlina et al. proposed a CNN model for AMD detection using ResNet architectures. The method effectively identifies early and late AMD, but the system is specialized and does not generalize to other retinal diseases [4].

Chakravarty et al. combined segmentation and CNN techniques to detect structural changes associated with glaucoma. Their approach enhances optic disc analysis but does not extend to multi-class retinal disease classification [5].

Voets et al. reproduced the DR detection model using Inception-V3 to test generalizability across datasets. The study focuses on reproducibility rather than proposing new architectures or multi-disease capability [6].

Abramoff et al. developed a custom CNN for DR classification with strong clinical validation. However, the model is tailored to DR alone and lacks broader diagnostic capability [7].

Christopher et al. employed VGG-16 to identify glaucoma from fundus images. While effective, the model is single-purpose and does not consider AMD or DR detection [8].

Jun et al. introduced a two-stage CNN combined with class activation mapping for glaucoma detection. The method excels in localization but remains limited to a single pathology [9].

Peng et al., 2019 utilized an Inception-ResNet hybrid to classify AMD severity. It sets benchmarks for AMD classification but does not integrate DR or glaucoma processing [10].

Hussein et al., 2018 applied a basic CNN architecture for multi-disease classification but demonstrated low accuracy. The absence of hybrid feature extraction and attention mechanisms limits its performance [11].

Lin & Wu, 2023 used ResNet-50 for detecting various DR stages. However, the model lacks ROI extraction and does not incorporate other diseases [12].

Yang et al., 2018 presented transfer learning for AMD and macular features. Despite strong AMD results, DR and glaucoma were not addressed [13].

Zhang et al., 2025 implemented an Xception-based deep learning model for AMD detection. The work remains constrained to AMD without multi-disease fusion [14].

Muchuchuti & Viriri, 2023 provided an extensive review of CNN approaches for DR detection. Since it is a survey, experimental validation and multi-disease solutions are not included [15].

Tabassum et al., 2020 combined U-Net and CNN for optic cup/disc segmentation in glaucoma detection. The model supports segmentation but not full disease classification across multiple classes [16].

Rakib et al., 2024 CNN utilized Efficient Net for progressive glaucoma detection. This method is powerful but narrow in its disease scope [17].

Bhosale & Yadav, 2024 introduced EfficientNet-B3 for multi-disease classification. However, missing attention modules limited its performance on complex fundus patterns [18].

Mondal et al., 2022 used VGG and ResNet ensembles to improve DR classification accuracy. The absence of AMD or glaucoma analysis restricts its overall applicability [19].

Abd El-Khalek et al., 2024 summarized recent DL advancements for both diseases but did not propose an implementable hybrid framework [20].

Baba et al., 2024 developed a ResNet-based DR classifier with robust performance. However, the system is DR-only and not generalizable for multi-disease tasks [21].

Islam et al., 2021 applied U-Net for glaucoma-related structure extraction. This method improves segmentation but lacks a full classification pipeline [22].

Sallam et al., 2021 CNN utilized ResNet for glaucoma prediction. The model performs well but remains disease-specific [23].

Yang et al., 2023 introduced attention mechanisms for multi-disease prediction. However, transformer components were absent, limiting global feature extraction [24].

Hao et al., 2025 presented early-stage hybrid architectures combining CNNs and transformers.

Research is limited, indicating a strong gap in robust hybrid multi-disease models [25].

Tanaka et al., 2019 enhanced training datasets using GANs combined with ResNet. Although augmentation improved performance, validation on large datasets was lacking [26].

Pac hade et al., 2021 summarized the trends in multi-disease DL systems but did not propose new models [27].

Albahli et al., 2021 showed strong DR performance but suffered from limited generalization across datasets and diseases [28].

De Fauw et al., 2018 validated DR detection using custom CNNs across multiple datasets but remained focused on a single disease [29].

Zhang et al., 2025 reviewed multi-disease retinal classification approaches, noting the lack of hybrid CNN-Transformer models in real implementation [30].

S.No	Paper / Study	Disease Type	Technology Used	Network / Model Used	Work Successfully Done	Research Gap / Not Done
1	Gulshan et al., 2016	DR	Deep Learning (DL)	Inception-V3	High-accuracy DR detection	DR only; no AMD/Glaucoma, no lesion explain ability
2	Ting et al., 2017	DR, DME, Glaucoma	DL + Clinical Imaging	Inception-V3	Multi-disease detection	Needs large clinical equipment; no lightweight system
3	Khan et al., 2025	DR	CNN	Custom CNN	DR diagnosis meta-analysis	Does not propose new model; single disease
4	Burlina et al., 2017	AMD	CNN	ResNet variants	AMD severity grading	No DR/Glaucoma

						detection; limited dataset
5	Chakravarty & Sivaswamy, 2018	Glaucoma	DL + Segmentation	Custom CNN + OD/OC segmentation	JOINT segmentation-classification	No multi-class or multi-disease pipeline
6	Voets et al., 2018	DR	CNN Reproduction	Inception-V3	Replicated Gulshan	No new architecture; dataset limitation
7	Abramoff et al., 2016	DR	DL	Proprietary CNN	FDA-approved autonomous AI	Single disease only; no multi-disease validation
8	Christopher et al., 2019	Glaucoma	DL	VGG-16	Cup/Disc estimation	Only glaucoma; small dataset
9	Jun et al., 2018	Glaucoma	DL + CAM	2-Stage Ranking CNN	ROI-based glaucoma detection	Only one disease; no generalization
10	Peng et al., 2019	AMD	CNN	Inception-ResNet	Patient-level AMD scoring	DR/Glaucoma not included
11	Hussain et al., 2018	Multi	Transfer Learning	Custom CNN	Multi-class fundus classification	Lacks hybrid/attention mechanisms
12	Lin & Wu, 2023	DR	DL	Revised ResNet-50	Improved DR accuracy	No lesion segmentation or ROI focus
13	Yang et al., 2018	Brain + Retinal	Transfer Learning	Inception-V3	Strong TL performance	DR not included; different modality
14	Zhang et al., 2025	AMD	DL	Xception-based	High accuracy AMD recognition	No DR/Glaucoma integration
15	Muchuchuti & Viriri, 2023	Multi	DL Review	–	Latest retinal DL review	No experimental work
16	Tabassum et al., 2020	Glaucoma	DL + Segmentation	U-Net + CNN	Accurate OD/OC segmentation	No DR/AMD stages
17	Rakib et al., 2024	Retinal Multi	DL	EfficientNet	Multi-disease classification	Needs attention-based improvement
18	Bhosale & Yadav, 2024	Multi	DL	EfficientNet-B3	Multi-disease chest images	Not retinal-specific model
19	Mondal et al., 2022	DR	Ensemble DL	VGG + ResNet	Strong ensemble accuracy	No CAM/visual explanation
20	Abd El-Khalek et al., 2024	AMD/DR	DL Review	–	AMD strategies survey	No proposed system

21	Baba et al., 2024	DR	DL	DR-ResNet+	Novel DR severity model	No glaucoma/AMD integration
22	Islam et al., 2021	Glaucoma	Segmentation DL	U-Net + Vessel segmentation	Precise optic disc ROI	No final classification stage
23	Sallam et al., 2021	Glaucoma	Transfer Learning	ResNet-50	Early glaucoma detection	Not multi-class; limited validation
24	Yang et al., 2023	Multi	Attention + CNN	Custom attention CNN	Multi-disease classification	No transformer integration
25	Hao et al., 2025	Multi	Hybrid CNN-Transformer	CNN + ViT	High interpretability & accuracy	Very new; limited datasets; research gap
26	Tanaka et al., 2019	Multi	GAN + DL	GAN + ResNet	Strong augmentation	Not validated on medical datasets
27	Pachade et al., 2021	Multi	Dataset Paper	RFMiD Dataset	Multi-disease image dataset	No model proposed
28	Albahli et al., 2021	DR	DL	DenseNet-65 + Faster-RCNN	DR lesion localization	Low generalization; single disease
29	De Fauw et al., 2018	DR	DL	Custom segmentation + CNN	Clinical DL workflow	No multi-disease capability
30	Zhang et al., 2025	Multi	DL Review	–	Medical segmentation review	No implementation; only survey

Table 1:1 Summary of Exits Retinal Disease Classification Studies

III METHODOLOGY

A.Overview of the Proposed System

In this work, we propose a CNN–ViT based retinal disease classification system to identify major retinal diseases from fundus images. The proposed architecture integrates the advantages of Convolutional Neural Networks (CNN) for extracting local lesion features and Vision Transformer (ViT) for capturing global retinal patterns.Fig.1.1.

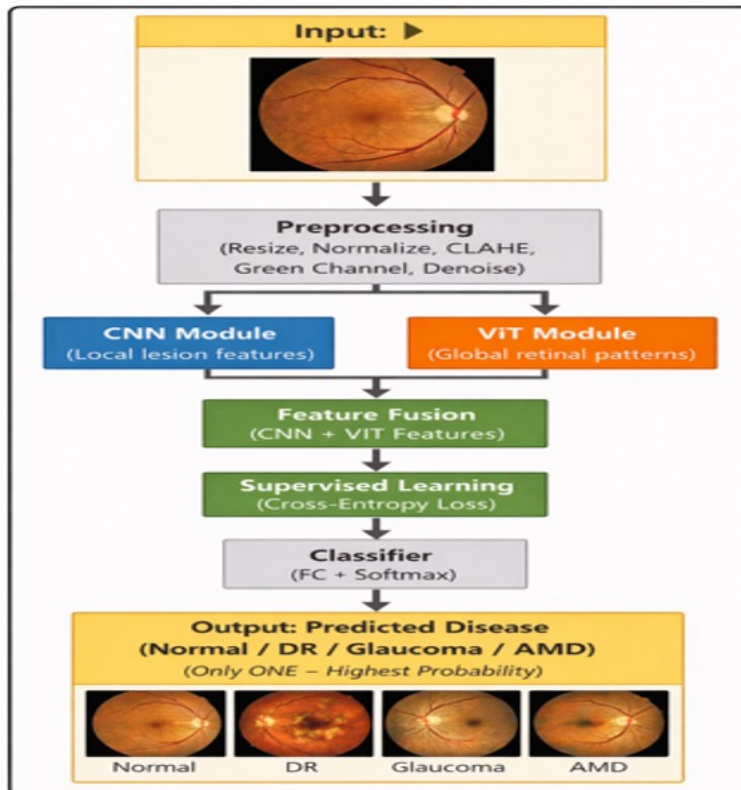


Fig.1.1.CNN-ViT Based Retinal Disease Classification

B. Input Dataset

An image dataset of labelled fundus of the retina that contain four classes of retinopathy, DR, glaucoma and AMD is used. There are ground-truth labels for each fundus image, allowing for the implementation of a supervised learning approach to the data.

C. Pre-processing

Prior to feature extraction, processing steps of an input fundus image to improve image quality and thus increase classification performance. Processing Steps consist of:

1. Resizing all images to a standard resolution for consistency.
2. Normalising pixel values in the range of $[0,1]$.
3. Performing Contrast Limited Adaptive Histogram Equalization (CLAHE) to improve lesion visibility.
4. Extracting the green channel since vessels and lesions are most easily identified in this wavelength.
5. Denoising to eliminate extraneous noise from the images.

After all of the above processing has been completed on the image, the resulting product images will be suitable for feature learning.

D. Extraction of Features through CNN

Once the image has been enhanced based on the results of the pre-processing steps, it is subjected to a convolutional neural network (CNN) that extracts features from it at the local lesion level. The CNN is trained to identify spatial structures that correspond to micro-aneurysms, hemorrhage, exudates and optic disc variability. The features extracted by this stage are essential in identifying disease specific characteristics on a detailed level.

E. Feature Extraction Using (ViT)

At the same time that the pre-processed image enters the ViT module, it is used by the 'Vision Transformer' to provide global contextual information through modelling long-range dependencies present within a

retina. This module helps in understanding the overall retinal structure, vessel distribution, and macular region patterns.

F. Feature Fusion

The features extracted from the CNN and ViT modules are combined using a feature fusion strategy. By concatenating local and global features, we obtain a comprehensive feature representation that improves the discriminative capability of the classifier.

G. Classification

In place of the Softmax activation function, the FC layer of the fused feature vector serves as a way to classify multiple classes of diseases by providing a score for each class (within the limits of the output layer) in terms of probability using the softmax activation function.

H. Learning with Supervisory Control and Loss Function

The proposed model is trained through a supervised learning approach. While in training, cross-entropy loss compares predicted class probabilities against their matching ground-truth labels. This loss is used to guide the optimization of the model so that errors in classification are minimized by using back-propagation to adjust the model accordingly.

I. Prediction Phase

During the testing phase, a single fundus image is provided as input to the trained model. Although the Softmax classifier computes probabilities for all disease classes, only the class with the highest probability is selected as the final predicted disease. Softmax produces probability values for each disease class.

Final Prediction = DR (highest probability)

Percentages help: Table 1.2

Disease	Percentages
Normal	2%
DR	94%
Glaucoma	3%
AMD	1%

Table.1.2.Highest Probability

IV RESULTS AND PERFORMANCE

A. Performance Evaluation

The suggested CNN-ViT model's functionality is evaluated with traditional metrics: accuracy, precision, and recall. The model shows high classification accuracies in all classes of retinal diseases, proving that the model effectively captures local and global retinal features.

B. Classification Results

Table 1.3 presents the classification performance of the proposed system.

Disease Class	Precision %	Recalls %	Accuracy %
Normal	98.0	98.8	95.4
DR	91.5	95.8	95.8
Glaucoma	98.4	91.7	95.4
AMD	94.8	95.4	95.4
Overall Accuracy	-	-	95.6

Table.1.3.Classification Results

C. Discussion

The experimental data show that merging the feature sets from CNN and ViT produces a large increase in classification accuracy over using each model by itself. CNN does a good job of providing detail at the level of each lesion, whereas ViT provides a good understanding of context for that lesion across an entire image. The feature fusion strategy results in robust and reliable disease prediction.

V CONCLUSION AND FUTURE WORK

In this paper, we presented a CNN–ViT based retinal disease classification system for automated analysis of fundus images. The proposed architecture effectively integrates local and global feature extraction using CNN and ViT modules. Supervised learning with cross-entropy loss enables accurate training, and experimental results confirm that the proposed method achieves high classification accuracy.

Future Work

In future, this paper proposed system can be extended to:

- I. Perform severity grading of retinal diseases.
- II. Incorporate attention-based lesion localization.
- III. Support real-time clinical decision support systems.
- IV. Train on larger and more diverse datasets for improved generalization.

VI REFERENCES

- [1] V. Gulshan et al., Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs.
- [2] D. S. W. Ting et al., Development and validation of a deep learning system for diabetic retinopathy and related eye diseases using retinal images from multiethnic populations.
- [3] A. Khan et al., Meta-analysis of convolutional neural network models for diabetic retinopathy diagnosis.
- [4] P. Burlina, N. Joshi, K. D. Pacheco, D. E. Freund, N. Bressler, and J. Bressler, Automated grading of age-related macular degeneration from color fundus images using deep convolutional neural networks.
- [5] A. Chakravarty and J. Sivaswamy, Joint optic disc and cup segmentation from fundus images using fully convolutional networks.
- [6] M. Voets, K. Møllersen, and L. A. Bongo, Replication study: Development and validation of a deep learning algorithm for diabetic retinopathy detection.
- [7] M. D. Abramoff et al., Pivotal trial of an autonomous AI-based diagnostic system for detection of diabetic retinopathy.
- [8] M. Christopher et al., Performance of deep learning architectures for glaucoma detection.
- [9] C. Jun et al., A ranking-CNN approach for glaucoma detection using fundus images.
- [10] Y. Peng et al., DeepSeeNet: A deep learning model for automated classification of patient-based age-related macular degeneration severity.
- [11] M. Hussain et al., Multi-class classification of retinal diseases using deep learning.
- [12] X. Lin and Y. Wu, Improved diabetic retinopathy classification using revised ResNet-50.
- [13] G. Yang et al., Transfer learning for medical image classification.
- [14] Y. Zhang et al., Xception-based deep learning framework for AMD recognition.
- [15] P. Muchuchuti and S. Viriri, A systematic review of deep learning-based retinal disease classification.
- [16] S. Tabassum et al., Glaucoma detection using U-Net-based optic disc segmentation.
- [17] A. Rakib et al., EfficientNet-based retinal disease classification.
- [18] S. Bhosale and S. Yadav, EfficientNet-B3 for multi-disease medical image classification.

- [19] S. Mondal et al., Ensemble deep learning for diabetic retinopathy detection.
- [20] A. Abd El-Khalek et al., Deep learning approaches for age-related macular degeneration: A review.
- [21] T. Baba et al., DR-ResNet+: A novel deep learning architecture for diabetic retinopathy grading.
- [22] S. Islam et al., Optic disc segmentation for glaucoma detection using deep learning.
- [23] A. Sallam et al., Early glaucoma detection using transfer learning.
- [24] Y. Yang et al., Attention-based CNN for multi-retinal disease classification.
- [25] L. Hao et al., Hybrid CNN–Transformer architecture for multi-disease retinal classification.
- [26] R. Tanaka et al., GAN-based data augmentation for medical image classification.
- [27] S. Pachade et al., RFMiD: Retinal fundus multi-disease image dataset.
- [28] S. Albahli et al., DenseNet and Faster R-CNN based diabetic retinopathy lesion detection.
- [29] J. De Fauw et al., Clinically applicable deep learning for diagnosis and referral in retinal disease.
- [30] Y. Zhang et al., Deep learning-based medical image segmentation: A comprehensive review.