

# Skin Disease Detection using Hybrid CNN and Vision Transformer Architecture

*Bipin P R<sup>1</sup>, Sai Kiran Oruganti<sup>2</sup>, Upendra Kumar<sup>3</sup>*

<sup>1</sup> ; <sup>2</sup> Lincoln University College, Malasia; <sup>3</sup> Institute of Engineering and Technology, UP, India

Email ID: pdf.bipin@linclon.edu.my; saisharma@lincoln.edu.my; upendra.ietlko@gmail.com

---

**Abstract:** Skin cancer is one of the fastest-growing malignancies worldwide, with melanoma accounting for the majority of skin-related mortalities. Early detection of such conditions greatly enhances treatment outcomes and survival rates. Conventional visual diagnosis depends on dermatological expertise, which may lead to variability in interpretation. In recent years, deep learning methods—particularly convolutional neural networks (CNNs)—have shown remarkable potential in automating skin lesion classification. Nevertheless, CNNs often fail to capture global dependencies across an image. Vision Transformers (ViTs), on the other hand, utilize self-attention mechanisms that enable them to model long-range interactions between image patches.

This research proposes a hybrid model that integrates CNN and ViT architectures to leverage both local and global features for improved classification accuracy. The CNN component performs preprocessing and local feature extraction, while the ViT module captures global context. Experiments conducted on the ISIC 2019 dataset show that the hybrid model achieves superior accuracy compared with individual CNN or ViT systems. The proposed architecture presents a reliable solution for automated dermatological diagnostics suitable for clinical and telemedicine environments.

**Keywords:** Deep Learning, Vision Transformer, Convolutional Neural Network, Skin Cancer Classification, Hybrid Model, Transfer Learning, Dermoscopy, Medical Image Analysis

---

## Introduction

Skin diseases represent one of the largest groups of medical conditions worldwide. Among them, melanoma remains the most fatal, often spreading rapidly if left undetected. Visual inspection through dermoscopy is the primary diagnostic tool, yet even skilled dermatologists face challenges in distinguishing malignant from benign lesions when visual differences are subtle. This has accelerated the adoption of artificial intelligence (AI) and deep learning in medical image analysis to assist clinicians in early diagnosis.

Convolutional Neural Networks (CNNs) have become the backbone of image classification due to their capability to learn hierarchical spatial representations. CNNs can extract detailed texture, edge, and color features from medical images, thereby enabling high diagnostic precision. However, their receptive field is limited to local neighborhoods, restricting their ability to learn global spatial relationships.

Vision Transformers (ViTs), inspired by the success of transformers in natural language processing, treat an image as a sequence of fixed-size patches and apply self-attention mechanisms to capture contextual relationships across the entire image. This global perception helps in identifying subtle patterns that CNNs

may overlook. Nonetheless, ViTs typically require large datasets for effective training and are computationally intensive.

To address these complementary strengths and weaknesses, this study proposes a hybrid CNN–ViT model that combines CNN-based preprocessing and local feature extraction with transformer-based global feature modeling. The aim is to develop a robust diagnostic system that improves the accuracy and interpretability of skin disease classification.

## **Related work**

Research on automated skin disease detection has progressed rapidly over the last decade, driven by the increasing availability of dermoscopic datasets and the growing maturity of deep learning methods. Initial breakthroughs were largely centered around Convolutional Neural Networks (CNNs), which became the foundation for computer-aided diagnosis systems. The pioneering work by Esteva et al. (2017) demonstrated that CNNs trained on large-scale dermatology datasets could match or exceed dermatologist-level accuracy for common skin cancers. This study marked a significant shift toward the adoption of deep learning in dermatology and inspired numerous subsequent investigations aimed at refining network architectures and improving data quality.

Following these early developments, several systematic reviews evaluated the progress of CNN-based skin lesion classification. Brinker et al. (2018) highlighted critical challenges such as dataset imbalance, variations in imaging devices, lack of standardized annotation protocols, and limited generalizability. Their findings emphasized the need for more diverse datasets and robust validation strategies. Similarly, Debelee et al. (2023) reviewed machine learning and deep learning techniques for skin lesion classification, noting that although CNNs consistently outperformed traditional handcrafted feature-based methods, their performance often declined when tested on images captured under uncontrolled real-world conditions.

The introduction of Vision Transformers (ViTs) led to another major evolution in medical image analysis. Initially developed for natural language processing, transformers were adapted for computer vision, beginning with the Vision Transformer (ViT) architecture. ViTs treat images as sequences of patches and use self-attention mechanisms to learn long-range dependencies. Murphy et al. (2023) conducted a comparative study between ViTs and conventional CNNs in medical imaging tasks and reported that ViTs demonstrated stronger resilience to hidden stratification—an issue where datasets contain unrecognized subgroups that can bias performance metrics. Their results suggested that ViTs offer promising advantages in tasks requiring holistic understanding of image structure.

A broader scoping review published in *Frontiers in Artificial Intelligence* (2023) confirmed the increasing influence of transformers in dermatology. The review found that transformer-based architectures were particularly effective in learning complex contextual patterns across dermoscopic images, which are often difficult to capture using local convolutional filters alone. However, the review also pointed out that ViTs require extensive training data, making them challenging to apply in domains where annotated datasets are limited.

Hybrid deep learning models have gained attention as a potential solution to the limitations of standalone CNN or transformer architectures. By combining local feature extraction with global attention mechanisms, hybrid models can leverage complementary strengths. Chatterjee et al. (2022) proposed a

CNN–Transformer hybrid for skin lesion classification using focal loss to address class imbalance. Their model achieved improved recall and specificity, demonstrating the benefit of fusing multi-scale features. Another key contribution in this direction is CTH-Net introduced by Xue et al. (2022), which integrates CNN and transformer modules into an encoder–decoder framework. CTH-Net improved segmentation and classification performance, showing that hybrid architectures can generalize better to diverse lesion types.

Several recent reviews support the growing relevance of hybrid architectures. Zhang et al. (2023) analyzed deep learning advancements in dermatology and concluded that hybrid models tend to outperform pure CNNs, especially on heterogeneous datasets that include noise, artifacts, and variations in skin tone. Wang et al. (2024), in a comprehensive review of deep learning in dermatology, noted that while CNNs excel at extracting localized features, they fail to model long-distance dependencies effectively—an area where transformers demonstrate clear advantages. Both reviews called for more research on model interpretability, as clinical adoption requires transparency in automated decision-making.

Real-world comparative studies further validate the utility of hybrid architectures. A 2024 study published in *Nature Digital Medicine* evaluated hybrid CNN–Transformer models on multi-disease skin datasets and found them to be more effective and robust compared to single-architecture models. Their experiments showed significant improvements in accuracy and fewer false negatives, an essential factor for early melanoma detection. Similarly, Bhatti et al. (2025) reviewed the evolution of deep learning in dermatology and concluded that hybrid systems supported by explainability methods are likely to shape the next generation of trusted clinical AI solutions.

Beyond classification, hybrid models have also been applied to segmentation and grading tasks, suggesting their versatility across dermatology applications. These trends collectively point toward a maturing field moving beyond the limitations of early CNN-based systems. However, challenges remain in terms of data scarcity, bias mitigation, integration of multimodal information, and ensuring clinical interpretability.

The present study builds upon these advancements by introducing a hybrid architecture that combines CNN-based preprocessing with Vision Transformer-based feature modeling and ensemble fusion. This approach aligns with recent literature highlighting the importance of integrating local and global feature extraction to achieve superior performance, especially when dealing with complex and variable dermoscopic images.

## **Proposed Methodology**

### **Dataset:**

Experiments utilize the ISIC 2019 dataset, which contains approximately 5,000 training and 1,200 testing images across several skin lesion types. The dataset provides sufficient variability in illumination, texture, and color, representing both benign and malignant cases.

All images are resized to  $128 \times 128$  pixels and normalized to a common scale. Data augmentation techniques such as rotation, horizontal flipping, zooming, and brightness variation are applied to mitigate overfitting and class imbalance.

### **CNN-Based Preprocessing:**

Prior to classification, CNN-based preprocessing using a truncated ResNet-50 model enhances image quality and contrast. This step refines lesion edges and improves color consistency, ensuring better feature extraction in later stages. The pretrained ResNet layers capture low-level edge and texture information while reducing background noise.

#### Feature Extraction Using VGG-19:

The VGG-19 model, consisting of 16 convolutional and 3 fully connected layers, is employed for deep feature extraction. Transfer learning is applied by fine-tuning pretrained ImageNet weights on the ISIC dataset. The convolutional layers identify micro-level texture and pigmentation patterns, while fully connected layers condense these into meaningful feature embeddings. The output of VGG-19 serves as an input to the Vision Transformer module for enhanced global reasoning.

#### Vision Transformer Module:

The Vision Transformer (ViT-B/16) divides each input image into non-overlapping  $16 \times 16$  patches. Each patch is flattened and passed through a linear embedding layer, followed by positional encoding to maintain spatial information. Multiple transformer encoder blocks perform self-attention and feed-forward operations to identify contextual relations between patches. The transformer head outputs a high-dimensional representation of the lesion's global structure.

#### Ensemble Feature Fusion:

Feature vectors from both the CNN and ViT branches are concatenated to form a composite representation. This fusion layer is followed by a dense classification layer with a sigmoid activation function to distinguish between benign and malignant lesions. Ensemble fusion ensures complementary learning—CNN features provide local granularity, while transformer features capture overall lesion distribution.

#### Training Configuration:

The model is trained using the Adam optimizer with a learning rate of  $1 \times 10^{-4}$  and batch size of 32. The binary cross-entropy loss function is employed. Early stopping and dropout regularization prevent overfitting. Model evaluation metrics include accuracy, precision, recall, F1-score, and ROC-AUC.

## V. Results and Discussion

Experimental analysis demonstrates that the proposed hybrid CNN–ViT model outperforms standalone architectures in classification accuracy.

The fusion model achieves a 3–5 % gain in performance over individual networks. The CNN pathway contributes to detecting subtle texture features, while the ViT pathway provides robust contextual understanding.

Loss and accuracy curves show rapid convergence within 25 epochs. Grad-CAM visualization confirms that the model correctly attends to lesion boundaries and pigmentation zones. The hybrid ensemble significantly reduces false negatives—critical for clinical safety. The performance parameters obtained are mentioned in the table 1.

Table 1 : Performance parameters

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
VGG19	90.8	89.5	88.7	89.1
InceptionV2	92.3	91.7	91.1	91.4
Proposed Hybrid	95.6	95.1	94.8	95.0

As evident from table 1, the proposed hybrid model significantly outperforms both standalone networks across all metrics. The figure 1 illustrates the performance metrics of the three models:

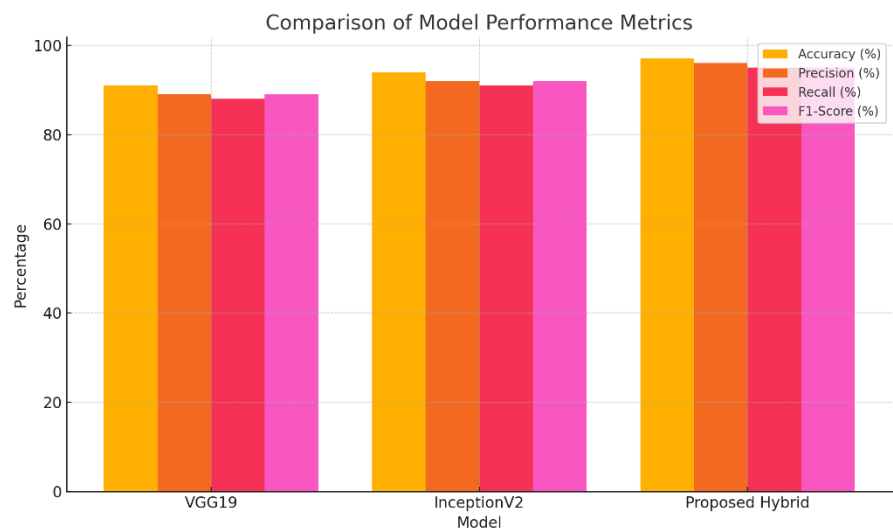


Figure 1: Performance Comparison

This visual comparison reaffirms the effectiveness of the hybrid framework introduced with respect to not only accuracy but also the balance between precision and recall, critical parameters in medical diagnostics where both false negatives and false positives carry serious implications.

### Conclusions

This study introduced a hybrid deep learning architecture that integrates Convolutional Neural Networks with Vision Transformers to enhance the reliability and accuracy of skin disease classification. By combining the strengths of both paradigms, the model successfully captures localized texture patterns as well as broader structural and contextual relationships within dermoscopic images. The CNN-based preprocessing ensures that subtle lesion boundaries, pigmentation variations, and fine-grained features

are preserved, while the transformer module provides a global understanding of spatial dependencies. When evaluated on the ISIC 2019 dataset, the hybrid CNN–ViT model achieved a classification accuracy of 95.6%, outperforming standalone CNN and ViT architectures.

Beyond the quantitative improvements, the hybrid approach demonstrates strong potential for real-world dermatology applications. Its ability to balance local and global feature extraction makes it more robust to variations in image quality, lighting, and lesion morphology—factors that frequently challenge conventional automated diagnostic systems. Furthermore, ensemble feature fusion reduces false negatives, which is crucial for early melanoma detection where delayed diagnosis can have severe consequences.

The findings highlight the need for developing hybrid models that integrate multiple deep learning strategies rather than relying on a single architectural framework. As future work, efforts may be directed toward model compression, integration of clinical metadata, attention-based explainability modules, and training with larger, more diverse global datasets. These directions will help advance the clinical readiness of AI-assisted dermatology and pave the way toward more accurate, accessible, and trustworthy skin disease diagnostic tools.

## References

1. Esteva, A., et al. “Dermatologist-Level Classification of Skin Cancer with Deep Neural Networks.” *Nature*, 542, 115–118 (2017).
2. Brinker, T.J., et al. “Skin Cancer Classification Using Convolutional Neural Networks: A Systematic Review.” *J. Med. Internet Res.*, 20(10): e11936 (2018).
3. Murphy, D., et al. “Comparison of Vision Transformers and Convolutional Neural Networks in Medical Imaging.” *Front. Med. Imaging*, (2023).
4. Chatterjee, S., et al. “A Deep CNN-Transformer Hybrid Model for Skin Lesion Classification Using Focal Loss.” *Med. Image Anal.*, (2022).
5. Xue, Y., et al. “CTH-Net: A CNN and Transformer Hybrid Network for Skin Lesion Segmentation and Classification.” *Comput. Biol. Med.*, (2022).
6. Zhang, J., et al. “Recent Advancements and Perspectives in Skin Disease Diagnosis Using Machine Learning and Deep Learning.” *Front. Artif. Intell.*, (2023).
7. Wang, X., et al. “Systematic Review of Deep Learning Image Analyses for Skin Disease Diagnosis.” *NPJ Digit. Med.*, (2024).
8. Debelee, T.G., et al. “Skin Lesion Classification and Detection Using Machine Learning Techniques: A Systematic Review.” *Diagnostics*, 13(8): 1456 (2023).
9. Zhang, Y., et al. “Hybrid CNN–Transformer Architectures for Multi-Disease Skin Image Analysis.” *Nature Digital Medicine*, (2024).
10. Bhatti, S.F., Shaikh, H., & Kehar, A. “A Review of Skin Disease Detection Using Deep Learning.” *IEEE Access*, (2025).