

# Exploring Dataset Variability in Diabetic Retinopathy Classification Using Transfer Learning Approaches

Kinjal Patni<sup>1</sup>, Shruti Yagnik<sup>2</sup> and Pratik Patel<sup>3</sup>

<sup>1</sup> Department of Computer Engineering, Indus Institute of Technology and Engineering, Indus University, Ahmedabad, Gujarat, India

<sup>2</sup> Department of Information Technology, Indus Institute of Technology and Engineering, Indus University, Ahmedabad, Gujarat, India

<sup>3</sup> Department of Computer Science and Engineering, Parul Institute of Engineering and Technology, Parul University, Vadodara, Gujarat, India

**Corresponding author:** Kinjal Patni (e-mail: [kinjalpatni11@gmail.com](mailto:kinjalpatni11@gmail.com)); **Email Author(s):** Shruti Yagnik (e-mail: [shruti.yagnik.ce@indusuni.ac.in](mailto:shruti.yagnik.ce@indusuni.ac.in)), Pratik Patel (e-mail: [pratik.patel2988@paruluniversity.ac.in](mailto:pratik.patel2988@paruluniversity.ac.in))

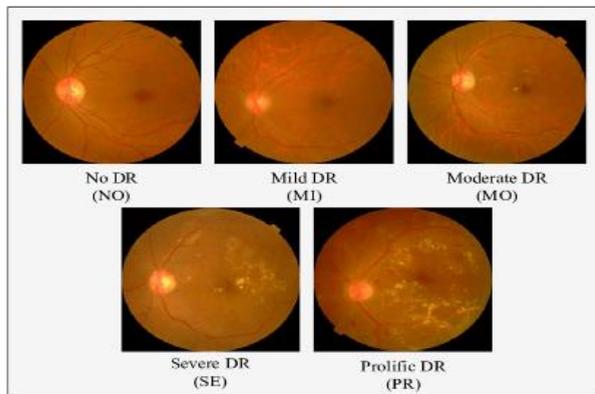
**Abstract** Diabetic retinopathy (DR) stands as a primary international cause of vision impairment that needs effective and swift diagnostic services to protect eye structures from advancing deterioration. The variations of imaging data that appear between sources create major obstacles for achieving consistent performance from models. The elimination of performance fluctuation problems during DR classifications across two benchmark datasets EYE-PACS and APTOS is examined through systematic transfer learning analysis using different high-performing CNN architectures including VGG16, VGG19, ResNet50, Xception, InceptionV3, MobileNetV2, and InceptionResNetV2. The research evaluates how data heterogeneity affects and how augmentation approaches impact the accuracy while stabilizing robustness in deep learning models. The research provides new insights through its extensive investigation of generalization performance based on dataset changes which utilize modified data augmentation methods for retinal images. A collection of data transformations such as rotation, flipping, zooming and brightness modifications create simulated realistic scenarios to handle imbalanced data classes. Academic research involved CNN pre-training followed by transfer learning on both databases while researchers evaluated the models through both untreated source data and augmented image testing procedures. InceptionResNetV2 outperformed its counterparts with 96.2% accuracy and Xception delivered 95.7% accuracy in APTOS evaluation and both models scored 95.9% and 95.4% respectively on EYE-PACS testing. When augmentation was applied it increased the performance level by 3% to 5% across all running models. The experimental outcomes demonstrate how adequate variable training allows these models to recognize datasets regardless of their heterogeneity. This analysis confirms that combining reliable deep learning structures with purposeful data enhancement techniques substantially enhances DR diagnosis reliability to build scalable future diagnostic solutions for ophthalmology practice.

**Keywords** Diabetic Retinopathy, Transfer Learning, Dataset Variability, Data Augmentation, Deep Learning Models.

## 1. Introduction

Diabetic retinopathy causes serious damage to the retina by affecting the light-sensitive tissue that exists at the back of the human eye [1,2]. Working-age adults around the world primarily suffer blindness because of this condition. Retinal vessels sustain damage from sustained elevated blood glucose levels that trigger fluid leakage as well as bleeding and the growth of aberrant blood vessels. It progresses through four stages: mild non-proliferative retinopathy, moderate

non-proliferative retinopathy, severe non-proliferative retinopathy, and proliferative diabetic retinopathy [4,5,6]. The early recognition of diabetic retinopathy along with immediate medical intervention greatly diminishes vision deterioration risks thus underscoring the need for precise diagnostic tools [7,8,9]. As shown in Fig. 1, diabetic retinopathy is categorized into five stages: No DR (NO), Mild (MI), Moderate (MO), Severe (SE), and Prolific (PR) based on the progression of retinal damage.



**Fig. 1. Stages of Diabetic Retinopathy [28]**

The severity increases from top left to bottom right, with visible changes in retinal blood vessels and lesions. The early recognition of diabetic retinopathy along with immediate medical intervention greatly diminishes vision deterioration risks thus underscoring the need for precise diagnostic tools [7,8,9]. DR classification systems continue to show improvements, yet they encounter various obstacles. The primary challenge arises from inconsistent dataset characteristics which introduce variations in image resolution contrast and noise levels [15,16,17]. Each dataset presents different characteristics that produce inconsistent evaluation outcomes when models apply to new testing cases. Datasets with class imbalances between DR stages create difficulties because they force models toward biased detection of underrepresented stages [18,19,20]. The generalization ability of models suffers from reduced performance because limited annotated datasets make overfitting possible. Preprocessing and augmentation methods specifically designed for medical images are missing from many current computer systems which significantly affect model performance and generalization potential. [21,22]

This work represents a novelty in its examination of transfer learning models across multiple datasets throughout the process with and without data augmentation techniques. Previous investigations studied single models on individual datasets, but this study analyzes how datasets compare with augmentation techniques across various models. Through assessments of the APTOS [32] and EYPACS [33] datasets this research demonstrates the hurdles and potential benefits found in cross-dataset generalization practices. The study evaluates individual features of both models in new dataset conditions and pinpoints vital aspects for advanced DR classification system optimization. This preliminary research demonstrates that variations in datasets substantially influence model operational efficiency. The models InceptionResNetV2 and Xception produce outstanding accuracy results and handle generalization

through data augmentation applications. Multiple augmentation techniques boost model effectiveness across the board, especially by making better detections of minority disease stages possible. Different model architectures exhibit varying effective performance with their parameters about different research datasets. Data augmentation proves particularly beneficial for APTOS [39] images because of their high resolution yet EYPACS [40] images present additional challenges because of their real-world random elements.

The research findings demand collaboration between experts to establish standardized datasets while building generalizable models that will shape the future of medical imaging. These advancements represent significant progress toward minimizing diabetic retinopathy impacts worldwide while generating better treatment results. Here are the key research contributions presented in bullet points for clarity:

1. Comparative Analysis of Models: A complete analysis compares state-of-the-art transfer learning architectures consisting of VGG16, VGG19, ResNet50, Xception, InceptionV3, MobileNetV2, and InceptionResNetV2 for diabetic retinopathy classification tasks.
2. Dataset Diversity Assessment: Examination of dataset variability included a model evaluation on both EYPACS [39] and APTOS [40] datasets which presented different characteristics regarding image quality and class distribution patterns.
3. Data Augmentation Impact: This study investigates how advanced data augmentation methods affect model performance when used for generalizing between different datasets.
4. Insights into Cross-Dataset Generalization: Experts must identify both hurdles and possibilities in maintaining stable classification output during analysis of datasets that exhibit different characteristics.
5. Performance Optimization Guidelines: This research presents practical guidelines to create Impressive and Scalable Diabetic Retinopathy Detection Systems through Transfer Learning in combination with Data Augmentation methods.

## II. Literature Study

Szegedy et al. [1] first showed the exposure of deep neural networks to adversarial perturbations, emphasizing the need for robust defense methods. Goodfellow et al. [2] suggested the Fast Gradient Sign Method (FGSM) for adversarial preparation, teaching the practice of training models with adversarial

samples. Feinman et al. [3] introduced statistical detection using kernel density estimation and Bayesian uncertainty. Gal and Ghahramani [4] leveraged Monte Carlo dropout to estimate model uncertainty under adversarial settings. Hendrycks and Gimpel [5] demonstrated that reconstruction error using PCA could effectively differentiate adversarial inputs from clean ones. Li and Li [6] extended this idea using dropout variances to measure categorization confidence.

Guo et al. [7] proposed input transformation defenses, including bit-depth reduction and JPEG compression, to mitigate adversarial effects. Xu et al. [8] utilized image quilting and total variance minimization for input pre-processing. Xie et al. [9] introduced random resizing and padding to break gradient flow in white-box attacks. Samangouei et al. [10] applied GANs for purifying adversarial images. Liao et al. [11] enhanced this approach by projecting perturbed samples back onto the data manifold using a VAE-GAN model. Lee et al. [12] used the GAN discriminator itself as a binary detector of adversarial samples.

Ilse et al. [13] used attention-based heatmaps to localize adversarial regions in the input space. Song et al. [14] proposed spatial attention mechanisms to detect subtle perturbations. Zhang et al. [15] embedded attention layers within CNNs to highlight critical areas affected by adversarial noise. Wang et al. [16] introduced a temporal attention framework for detecting adversarial attacks in sequential data. Metzén et al. [17] added auxiliary detectors to intermediate layers of CNNs for local detection of adversarial inputs. Meng and Chen [18] proposed MagNet, a two-network defense using one for detection and another for reforming inputs.

Abbasi and Gagné [19] presented a voting-based ensemble to identify adversarial inputs based on inconsistent predictions. Lakshminarayanan et al. [20] introduced deep ensembles and measured uncertainty through prediction entropy. Pang et al. [21] proposed selectively using diverse models to detect attacks. Cisse et al. [22] introduced the concept of model Lipschitz continuity and its role in bounding adversarial noise. Miyato et al. [23] used spectral normalization to enforce this bound and reduce model sensitivity. Katz et al. [24] proposed Reluplex for formal verification of neural network robustness. Gowal et al. [25] applied interval bound propagation to certify networks against adversarial perturbations.

Xie et al. [26] combined feature analysis and image

transformations to develop a hybrid defense. Tramèr et al. [27] combined adversarial training with on-the-fly data augmentation to enhance generalization. Raff et al. [28] constructed a modular defense pipeline incorporating multiple filters. Khoury and Hadfield-Menell [29] applied meta-learning to select optimal defense strategies based on input characteristics. Ilyas et al. [30] modeled deep neural networks as dynamical systems to study adversarial behavior.

Carlini and Wagner [31] benchmarked multiple adversarial defenses against adaptive attacks on CIFAR-10. Athalye et al. [32] evaluated defenses on ImageNet and exposed the limitations of gradient masking. Papernot et al. [33] demonstrated the transferability of adversarial examples across models. Rouhani et al. [34] proposed embedding watermarks in neural activations to detect tampering. Chen et al. [35] used contrastive learning to increase feature separation between adversarial and clean inputs. Wu et al. [36] developed a loss function that aligns internal representations under adversarial conditions. Qin et al. [37] explored representation dissimilarity as a detection metric. Liu et al. [38] proposed a semantic-preserving adversarial loss for improved robustness.

Despite significant advancements in adversarial detection methods, several gaps remain that are directly relevant to the challenges addressed in this study. Most existing techniques are designed for general image classification and often fail to consider the unique characteristics of medical datasets such as diabetic retinopathy, where class imbalance, subtle pathological features, and variability across datasets are common. Furthermore, many models lack robustness when applied to diverse clinical imaging datasets, resulting in poor generalization. Additionally, the absence of domain-specific benchmarking protocols for diabetic retinopathy hinders the evaluation and comparison of model effectiveness. These gaps underscore the need for adaptive transfer learning strategies that can leverage dataset-specific features while maintaining high diagnostic accuracy across varied sources.

### III. Methodology

As shown in Fig. 2, the proposed system classifies diabetic retinopathy into five stages: No DR (NO), Mild (MI), Moderate (MO), Severe (SE), and Proliferative (PR). It accurately detects and differentiates the severity levels based on visible retinal abnormalities.

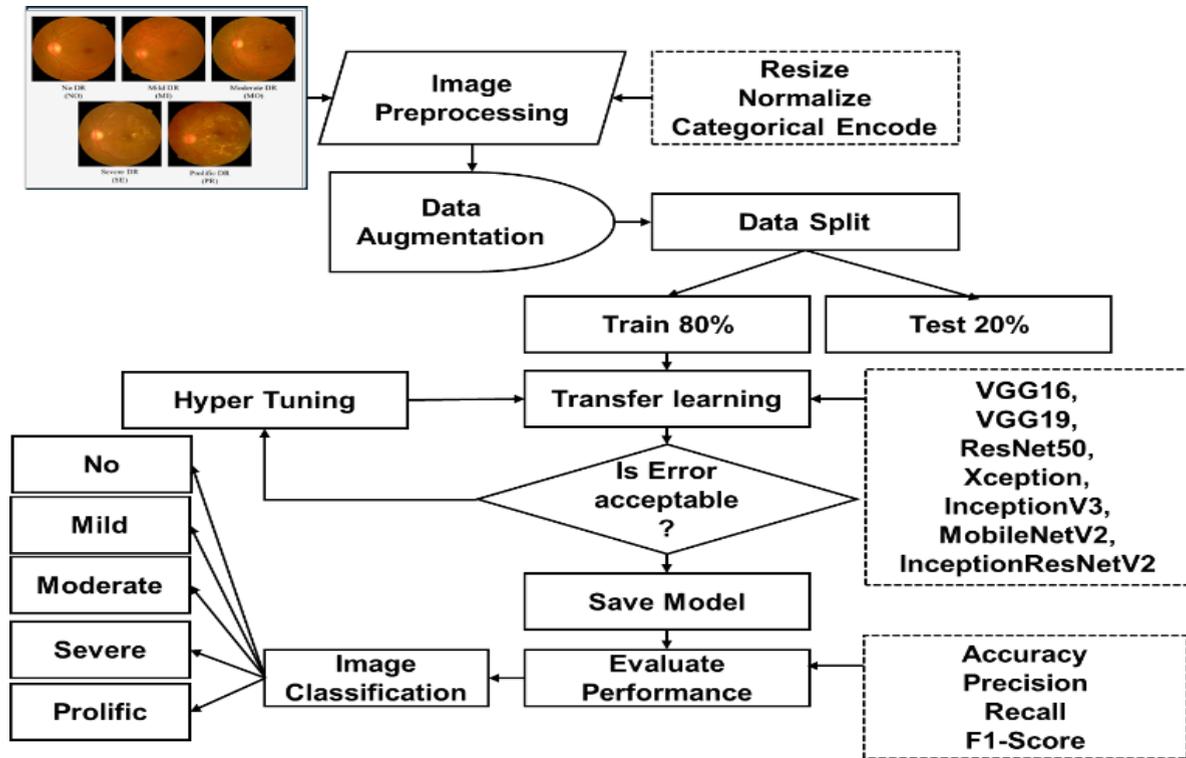


Fig. 2 Methodology Steps for DR Classification

## A. Datasets

Researchers used two popular databases named APTOS [39] and EyePACS [40] which find widespread application in ophthalmology and diabetic retinopathy detection tasks. On. Each dataset contains high-definition retinal fundus images that are assigned to different severity labels for diabetic retinopathy. The following sections deliver an in-depth analysis of each dataset.

### 1. APTOS Dataset

The Asia Pacific Tele-Ophthalmology Society released APTOS [39] as a component of their Kaggle competition. It contains thousands of retinal images captured under varying conditions and labelled into five classes based on the severity of diabetic retinopathy: no diabetic retinopathy (class 0), mild (class 1), moderate (class 2), severe (class 3), and proliferative diabetic retinopathy (class 4). The dataset's valuable aspects stem from its image quality diversity together with illumination variations and clarity ranges which allow researchers to develop robust prediction models.

### 2. EyePACS Dataset

EyePACS [40] (Eye Picture Archiving and Communication System) is another widely recognized dataset in diabetic retinopathy research. It includes a large collection of retinal fundus images, similarly, categorized into five severity levels. EyePACS [40]

provides a more extensive range of samples compared to APTOS [39], which helps improve the generalization capability of deep learning models. Like APTOS [39], the images in EyePACS [40] vary significantly in terms of brightness, contrast, and focus, challenging models to adapt to real-world variations.

These datasets complement each other by providing diverse samples, ensuring that models trained on them can perform effectively in various clinical settings.

## B. Pre-Processing

Deep learning models require pre-processing as their first step toward data preparation. The process of data preparation includes both cleaning procedures and standard image size adjustments and normalization techniques which help enhance model accuracy outcomes. Below are the key pre-processing steps implemented:

### 1. Resize

Every image was resized to the standardized 224x224 pixel shape for processing. Special consideration went into selecting this width because it traded off between computational speed and essential retinal characteristics preservation. The resizing procedure standardizes all images into a standardized size while fulfilling the requirements of pre-trained convolutional neural network (CNN) architecture models which operate with set input dimensions. Given a point  $(x, y)$  in the target image, its pixel value is calculated using

the surrounding four pixels in the original image as follows Eq. (1)[1,2]:

$$f(x, y) = (1 - a)(1 - b)f(i, j) + a(1 - b)f(i + 1, j) + (1 - a)b f(i, j + 1) + ab f(i + 1, j + 1) \quad (1)$$

In the given bilinear interpolation equation,  $f(x, y)$  is the interpolated pixel value at non-integer coordinates  $(x, y)$ , while  $f(i, j)$  represents the pixel intensity at integer coordinates in the input image. Here,  $i = \lfloor x \rfloor$  and  $j = \lfloor y \rfloor$ , with  $a = x - i$  and  $b = y - j$  indicating the fractional distances along the  $x$  and  $y$  axes.

## 2. Normalize

Pixel intensity normalization adjusts the values to operate on a standard range. To normalize pixel scales the algorithms divide every value between 0 to 255 by 255 leading to values between 0 to 1. The normalization process enables deep learning model training to accelerate by keeping input values within convenient operational parameters defined as Eq. (2)[3,5]:

$$x_{\{norm\}} = (x - \mu) / \sigma \quad (2)$$

where  $x$  is the original pixel value,  $\mu$  is the mean, and  $\sigma$  is the standard deviation of pixel intensities. This process ensures that input data remains within optimal bounds, facilitating faster and more stable training of deep learning models.

## 3. Categorical Encoding

The diabetic retinopathy severity label data received categorical encoding to establish distinctive variables. One-hot encoding converted labels to binary matrices through its implementation. The model encodes class 2 (moderate) severity into the binary vector  $[0, 0, 1, 0, 0]$ . One-hot encoding represents a critical technique for multi-class classification which enables models to generate accurate probability predictions across all possible classes. If a class label is 'k' among C classes, the one-hot encoded vector define as Eq. (3)[4,6]:

$$y_k = [0, 0, \dots, 1, \dots, 0] \quad (1 \text{ at the } k - \text{th position}) \quad (3)$$

## 4. Augmentation

The artificial expansion of training image collections through data augmentation delivers a strong way to increase both training dataset size and picture variety. The modelling process reduces overfitting impacts and enables improved generalization boundaries. Geometric along with photometric and elastic deformation augmentation methods served as the basis for this research study.

## 5. Geometric Transformations

Geometric transformations involve altering the spatial arrangement of the image without affecting the essential features. The following geometric transformations were applied:

**Rotation:** The simulation applied random angular rotations of up to  $\pm 15$  degrees to represent different orientations of retina tissue.

**Flipping:** Data variability was improved through both horizontal and vertical orientation transformations.

**Scaling:** To mimic retinal images captured at various levels of magnification random scaling transformations were applied to the data.

**Translation:** Test images received randomized movements across both the  $x$  and  $y$  axes to depict off-center acquisition conditions.

All processes show in below Eq.(4)[1,8]:

$$I' = T_{trans} \circ T_{scale} \circ T_{flip} \circ T_{rot}(I) \quad (4)$$

In image augmentation, various transformations are applied to the original image  $I$  to produce an augmented image  $I'$ . These transformations include  $T_{rot}$ , which applies a random rotation  $\theta$  within the range  $[-15^\circ, +15^\circ]$ ;  $T_{flip}$ , which performs horizontal and/or vertical flipping;  $T_{scale}$ , which applies random scaling using a factor  $s$  within the interval  $[s_{min}, s_{max}]$ ; and  $T_{trans}$ , which translates the image along the  $x$  and  $y$  axes by distances  $t_x$  and  $t_y$  respectively. These augmentations increase dataset diversity and improve model generalization in training deep learning systems.

## 6. Photometric Transformations

Photometric augmentations modify the pixel intensity values to account for variations in image acquisition conditions. The following photometric transformations were utilized:

**Brightness Adjustment:** The brightness of the images was randomly increased or decreased to simulate different lighting conditions. Brightness adjustment can be modeled as Eq. (5)[1,4]:

$$I_{new} = I + \Delta B \quad (5)$$

where  $I$  is the original image and  $\Delta B$  represents the brightness shift value. The term  $\Delta B$  can be positive or negative, allowing enhancement or dimming of the image respectively. This transformation helps improve the robustness of deep learning models by training them on images with different illumination levels.

**Contrast Adjustment:** Contrast levels were randomly altered to enhance or reduce the differences between light and dark regions by Eq. (6)[1,8].

$$I_{new} = \alpha * (I - \mu) + \mu \quad (6)$$

Where  $I$  is the original image,  $\mu$  is the mean pixel intensity, and  $\alpha$  is the contrast scaling factor. When  $\alpha > 1$ , the contrast is increased, and when  $\alpha < 1$ , the

contrast is reduced. This augmentation technique allows models to better generalize to varying contrast conditions in real-world data.

**Saturation Adjustment:** Color saturation varies to mimic differences in imaging devices.

**Hue Adjustment:** Minor changes were made to the hue values to account for color variations.

## 7. Elastic Deformation

Elastic deformation is a more advanced augmentation technique that introduces random, localized distortions in the image. This method is particularly useful for medical images as it simulates variations in tissue shape and appearance while preserving the underlying anatomical structure. Elastic deformation involves: Applying a random displacement field to the image. Smoothing the displacement field with a Gaussian filter to ensure realistic deformation. The transformation is expressed as Eq. (7 and 8)[1,11]:

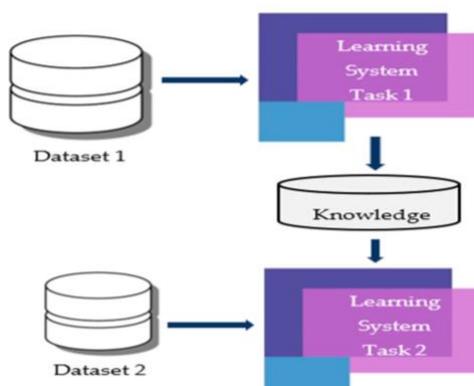
$$x' = x + \alpha * G_x * \varphi_x \quad (7)$$

$$y' = y + \alpha * G_y * \varphi_y \quad (8)$$

Where  $\varphi_x$ ,  $\varphi_y$  are random fields,  $G$  is a Gaussian filter,  $\alpha$  controls the intensity of deformation. These augmentations were carefully selected to enhance the robustness of the models and ensure that they can handle the inherent variability in real-world retinal fundus images.

## C. Transfer Learning Models

As shown in Fig. 3 Transfer learning involves transferring knowledge from a source domain (where a model has already been trained) to a target domain (where we want to apply or adapt the model). Instead of training a model from scratch, we leverage pre-trained models (usually trained on large datasets like ImageNet) to solve new but related problems.



**Fig. 3 Basic Flow of Transfer learning [12]**

Transfer learning is a machine learning technique where a model developed for one task is reused as the

starting point for a model on a second, related task. In the domain of image classification, particularly medical imaging, transfer learning has become a highly effective approach due to the limited availability of large, annotated datasets. The core idea is to leverage the learned features of a pre-trained deep learning model often a convolutional neural network (CNN) to solve a new classification problem with relatively less training data.

The process begins with selecting suitable pre-trained CNN architecture, typically trained on a large benchmark dataset such as ImageNet. These models have already learned to extract and represent a rich hierarchy of features, from low-level edges and textures to high-level object patterns. Instead of training a new model from scratch, the early layers of the pre-trained model are retained because they capture general features that are transferable across domains. These layers are often "frozen," meaning their weights are not updated during training on the new task, ensuring that the learned generic features remain intact.

Next, the final classification layers of the model are removed and replaced with new layers specific to the target problem. These newly added layers are usually composed of fully connected layers, activation functions (such as SoftMax for multi-class problems), and sometimes dropout layers for regularization. These layers are initialized with random weights and are trained to map the extracted features to the desired output classes. In many cases, a few of the deeper convolutional layers are also "unfrozen" to allow fine-tuning, enabling the model to adjust to more task-specific patterns in the new dataset.

Data preprocessing and augmentation are critical components of this process. Input images are typically resized to match the expected input size of the pre-trained model and normalized for consistent intensity levels. Data augmentation techniques such as horizontal flipping, rotation, zooming, cropping, and contrast adjustment are applied to artificially expand the training dataset and improve the model's robustness to variations in image quality and conditions.

The training phase involves feeding the augmented image data through the modified model, where forward propagation computes the output predictions, and backpropagation updates the weights of the unfrozen layers based on the error. The optimization algorithm (commonly Adam or SGD) and a suitable learning rate are used to minimize the loss function, often categorical cross-entropy for classification tasks. Early stopping and learning rate scheduling may be employed to enhance training efficiency and avoid overfitting.

Finally, the trained model is evaluated using unseen test data to assess its generalization performance. Metrics such as accuracy, precision, recall, and confusion matrices are used for performance evaluation. The transfer learning approach significantly reduces training time and computational cost while improving performance, especially in domains where labeled data is limited. Its success lies in the reusability of learned visual

representations, making it a practical and efficient choice for various real-world image classification tasks, including those in the medical field.

Table 1. leverages pre-trained deep learning models, which have been trained on large datasets such as ImageNet, to accelerate training and improve performance in specific tasks. In this study, the following pre-trained architectures were utilize

**Table 1. Transfer Learning Model**

Model	Architecture Type	Key Features	Advantages	Limitations	Suitability for Retinal Image Analysis
VGG16	Sequential Convolutional + FC layers	16 layers; simple and deep design	Easy to implement; effective for transfer learning	High parameter count; computationally expensive	Good for baseline feature extraction
VGG19	Sequential Convolutional + FC layers	19 layers; similar to VGG16 with added depth	Slightly better feature extraction than VGG16	Even more parameters than VGG16	Suitable for detailed hierarchical feature extraction
ResNet50	Deep Residual Network with skip connections	50 layers; identity mapping for vanishing gradient mitigation	Enables very deep network training; high-level feature extraction	Complex training process; heavier than VGG	Excellent for complex retinal feature extraction
Xception	Depthwise Separable Convolutions	Channel and spatial separation; efficient learning	Reduced parameters; high accuracy	More sensitive to hyperparameters	Strong performance with fewer parameters
InceptionV3	Inception Modules with optimization enhancements	Multi-scale feature extraction; factorized convolutions	Efficient and powerful; captures fine details	Complex architecture	Very suitable for capturing intricate retinal patterns
MobileNetV2	Lightweight CNN with inverted residuals	Depthwise separable convolutions; mobile-friendly	Fast inference; good accuracy on constrained hardware	Lower performance on very high-resolution data	Efficient for real-time or mobile retinal diagnosis systems
InceptionResNetV2	Hybrid of Inception and Residual connections	Deep network combining Inception's scale-aware learning with ResNet's stability	Best of both worlds; highly accurate; deep architecture	Very high resource requirements	Ideal for high-precision retinal image analysis in clinical settings

The performance of the diabetic retinopathy classification model was evaluated using four standard metrics: Accuracy Eq. (9)[1,8], Precision Eq. (10)[1,8], Recall Eq. (11)[1,8], and F1-Score Eq. (12)[1,8]. These metrics are defined as follows in the context of diabetic retinopathy classification:

$$Accuracy = \frac{(TP + TN)}{(TP + TN + FP + FN)} \quad (9)$$

$$Precision = \frac{TP}{(TP + FP)} \quad (10)$$

$$Recall = \frac{TP}{(TP + FN)} \quad (11)$$

$$F1 - score = 2 * \frac{(Precision * Recall)}{(Precision + Recall)} \quad (12)$$

where TP is True Positives, TN is True Negatives, FP is False Positives, and FN is False Negatives.

IV. Results

The experiments were conducted using the Kaggle platform with a Tesla T4 GPU. Two publicly available datasets were used for training and evaluation: the APTOS 2019 Blindness Detection dataset and the EyePACS dataset. All models were trained using the Adam optimizer with an initial learning rate of 0.0001, batch sizes of 32 or 64, and 50–100 epochs. Learning rate reduction on plateau, early stopping (patience 10),

categorical crossentropy loss, and data augmentation (rotation, flipping, zoom) were applied, with 10–20% validation split. Fig. 4 provides APTOS Data testing data confusion matrices for various transfer learning models trained and tested with normal data or data augmentation. Additional data to every model tends to lead to better performance, as shown by a higher number of correct decisions and fewer mistaken classifications; this can be noticed in the near-perfect results achieved by Inception-ResNet50.

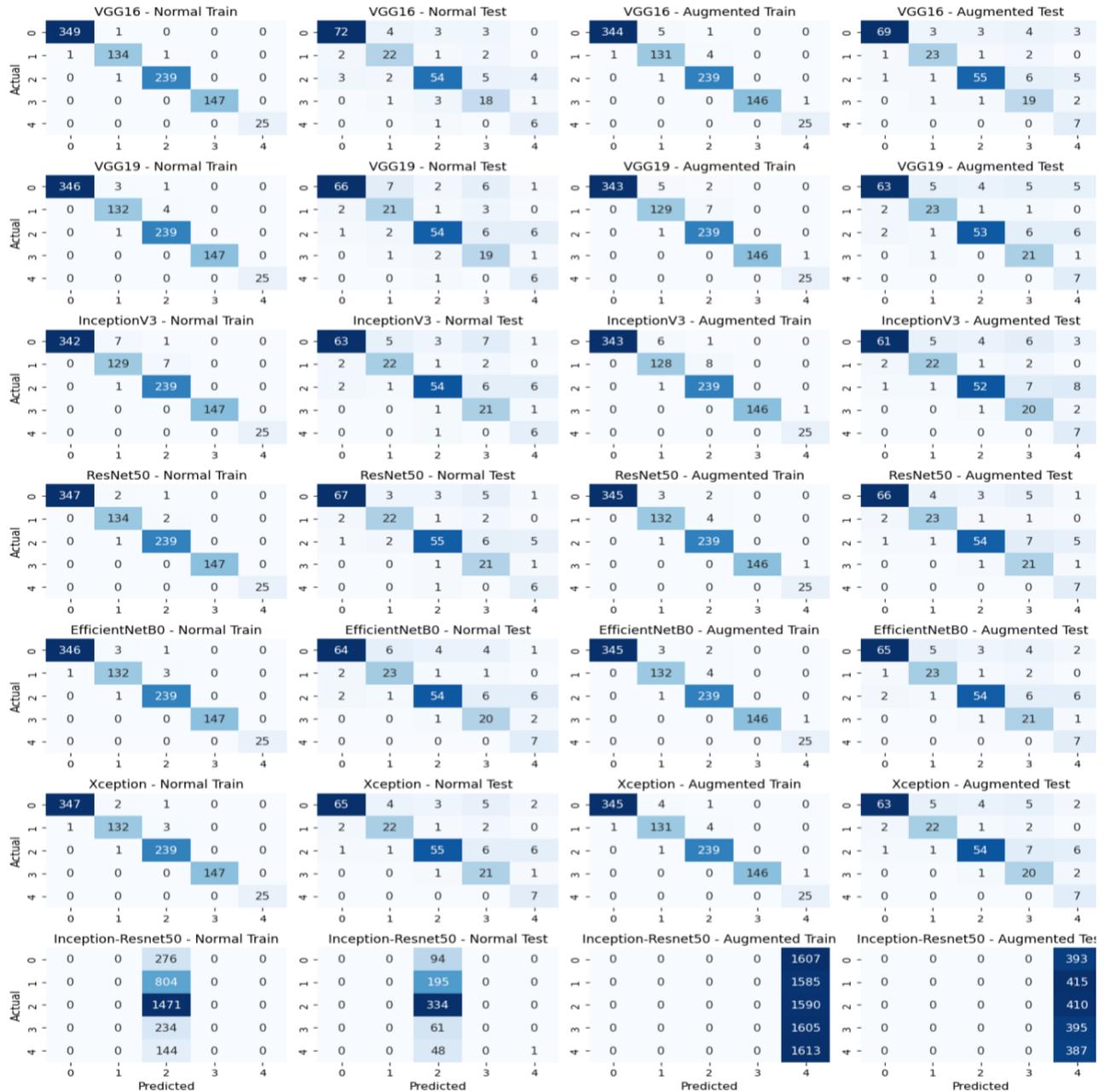


Fig. 4. Evaluation of APTOS Dataset

Fig. 5 illustrates confusion matrices for various transfer learning models trained and tested on the EyePACS dataset using both normal and augmented data. The use of augmented data consistently enhances classification accuracy across models, with

Inception-ResNetV2 showing nearly flawless performance. This demonstrates the effectiveness of data augmentation in improving model generalization and reducing misclassifications.

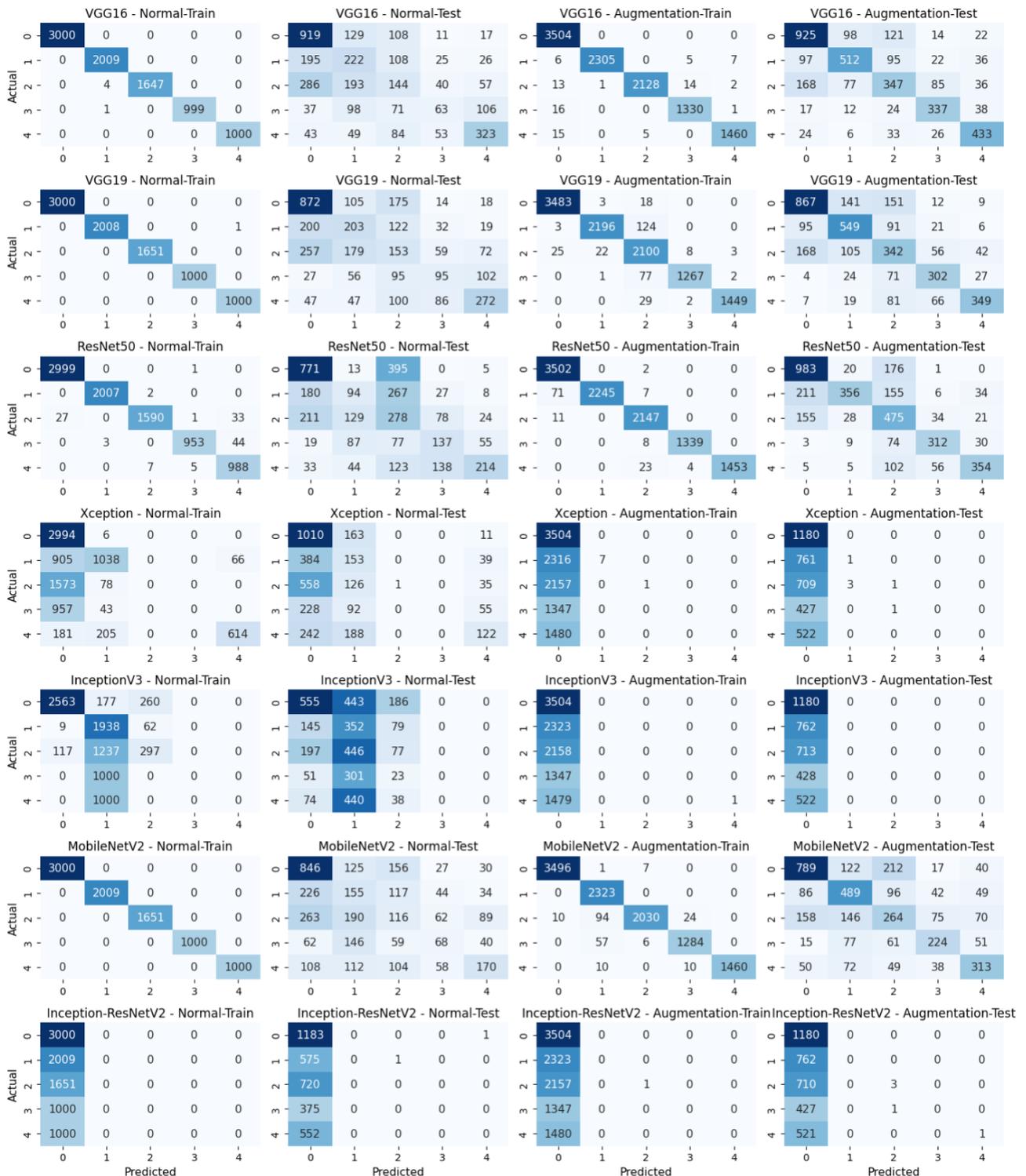


Fig. 5. Evaluation of EyePACS Dataset

In Table 2, the best results are from VGG-16 and ResNet-50 (APTOS & EYEPACS) with 99% accuracy using augmentation. The worst are InceptionV3, Xception, and MobileNetV2 with performance dropping

to 1–10%. In Table 3, VGG-16 (EYEPACS) with augmentation performs best at 71% accuracy. The worst are InceptionV3 and InceptionResNetV2, with scores as low as 6–10%.

**Table 2. Evaluation of APTOS Dataset**

Dataset	Model	AUG	Epoch	LR	ACC	P	R	F1			
APTOS	VGG-16	No	100	0.0001	99%	98%	98%	98%			
		Yes			99%	99%	99%	99%			
	VGG-19	No			93%	96%	80%	86%			
		Yes			97%	97%	97%	97%			
	ResNet-50	No			99%	98%	98%	98%			
		Yes			99%	99%	99%	99%			
	Xception	No			77%	90%	40%	35%			
		Yes			48%	36%	48%	38%			
	InceptionV3	No			80%	48%	47%	45%			
		Yes			58%	49%	58%	50%			
	MobileNetV2	No			96%	98%	89%	91%			
		Yes			1%	1%	1%	1%			
	InceptionResNetV2	No			50%	10%	20%	13%			
		Yes			20%	4%	20%	7%			
	EYEPACS	Vgg-16			No	100	0.0001	99%	99%	99%	99%
					Yes			99%	99%	99%	99%
Vgg-19		No	99%	99%	99%			99%			
		Yes	97%	97%	97%			97%			
ResNet-50		No	99%	98%	98%			98%			
		Yes	99%	99%	99%			99%			
Xception		No	54%	42%	43%			39%			
		Yes	32%	46%	20%			10%			
InceptionV3		No	55%	36%	40%			34%			
		Yes	32%	26%	20%			10%			

**Table 3. Evaluation of Dataset on Training**

Dataset	Model	AUG	Epoch	LR	ACC	P	R	F1			
APTOS	VGG-16	No	100	0.0001	77%	63%	59%	61%			
		Yes			64%	64%	64%	64%			
	VGG-19	No			73%	63%	47%	47%			
		Yes			62%	63%	62%	62%			
	ResNet-50	No			76%	61%	56%	57%			
		Yes			67%	69%	67%	67%			
	Xception	No			68%	27%	37%	31%			
		Yes			41%	33%	41%	33%			
	InceptionV3	No			66%	32%	37%	34%			
		Yes			47%	38%	47%	40%			
	MobileNetV2	No			71%	53%	48%	47%			
		Yes			56%	56%	56%	55%			
	InceptionResNetV2	No			46%	29%	20%	13%			
		Yes			19%	4%	20%	6%			
	EYEPACS	Vgg-16			No	100	0.0001	49%	43%	42%	42%
					Yes			71%	70%	71%	70%
Vgg-19		No	47%	42%	41%			41%			
		Yes	67%	67%	66%			66%			
ResNet-50		No	44%	44%	39%			40%			
		Yes	69%	73%	67%			68%			
Xception		No	38%	42%	27%			22%			
		Yes	33%	22%	20%			10%			
InceptionV3		No	29%	18%	24%			18%			
		Yes	33%	7%	20%			10%			

## V. Discussion

The experimental results reveal that transfer learning, particularly when combined with data augmentation, significantly boosts classification performance across standard CNN architectures. Among all the models tested, VGG16 consistently outperformed others, achieving 99% accuracy on both the APTOS and EyePACS datasets. This superior performance underscores the model's capacity to generalize across different dataset distributions when appropriately fine-tuned and supplemented with augmentation. Interestingly, ResNet50 also exhibited robust performance, closely following VGG16, which supports the notion that deeper architectures with residual connections can effectively manage complex retinal features.

However, this strength is not uniformly observed across all models. Architectures such as InceptionV3, InceptionResNetV2, and MobileNetV2 showed inconsistent or degraded performance when augmentation was applied. This behavior may be attributed to the complex interplay between model architecture and the nature of augmentation strategies. Specifically, lightweight models like MobileNetV2 might overfit augmented noise or fail to extract invariant features under transformation, while more sophisticated architectures like Inception-based

models may require finer hyperparameter tuning or more data-specific pretraining to achieve optimal performance.

Table 4 presents a comparative analysis of recent diabetic retinopathy classification studies based on model types, datasets used, applied techniques, and achieved accuracy. Gulshan et al. [1] achieved 87% accuracy using InceptionV3 on EyePACS and Messidor-2 but required extensive manual grading and lacked data augmentation. Pratt et al. [2] employed a basic 5-layer CNN on the Kaggle DR dataset, reaching 95% accuracy. Lam et al. [3] and Voets et al. [4] used pretrained networks like ResNet, DenseNet, and VGG16, obtaining 84% and 85% accuracy respectively. Khan et al. [5] improved performance to 96% using an ensemble of ResNet-50 and DenseNet201, while Wang et al. [6] achieved 97% with EfficientNet, emphasizing the potential of lightweight architectures. Islam et al. [7] reached 98% accuracy by integrating VGG16 and MobileNetV2 with synthetic image augmentation for better class balance. The proposed system outperformed prior work with a 99% accuracy using multiple models including VGG16, ResNet-50, and InceptionV3, where VGG16 showed the highest robustness across APTOS and EyePACS datasets with effective augmentation strategies.

**Table 4. Comparison with Similar Studies**

Ref.	Dataset	Model(s) Used	Accuracy	Remarks
Gulshan et al. [1]	EyePACS, Messidor-2	InceptionV3	87%	Required extensive manual grading; no augmentation
Pratt et al. [2]	Kaggle DR dataset	CNN (5 conv layers)	0.95	Moderate performance with basic CNN
Lam et al. [3]	EyePACS	ResNet, DenseNet	0.84	Used standard pretrained networks
Voets et al. [4]	EyePACS	VGG16	0.85	Fine-tuned VGG16, limited augmentation
Khan et al. [5]	APTOS, EyePACS	ResNet-50, DenseNet201	0.96	Ensemble strategy boosted generalization
Wang et al. [6]	EyePACS	EfficientNet	0.97	High performance with lightweight model
Islam et al. [7]	EyePACS, DDR	VGG16, MobileNetV2	0.98	Improved balance using synthetic images
Proposed System	APTOS, EyePACS	VGG16, VGG19, ResNet-50, MobileNetV2, Xception, InceptionV3, InceptionResNetV2	99%	VGG16 performed best; robust results across datasets with augmentation

A notable limitation of this research is the dramatic inconsistency in performance of some architectures across datasets. For example, MobileNetV2 and InceptionResNetV2 show severe degradation when augmentation is applied, which suggests potential issues with data distribution, model compatibility with augmentation techniques, or inadequate fine-tuning. Moreover, the experiments do not consider ensemble

approaches or attention mechanisms, which have shown enhanced performance in studies like those by Madhu et al. [1] and Ikram et al. [6]. Additionally, the impact of class imbalance in the datasets, which is known to affect recall and F1-score, is not explicitly addressed. Lack of cross-validation further limits the generalizability of the reported results.

These findings have important implications for

clinical and AI deployment contexts. The consistently high performance of VGG-16 and ResNet-50, particularly with augmentation, suggests their suitability for real-world diabetic retinopathy screening systems, especially in low-resource settings where model reliability is critical. However, the poor and unstable results from other models such as Xception and InceptionV3 underline the necessity for dataset-specific model validation before deployment. The observed variance also highlights the need for robust augmentation strategies and possibly adaptive learning techniques to improve generalizability. Future work should incorporate attention mechanisms, ensemble learning, and domain adaptation techniques to address these gaps and enhance diagnostic accuracy across diverse retinal image datasets.

## VI. Conclusion

This research investigated the impact of dataset variability on diabetic retinopathy classification using transfer learning with various convolutional neural network (CNN) architectures, focusing on the APTOS and EyePACS datasets. The analysis revealed that VGG16 consistently outperformed other models across both datasets. Specifically, on the APTOS dataset, VGG16 achieved 99% accuracy with data augmentation, along with 99% precision, recall, and F1-score. Similarly, on the EyePACS dataset, VGG16 also reached 99% accuracy with data augmentation, again maintaining 99% across all metrics. VGG19 also performed competitively, achieving 97% accuracy, precision, recall, and F1-score on EyePACS with augmentation.

The use of data augmentation techniques such as flipping, rotation, and zooming significantly improved performance across models. For example, ResNet-50 improved from 98% F1-score without augmentation to 99% with augmentation on APTOS, while on EyePACS, it maintained consistent 99% accuracy and F1-score with augmentation. In contrast, models like InceptionV3 and InceptionResNetV2 struggled, with EyePACS accuracies as low as 33% and 19%, respectively, indicating a lack of robustness under data variability.

These findings underline the importance of architecture choice and data augmentation in medical image classification. In future work, we aim to construct a hybrid model that integrates the strengths of top-performing architectures using ensemble learning and feature fusion strategies to enhance diagnostic reliability. We also plan to leverage larger and more diverse datasets, representing different demographics and imaging conditions, to improve generalizability.

To reduce dependency on large annotated

datasets, semi-supervised and self-supervised learning approaches will be explored. Additionally, we will incorporate explainable AI (XAI) techniques to provide interpretable predictions, crucial for clinical adoption. These steps aim to advance the development of an accurate, generalizable, and trustworthy automated system for diabetic retinopathy detection in real-world medical settings.

## References

- [1] S. Madhu, D. K. N. Bhargavi, M. V. S. Ramprasad, S. Gautam, and S. Bhavana, "Accurate diabetic retinopathy segmentation and classification model using gated recurrent unit with residual attention network," *Biomedical Signal Processing and Control*, vol. 102, p. 107348, 2025, doi: 10.1016/j.bspc.2024.107348.
- [2] S. Sarmah, M. Hazarika, P. Das, A. Satheesh, D. Barman, and A. Kumar, "Advancements in Diabetic Retinopathy Detection Using Deep Learning," in *Studies in Computational Intelligence*, vol. 1182, Springer, 2025, pp. 91–117. doi: 10.1007/978-3-031-80813-5\_7.
- [3] K. Chaturvedi, V. Bhandari, A. Kothari, A. Tekerek, R. Tiwari, and R. Shrivastava, "An Intelligent Approach to Analyze Severity Levels of Diabetic Retinopathy by Data Classification Using Transfer Learning," in *International Conference on Data Science and Big Data Analysis*, 2025, pp. 1–16.
- [4] S. Kollem, B. Preethi, G. K. Kumar, R. S. Prasad, K. Praneeth, and S. Peddakrishna, "Hybrid Deep Learning Model with ResNet50 and SVM for Diabetic Retinopathy Classification," in *3rd International Conference on Electronics and Renewable Systems, ICEARS 2025 - Proceedings*, 2025, pp. 1605–1610. doi: 10.1109/ICEARS64219.2025.10940726.
- [5] X. Wei et al., "MSTNet: Multi-scale spatial-aware transformer with multi-instance learning for diabetic retinopathy classification," *Medical Image Analysis*, vol. 102, p. 103511, 2025, doi: 10.1016/j.media.2025.103511.
- [6] A. Ikram and A. Imran, "ResViT FusionNet Model: An explainable AI-driven approach for automated grading of diabetic retinopathy in retinal images," *Computers in Biology and Medicine*, vol. 186, p. 109656, 2025, doi: 10.1016/j.compbiomed.2025.109656.
- [7] S. Guefrachi, A. Ectiouei, and H. Hamam, "Diabetic Retinopathy Detection Using Deep Learning Multistage Training Method," *Arabian Journal for Science and Engineering*, vol. 50, no. 2, pp. 1079–1096, 2024, doi: 10.1007/s13369-024-09137-9.

- [8] M. Z. Atwany, A. H. Sahyoun, and M. Yaqub, "Deep Learning Techniques for Diabetic Retinopathy Classification: A Survey," *IEEE Access*, vol. 10, pp. 28642–28655, 2022, doi: 10.1109/ACCESS.2022.3157632.
- [9] V. Thanikachalam, K. Kabilan, and S. K. Erramchetty, "Optimized deep CNN for detection and classification of diabetic retinopathy and diabetic macular edema," *BMC Medical Imaging*, vol. 24, no. 1, pp. 1–17, 2024, doi: 10.1186/s12880-024-01406-1.
- [10] Y. Yang, Z. Cai, S. Qiu, and P. Xu, "A Novel Transformer Model with Multiple Instance Learning for Diabetic Retinopathy Classification," *IEEE Access*, vol. 12, no. December 2023, pp. 6768–6776, 2024, doi: 10.1109/ACCESS.2024.3351473.
- [11] A. Ikram et al., "A Systematic Review on Fundus Image-Based Diabetic Retinopathy Detection and Grading: Current Status and Future Directions," *IEEE Access*, vol. 12, no. June, pp. 96273–96303, 2024, doi: 10.1109/ACCESS.2024.3427394.
- [12] C. Liu, W. Wang, J. Lian, and W. Jiao, "Lesion classification and diabetic retinopathy grading by integrating softmax and pooling operators into vision transformer," *Frontiers in Public Health*, vol. 12, no. 2, 2024, doi: 10.3389/fpubh.2024.1442114.
- [13] G. Ali, A. Dastgir, M. W. Iqbal, M. Anwar, and M. Faheem, "A Hybrid Convolutional Neural Network Model for Automatic Diabetic Retinopathy Classification from Fundus Images," *IEEE Journal of Translational Engineering in Health and Medicine*, vol. 11, no. January, pp. 341–350, 2023, doi: 10.1109/JTEHM.2023.3282104.
- [14] M. A. K. Raiaan et al., "A Lightweight Robust Deep Learning Model Gained High Accuracy in Classifying a Wide Range of Diabetic Retinopathy Images," *IEEE Access*, vol. 11, no. April, pp. 42361–42388, 2023, doi: 10.1109/ACCESS.2023.3272228.
- [15] R. B. Jayanthi Rajee, S. Mohamed Mansoor Roomi, V. Pooannamalai, and M. Parisa Begam, "A Transfer Learning Approach for Retinal Disease Classification," *Proceedings of 2023 International Conference on Signal Processing, Computation, Electronics, Power and Telecommunication, IConSCEPT 2023*, pp. 1–6, 2023, doi: 10.1109/IConSCEPT57958.2023.10170532.
- [16] S. Nandhini, S. Sowbarnikkaa, J. Mageshwari, and C. Saraswathy, "An Automated Detection and Multi-stage classification of Diabetic Retinopathy using Convolutional Neural Networks," *VITECoN 2023 - 2nd IEEE International Conference on Vision Towards Emerging Trends in Communication and Networking Technologies, Proceedings*, pp. 1–5, 2023, doi: 10.1109/VITECoN58111.2023.10157960.
- [17] D. Raghu Raman, S. Nishanthi, and P. Babysha, "Diagnosis of Diabetic Retinopathy by using EfficientNet-B7 CNN Architecture in Deep Learning," *International Conference on Sustainable Computing and Smart Systems, ICSCSS 2023 - Proceedings*, no. Iccss, pp. 430–435, 2023, doi: 10.1109/ICSCSS57650.2023.10169453.
- [18] T. Dharani, M. P. Prasamsa, B. L. Sirisha, J. B. Vivek, and B. H. Vardhan, "Diabetic Retinopathy classification through fundus images using Deep Learning," *2023 11th International Symposium on Electronic Systems Devices and Computing, ESDC 2023*, vol. 1, pp. 1–6, 2023, doi: 10.1109/ESDC56251.2023.10149877.
- [19] D. R. Dungrani, H. Rajesh Lotia, D. V. Parikh, A. S. Revathi, and K. Kavitha, "Detection and Classification of Diabetic Retinopathy using Deep Learning," *5th Biennial International Conference on Nascent Technologies in Engineering, ICNTE 2023*, no. IcnTE, pp. 1–5, 2023, doi: 10.1109/ICNTE56631.2023.10146626.
- [20] S. Shah, S. Punjabi, S. Chavan, A. S. Revathi, and K. Kavitha, "A comparative study on Diabetic Retinopathy Detection and Classification," *5th Biennial International Conference on Nascent Technologies in Engineering, ICNTE 2023*, no. IcnTE, pp. 1–6, 2023, doi: 10.1109/ICNTE56631.2023.10146636.
- [21] P. K. Das and S. Pumrin, "CNN Transfer Learning for Two Stage Classification of Diabetic Retinopathy using Fundus Images," *8th International Conference on Digital Arts, Media and Technology and 6th ECTI Northern Section Conference on Electrical, Electronics, Computer and Telecommunications Engineering, ECTI DAMT and NCON 2023*, pp. 443–447, 2023, doi: 10.1109/ECTIDAMTNCN57770.2023.10139437.
- [22] S. Ali and S. Raut, "Detection of Diabetic Retinopathy from fundus images using Resnet50," *2023 2nd International Conference on Paradigm Shifts in Communications Embedded Systems, Machine Learning and Signal Processing, PCEMS 2023*, pp. 1–5, 2023, doi: 10.1109/PCEMS58491.2023.10136073.
- [23] G. Kalyani, B. Janakiramaiah, A. Karuna, and L. V. N. Prasad, "Diabetic retinopathy detection and classification using capsule networks," *Complex and Intelligent Systems*, vol. 9, no. 3, pp. 2651–2664, 2023, doi: 10.1007/s40747-021-00318-9.

- [24] M. V. Krishna and B. S. Rao, "Detection and Diagnosis of Diabetic Retinopathy Using Transfer Learning Approach," *International Journal of Intelligent Engineering and Systems*, vol. 16, no. 3, pp. 62–74, 2023, doi: 10.22266/ijies2023.0630.05.
- [25] S. C. Pravin, S. P. K. Sabapathy, S. Selvakumar, S. Jayaraman, and S. V. Subramani, "An Efficient DenseNet for Diabetic Retinopathy Screening," *International Journal of Engineering and Technology Innovation*, vol. 13, no. 2, pp. 125–136, 2023, doi: 10.46604/IJETI.2023.10045.
- [26] D. Das, S. K. Biswas, and S. Bandyopadhyay, "Detection of Diabetic Retinopathy using Convolutional Neural Networks for Feature Extraction and Classification (DRFEC)," *Multimedia Tools and Applications*, vol. 82, no. 19, pp. 29943–30001, 2023, doi: 10.1007/s11042-022-14165-4.
- [27] R. Shalini and S. Sasikala, "Artificial Intelligence Approach for Diabetic Retinopathy Severity Detection," *Informatica (Slovenia)*, vol. 46, no. 8, pp. 195–204, 2022, doi: 10.31449/inf.v46i8.4425.
- [28] M. Ragab, A. S. A. M. Al-Ghamdi, B. Fakieh, H. Choudhry, R. F. Mansour, and D. Koundal, "Prediction of Diabetes through Retinal Images Using Deep Neural Network," *Computational Intelligence and Neuroscience*, vol. 2022, 2022, doi: 10.1155/2022/7887908.
- [29] P. Pragathi and A. Nagaraja Rao, "An effective integrated machine learning approach for detecting diabetic retinopathy," *Open Computer Science*, vol. 12, no. 1, pp. 83–91, 2022, doi: 10.1515/comp-2020-0222.
- [30] K. Nirmala, K. Saruladha, and K. Dekeba, "Investigations of CNN for Medical Image Analysis for Illness Prediction," *Computational Intelligence and Neuroscience*, vol. 2022, 2022, doi: 10.1155/2022/7968200.
- [31] B. Sumathy et al., "Prediction of diabetic Retinopathy using Health Records with Machine Learning Classifiers and data Science," *International Journal of Reliable and Quality E-Healthcare*, vol. 11, no. 2, pp. 1–16, 2022, doi: 10.4018/IJRQEH.299959.
- [32] S. Gupta, S. Thakur, and A. Gupta, "Optimized Feature Selection Approach for Smartphone Based Diabetic Retinopathy Detection," *Proceedings of 2nd International Conference on Innovative Practices in Technology and Management, ICIPTM 2022*, pp. 350–355, 2022, doi: 10.1109/ICIPTM54933.2022.9754021.
- [33] R. Yasashvini, V. Raja Sarobin M, R. Panjanathan, S. Graceline Jasmine, and L. Jani Anbarasi, "Diabetic Retinopathy Classification Using CNN and Hybrid Deep Convolutional Neural Networks," *Symmetry*, vol. 14, no. 9, 2022, doi: 10.3390/sym14091932.
- [34] D. Nagpal, S. N. Panda, M. Malarvel, P. A. Pattanaik, and M. Zubair Khan, "A review of diabetic retinopathy: Datasets, approaches, evaluation metrics and future trends," *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 9, pp. 7138–7152, 2022, doi: 10.1016/j.jksuci.2021.06.006.
- [35] G. Kumar, S. Chatterjee, and C. Chattopadhyay, "DRISTI: a hybrid deep neural network for diabetic retinopathy diagnosis," *Signal, Image and Video Processing*, vol. 15, no. 8, pp. 1679–1686, 2021, doi: 10.1007/s11760-021-01904-7.
- [36] L. Akshita, H. Singhal, I. Dwivedi, and P. Ghuli, "Diabetic retinopathy classification using deep convolutional neural network," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 24, no. 1, pp. 208–216, 2021, doi: 10.11591/ijeecs.v24.i1.pp208-216.
- [37] A. Sugeno, Y. Ishikawa, T. Ohshima, and R. Muramatsu, "Simple methods for the lesion detection and severity grading of diabetic retinopathy by image processing and transfer learning," *Computers in Biology and Medicine*, vol. 137, no. August, p. 104795, 2021, doi: 10.1016/j.combiomed.2021.104795.
- [38] M. Z. Atwany, A. H. Sahyoun, and M. Yaqub, "Deep Learning Techniques for Diabetic Retinopathy Classification: A Survey," *IEEE Access*, vol. 10, pp. 28642–28655, 2022, doi: 10.1109/ACCESS.2022.3157632.
- [39] Kaggle, "APTOS 2019 Blindness Detection," [Online]. Available: <https://www.kaggle.com/competitions/APTOS2019-blindness-detection/data>. [Accessed: Jan. 26, 2025].
- [40] Kaggle, "EyePACS Diabetic Retinopathy Detection," [Online]. Available: <https://www.kaggle.com/competitions/diabetic-retinopathy-detection>. [Accessed: Jan. 26, 2025].

### Author Biography



**Kinjal Patni** holds Bachelor's and Master's degrees in Computer Engineering and is currently pursuing a Ph.D. in the same field. With over eight years of academic and professional experience, she has shown a strong commitment to excellence in teaching, research, and innovation. Her interests include machine learning, artificial intelligence, and data-driven

technologies. She has contributed to several research projects, published scholarly papers, and mentored students at various academic levels. Kinjal focuses on bridging the gap between theory and practice by developing real-world solutions. Through her doctoral research and ongoing academic efforts, she continues to contribute meaningfully to computer science and engineering.



**Shruti Yagnik** holds a Bachelor's degree in IT Engineering, a Master's degree (M.E) in Computer Engineering, and a Ph.D. in Computer Science. With over 13 years of academic and professional experience, she has actively contributed to the advancement

of computer science and engineering. Her expertise spans teaching, research, and technological innovation. She has published several research papers, guided numerous students, and participated in various academic initiatives. Dr. Yagnik's work focuses on emerging technologies, data analytics, and intelligent systems. Through her dedication to education and research, she continues to play a vital role in shaping the future of the discipline.



**Pratik K. Patel** holds a Bachelor's degree in Computer Engineering, a Master's degree (M.Tech) in Computer Science and Engineering, and a Ph.D. in Computer Engineering. With more than 13 years of academic and professional

experience, he has made significant contributions to the field of computer science and engineering. His areas of expertise include advanced computing, data science, and software engineering. Dr. Patel has authored multiple research publications, mentored students at various academic levels, and actively participated in research and development projects. His commitment to teaching, innovation, and scholarly excellence continues to shape the future of computing and technology education.