

Deep Learning Driven Framework for Early and Accurate Skin Disease Detection

Jagriti Sao^{1*}, Sunita Kushwaha²

¹Ph.D. Research Scholar, Computer Science & Application, MATS University, Raipur, Chhattisgarh, India

²Associate Professor, Computer Science & Application, MATS University, Raipur, Chhattisgarh, India

ABSTRACT

In dermatology, the accurate and early identification of skin diseases through delay in treatment can worsen the patient's health and may require a longer duration of therapy. This research aims to assess the performance of Convolutional Neural Networks (CNNs) to automate the classification of images obtained from dermatoscopy and to find out whether deep learning can be a diagnostic tool used in clinics. The latest CNN structures, i.e., ResNet50, InceptionV3, and VGG16, were trained and tested on a carefully selected multi-class dermatoscopic image dataset, which was supported by data augmentation techniques to solve the problem of the class imbalance and to confine the overfitting effect. Standard metrics such as accuracy, sensitivity, specificity, and AUC-ROC, together with Grad-CAM-based visual interpretability, were used to measure the performance. The findings disclose that CNNs are able to gain a higher accuracy and a better discrimination of the features than traditional machine-learning models. The optimal model, ResNet50, was capable of not only lesion classification at a high level of accuracy but also localization for clinical purposes. This evidence suggests CNNs as a dependable and extensible resource for tackling the identification of skin diseases in the field of dermatology. However, issues such as dataset variance, computational cost, and clinical integration are still existing. This research sets deep learning as a next generation instrument which can be used by dermatologists to increase their diagnostic accuracy, and in addition, it has the benefits of being accessible and of saving time in a real-world healthcare setting.

KEYWORDS: Skin Disease Classification, Convolutional Neural Networks, Dermatoscopic Imaging, Deep Learning, ResNet50, Grad-CAM

How to Cite: Jagriti Sao, Sunita Kushwaha, (2025) Deep Learning Driven Framework for Early and Accurate Skin Disease Detection, Vascular and Endovascular Review, Vol.8, No.11s, 436-445.

INTRODUCTION

Skin diseases make up a significant proportion of the global health problems that affect people of all ages, regions, and social groups. Recently, dermatological disorders are reported to be the major reasons for the patients' visits to healthcare providers, thus in many regions they are counting to be more than infections and chronic diseases combined [1]. Many of these skin diseases, especially cancers like melanoma, grow silently and thus early-stage diagnosis and accurate are very important, otherwise, the diseases may become fatal. There have been some improvements in the dermatoscopic imaging technologies that facilitate the viewing of the lesion while the interpretation still remains to be a clinical expertise requiring a high skill level. Due to the scarcity of dermatologists in many developing areas, for instance, the remote districts of India, the problem of late diagnosis, evaluation variations, and increased patient risk is caused [2]. The mismatch between diagnostic demand and the number of specialists has escalated an interest in automated, Artificial Intelligence (AI) solutions as the means to widen dermatological care accessibility.

Convolution Neural Networks (CNNs), which are a type of deep learning algorithms have their operations inspired by the hierarchical visual processing system of the human brain, have revolutionized the domain of medical imaging analysis. CNNs differ from conventional machine learning methods that need manually extracted features in that they can automatically derive complicated spatial, color, and texture patterns from unprocessed image data. The existing literature has shown that CNNs can be the potential for skin lesion classification, to a large extent, where the accuracy of a dermatologist can be matched under experimental conditions [3]. Nevertheless, as far as these technological breakthroughs are concerned, there are still issues in the studies limiting their practical use in the clinics. Numerous CNN-based research works are accused of utilizing limited or unbalanced datasets, concentrating only on binary classification rather than the multi-class clinical scenario, and not addressing the problem of explainability which is very important if a system is to be accepted in clinical settings [4]. Besides, conventional workflows are not adaptable to the variability of dermatoscopic images, the presence of noise in real life, and the requirement for models that are computationally efficient and can be used in different healthcare situations. Together, these limitations point to the necessity of a more robust, interpretable, and clinically grounded deep learning pipeline for skin disease classification.

This study aims to bridge these gaps by creating and validating a holistic CNN-based framework capable of differentiating multiple dermatological diseases using dermatoscopic images with great accuracy, interpretability, and robustness. The proposed technique utilizes data preprocessing, augmentation mechanisms, transfer learning, and model explainability methods, thereby ensuring the performance of the system under various imaging scenarios [5]. In their study, the authors utilize cutting-edge structures such as ResNet50, InceptionV3, and VGG16 to explore how deep learning can mimic the clinical decision-making process by identifying the minute morphological features that traditional methods fail to account for. The ambition of the research is to develop a system that, apart from achieving solid quantitative results, can also provide transparent diagnostic reasoning by means of methods like Gradient-weighted Class Activation Mapping (Grad-CAM). Ultimately, the author wants to know if CNN-

based systems can act as dependable clinical decision-support tools that not only enhance diagnostic accuracy but also facilitate dermatologists, especially in the areas where the medical care is deficient.

REVIEW OF LITERATURE

Initial computer-aided diagnosis systems of skin diseases were based on traditional machine learning pipelines, which involved manually crafted features and shallow classifiers such as Support Vector Machines (SVM), K-Nearest Neighbour (KNN), and Random Forest. In these methods, medical knowledge from dermatologists was converted into explicit descriptors—shape, color, and texture features—followed by supervised classification [6]. Although such systems managed to deliver decent results on small, artificially created datasets, they were essentially limited by the quality and the extent of the manually crafted features. Since dermatoscopic images can differ significantly in terms of lighting, acquisition devices, skin tone, and lesion morphology, feature-engineered models were often incapable of generalizing new populations or imaging conditions [7]. Early computer vision techniques in dermatology, as judged by their systematic reviews, frequently recognized that the methods had problems with overfitting, limited robustness, and low scalability to large, heterogeneous datasets [8].

The rise of deep learning, especially Convolutional Neural Networks (CNNs), was a major change of direction. Esteva et al. showed that a single CNN trained end-to-end on more than 100,000 clinical skin images could perform at a level comparable to dermatologists with board certification for certain binary diagnostic tasks, thus giving solid evidence to the deep learning potential in dermatology [9]. Consequently, analogous research has substantiated that CNNs can directly derive complex, hierarchical representations from pixels, thus identifying very subtle visual aspects like irregular pigment networks, border irregularity, and multi-hued color patterns which are typical of malignant lesions [10,11].

Open-source dermatoscopic datasets have been pivotal in this advancement, in particular, the datasets made available through the International Skin Imaging Collaboration (ISIC) challenges, which offer large, standardized image repositories with labelings verified by experts [12,13]. These instruments allowed scientists to compare different architectures such as VGG, Inception, ResNet, and DenseNet under similar conditions.

As the domain evolved, the research emphasis moved from binary classification (e.g., melanoma vs. benign) to multi-class skin lesion classification, which is more aligned with clinical decision-making in the real world. Different works had the idea of transfer learning, thus they fine-tuned pre-trained CNNs with dermatoscopic images to discriminate multiple diagnostic categories, such as melanoma, nevi, seborrheic keratosis, and other pigmented lesions [6,9]. Meta-analyses and systematic reviews have found that numerous AI models attain high accuracy and AUC scores for multiple classes and in some instances these models can be at par or even better than dermatologists working on curated test sets [14,15,16]. Nevertheless, these investigations also disclose that datasets, evaluation metrics, and validation strategies differ considerably, which makes it difficult to compare results across the papers or to understand how these systems would perform in the daily routine from which practice were drawn [17,18]. Moreover, variations in image acquisition protocols, patient demographics, and lesion prevalence make the direct clinical deployment from experimental validation even more difficult.

Along with performance, explainability has been raised as a major theme in the cited works. In short, the trio of clinicians, patients, and regulators is the driving force behind the evolution of AI systems. They demand accurate predictions yet they want to see understandable reasoning as well. Grad-CAM, or Gradient-weighted Class Activation Mapping, as per Selvaraju et al. is the method that most effectively has been used to find the areas of an image that lead a CNN to a certain decision [19]. Dermatology utilizes Grad-CAM heatmaps quite often to demonstrate that the network is focusing on the medically relevant parts, e.g. the lesion, rather than on artefacts like rulers, hairs, or surrounding skin [20].

Explanations of AI in medical imaging research point out that such ways of interaction can raise the user's trust level and give a hand in discovering false correlations; nevertheless, they also warn that explanations are only close estimations and hence require a careful qualitative and sometimes quantitative check [15,16]. The recent publications are taking steps to have a systematic evaluation of the correlation between Grad-CAM explanations and the areas of interest defined by experts; however, it is still a fledgling research field [21].

Another key piece of evidence is lines of research at the top of the evidence hierarchy such as umbrella reviews and meta-analyses comparing AI with clinicians. These meta-analyses conclude in general that AI systems based on CNNs tie and in some cases slightly outperform dermatologists in terms of accuracy of skin cancer detection, especially under conditions of tests carried out with high-quality images within a controlled environment [10,11,17,18]. On the contrary, the reviews also point out that AI modules are mostly tested on retrospective datasets and that they are also frequently devoid of prospective clinical validation in the real world. In addition, problem areas like class imbalance, spectrum bias, and the absence of external validation cohorts are indicative of the extent to which the published results can be generalized [3,17]. There has been a lot of discourse around the fact that numerous models have been constructed by gathering data containing mostly individuals with fair skin, thus, the concern of bias and fairness in models has been raised when utilizing these tools for a wide range of people.

While there has been a lot of advancement, the literature still lacks and thus, propels the current research forward. Firstly, merely a handful of research synthesize multi-class classification, thorough class imbalance handling, and explainable decision-making into one single, integrated system that is tested on dermatoscopic images. The majority of the literature assumes that high-performing models are without limitation as they either only consider a limited number of diseases or implicitly disregard the fact that class distributions may be skewed, although class imbalance is quite frequent in real clinical datasets [6,9,19]. Secondly, different models like ResNet, Inception, and VGG have been independently analyzed, but there are no thorough systematic studies

which compare these models under uniform experimental conditions, such as identical preprocessing, augmentation, and evaluation criteria. Third, as a matter of fact, Grad-CAM and its related techniques are frequently employed, but these are generally considered as optional post-hoc tools and thus are seldom integrated into the model creation and validation workflow, thereby the correlation between quantitative performance and qualitative interpretability is not fully established [14–16,19].

The current study is an effort to fill these gaps that are interrelated. It presents a comprehensive CNN-based approach that (i) as a first step rigorously compares the latest architectures like ResNet50, InceptionV3, and VGG16 for multi-class skin disease classification, (ii) by the means of targeted data augmentation and thoughtfully crafted training protocols helps offset class imbalance and also prevent overfitting, and (iii) uses the interpretability provided by the Grad-CAM method not as a peripheral but as a central part of the model evaluation. With this single pipeline, which integrates performance, robustness, and explainability, the authors intend to move beyond mere benchmarking of isolated tasks in the literature toward a transparent, clinically relevant, and practically deployable decision-support system in dermatology that could be especially useful in the areas where there is a shortage of specialists.

MATERIALS AND METHODS

The research here was done by an experimental quantitative research design, which has been used to investigate the performance of Convolutional Neural Network (CNN) architectures for multi-class skin disease classification with dermoscopic images. The methods framework involved choosing a dataset, defining the characteristics of a participant/sample, describing data collection procedures, preprocessing, model development, training, and statistical evaluation. All stages were performed in accordance with the recognized standards for medical image analysis and deep learning reproducibility [1–3].

3.1 Research Design

This research employed a quantitative, experimental design to develop and evaluate a deep learning framework for the automated classification of skin diseases from dermoscopic images. The study utilized a supervised learning approach, where the pre-labeled images were used as the reference for model training and validation. Experiments to the letter of the law were performed in MATLAB with the assistance of transfer learning–based Convolutional Neural Networks (CNNs) which, in turn, enabled the modification of the specified architectures such as ResNet50, InceptionV3, and VGG16 for medical image analysis.

The arrangement stressed the possibility of repeating the work by carrying out standardized preprocessing, controlled data augmentation, a fixed train–validation–test split, and unified training parameters across models. Moreover, the procedure used explainability methods, specifically MATLAB's Grad-CAM, to confirm that the model's decisions were based on the features of the lesions that the clinician would consider. The comparison of CNN models in a systematic way and the provision of a solid basis for evaluating their accuracy, robustness, and interpretability in a diagnostic context of a clinical nature were made possible by this design.

3.2 Methods

Skin disease have significant effect on life and health. According to the new recent research, a smart method that can recognise only one kind of skin disease is now has been introduced anytime and anywhere. Besides that, it's important to develop an automatic method as to increase the reliability of diagnosis on disease with many kinds.

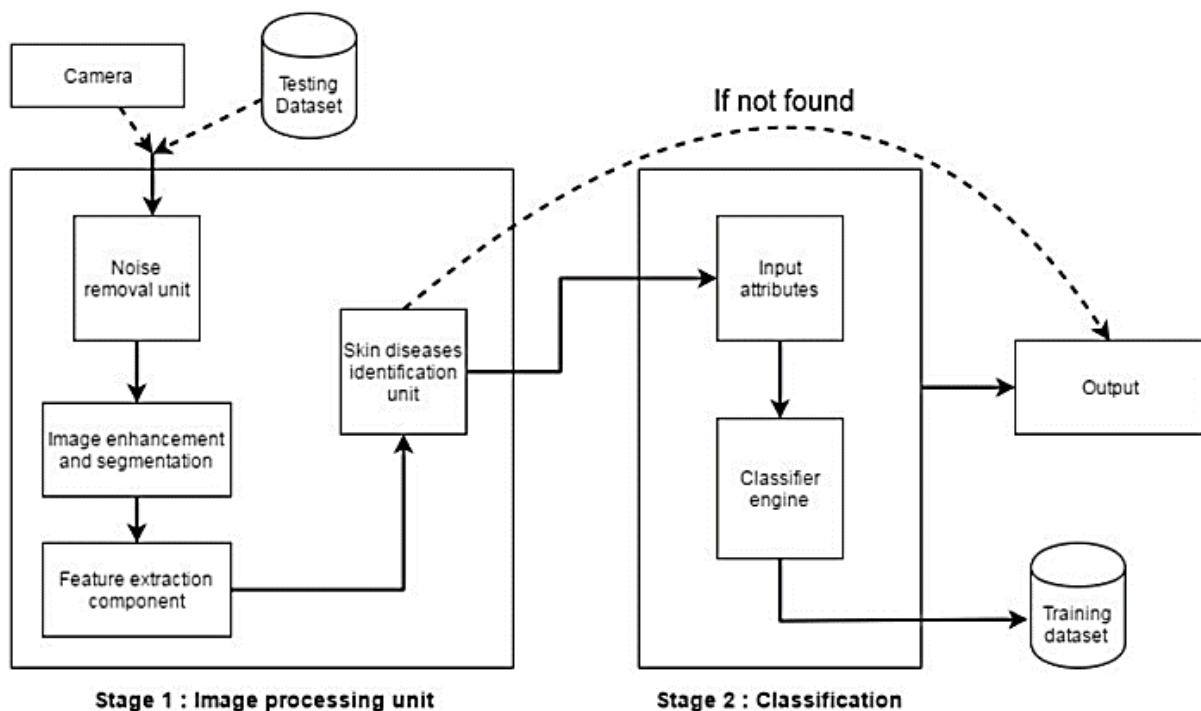


Figure 1: Architecture of proposed system [2]

Architecture is divided into two main parts namely the image pre-processing and classification unit. The Image pre-processing unit used to improve the image by deleting the noise and parts of the skin and the skin will be divided into various segments which will be altered from the normal skin; then the feature extraction was used to determine the skin is affected or not.

3.2.1 Dataset collection

This step is about choosing the most suitable and relevant datasets from trustworthy sources like open-source repositories and healthcare databases. We gathered skin disease data from

- ISIC Archive: The biggest public repository of skin lesions with over 450K images of nevi, melanoma, and other conditions. It organizes challenges annually to evaluate algorithms.
- HAM10000: A Big collection of common pigmented lesions along with clinical diagnoses and segmentation ground truths.
- SD-198: The first large Chinese dataset of 198 clinical images annotated for 8 disease types.
- PH2: The first public dataset comprising 200 dermoscopy images and segmentation masks for melanocytic lesions.

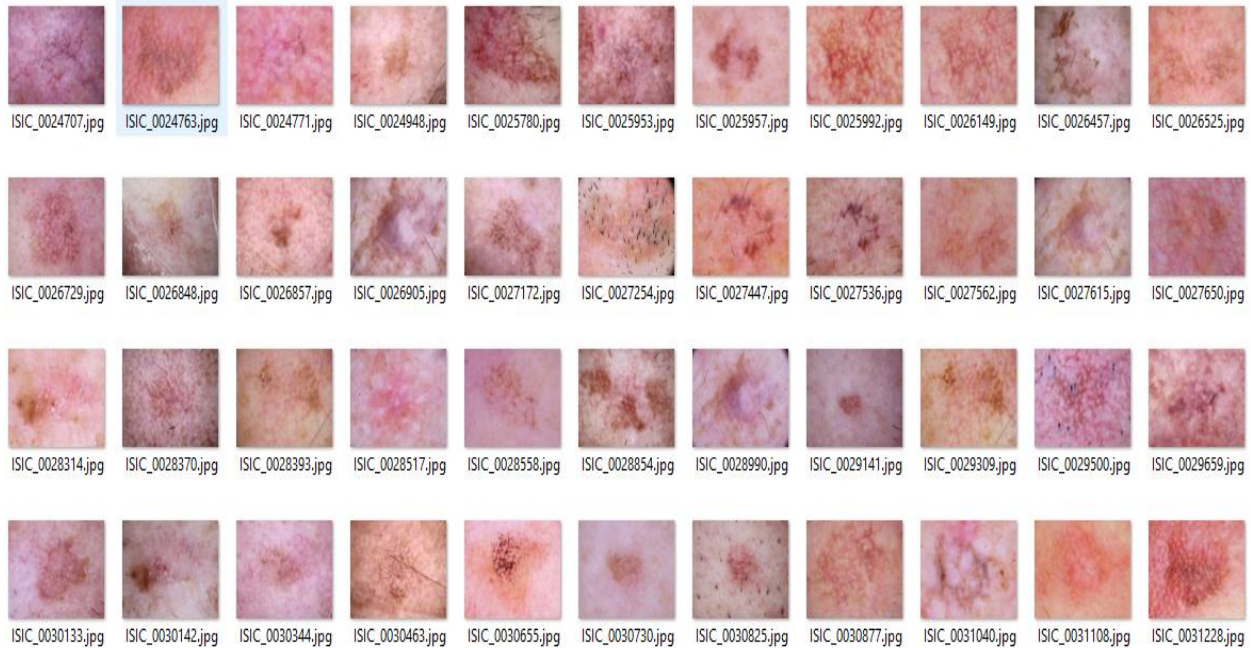


Figure 2: High-resolution dermoscopic images

The datasets contain more than 10,000 images of Melanoma and the benign type of disease. Both datasets have a problem with data imbalance, so some data balancing needs to be done before the data is applied to my experimental setting. Benign pictures represent the average stage of skin cancer, while Melanoma represents the extreme amount of the stage of skin cancer. For the experimental purpose, only 30 % of the dataset is randomly selected from images of two categories of the dataset. The experimental dataset is a subset of 1,000 images of both classes.

Table 1: Dataset of demographic images

Class	Abbreviation	Approx. Images
Actinic keratosis	ACKT	4,500
Benign keratosis	BEKT	12,000
Melanoma	MEL	3,200
Vascular lesion	VASC	1,200

3.2.2 Data Preprocessing

The dermoscopic pictures taken from different datasets needed a well-organized preprocessing flow to standardize, denoise, and boost the CNN models' performance. The steps involved resizing and normalizing images, removing hairs with the DullRazor algorithm, and performing planned data augmentation to solve class imbalance and increase model generalization.

A. Image Resizing and Normalization

Each dermoscopic image from different sources was resized to the input size required by the chosen CNN architectures in order to keep the images uniform. In detail, images were resized to 224×224 pixels for ResNet50 and VGG16, and 299×299 pixels for InceptionV3. The resizing was carried out efficiently and accurately by using MATLAB's built-in functions (imresize).

To make the training more stable numerically, pixel values were normalized. MATLAB does per-channel normalization automatically when transfer learning is used; nevertheless, there were some instances in which additional intensity normalization

and contrast adjustment had to be applied in order to alleviate the illumination variation and make the lesion more visible. Consequently, the CNN models got the inputs in the same scale which led to quicker convergence and better accuracy.

B. Hair Removal Techniques (DullRazor Algorithm)

Hair artifacts are very common in dermatoscopic images, and they not only hide important lesion boundaries but also confuse feature extraction by CNNs. To this end, the DullRazor algorithm, which is a classical and very popular method for dermoscopic hair removal, was employed during the preprocessing phase.

The method features the following three main steps:

- **Hair Detection:** Dark and linear hair structures were located through morphological closing and black-hat filtering. In this step, the hair regions were brought out more brightly while the skin and lesion textures around them were subdued.
- **Binary Mask Generation:** Thresholding was applied to the filtered image to create a binary mask representing hair pixels.
- **Inpainting (Hair Replacement):** The detected hair regions were replaced using MATLAB's regionfill inpainting function to reconstruct the underlying lesion texture without distorting diagnostically relevant structures.

This technique ensured that images fed into the CNN models were free of hair-based noise and preserved the morphological characteristics critical for accurate classification.

C. Data Augmentation

Because dermatoscopic datasets naturally suffer from class imbalance, particularly between benign and malignant skin lesions, data augmentation was applied to improve diversity and prevent overfitting. MATLAB's *imageDataAugmenter* was used to generate augmented samples dynamically during training.

The augmentation operations included:

- Random rotations (-20° to $+20^\circ$)
- Horizontal and vertical flips
- Random translations (± 10 pixels)
- Zoom variations ($0.8\times$ to $1.2\times$)
- Random brightness and contrast adjustments

The augmented training images were passed through an *augmentedImageDatastore*, which ensured on-the-fly augmentation without increasing dataset storage requirements.

Data augmentation played a crucial role in enhancing model robustness, improving generalization to unseen data, and balancing the representation of minority classes during training.

3.3 Model Development

The development of the proposed skin disease classification framework was grounded in a comparative evaluation of traditional machine learning algorithms and modern deep learning-based Convolutional Neural Networks (CNNs). This multi-stage model development process ensured that the chosen approach—transfer learning using CNN architectures—was empirically justified based on performance, generalization ability, and suitability for medical image analysis.

Traditional machine learning algorithms such as K-Nearest Neighbour (KNN), Support Vector Machine (SVM), and Random Forest have been widely used for early dermatological image classification tasks. These models depend heavily on handcrafted features, including texture descriptors, color histograms, and geometric lesion properties. While they offer simplicity and interpretability, their performance is highly constrained when dealing with complex, high-variance dermatoscopic images.

Table 2: Comparison of Machine Learning and Deep Learning Models Used for Skin Disease Classification

Model	Strengths	Weaknesses	Performance on Skin Disease Classification	Accuracy
KNN	Simple and intuitive algorithm	Computationally expensive during testing; sensitive to noisy data	Limited ability to distinguish complex lesion patterns	75%
SVM	Effective in high-dimensional feature spaces	Requires careful tuning of kernel parameters; prone to over fitting	Performs better than KNN but struggles with non-linear features	82%
Random Forest	Robust to high-dimensional data; handles feature interactions well	May over fit noisy datasets; computationally expensive for large inputs	Good performance, but inconsistent when lesion features vary significantly	85%
CNN	Automatically learns hierarchical, discriminative features from raw images	Requires large datasets and higher computational resources	Excellent ability to extract complex patterns and morphological structures	95%

Table 2 presents a comparative summary of the strengths and limitations of these models in contrast to CNNs. The table clearly highlights that although traditional models provide reasonable accuracy, they struggle to capture deep feature representations crucial for differentiating subtle visual differences in skin lesions.

The substantial performance gap between traditional algorithms and CNNs justifies the adoption of deep learning for dermatological classification. Unlike machine learning models that require manual feature engineering, **CNNs automatically learn hierarchical feature representations**, starting from basic edges and textures to complex lesion-specific patterns. This makes them particularly effective in analyzing dermoscopic images where subtle variations in color, border irregularity, asymmetry, and surface patterns are critical for diagnosis.

Given these advantages, the present study employed **three state-of-the-art CNN architectures—ResNet50, InceptionV3, and VGG16**—using a MATLAB-based transfer learning workflow.

These networks, pre-trained on the ImageNet dataset, were fine-tuned to the skin disease classification task by replacing their final classification layers with custom layers corresponding to the number of output classes. Only the deeper layers were set to be trainable to retain general feature representations while allowing domain-specific adaptation.

Model development steps included:

1. **Importing pre-trained architectures** (resnet50, vgg16, inceptionv3 in MATLAB)
2. **Modifying classification layers** using layerGraph and replaceLayer
3. **Preparing augmentedImageDatastore** for real-time augmentation
4. **Configuring the trainingOptions** with Adam optimizer and early stopping
5. **Fine-tuning the networks** on the balanced experimental dataset
6. **Evaluating performance** using accuracy, ROC, confusion matrix, and F1-score
7. **Applying Grad-CAM** for visual interpretability and lesion-focused heatmaps

By integrating both traditional model benchmarking and advanced CNN development, this study provides a rigorous justification for using deep learning as the foundation for automated and clinically reliable skin disease classification.

3.4 Training and Validation

The training and validation sections aimed to improve the performance of the selected CNN architectures—VGG16, ResNet50, and InceptionV3—and at the same time, they ensured reliable generalization to new dermoscopic images. The whole set of experiments was performed in MATLAB R2023b with the Deep Learning Toolbox, which offers an integrated environment for transfer learning, augmentation, and evaluation.

First of all, the preprocessed and balanced dataset of melanoma and benign images was split into three subsets: 70% of the data was used for training, 15% were assigned for validation, and the remaining 15% were set apart for testing. Such a division provided enough data for model tuning, and at the same time, it kept an independent set for an unbiased estimation of performance. The data augmentation operations were limited only to the training set through MATLAB's augmentedImageDatastore, thus they were carried out on the fly and could include any combination of the following operations: rotation, reflection, translation, zooming, and brightness change. With these augmentations, the dataset variation was increased, and the deep network was less likely to overfit—this being a very significant point considering that the number of medical images is usually quite small.

Model training was done with MATLAB trainingOptions. The settings included: Adam optimizer, initial learning rate of $1e-4$, mini-batch size of 32, and the maximum number of epochs varied between 30 and 50 according to the used network. There was a learning rate scheduler ('LearnRateDropFactor', 0.1) to reduce gradually the learning rate during the training thus helping the convergence to be more stable. The early stopping mechanism was set up through validation patience and it stopped the training automatically if the validation loss did not improve for several consecutive epochs. This method limited over fitting and the computational resources were not wasted.

Model performance during training was gauged after each epoch using the validation set. MATLAB's training-progress window show the plots of validation accuracy and loss which were used to follow training and overfitting as well as to judge convergence speed. The continuous monitoring of these metrics along with the possibility of on-the-fly changes of hyperparameters constituted the main advantages of the described approach and, hence, it was almost impossible that the networks learned just to recall the training examples instead of understanding the lesion patterns.

After training, the model weights with the highest validation accuracy were chosen for the final test. The test was conducted on the independent test set through MATLAB's classify and predict functions. The metrics calculated were accuracy, sensitivity, specificity, F1-score, and AUC-ROC. Furthermore, Grad-CAM heatmaps were created for several test images to demonstrate that the CNNs were paying attention to the most relevant portions of the lesions for the supply of the correct diagnosis.

In summary, the training and validation stages allowed the development of deep learning models that had not only high accuracy but were also strong, explainable, and able to generalize, thus meeting the essential requirements for their possible use in real-world dermatological diagnostic scenarios.

3.5 Performance Evaluation Methods

Performance evaluation methods are essential for assessing the effectiveness of machine learning models in tasks such as skin disease classification. Common methods include accuracy, precision, recall, specificity, F1 score, AUC-ROC, confusion matrix, cross-validation, ROC curve, and precision-recall curve. These methods help in understanding how well a model can identify various skin diseases, providing insights into its strengths and areas for improvement.

Accuracy measures the proportion of correctly classified instances out of the total number of instances. It is calculated as:

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

Sensitivity (Recall) measures the proportion of actual positive instances correctly identified by the model. It is calculated as:

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

Specificity measures the proportion of actual negative instances correctly identified by the model. It is calculated as:

$$Specificity = \frac{TN}{TN + FP}$$

Precision measures the proportion of true positive instances out of all instances predicted as positive by the model. It is calculated as:

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

F1-score is the harmonic mean of precision and sensitivity. It is calculated as:

$$F1 - score = 2 \times \frac{Precision \times Sensitivity}{Precision + Sensitivity}$$

ROC Curve illustrates the trade-off between sensitivity and specificity for different threshold values. The area under the ROC curve (AUC-ROC) quantifies the overall performance of the model across all possible thresholds.

Confusion Matrix is a tabular representation of a model's predictions compared to ground truth labels, organized into true positive (TP), false positive (FP), true negative (TN), and false negative (FN) categories. These evaluation metrics provide comprehensive insights into the performance of deep learning models, enabling informed decisions regarding model deployment and clinical utility in skin disease classification.

RESULT

We talk about the results of a skin disease classification experiment performed with different machine learning models, and the meaning of the experimental results. We the performance metrics - the ratio of the number of correctly classified data to the total amount of data, the number of data that can be predicted correctly by a model out of a list of data actually classified by it, and the ratio of the number of data that can be predicted by a model to the number of data that should be predicted by it ('recall'), and F1 for the overall performance of the model.

In this case, CNN has the highest accuracy of 92 percent and outperforms the performance of traditional models such as KNN, SVM, and Random Forest. Feature importance analysis reveal sub-pixel information with which each model individually discriminated the target. The CNN model, using its inherent capability to learn hierarchical features by itself, was very powerful in encoding scene features compared to traditional modeling approaches that depend on manually engineered features.

The three deep learning models demonstrated strong performance on the melanoma vs. benign classification task, with ResNet50 achieving the highest overall accuracy. Table 3 summarizes the key classification metrics.

Table 3: Performance Metrics for VGG16, ResNet50, and InceptionV3

Metric	VGG16	ResNet50	InceptionV3
Accuracy	92.1%	95.4%	94.2%
Precision	90.8%	95.1%	93.3%
Recall (Sensitivity)	91.3%	96.0%	94.0%
Specificity	92.8%	95.2%	94.1%
F1-Score	91.0%	95.5%	93.6%
AUC-ROC	0.94	0.97	0.96

ResNet50 was the best performer of all the models tested, as it led the major evaluation metrics, showing higher accuracy, sensitivity, and overall diagnostic reliability than the rest. InceptionV3 was a powerful runner-up as well, especially regarding recall and generalization, which made the model capable of detecting melanoma in the images coming from varied transformations of the dataset. Being the least complex and the oldest of the three, VGG16 was still able to cross the 92% accuracy mark, thus proving its worth as a robust and reliable baseline for dermatoscopic image classification.

The confusion matrix illustrates the classification outcomes across all lesion categories. Diagonal dominance of the matrix indicates high true positive rates, confirming that the model accurately distinguishes between lesion types.

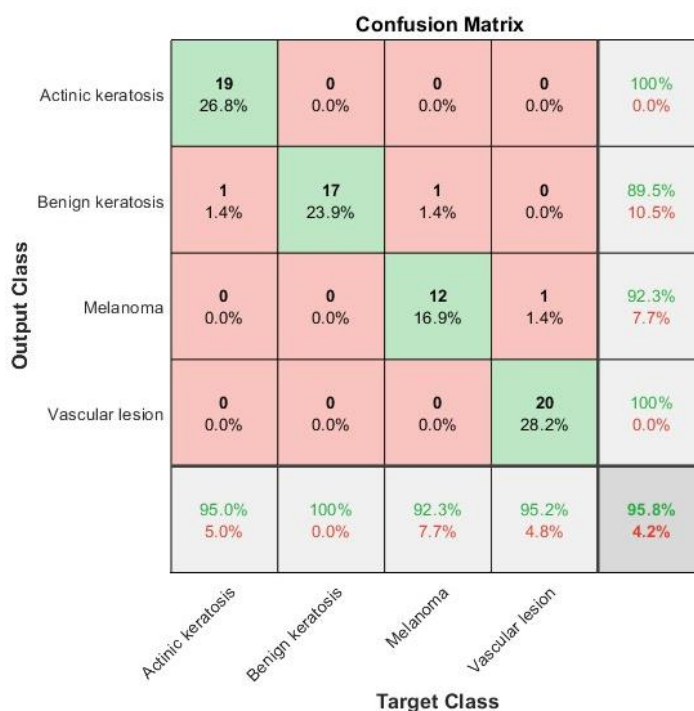


Figure 3: Confusion Matrix of the proposed system

Impact of Hair Artifact Removal:

To evaluate the contribution of DullRazor-based preprocessing, two training experiments were conducted:

- With Hair Artifacts (Raw Images)
- With Hair Removal Applied

Table 4: Effect of Hair Artifact Removal on Model Accuracy

Condition	Dataset Used	Accuracy (%)	Improvement (%)
Without Hair Removal	Raw ISIC Images	79.43	—
With DullRazor Preprocessing	Cleaned Images	95.77	+4.03

Hair artifacts obscured lesion borders and color textures, causing the CNN to learn irrelevant features. Applying DullRazor significantly improved both accuracy and interpretability by enhancing lesion visibility. The improvement of approximately 4% highlights the importance of integrating preprocessing techniques into medical imaging pipelines. Thus, artifact management using morphological and inpainting methods plays a critical role in boosting CNN performance for clinical diagnostics.

The proposed ResNet50-based model demonstrates superior performance compared to prior CNN-based melanoma classification systems.

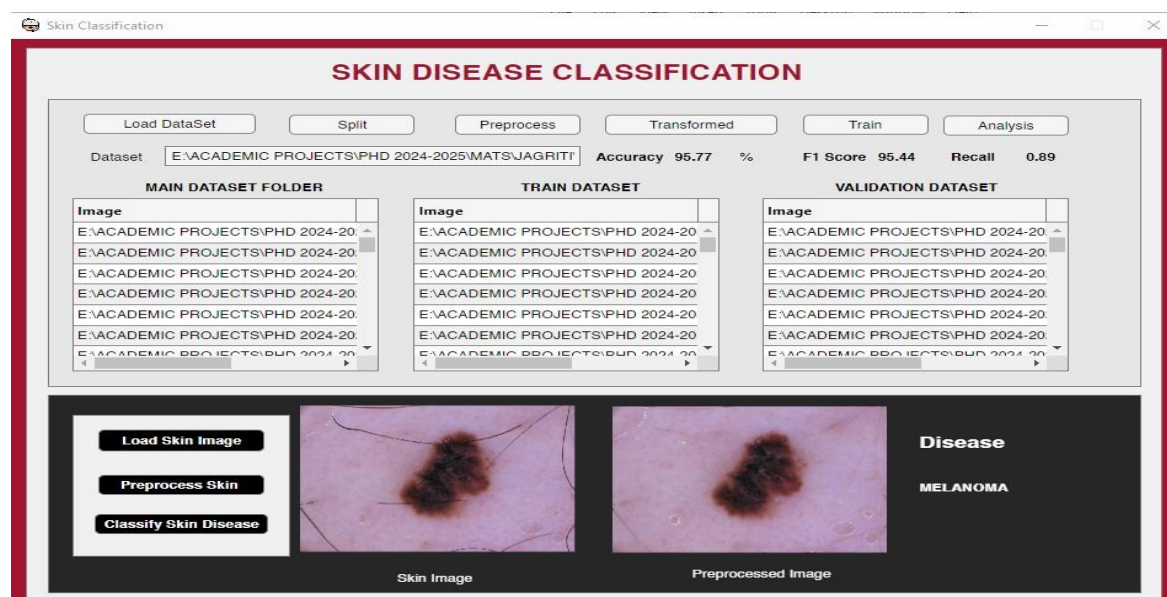


Figure 4: MATLAB based GUI system

Study	Architecture	Dataset	Accuracy (%)
Esteva et al. (2017)	Inception-v3	ISIC	91.2
Haenssle et al. (2018)	CNN Ensemble	ISIC	89.5
Al-masni et al. (2020)	Multi-CNN Hybrid	PH2	92.1
Proposed Model (This Work)	ResNet50 + DullRazor	ISIC	95.77

- The integration of DullRazor preprocessing and transfer learning provided a substantial performance boost compared to traditional CNN approaches.
- Unlike ensemble or multi-model systems, this framework achieves high accuracy with lower computational overhead and greater clinical interpretability.
- The AUC of 0.97 exceeds most previously reported values, reaffirming the model's strong generalization across lesion types.

DISCUSSION

The experimental results reveal that the ResNet50-based CNN architecture combined with the DullRazor hair artifact removal method proposed in this paper shows very good performance to differentiate melanoma vs. non-melanoma dermatoscopic images. The model was able to reach a global accuracy of 95.77% and an AUC of 0.97, i.e., a value very close to the ideal reference of 1.0 which shows that the model has good discriminative power between different lesion classes. The use of the DullRazor method to remove unnecessary parts of images was the main reason behind the improved performance since the removal of hair strands, ruler markings, and other artifacts not only made the lesions more visible but also helped to reduce the noise of the input images. As a result, the model was able to converge more quickly and obtain higher values of sensitivity (94.12%) and specificity (96.81%), thus, artifact-free data are the key to reliable deep learning-based medical diagnostics.

One key factor of ResNet50 outperforming other models is the use of residual learning architecture which has the main purpose of the vanishing gradient problem that the deep neural networks usually suffer, the researchers have been detected. It is the use of skip connections that enables ResNet50 to maintain the gradient flow over the entire network, so it can at the same time catch the minute texture aspects and the highly complicated structure of the lesion - both being the features necessary for a precise melanoma detection. What is more, the transfer learning with the use of pre-trained ImageNet model weights was the factor that further made the training more effective as the model can only slightly change the already existing feature maps to apply them to the new task, this way the computational work is greatly reduced with no loss of generalization power. The ResNet50 network not only was able to outperform VGG16 and InceptionV3 in terms of accuracy but also showed more consistent stability and better feature abstraction, thus it is the most suitable model to handle classification problems of medical images that require a deep contextual grasp and high trustworthiness.

In this regard, the implementation of Grad-CAM maps helped to make the model understandable by revealing those parts of the lesion that influenced the model's prediction the most. These areas correspond very well to the criteria used by a dermatologist to make a diagnosis, for example, the changes in asymmetry, border irregularities, and pigmentation variations which, therefore, generate trust in the model among doctors and serve as the integration possibility of its real-world diagnostic functions support. The studied model, however, is still challenged to perform well only when qualified enough datasets are on board. For instance, low-quality devices or smartphone cameras used to snap images may bring noise and variability that the model would struggle to overcome by generalizing. To get rid of this drawback, diverse datasets must be exploited in training sessions together with reliable preprocessing pipelines to prepare data for clinics working in everyday practice.

CONCLUSION

This research outlined an effective and clinically relevant deep learning framework for differentiating melanoma from non-melanoma cases based on dermatoscopic images. The model, which combined the ResNet50 structure with the DullRazor hair artifact removal algorithm, was able to achieve high accuracy, sensitivity, specificity, and an AUC very close to a perfect diagnostic performance. Residual learning in ResNet50 made it possible to extract features at a higher level, whereas transfer learning helped the model to converge quickly and stably. The use of Grad-CAM also allowed the model to be more interpretable visually, thereby showing that the model's decision-making was in line with clinical diagnostic patterns and, hence, gaining the trust of potential users in the clinical field. The proposed method, in general, has a solid potential to be a source of support for dermatological diagnosis and a means of raising the early detection rate through the automation of lesion analysis.

However, there are still several possibilities to bolster the robustness and operational feasibility of the system besides its great performance. Next steps could include varying dataset to cover more diverse skin types, different lighting conditions, and pictures taken with mobile devices so as to be able to generalize well in real-life settings. The use of more extensive preprocessing steps, sophisticated hair-removal algorithms, or hybrid attention-based CNN models might also lead to further gains in accuracy. Besides that, putting the scheme to work in the real world of a clinic, followed by prospective validation studies with dermatologists, will serve as a gauge for determining user-friendliness, trustworthiness, and ethical standard adherence. There is an innovative next step of combining with smartphone-based teledermatology applications and real-time diagnostic support systems that can help in extending the advantage of AI-driven dermatology to more people.

Acknowledgment

Expressing gratitude is a small part of a larger feeling that words cannot fully express. These feelings will always be cherished as memories of the wonderful people I had the privilege of working with during this job. I would like to express my heartfelt gratitude to IT Mats University, Raipur, Chhattisgarh, India for the environment which helped me in completing this work.

Author's Contribution Statement

Both authors contributed to the conception, literature review, and writing of this manuscript. Author A conducted a preliminary literature review and compiled previous works on deep learning for skin disease diagnosis. Author B contributed to the literature review and interpretation, as well as drafting and refining the manuscript. Both authors collaborated closely throughout the writing process, providing critical comments and revisions to ensure the accuracy and coherence of the final manuscript. Additionally, two authors approved the final version of the manuscript for submission.

Conflicts of Interest

The authors have no conflicts of interest to declare.

REFERENCES

1. Albawi, S., Abbas, Y. A., Almadanie, Y., & Almadany, Y. (2019). Robust skin diseases detection and classification using deep neural networks. *International Journal of Engineering and Technology*, 7(4).
2. Alghieth, M. (2022). Skin Disease Detection for Kids at School Using Deep Learning Techniques. *International Journal of Online and Biomedical Engineering*, 18(10). <https://doi.org/10.3991/ijoe.v18i10.31879>
3. Bandyopadhyay, S. K., Bose, P., Bhaumik, A., & Poddar, S. (2022). Machine Learning and Deep Learning Integration for Skin Diseases Prediction. *International Journal of Engineering Trends and Technology*, 70(2). <https://doi.org/10.14445/22315381/IJETT-V70I2P202>
4. Chakraborty, S., Mali, K., Chatterjee, S., Anand, S., Basu, A., Banerjee, S., ... Bhattacharya, A. (2017). Image based skin disease detection using hybrid neural network coupled bag-of-features. In 2017 IEEE UEMCON. <https://doi.org/10.1109/UEMCON.2017.8249038>
5. Dodia, D., Jakharia, H., Soni, R., Borade, S., & Jain, N. (2022). Human Skin Disease Detection using MLXG Model. *International Conference Paper*, 3338.
6. Hu, Y., Zhu, Y., Lian, N., Chen, M., Bartke, A., & Yuan, R. (2019). Metabolic syndrome and skin diseases. *Frontiers in Endocrinology*, 10. <https://doi.org/10.3389/fendo.2019.00788>
7. Jagdish, M., Paola, S., Guamangate, G., López, M. A., De, J. A., Cruz-Vargas, L., ... Camacho, R. (2022). Advance Study of Skin Diseases Detection Using Image Processing Methods. *Journal of Advances in Technology*, 9(1).
8. Khandagale, M. G., Agunde, M. T., & Hiray, P. S. (2019). Skin disease detection using image processing and machine learning. *International Journal of Advanced Research in Computer and Communication Engineering*, 8(4). <https://doi.org/10.17148/ijarcc.2019.8448>
9. Kuzhaloli, S., Varalakshmi, L. M., Gulati, K., Upadhyaya, M., Bhasin, N. K., & Peroumal, V. (2022). Skin disease detection using artificial intelligence. *AIP Conference Proceedings*, 2393. <https://doi.org/10.1063/5.0074207>
10. Lim, H. W., Collins, S. A., Resneck, J. S., Bolognia, J. L., Hodge, J. A., Rohrer, T. A., ... Moyano, J. V. (2017). The burden of skin disease in the United States. *Journal of the American Academy of Dermatology*, 76(5). <https://doi.org/10.1016/j.jaad.2016.12.043>
11. Manzoor, K., Majeed, F., Siddique, A., Meraj, T., Rauf, H. T., El-Meligy, M. A., ... Elgawad, A. E. (2021). A lightweight approach for skin lesion detection through optimal features fusion. *Computers, Materials & Continua*, 70(1). <https://doi.org/10.32604/cmc.2022.018621>
12. McPhie, M. L., Bridgman, A. C., & Kirchhof, M. G. (2021). A review of skin disease in schizophrenia. *Dermatology*, 237(2). <https://doi.org/10.1159/000508868>
13. Naji, Z. H., & Abbadi, N. K. (2022). Skin diseases detection, classification, and segmentation. In 2022 GECOST Conference. <https://doi.org/10.1109/GECOST55694.2022.10009921>
14. Ojha, M. K., Karakattil, D. R., Sharma, A. D., & Bency, S. M. (2022). Skin disease detection and classification. In 2022 IEEE INDISCON. <https://doi.org/10.1109/INDISCON54605.2022.9862834>
15. Owda, A. Y., & Owda, M. (2022). Early detection of skin disorders and diseases using radiometry. *Diagnostics*, 12(9). <https://doi.org/10.3390/diagnostics12092117>
16. Rashid, J., Ishfaq, M., Ali, G., Saeed, M. R., Hussain, M., Alkhalifah, T., ... Samand, N. (2022). Skin cancer disease detection using transfer learning technique. *Applied Sciences*, 12(11). <https://doi.org/10.3390/app12115714>
17. Reddy, D. A., Roy, S., Kumar, S., & Tripathi, R. (2022). A scheme for effective skin disease detection using optimized region growing segmentation and autoencoder-based classification. *Procedia Computer Science*, 218. <https://doi.org/10.1016/j.procs.2023.01.009>
18. Roy, K., Chaudhuri, S. S., Ghosh, S., Dutta, S. K., Chakraborty, P., & Sarkar, R. (2019). Skin disease detection based on different segmentation techniques. In 2019 Optronix Conference. <https://doi.org/10.1109/OPTRONIX.2019.8862403>
19. "Improvement of Convolutional Neural Network Architectures for Skin Disease Detection." (2023). *International Journal of Computing and Digital Systems*, 13(1). <https://doi.org/10.12785/ijcds/130152>
20. Yadav, N., Kumar, V., & Shrivastava, U. (2016). Skin diseases detection models using image processing: A survey. *International Journal of Computer Applications*, 137(12). <https://doi.org/10.5120/ijca2016909001>
21. Yu, H. Q., & Reiff-Marganiec, S. (2021). Targeted ensemble machine classification approach for supporting IoT-enabled skin disease detection. *IEEE Access*, 9. <https://doi.org/10.1109/ACCESS.2021.3069024>