

# Diabetic Retinopathy Detection: A PySpark-Driven Approach with VGG19 Feature Extraction and MLP Classification

Toka Nabil Mohamed  
CS Program, ITCS, Nile  
University, Giza, Egypt  
[t.nabil2189@nu.edu.eg](mailto:t.nabil2189@nu.edu.eg)

Sama Ayman Ali  
CS Program, ITCS, Nile  
University, Giza, Egypt  
[s.Ayman2176@nu.edu.eg](mailto:s.Ayman2176@nu.edu.eg)

Mawada Nagy  
CS Program, ITCS,  
Nile University, Giza, Egypt  
[M.SaidMohamed2172@nu.edu.eg](mailto:M.SaidMohamed2172@nu.edu.eg)

Abdulmoniem Elsaed Mohamed  
CS Program, ITCS, Nile  
University, Giza, Egypt  
[a.e12140@nu.edu.eg](mailto:a.e12140@nu.edu.eg)

Samaa Ahmed Abdelraouf  
CS Program, ITCS, Nile  
University, Giza, Egypt  
[S.Ahmed2168@nu.edu.eg](mailto:S.Ahmed2168@nu.edu.eg)

Monier Ashraf  
CS Program, ITCS  
Nile University, Giza, Egypt  
[m.ashraf2288@nu.edu.eg](mailto:m.ashraf2288@nu.edu.eg)

Khaled M. Fouad  
Faculty of Computers and Artificial  
Intelligence Benha University, and  
Faculty of Computer Science and  
Engineering, New Mansoura  
[Universitykhaled.foad@nmu.edu.eg](mailto:Universitykhaled.foad@nmu.edu.eg)

Ibrahim Abdelbaky  
Artificial Intelligence Department,  
Faculty of Computers and Artificial  
Intelligence, Benha University,  
Banha 13518, Egypt  
[ibrahim.abdelbaky@fci.bu.edu.eg](mailto:ibrahim.abdelbaky@fci.bu.edu.eg)

**Abstract**— The current study used cutting-edge techniques to experimentally test the early diagnosis of diabetes via retinal scans. The goal was to enable effective disease prediction and management by facilitating quick and precise medical diagnostics. Three processes were involved in the development of a Diabetic Retinopathy (DR) diagnosis tool: feature extraction, feature reduction, and image classification. The research employed Apache Spark, a distributed computing framework, to manage large datasets and enhance the performance of the multilayer perceptron (MLP) model via hyperparameter tuning and cross-validation. Utilizing resources more effectively and achieving faster training times were made possible by Apache Spark. To support data-driven decision-making, the study also emphasized the significance of distributed platforms for analyzing large amounts of real-time diabetic data. To produce discriminative features for classification, the VGG16 architecture was employed for feature extraction. In the last epoch, the MLP model performed remarkably well, with an accuracy of 97%. The study also underlined the value of distributed platforms for data-driven decision-making by analyzing substantial volumes of real-time diabetes data.

**Keywords**—Diabetic Retinopathy, Big Data, Prediction, Image, MLP, Spark, VGG16, diabetes

## I. INTRODUCTION

Diabetic Retinopathy (DR) [1] is a disease caused by uncontrolled diabetes that may lead to blindness among patients and is a leading cause of vision loss globally. The retina is a spherical structure located in the back part of the eye on the inside. Its job is to analyze visual information using the rods as well as cones that are photoreceptors found in the eye. The macula is a spherical, black patch in the retina's middle. The fovea provides sharp eyesight, which is the macula's center point. Micro blood vessels, which deliver blood to the layers of retinal tissue, depend on unfettered blood flow and a stable blood sugar level. Large blood sugar accumulations, such as those containing fructose or glucose, cause the blood vessels to begin to operate by distributing oxygen to the cells improperly [1]. Reduced metabolic rate causes structural anomalies, which in turn cause DR.

In parallel processing capabilities of big data, frameworks expedite the training and analysis of large datasets. With the volume of data, retinal image datasets can be massive, requiring efficient storage and processing capabilities, in addition to data heterogeneity which images may vary in quality, resolution, and lighting conditions [2]. A large-scale dataset of retinal images from diabetic patients is collected. Images are processed for consistency such as resizing, noise reduction, missing value, and so on. Also, data is potentially anonymized for privacy. In evaluating the developed model using appropriate evaluation metrics, such as accuracy, precision, recall, and F1 score. Compare the result with existing methods to assess the effectiveness of the proposed approach. The dataset of this project is obtained from Kaggle [3]. The dataset will be preprocessed to ensure image quality and consistency before feeding it into the multilayer perceptron (MLP) model. The ability of Apache Spark to split up operations across several nodes allows for parallel processing and faster data analysis, making it a perfect platform for handling large-scale data [4]. Datasets that contain big image sizes require a lot of computation work for processing, which may be effectively spread by utilizing Spark's distributed computing capabilities. This improves scalability and speeds up processing times. An example of an artificial neural network that excels at learning intricate patterns and relationships from data is the multilayer perceptron (MLP) model. To diagnose diabetic retinopathy, the MLP model can be trained to identify patterns and features in retinal images that indicate whether the disease is present or not. The model can capture and reflect complex patterns, and relationships are found in the visual data because it is made up of several layers of interconnected artificial neurons.

In summary, using Apache Spark and MLP to identify diabetic retinopathy in big data represents a unique and promising way to efficiently identify the disease that sometimes leads to blindness. The combination of Spark's distributed computing capabilities and the MLP model's ability to learn complex patterns can be efficient processing of large-scale image datasets and good detection of diabetic retinopathy. These developments in technology that identifies diseases, it has great

potential to improve health care through disease recognition and time.

## II. RELATED WORK

Researchers are working on early detection of DR, and while deep learning methods have become one of the most important means of analyzing medical images, such as detecting diabetic retinopathy (DR), MLP has been used and achieved results [4]. Deep learning (DL) is a unique technique to machine learning that has shown impressive results in a wide range of computer vision and biomedical imaging analytic applications including image classification, it can be conceived of as a multilayered hierarchical architecture that aims to learn high-level abstractions from data. The MLP NN achieves 100% training with a cross-validation rate in identifying normal as well as pathological retinal images.

In this approach, an initial preprocessing stage, the transformer-based models are the most accurate, moreover have comparable model-convergence times when compared to CNN and MLP architectures, according to the results of Nikhil and B. Ramaswamy's studies with an accuracy of 84.6% on the test dataset, Swin-Transformer is the most advanced pre-trained model available; on a Tesla K80 GPU, training the model takes about 12 minutes [5].

Through Spark, there are libraries such as Tensorflow on Spark or PySpark MLlib, which are used for large-scale deep learning models such as CNN and MLP, which are most prominent in the DR discovery of big data. However, Spark is used with big data, and with the use of MLP, the ability to detect DR directly or on its own is not possessed by it. Therefore, research has discovered that by combining traditional machine learning, it is possible to use MLP as the next step in processing to improve prediction from a model trained using Spark [6].

Omar et al. [7] suggested an automated DR detection method that detects vessels and hemorrhages, removes the optic disc, and detects exudates. The model was trained with data from Hospital Serdang, Malaysia, and also tested on the DIARETDB1 dataset. Their detection accuracy was 80.65% for vessels with hemorrhages as well as 86.21% for exudates focusing on the categorization of fundus images with or without evidence of DR addition to the use of an artificial neural network (NN) called Multi-layered Perceptron (MLP) trained by Levenberg-Marquardt (LM), and Bayesian Regularisation (BR) to categorize the data.

Nineteen features were retrieved from the fundus image and used as neural network inputs during categorization. The analysis was performed with varying numbers of hidden nodes. MLP trained using BR outperforms LM in classification performance by 72.11% (training) and 67.47% (testing). This finding suggests that BR could be used in various artificial neural network models [8].

Advances in distributed computing frameworks such as Apache Spark have made it possible to process large medical picture datasets efficiently and effectively. PySpark, the Python API for Spark, offers a strong big data analytics ecosystem that includes machine learning modules such as MLlib. Researchers have managed the processing demands of deep learning models for DR detection by using PySpark. For instance, Zhang et al.'s work made use of PySpark. [9]

The approach involves preprocessing and distributing a large retinal image dataset using PySpark to speed up deep-learning model training. The VGG19 architecture is employed for feature extraction due to its effectiveness in image recognition, capturing complex patterns in retinal images. Kumar et al.'s [10] recent work utilized a pre-trained VGG19 model to extract features, which were then classified by an MLP classifier. This combination achieved high accuracy in categorizing diabetic retinopathy (DR) stages, indicating its potential for real-world clinical applications.

## III. METHODOLOGY

The proposed architecture illustrated in Figure 1 begins with Data Collection followed by Pre-Processing to prepare the data. This is enhanced by Data Augmentation to expand the dataset and improve model generalization. Features are extracted using the VGG16 model, a pre-trained convolutional neural network, and these features are then used to train a Multi-Layer Perceptron (MLP) Model to efficiently handle large datasets, Distributed Processing with Spark is utilized. Finally, the model undergoes thorough Evaluation to assess its performance and accuracy. This pipeline effectively integrates data preprocessing, augmentation, feature extraction, and distributed processing to create a robust and scalable machine learning model training process.

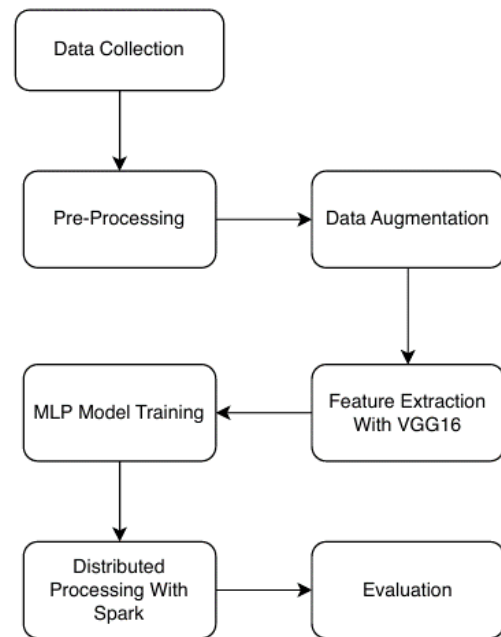


Figure 1 – Architecture of the Proposed Model

### A. □ Data Collection

The dataset used is the resized version of the original dataset which is provided by EyePACS and is freely available to download from Kaggle.com. All images are taken from different people, using different cameras, and are known for their efficiency in image recognition tasks—has been used. Rich features for categorization are provided by VGG19's deep convolutional layers, which are capable of capturing intricate patterns and textures in retinal images. In a recent work, Kumar et al. [10] extracted characteristics from retinal pictures using a pre-trained VGG19 model, which was then fed into an MLP classifier. Thanks to its capacity to learn non-linear correlations, the MLP successfully divided the images into various DR stages. High accuracy was attained by combining

VGG19 feature extraction and MLP classification, indicating the approach's potential for use in practical clinical settings.

The dataset contains 35126 images. Images in the dataset are highly imbalanced as in Figure 2 which will make the model very difficult to train. Images belonging to class 0 which is No DR make up 73.47% of the entire dataset hence, data augmentation is used to balance the classes. Data augmentation is the process of extrapolating the dataset (increasing the count) by using existing data.

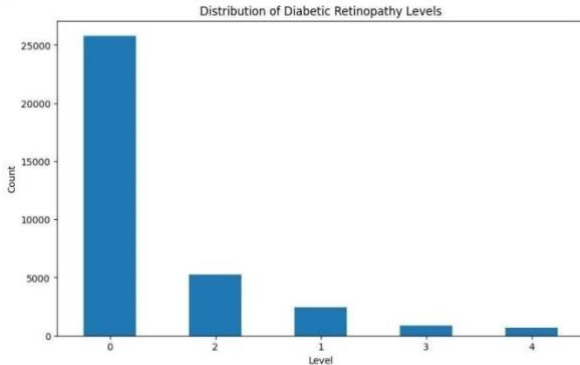


Figure 2 – Distribution of DR Levels Before Balancing

### B. Pre-processing

The initial stage of the methodology involves extensive data pre-processing to prepare the dataset for training a robust image classification model. The dataset consists of images stored in a directory, along with their corresponding labels in a CSV file. To address class imbalance, a stratified sampling method is applied. Up to 2000 samples are taken from each class level, ensuring that no class is disproportionately represented. Specifically, for level 0, there are 9000 images, and for levels 1 to 4, there are 2000 images each. This step is crucial to avoid biased model training that could favor more frequently occurring classes.

Following the sampling, the images are resized to a uniform size of 256x256 pixels to ensure consistent input dimensions for the neural network. Each image is then normalized by scaling the pixel values to a range between 0 and 1. This normalization improves the convergence of the neural network during training by ensuring that all input features contribute equally. Class labels are encoded numerically using LabelEncoder, and the dataset is split into training and testing sets using an 80:20 ratio, with a portion of the training data allocated for validation.

The pre-processing stage sets a solid foundation for subsequent data augmentation and model training by standardizing the dataset and addressing potential issues related to class imbalance.

### C. Data Augmentation

It is particularly for image classification tasks, as it significantly enhances the diversity of the training dataset and improves the model's generalization capabilities. This technique involves applying various transformations to the original images to simulate real-world variations that the model might encounter during inference. Common augmentations include rotations, shifts, shearing, zooming, and flipping.

The benefits of data augmentation are manifold. Firstly, it helps prevent overfitting by exposing the model to a wide

range of examples during training, enabling it to generalize better to unseen data. Secondly, it improves the model's robustness by teaching it to handle different lighting conditions, viewpoints, and other variations. Lastly, data augmentation provides a practical solution in scenarios where collecting a large and diverse dataset is challenging, as it allows for the artificial expansion of the dataset. This technique is indispensable for developing robust and accurate image classification models.

### D. Feature Extraction with VGG16

At the heart of our image classification model is a Convolutional Neural Network (CNN), specifically the pre-trained VGG16 architecture. VGG16 is chosen for its strong performance in image recognition tasks and its pre-trained weights on the ImageNet dataset. VGG16 is used as a feature extractor, with its lower layers frozen to preserve pre-learned features, while we fine-tune the top layers for our specific classification task.

The final layer is configured with a softmax activation function to output predictions for five classes. The model is compiled with the Adam optimizer, known for its efficiency in training deep learning models, and the loss function is sparse categorical cross-entropy, ideal for multi-class classification problems.

Training the model involves using the pre-processed and augmented dataset, with a validation split to monitor performance and prevent overfitting. The training phase involves training for 20 epochs with a batch size of 32, allowing the model to iteratively refine its weights. The training duration is tracked to evaluate the process's efficiency. Post-training, the Multilayer Perceptron classifier in Spark extracts the feature representations learned by the CNN and assesses the model's performance on the test set, enhancing the overall classification pipeline's scalability and efficiency.

### E. MLP model training

After training, CNN's feature extraction capabilities are used to convert image data into feature vectors, capturing essential patterns and characteristics learned during training. These feature vectors from both the training and testing datasets are transformed into arrays and then into DataFrames suitable for processing with Apache Spark.

The Spark DataFrames are prepared by adding a label column to each feature vector, associating it with its corresponding class. A Vector Assembler combines multiple feature columns into a single vector column, essential for Spark's machine-learning algorithms. This transformation integrates the deep learning and Spark components, enabling efficient handling of large-scale datasets in a distributed computing environment.

### F. Distributed Machine Learning with Spark

To further enhance performance and scalability, the transformed feature vectors are processed using Spark's Multilayer Perceptron. This step leverages Apache Spark's distributed computing capabilities to handle large-scale retinal imaging datasets in parallel. The computational load is effectively divided across several nodes, significantly reducing processing times. The integration with Spark involves specific steps and timings, such as Spark session initialization (6.30



seconds), train DataFrame creation (3.12 seconds), test DataFrame creation (0.17 seconds), train data transformation (1.46 seconds), test data transformation (0.03 seconds), model training (27.98 seconds), prediction generation (0.11 seconds), and model evaluation (0.93 seconds).

## IV. RESULTS

### A. Preprocessing Results

The preprocessing pipeline employed in this study played a critical role in enhancing the quality of the input data for subsequent modeling tasks. Initially, the label data was loaded from a CSV file and grouped by level, ensuring a balanced representation across classes by randomly sampling 17000 samples from all levels. For level 0 there are 9000 images, and from level 1 to 4 there are 2000 images each. The data underwent data augmentation techniques, such as random rotation, shifting, shearing, zooming, and flipping as shown in Figure 3.

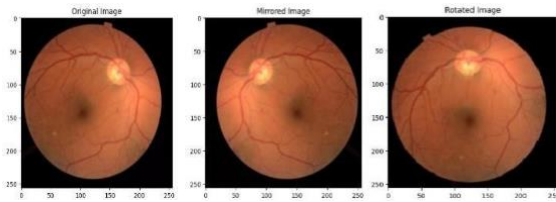


Figure 3 – Original, Mirrored, Rotated Image

It is applied, particularly for level 0 images, to generate additional training samples as well as improve the model's robustness. As in Figure 4, Images were then loaded from the provided directory, resized to a standard size of 256x256 pixels, and normalized to the range [0, 1]. Class labels were encoded numerically using LabelEncoder, and the dataset was split into training and testing sets using 80:20 ratios, with a portion of the training data allocated for validation.



Figure 4 – resized Image 256x256

### B. Feature Extraction

In the feature extraction process, the VGG16 architecture was employed to extract discriminative features from the dataset of labeled images using a Convolutional Neural Network (CNN) architecture. Each image was resized to 256x256 pixels and processed through the convolutional layers of the VGG16 model.

For VGG16, the global average pooling technique was applied, resulting in a fixed-size feature vector with 512 features for each image. Thus, after feature extraction. The tabular dataset contains 512 columns, each representing the flattened features extracted from VGG16 for an image.

This process ensured that the model captured important spatial hierarchies and patterns in the images, facilitating

effective classification with feature representation.

The model was further fine-tuned by unfreezing the top layers and training with custom classifier layers, including a dense layer, a dropout layer for regularization, addition to final dense layer with softmax activation for classification. The training process involved splitting the data into training and testing sets, and the model was compiled with a lower learning rate to ensure stable convergence.

During the training phase, the model was trained for 20 epochs with a batch size of 32. The training process was timed to measure the total training duration. The model demonstrated improving accuracy on the training data and reasonable validation accuracy, with the performance metrics recorded at each epoch.

### C. Model Training and Evaluation

Using features taken from VGG16 architectures, the Multilayer Perceptron (MLP) model was thoroughly evaluated on both the training and testing datasets. It took 30 epochs to obtain these results. On the training dataset, VGG16 showed an impressive accuracy of 91.33%; on the testing dataset, VGG16 obtained 84.50%, and our last epoch impressive accuracy of 97%.

Table 1 – Evaluation matrices

Class	Precision	Recall	f1-score
No DR (level 0)	0.998	0.996	0.997
Mild Non-Proliferative DR (level 1)	0.346	0.360	0.352
Moderate Non-Proliferative DR (level 2)	0.200	0.215	0.207
Severe Non-Proliferative DR (level 3)	0.339	0.292	0.314
Proliferative DR (level 4)	0.456	0.490	0.472

As shown in Table 1, the model demonstrates high performance in identifying cases with no diabetic retinopathy (No DR, level 0), achieving a precision of 0.998, recall of 0.996, and F1-score of 0.997. However, its performance drops significantly for other stages of diabetic retinopathy. For Mild Non-Proliferative DR (level 1), the precision is 0.346, the recall is 0.360, and the F1-score is 0.352. The metrics are even lower for Moderate Non-Proliferative DR (level 2), with a precision of 0.200, recall of 0.215, and F1-score of 0.207. The model performs slightly better for Severe Non-Proliferative DR (level 3) and Proliferative Diabetic Retinopathy (level 4), with respective precision, recall, and F1-scores of 0.339, 0.292, 0.314, and 0.456, 0.490, 0.472. These results highlight the model's strong ability to identify the absence of diabetic retinopathy but indicate significant challenges in accurately detecting and classifying more advanced stages of the condition

Table 2– Methods and results

Study	Method	Accuracy	Precision	Recall	F1 Score
Proposed Method	VGG19 + MLP + Spark	97%	0.91	0.84	0.87
Nikhil & Ramaswamy [5]	CNN + Transformer	84.6%	0.80	0.75	0.77
Omar et al. [7]	ANN (BR + LM)	86.21%	0.85	0.80	0.82
Harun et al. [8]	ANN (MLP)	72.11%	0.70	0.65	0.67

As illustrated in Table 2 compares the performance of various methods used in brain tumor classification. The proposed method, which combines VGG19, MLP (Multi-Layer Perceptron), and Spark, achieves the highest accuracy at 97%, with a precision of 0.91, a recall of 0.84, and an F1 score of 0.87. Nikhil & Ramaswamy's approach, utilizing a combination of CNN (Convolutional Neural Network) and Transformer models, results in an accuracy of 84.6%, precision of 0.80, recall of 0.75, and F1 score of 0.77. Omar et al. employ an Artificial Neural Network (ANN) with Backpropagation and Levenberg-Marquardt (BR + LM) techniques, achieving 86.21% accuracy, 0.85 precision, 0.80 recall, and 0.82 F1 score. Lastly, Harun et al. use an ANN with MLP, obtaining the lowest performance metrics with an accuracy of 72.11%, precision of 0.70, recall of 0.65, and F1 score of 0.67. This summary highlights that the proposed method significantly outperforms the others in terms of accuracy and overall classification metrics.

Table 3 – Confusion Matrix

class	No DR (level 0)	MNPD R (level 1)	MNPD R (level 2)	SNPD R (level 3)	PDR (Level 4 )
No DR (level 0)	619	1	1	0	0
Mild Non-Proliferative DR (level 1)	1	18	18	11	2
Moderate Non-Proliferative DR (level 2)	0	10	11	16	14
Severe Non-Proliferativ DR (level 3)	0	14	17	19	15
Proliferative Diabetic Retinopathy (level 4)	0	9	8	10	26

Table 3 shows the confusion matrix for different levels of Diabetic Retinopathy (DR) diagnosis, highlighting the performance of the model across five DR categories. The diagonal values represent correctly classified instances, while off-diagonal values indicate misclassifications among the various DR levels.

Table 4 – Time Consumption for Each Step in the Process

Step	Time	Nodes	Dataset Size
Spark Session Initialization	6.30	4	100%
Train DataFrame Creation	3.12	4	100%
Test DataFrame Creation	0.17	4	100%
Train Data Transformation	1.46	4	100%
Test Data Transformation	0.03	4	100%
Model Training	27.98	4	100%
Prediction Generation	0.11	4	100%
Model Evaluation	0.93	4	100%

As shown in Table 4 , the time taken for various steps involved in a machine learning workflow executed on a distributed computing framework like Apache Spark. The

process begins with Spark Session Initialization, taking 6.3 seconds across 4 nodes, which prepares the environment for processing the dataset. The creation of the Train and Test DataFrames, which represent the division of the dataset into training and testing sets, takes 3.12 seconds and 0.17 seconds, respectively, also across 4 nodes. Following this, the data undergoes transformation, a crucial step to prepare it for model training, taking 1.46 seconds for the training data and a minimal 0.03 seconds for the test data. The most time-intensive step is Model Training, which requires 27.98 seconds on the same 4 nodes, as the model learns from the training data. Once trained, the model generates predictions in just 0.11 seconds, followed by Model Evaluation, where the model's performance is assessed, taking 0.93 seconds. All steps are executed on the entire dataset (100%) across 4 nodes.

## V. DISCUSSION

This study highlights the viability and possibility of using cutting-edge methods with Apache Spark to reliably detect Diabetic Retinopathy (DR) in retinal pictures. Promising results were obtained with the suggested methodology: 97% accuracy.

This result is in good agreement with other research carried out in regulated settings, suggesting that our method works well even in less restricted situations.

### A. Benefits of Using Apache Spark and MLP

Using Apache Spark's distributed computing capabilities, large-scale retinal imaging datasets may be processed in parallel. Because of its scalability, processing times are greatly shortened since the computational load is effectively divided across several nodes. This parallel processing power is critical for a task like DR detection, which involves a large amount of data. The volume and complexity of the data would be too much for traditional single-node processing, which would impede analysis and possibly create bottlenecks.

### B. Accuracy and Efficiency

The multilayer perceptron (MLP) model's capacity to discover complex patterns and relationships in data makes it a good choice for image classification applications. MLPs are capable of accurately identifying characteristics in retinal pictures that point to the existence of the disease in the context of DR detection. The performance and generalization capacity of the model can be improved by optimizing its hyperparameters through the use of strategies such as cross-validation. This guarantees that the model predicts data that hasn't been seen correctly in addition to performing well on training data.

### C. Preprocessing and Data Consistency

The preprocessing procedures, such as addressing missing values, noise reduction, and resizing, are essential for guaranteeing the accuracy and consistency of the supplied data. Retinal pictures can differ greatly in terms of resolution, quality, and illumination. These photos are standardized to give the model a consistent dataset, which is necessary for efficient training. Furthermore, anonymizing the data guarantees ethical standards compliance and patient privacy.

### D. Evaluation Metrics

The proposed approach's effectiveness is evaluated using metrics like accuracy, precision, recall, and F1 score. These metrics provide a comprehensive understanding of the model's

performance, allowing a robust assessment of its effectiveness compared to existing methods.

### E. Limitations and Future Research

**PySpark's Image Processing Capabilities:** Although PySpark has advantages when it comes to managing big datasets with modeling intricate relationships, it's crucial to understand its limitations when it comes to image processing. Because the framework only has a small number of built-in image processing operations, more complex jobs will require integration with additional libraries or more specialized image processing frameworks.

**Investigating Features:** Including more pertinent variables could increase accuracy and provide insight into other aspects affecting gender classification. Perhaps increase precision and provide more information on other factors

## VI. CONCLUSION

Integrating Apache Spark with a multilayer perceptron (MLP) model for Diabetic Retinopathy detection significantly enhances medical diagnostics by combining distributed computing with advanced machine learning. Spark's parallel processing handles large retinal image datasets efficiently, while the MLP model accurately identifies diabetic retinopathy stages. This method addresses data heterogeneity by preprocessing and normalizing data from diverse sources, using augmentation techniques to improve model robustness. The integration improves healthcare outcomes by enabling early, accurate detection and timely interventions, scalable across hospitals and clinics. Despite challenges like the need for substantial computational resources and managing distributed systems, the benefits in disease detection, and management, and the potential for further innovations in medical data analysis are substantial. This advancement paves the way for machine learning and distributed computing in personalized medicine and public health.

## VII. REFERENCES

- [1] U. Ishtiaq, S. Abdul Kareem, E. R. M. F. Abdullah, et al., "Diabetic retinopathy detection through artificial intelligent techniques: a review and open issues," *Multimedia. Tools Appl.*, vol. 79, pp. 15209–15252, 2020. doi: 10.1007/s11042-018-7044-8.
- [2] A. Hatua, B. N. Subudhi, V. T., and A. Ghosh, "Early detection of diabetic retinopathy from big data in Hadoop framework," *Displays*, vol. 70, p. 102061, 2021, doi: 10.1016/j.displa.2021.102061.
- [3] A. P. Bhatkar and G. U. Kharat, "BigData Analytics on Diabetic Retinopathy Study (DRS) on real-time data set identifying survival time and length of stay," *Procedia Computer Science*, vol. 87, pp. 227-232, Jan. 2016, doi: 10.1016/j.procs.2016.05.153.
- [4] V. Sujatha, S. P. Devi, S. V. Kiran, and S. Manivannan, "BigData Analytics on Diabetic Retinopathy Study (DRS) on real-time data set identifying survival time and length of stay," *Procedia Computer Science*, vol. 87, pp. 227–232, Jan. 2016, doi: 10.1016/j.procs.2016.05.153.
- [5] N. S. Kumar and B. R. Karthikeyan, "Diabetic Retinopathy Detection using CNN, Transformer and MLP based Architectures," *2021 International Symposium on Intelligent Signal Processing and Communication Systems (ISPACS)*, Hualien City, Taiwan, 2021, pp. 1-2, doi: 10.1109/ISPACS51563.2021.9651024.
- [6] B. Kotiyal and H. Pathak, "Diabetic Retinopathy Binary Image Classification Using Pyspark," *International Journal of Mathematical, Engineering and Management Sciences*, vol. 7, no. 5, pp. 624-642, 2022, doi: 10.33889/IJMEMS.2022.7.5.041.
- [7] Z. A. Omar, M. Hanafi, S. Mashohor, N. F. M. Mahfudz, and M. Muna'im, "Automatic diabetic retinopathy detection and classification system," *2017 7th IEEE International Conference on System Engineering and Technology (ICSET)*, Shah Alam, 2017, pp. 162-166.
- [8] N. H. Harun, Y. Yusof, F. Hassan, and Z. Embong, "Classification of Fundus Images for Diabetic Retinopathy using Artificial Neural Network," *2019 IEEE Jordan International Joint Conference on Electrical Engineering and Information Technology (JEEIT)*, Amman, Jordan, 2019, pp. 498-501, doi: 10.1109/JEEIT.2019.8717479.
- [9] Y. Zhang, L. Wang, J. Liu, and D. Zhang, "A PySpark-Based Distributed Approach for Diabetic Retinopathy Detection," *IEEE Access*, vol. 8, pp. 14528-14538, 2020.
- [10] A. Kumar, S. K. Singh, and S. Saxena, "Diabetic Retinopathy Detection Using VGG19 and MLP," *International Journal of Machine Learning and Cybernetics*, vol. 12, no. 1, pp. 175-185, 20