

SWIN-TRANSFORMER BASED RECOGNITION OF DIABETIC RETINOPATHY GRADE

By

SANJAY GANDHI GUNDABATINI *

SAI SINDHU MANNE **

SUNKARA LIKHIT BABU ***

VANGAPANDU BHARGAVA RAO ****

SANKA TEJASWI *****

*-***** Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India.

<https://doi.org/10.26634/jpr.12.1.21927>

Date Received: 24/04/2025

Date Revised: 28/05/2025

Date Accepted: 10/07/2025

ABSTRACT

Diabetic Retinopathy (DR), a common diabetes-related disorder, is a leading driver of blindness worldwide. Quick detection and precise staging are essential for effective management and vision preservation. This study explores the Swin Transformer, an advanced deep learning framework with a multi-layered setup and a unique sliding window method, to create an automated tool for DR stage assessment. Utilizing the APTOS 2019 Blindness Detection dataset, the system accurately identifies small retinal signs like microaneurysms and more pronounced features such as hemorrhages, achieving high precision. Improved preprocessing, including image enrichment and calibration, enhances its versatility. Results indicate that this approach outperforms traditional Convolutional Neural Networks (CNNs) in precision, computational thrift, and growth potential, with a test accuracy of 99.57% and a test loss of 0.0220.

Keywords: Diabetic Retinopathy, Swin Transformer, Automated Staging, Retinal Analysis, Deep Learning Technology.

INTRODUCTION

The global increase in diabetes has spotlighted Diabetic Retinopathy (DR) as a major health threat, capable of causing blindness if not addressed early (World Health Organization, 2025). This condition results from harm to the retina's blood vessels due to prolonged high glucose levels, requiring swift intervention to prevent significant vision deterioration. Figure 1 shows an image of the normal retina and the effects of diabetes on the retina, highlighting the progressive damage that occurs with disease advancement. Current diagnostic methods typically involve expert ophthalmologists examining retinal photographs, a practice that, while accurate, is resource-heavy and time-consuming. In regions lacking

such specialists, diagnosis delays amplify the risk of disease worsening, and discrepancies in human analysis can lead to uneven results (Coan et al., 2023). These issues emphasize the pressing need for automated, reliable diagnostic alternatives. Table 1 shows the severity levels of Diabetic Retinopathy and associated retinal lesions, outlining how the disease progresses from mild symptoms to severe complications, including neovascularization and significant hemorrhages.

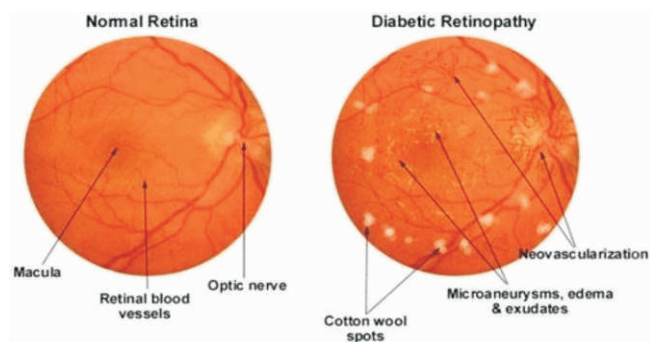


Figure 1. Image of the Normal Retina and the Effects of Diabetes on the Retina



This paper has objectives related to SDGs



DR Severity	Level Lesions
No DR	No observable retinal damage.
Mild	Limited to microaneurysms.
Moderate	Lesions beyond mild but short of severe.
Severe	Over 20 hemorrhages per retinal quadrant.
Proliferative DR	Features neovascularization or significant hemorrhages.

Table 1. Severity Levels of Diabetic Retinopathy and Associated Retinal Lesions

Innovations in deep learning (DL) and machine learning (ML) have transformed healthcare by enabling fast, dependable analysis of complex datasets (Beam & Kohane, 2018; Gulshan et al., 2016; Topol, 2019). These systems leverage vast data to detect subtle patterns, refine diagnostic accuracy, optimize treatment plans, and personalize patient care. Their capacity to process extensive information makes them highly promising for automating DR detection. Among various DL models, the Swin Transformer stands out for its superior accuracy and efficiency in medical image processing, thanks to its layered feature extraction and adaptive window technique (Liu et al., 2021). This research employs the Swin Transformer to tackle the demand for robust, scalable DR screening tools, capturing both minute details (such as microaneurysms) and larger anomalies (such as hemorrhages) more effectively than CNNs, which focus mainly on local patterns.

The Swin Transformer, distinguished by its tiered design and innovative window-shifting method, excels in this domain. Unlike CNNs, which prioritize localized features, it captures both detailed and comprehensive retinal traits, facilitating the identification of early markers like

microaneurysms and exudates. This study taps into its strengths to address the urgent call for effective, scalable DR screening solutions.

1. System Model

Automating DR stage evaluation marks a leap forward in medical diagnostics, cutting reliance on slow manual reviews and promoting consistent, rapid decisions. The Vision Transformer (ViT) pioneered this shift by adapting language-processing transformers for image tasks, offering scalability and strong local feature learning beyond CNN limitations (Dosovitskiy et al., 2020). The Swin Transformer enhances this with its sliding window innovation, structuring images in a tiered manner across linked processing units (Liu et al., 2021). This design excels in applications like image categorization, object spotting, and area segmentation (Hatamizadeh et al., 2021; Xu et al., 2021).

Unlike ViT's resource-heavy $O(N^2)$ complexity, the Swin Transformer uses a leaner $O(M * N)$ model with non-overlapping window calculations, merging patches to form multi-scale image views (Liu et al., 2021). This adaptability suits retinal scans with diverse feature sizes. Its window-sliding process shifts positions between layers, boosting spatial connectivity and context comprehension compared to fixed-window methods. Figure 2 shows an illustration of the Swin transformer's architecture.

Four essential steps make up the Swin Transformer architecture, which is intended to handle and display image data in a hierarchical fashion. The model can

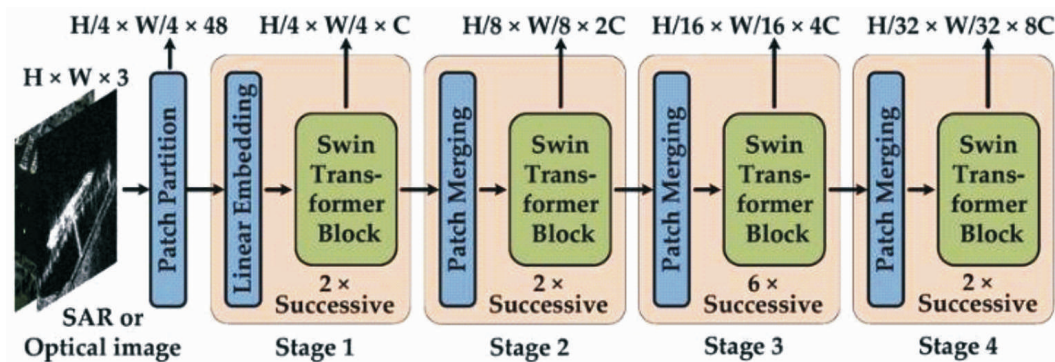


Figure 2. Architecture of Swin Transformer

contain both local details and global context thanks to these stages that build upon one another (Liu et al., 2021).

1.1 Patch Partitioning

The initial step involves segmenting the input image into smaller, manageable units, typically 4x4 pixel patches. Figure 3 shows a patch partition, where the image is divided into uniform segments to enable localized feature extraction. Each patch, comprising 16 pixels with three RGB color channels, is flattened into a 48-channel vector representation, reducing the dimensionality of the data for subsequent processing. Figure 4 shows an image resulting from the patch partition, illustrating how the original input is transformed into a grid of encoded units. This technique lowers the dimensionality of the input data, which makes further processing easier (Brownlee, 2019).

The input image is efficiently divided into smaller areas, usually four by four pixels, after patch partitioning (Topol, 2019). This division simplifies the creation of debug symbols by converting the picture size to $(W/4, H/4, \times \text{channel})$, which changes to $(W/4, H/4, \times 3)$ when patch partitioning is finished. The analysis and manipulation of the image data are made easier by this process.

1.2 Linear Embedding

Following patch division, these segments undergo a

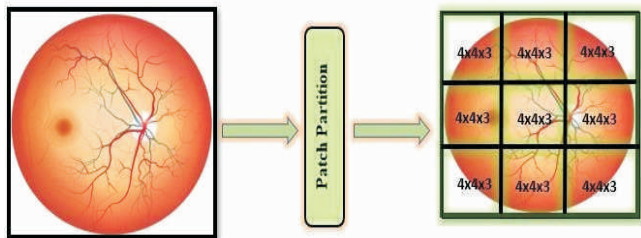


Figure 3. Patch Partition

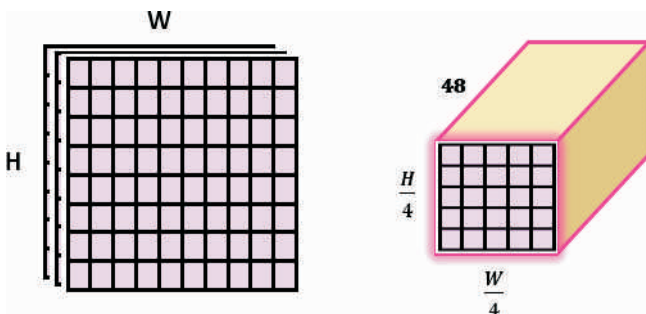


Figure 4. Image from Patch Partition

transformation into a higher-dimensional space through a linear embedding process, akin to applying a convolutional layer that maps 48 channels to 96 channels (Goodfellow et al., 2016). This step prepares the data for deeper analysis by enhancing its representational capacity. By utilizing convolution operations and kernel functions, the patch partitioning and linear embedding procedures together achieve patch embedding, much like convolutional neural networks (O'Shea & Nash, 2015).

1.3 Swin Transformer Block

At the heart of the model lies the Swin Transformer block, which replaces traditional multi-head self-attention with a window-based approach. Figure 5 shows the Swin transformer blocks, highlighting the architecture's hierarchical structure and local-to-global feature extraction mechanism. It employs two sub-modules: window-based multi-head self-attention (W-MSA) for localized processing and shifted-window multi-head self-attention (SW-MSA) for cross-window connectivity (Topol, 2019). These are followed by a multi-layer perceptron (MLP) with GELU nonlinearity, preceded by layer normalization, enabling the model to discern long-range dependencies within the image efficiently (Vaswani et al., 2017).

After linear embedding, the first set of patch tokens,

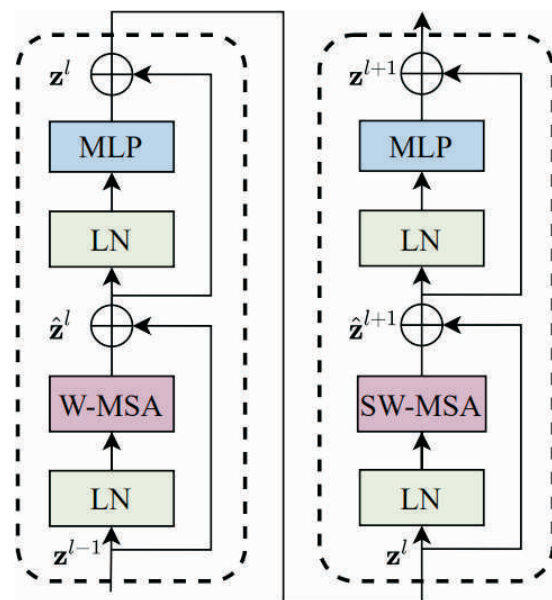


Figure 5. Swin Transformer Blocks

referred known as Stage 1, has dimensions of (W/4, H/4, C). Each following stage's dimensions change when fewer 2x2 patch tokens are produced: Stage 2 measures (W/8, H/8, 2C), Stage 3 measures (W/16, H/16, 4C), and Stage 4 measures (W/32, H/32, 8C).

The mathematical formulations for Window-based Multi-scale Attention (W-MSA) and Shifted Window Multi-scale Attention (SW-MSA) are expressed in Equations (1-4) (Hathot et al., 2021).

$$z \wedge l = W - (LM(zl - 1)) + zl - 1 \quad (1)$$

$$zl = M(LN(zl - 1)) + z \wedge l \quad (2)$$

$$z \wedge l + 1 = SW - M(LN(zl)) + zl \quad (3)$$

$$zl + 1 = M(LN(z \wedge l + 1)) + z \wedge l + 1 \quad (4)$$

In these equations, $(z \wedge l)$ represents the output of each block, and $(zl - 1)$ represents the input from the previous block. (LN) denotes layer normalization, (MLP) denotes the multi-layer perceptron, (W-MSA) denotes window-based multi-head self-attention, and (SW-MSA) denotes shifted window multi-head self-attention.

1.4 Patch Merging

This stage reduces the spatial resolution of the image by merging adjacent 2 x 2 patches into a single unit, concatenating their features, and adjusting the depth through fully connected layers. This down sampling process mirrors techniques used in CNNs, allowing the model to extract higher-level features and expand its receptive field progressively across stages.

Each (2x2) set of nearby pixels is combined into a patch during patch merging, and pixels with the same color are combined to create four feature maps. The depth dimension is then used to concatenate these feature maps. The output from this step passes through layer normalization (LN) and fully connected (FC) layers, which are intended to linearly modify the feature map's depth, changing it from C to C/2, as shown in Figure 6. The model can extract higher-level characteristics and capture bigger receptive fields thanks to this reduction in spatial dimensions (Dumoulin & Visin, 2016).

2. Experimental Settings

To evaluate the proposed system, this study utilizes the

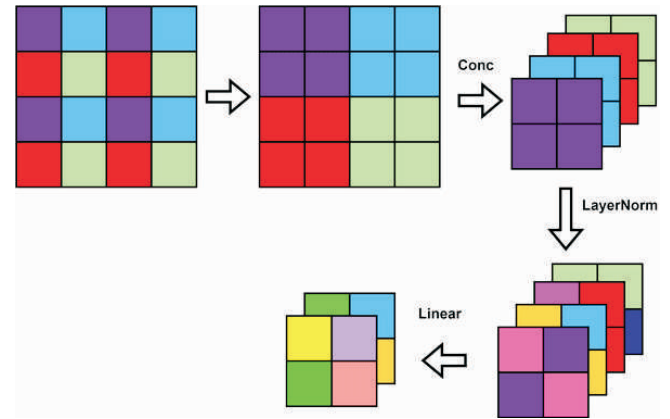


Figure 6. Patch Merging

APTOS 2019 Blindness Detection dataset, a widely recognized benchmark sourced from Kaggle and developed by the Asia Pacific Tele-Ophthalmology Society (Chen et al., 2018). The dataset comprises 3,662 retinal fundus photographs captured under a variety of imaging conditions, reflecting real-world diversity. Each image has been meticulously categorized by experts into one of five DR severity levels: No DR (Class 0), Mild (Class 1), Moderate (Class 2), Severe (Class 3), and Proliferative (Class 4). The distribution of images across these categories highlights the dataset's imbalance, with a significant concentration in the No DR class, posing a challenge that the model must address. The distribution of retinal images across each severity level within the dataset is shown in Table 2.

2.1 Implementation Details

The Swin Transformer was selected as the backbone of this system due to its proven efficiency in handling memory-intensive tasks and its flexibility in adapting to diverse image analysis challenges (Dosovitskiy et al., 2020). Input images were pre-processed to a uniform resolution of either 224x224 pixels or 160x160 pixels, ensuring

Severity level	Number of Images
Class 0 (Normal)	1805
Class 1 (Mild)	370
Class 2 (Moderate)	999
Class 3 (Severe)	193
Class 4 (Proliferative)	295
Total	3662

Table 2. Dataset Summary of the APTOS Dataset

compatibility with the model's architecture. Training and evaluation were conducted on a high-performance GPU equipped with 12.68 GB of storage capacity and 5.71 GB of dedicated memory, capable of supporting the computational demands of deep learning workflows. To achieve optimal performance, an extensive hyperparameter tuning process was undertaken, exploring a wide range of values to fine-tune the model's behavior (Kaggle, 2019). Figure 7 shows the training and validation loss and accuracy per epoch, illustrating the model's convergence trends and performance stability during training. The final configuration, shown in Table 3, reflects the settings that yielded the best results for DR severity grading.

3. Results and Discussion

3.1 Introduction

The application of the Swin Transformer for identifying and categorizing Diabetic Retinopathy represents a significant leap forward in the field of automated medical

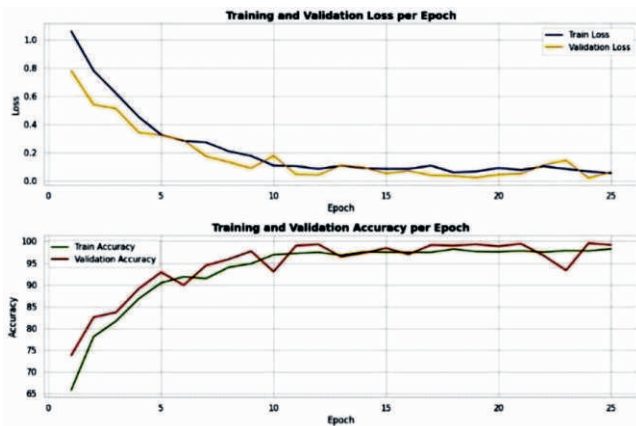


Figure 7. Training and Validation Loss and Accuracy per Epoch

Hyper-parameter	Value
Batch Size	32
Learning Rate	1e - 4
Size of shifting window	16
Size of attention window	16
Epoch	25
Weight decay	1e - 3
Optimizer	AdamW
Patch size	16×16
Embedded dimension	96

Table 3. Hyper-parameters in Swin Transformer Training (Bergstra & Bengio, 2012)

image analysis. A detailed examination of the experimental results assesses the system's performance across multiple dimensions and explores the broader implications of these findings for both research and clinical practice.

3.2 Dataset Overview

The APTOS 2019 Blindness Detection dataset formed the basis for both the training and testing phases of this study (Kaggle, 2019). Comprising 3,662 retinal fundus images, the dataset spans the full range of DR severity levels, from No DR to Proliferative DR, as previously outlined. Its diversity and expert annotations make it an ideal testbed for evaluating the model's diagnostic capabilities.

- No DR
- Mild DR
- Moderate DR
- Severe DR
- Proliferative DR

3.3 Performance Metrics

To comprehensively assess the model's effectiveness, several standard performance metrics were employed (Tharwat, 2021):

- *Accuracy (ACC)*: The percentage of images correctly classified out of the total sample.
- *Precision (P)*: The fraction of positive predictions that are true positives, indicating prediction reliability.
- *Recall (R)*: The proportion of actual positive cases correctly identified, reflecting detection sensitivity.
- *F1-score (F1)*: A harmonic average of precision and recall, offering a balanced measure of performance (Powers, 2020).
- *Balanced Accuracy*: The average of true positive and true negative rates, providing a robust metric for imbalanced datasets.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1 - Score = 2 \times \frac{Pr\ ecision \times Re\ call}{Pr\ ecision + Re\ call}$$

$$Balanced\ accuracy = \frac{TPR + TTNR}{2}$$

In these equations, TP represents true positives, TN represents true negatives, FP represents false positives, and FN represents false negatives.

3.4 Experimental Results

3.4.1 Model Performance

The Swin Transformer demonstrated exceptional performance, significantly outperforming established baseline models such as ResNet50 and the original vision Transformer (ViT). The detailed results from the test dataset are shown in Table 4.

3.4.2 Confusion Matrix

The confusion matrix for the Swin Transformer model is shown in Figure 8, illustrating the model's ability to discriminate between the five stages of diabetic retinopathy.

Label	Precision	Recall	F1- score	Support
0	1.00	1.00	1.00	355
1	0.97	1.00	0.99	72
2	1.00	0.98	0.99	181
3	0.93	1.00	0.96	37
4	1.00	0.98	0.99	59
	Accuracy		0.99	704
Macro Avg	0.98	0.99	0.99	704
Weighted Avg	0.99	0.99	0.99	704

Table 4. Performance Metrics of Swin Transformer

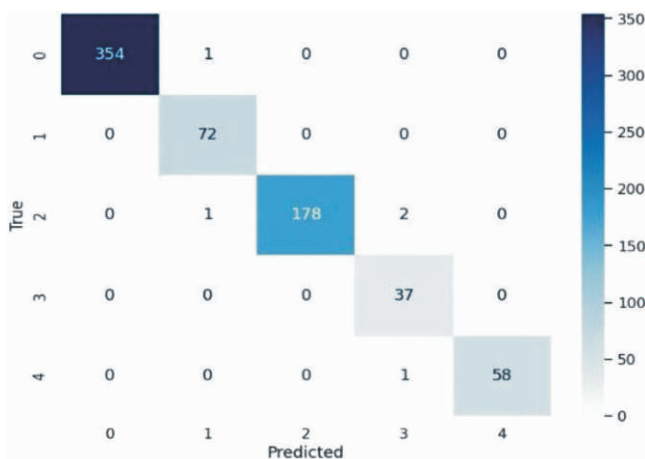


Figure 8. Confusion Matrix for Swin Transformer

3.4.3 Comparative Analysis

Compared to ResNet50 and Vision Transformer, the Swin Transformer exhibits:

- *Reduced Computational Overhead:* The window-shifting mechanism lowers the complexity of processing, enabling faster training times compared to the quadratic demands of ViT.
- *Improved Processing Speed:* Optimized attention computations enhance inference efficiency, making the model more practical for real-world applications.

3.4.4 Discussion

3.4.4.1 Strengths of Swin Transformer

- *Enhanced Feature Representation:* The adaptive window-shifting technique allows the model to analyze retinal images at multiple scales, improving its sensitivity to subtle DR indicators, such as microaneurysms, which are typically missed by less sophisticated models (Liu et al., 2021).
- *Scalability:* Its hierarchical processing framework enables efficient handling of high-resolution images while keeping memory consumption in check, a significant advantage over traditional CNNs that struggle with resource demands as image size increases (Liu et al., 2021).
- *Effects of Data Augmentation:* The incorporation of data augmentation techniques, such as random cropping and contrast adjustments, enhances the model's robustness, contributing to its high-test accuracy and ability to generalize across diverse retinal images (Shorten & Khoshgoffaar, 2019).

3.4.4.2 Challenges and Limitations

- *High Computational Requirements:* Despite its efficiency gains, the Swin Transformer still requires substantial computational resources, particularly GPU power, for training, which may limit its accessibility in resource-constrained settings (Liu et al., 2021).
- *Inter-class Similarity:* Distinguishing between DR stages with overlapping characteristics (such as mild versus moderate) remains a challenge due to the

subtle visual differences between these categories, occasionally leading to classification errors.

- *Interpretability:* Like many deep learning systems, the model's decision-making process lacks inherent transparency, necessitating additional tools like Grad-CAM to provide visual explanations that can foster trust among clinicians (Selvaraju et al., 2017).

3.4.4.3 Extended Insights

- *Feature Visualization:* Visualization of attention maps reveals that the model prioritizes key retinal regions, such as areas with hemorrhages or microaneurysms, during classification, offering valuable insight into its diagnostic focus.
- *Comparison of Hyperparameter Tuning:* Extensive hyperparameter tuning identified the AdamW optimizer, paired with a learning rate of $1e-4$, as the optimal configuration for achieving stable convergence and maximizing performance on the APTOS dataset (Défossez et al., 2020; Kingma, 2014).

3.4.4.4 Future Prospects

The success of this model opens several avenues for further exploration:

- *Multi-modal Enhancement:* Integrating retinal images with patient-specific data, such as age, blood sugar levels, or HbA1c values, could refine classification accuracy by providing additional context.
- *Real-time Application:* Adapting the model for deployment on edge devices could enable rapid DR screening in low-resource environments, improving accessibility and diagnostic turnaround time.
- *Cross-Dataset Validation:* Testing the model's performance on additional datasets beyond APTOS 2019 would validate its generalizability and robustness across different imaging conditions and populations.
- *Hybrid Model Development:* Combining the strengths of CNNs, which excel at local feature extraction, with the Swin Transformer's global processing capabilities could yield a hybrid architecture with even greater performance potential.

Conclusion

The Swin Transformer showcases extraordinary proficiency in evaluating the severity of Diabetic Retinopathy, achieving a test accuracy of 99.57% that markedly exceeds the capabilities of established models like ResNet50 and the original Vision Transformer. Its innovative window-shifting mechanism, coupled with a hierarchical processing structure, enables the model to adeptly capture both localized retinal details and broader contextual features, delivering unparalleled diagnostic precision. Advanced preprocessing strategies, including augmentation and normalization, further bolster its resilience against dataset imbalances, ensuring consistent performance across all DR stages. However, challenges such as high computational requirements and the need for greater interpretability remain, necessitating the integration of explainable AI techniques like Explainable AI (XAI) to enhance clinical acceptance and trust.

Beyond its application to DR, the Swin Transformer holds significant promise for broader medical imaging tasks, offering a versatile framework that could be adapted to other diagnostic challenges. Future research should prioritize the exploration of hybrid architectures that blend convolutional and transformer approaches, rigorous validation across diverse datasets to confirm generalizability, and optimization for real-time deployment on resource-limited devices. These efforts will further solidify the model's role as a transformative tool in healthcare, facilitating early detection of vision-threatening conditions like DR and ensuring timely interventions that improve patient outcomes. By addressing critical global health challenges through AI-driven innovation, this study underscores the potential of advanced deep learning systems to revolutionize diagnostic practices and enhance the quality of care worldwide.

References

- [1]. Beam, A. L., & Kohane, I. S. (2018). Big data and machine learning in health care. *JAMA*, 319(13), 1317-1318.

<https://doi.org/10.1001/jama.2017.18391>

[2]. Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *The Journal of Machine Learning Research*, 13(1), 281-305.

[3]. Brownlee, J. (2019). *Deep Learning for Computer Vision: Image Classification, Object Detection, and Face Recognition in Python*. Machine Learning Mastery.

[4]. Chen, L. C., Zhu, Y., Papandreou, G., Schroff, F., & Adam, H. (2018). Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 801-818).

[5]. Coan, L. J., Williams, B. M., Adithya, V. K., Upadhyaya, S., Alkafri, A., Czanner, S., & Czanner, G. (2023). Automatic detection of glaucoma via fundus imaging and artificial intelligence: A review. *Survey of Ophthalmology*, 68(1), 17-41.

<https://doi.org/10.1016/j.survophthal.2022.08.005>

[6]. Défossez, A., Bottou, L., Bach, F., & Usunier, N. (2020). A simple convergence proof of adam and adagrad. *arXiv preprint arXiv:2003.02395*.

<https://doi.org/10.48550/arXiv.2003.02395>

[7]. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., & Dehghani, M. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.

<https://doi.org/10.48550/arXiv.2010.11929>

[8]. Dumoulin, V., & Visin, F. (2016). A guide to convolution arithmetic for deep learning. *arXiv preprint arXiv:1603.07285*.

<https://doi.org/10.48550/arXiv.1603.07285>

[9]. Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). *Deep Learning*. MIT Press, Cambridge.

[10]. Gulshan, V., Peng, L., Coram, M., Stumpe, M. C., Wu, D., Narayanaswamy, A., & Webster, D. R. (2016). Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *JAMA*, 316(22), 2402-2410.

<https://doi.org/10.1001/jama.2016.17216>

[11]. Hatamizadeh, A., Nath, V., Tang, Y., Yang, D., Roth,

H. R., & Xu, D. (2021). Swin unetr: Swin transformers for semantic segmentation of brain tumors in MRI images. In *International MICCAI Brainlesion Workshop* (pp. 272-284). Springer International Publishing.

https://doi.org/10.1007/978-3-031-08999-2_22

[12]. Hathot, S. F., Jubier, N. J., Hassani, R. H., & Salim, A. A. (2021). Physical and elastic properties of TeO₂-Gd₂O₃ glasses: Role of zinc oxide contents variation. *Optik*, 247, 167941.

<https://doi.org/10.1016/j.ijleo.2021.167941>

[13]. Kaggle. (2019). *APTOS 2019 Blindness Detection*. Asia Pacific Tele-Ophthalmology Society (APTOS). Retrieved from

<https://www.kaggle.com/competitions/aptos2019-blindness-detection>

[14]. Kingma, D. P. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

<https://doi.org/10.48550/arXiv.1412.6980>

[15]. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., & Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 10012-10022).

[16]. O'shea, K., & Nash, R. (2015). An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458*.

<https://doi.org/10.48550/arXiv.1511.08458>

[17]. Powers, D. M. (2020). Evaluation: From precision, recall and F-measure to ROC, informedness, markedness and correlation. *arXiv preprint arXiv:2010.16061*.

<https://doi.org/10.48550/arXiv.2010.16061>

[18]. Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 618-626).

[19]. Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1), 1-48.

<https://doi.org/10.1186/s40537-019-0197-0>

- [20]. Tharwat, A. (2021). Classification assessment methods. *Applied Computing and Informatics*, 17(1), 168-192.
<https://doi.org/10.1016/j.aci.2018.08.003>
- [21]. Topol, E. J. (2019). High-performance medicine: The convergence of human and artificial intelligence. *Nature Medicine*, 25(1), 44-56.
<https://doi.org/10.1038/s41591-018-0300-7>
- [22]. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 1-11.
- [23]. World Health Organization. (2025). *Blindness and Vision Impairment*. Retrieved from
<https://www.who.int/news-room/fact-sheets/detail/blindness-and-visual-impairment>
- [24]. Xu, X., Feng, Z., Cao, C., Li, M., Wu, J., Wu, Z., & Ye, S. (2021). An improved swin transformer-based model for remote sensing object detection and instance segmentation. *Remote Sensing*, 13(23), 4779.
<https://doi.org/10.3390/rs13234779>

ABOUT THE AUTHORS

Sanjay Gandhi Gundabatini is a Professor in the Department of Computer Science and Engineering at Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India.

Sai Sindhu Manne is currently pursuing a B.Tech degree in the Department of Computer Science and Engineering at Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India.

Sunkara Likhit Babu is currently pursuing a B.Tech degree in the Department of Computer Science and Engineering at Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India.

Vangapandu Bhargava Rao is currently pursuing a B.Tech degree in the Department of Computer Science and Engineering at Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India.

Sanka Tejaswi is currently pursuing a B.Tech degree in the Department of Computer Science and Engineering at Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India.