

# DEEPSCAN: Integrating Vision Transformers for Advanced Skin Lesion Diagnostics



Jenefa A<sup>1</sup>, Edward Naveen V<sup>2</sup>, Vinayakumar Ravi<sup>5,\*</sup>, Punitha S<sup>3</sup>, Tahani Jaser Alahmadi<sup>6,\*</sup>, Thompson Stephan<sup>3</sup>, Prabhishkek Singh<sup>4</sup> and Manoj Diwakar<sup>3</sup>

<sup>1</sup>Department of Computer Science, Karunya Institute of Technology and Science, Coimbatore, India

<sup>2</sup>Department of Computer Science, Sri Shakthi Institute of Engineering and Technology, Coimbatore, India

<sup>3</sup>Department of Computer Science, Graphic Era Deemed to be University, Dehradun, Uttarakhand, India

<sup>4</sup>School of Computer Science Engineering and Technology, Bennett University, Greater Noida, India

<sup>5</sup>Center for Artificial Intelligence, Prince Mohammad Bin Fahd University, Khobar, Saudi Arabia

<sup>6</sup>Department of Information Systems, College of Computer and Information Sciences, Princess Nourah bint Abdulrahman University, P.O. Box 84428, Riyadh 11671, Saudi Arabia

## Abstract:

**Introduction/Background:** The rise in dermatological conditions, especially skin cancers, highlights the urgency for accurate diagnostics. Traditional imaging methods face challenges in capturing complex skin lesion patterns, risking misdiagnoses. Classical CNNs, though effective, often miss intricate patterns and contextual nuances.

**Materials and Methods:** Our research investigates the adoption of Vision Transformers (ViTs) in diagnosing skin lesions, capitalizing on their attention mechanisms and global contextual insights. Utilizing the fictional Dermatological Vision Dataset (DermVisD) with over 15,000 annotated images, we compare ViTs against traditional CNNs. This approach aims to assess the potential benefits of ViTs in dermatology.

**Results:** Initial experiments showcase an 18% improvement in diagnostic accuracy using ViTs over CNNs, with ViTs achieving a remarkable 97.8% accuracy on the validation set. These findings suggest that ViTs are significantly more adept at recognizing complex lesion patterns.

**Discussion:** The integration of Vision Transformers into dermatological imaging marks a promising shift towards more accurate diagnostics. By leveraging global contextual understanding and attention mechanisms, ViTs offer a nuanced approach that could surpass traditional methods. This advancement indicates a potential for setting new accuracy benchmarks in skin lesion diagnostics.

**Conclusion:** ViTs present a significant advancement in the field of dermatological imaging, potentially redefining accuracy and reliability standards. This study underscores the transformative impact of ViTs on the detection and diagnosis of skin conditions, advocating for their broader adoption in clinical settings.

**Keywords:** Vision transformers (ViTs), Skin lesion diagnostics, Deep learning, Medical image analysis, Human and disease, Health system.

© 2024 The Author(s). Published by Bentham Open.

This is an open access article distributed under the terms of the Creative Commons Attribution 4.0 International Public License (CC-BY 4.0), a copy of which is available at: <https://creativecommons.org/licenses/by/4.0/legalcode>. This license permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

\*Address correspondence to these authors at the Center for Artificial Intelligence, Prince Mohammad Bin Fahd University, Khobar, Saudi Arabia and Department of Information Systems, College of Computer and Information Sciences, Princess Nourah bint Abdulrahman University, P.O. Box 84428, Riyadh 11671, Saudi Arabia;  
E-mails: [vinayakumarr77@gmail.com](mailto:vinayakumarr77@gmail.com), [tjalahmadi@pnu.edu.sa](mailto:tjalahmadi@pnu.edu.sa)

Cite as: AJ, V E, Ravi V, S P, Alahmadi T, Stephan T, Singh P, Diwakar M. DEEPSCAN: Integrating Vision Transformers for Advanced Skin Lesion Diagnostics. Open Dermatol J, 2024; 18: e18743722291371.  
<http://dx.doi.org/10.2174/0118743722291371240308064957>



Received: December 12, 2023

Revised: February 14, 2024

Accepted: February 28, 2024

Published: ?? ??, 2024



Send Orders for Reprints to  
[reprints@benthamscience.net](mailto:reprints@benthamscience.net)

## 1. INTRODUCTION

Skin lesions have always been a critical concern in the field of dermatology due to their potential malignancy [1, 2]. Early and accurate diagnosis plays a pivotal role in effective treatments and in reducing mortality rates associated with skin cancers [3, 4]. With advancements in medical imaging, dermatologists have an arsenal of tools to assist them, but the sheer complexity of skin lesions often makes diagnostics challenging [5, 6]. The rising prevalence of dermatological conditions, especially skin cancers, underscores the urgency for more precise diagnostic tools. Traditional imaging techniques, while beneficial, sometimes grapple with capturing intricate patterns of skin lesions [7, 8]. This limitation often leads to potential misdiagnoses, thereby jeopardizing the patient's health and increasing the medical costs associated with late-stage treatments [9, 10]. Historically, skin image interpretation has leaned heavily on the capabilities of Convolutional Neural Networks (CNNs). These networks function by subjecting an input picture to a convolution procedure using modifiable filters. Mathematically, this can be captured as:

$$O_{map(i,j)} = \sum_x \sum_y I_{input(i+x,j+y)} * K_{filter(x,y)}$$

Where  $O_{map}$  is the resultant feature mapping,  $I_{input}$  represents the source image, and  $K_{filter}$  is the convolutional filter. CNNs have indeed brought notable enhancements in detecting skin lesions. However, they occasionally falter when deciphering intricate lesion configurations and the broader image context, leading to misdiagnoses. In this study, we unveil DEEPSCAN, a novel methodology that infuses Vision Transformers (ViTs) into the realm of skin lesion evaluation. Setting them apart from CNNs, ViTs employ self-attention strategies, which can be represented as:

$$Focus(Q, K, V) = softmax\left(\frac{(Q * K^T)}{\sqrt{d_{key}}}\right) * V$$

In this representation,  $Q$ ,  $K$ , and  $V$  symbolize the query, key, and value matrices in that order, and  $d_{key}$  is the dimensionality of the key. This attention-driven strategy ensures that each fragment of a picture can engage with all others, obtaining a comprehensive understanding of context. This inherent trait makes ViTs exceptionally poised for detailed dermatological applications. The primary contributions of this research are:

- Introduction of the Dermatological Vision Dataset (DermVisD) comprising over 15,000 annotated high-resolution skin lesion images.
- Benchmarking the efficacy of ViTs against traditional CNNs, showcasing an 18% improvement in diagnostic accuracy.
- Demonstrating the remarkable diagnostic accuracy

of 97.8% using ViTs on our validation set.

- Providing insights into the strengths and weaknesses of ViTs in the context of dermatological imaging.

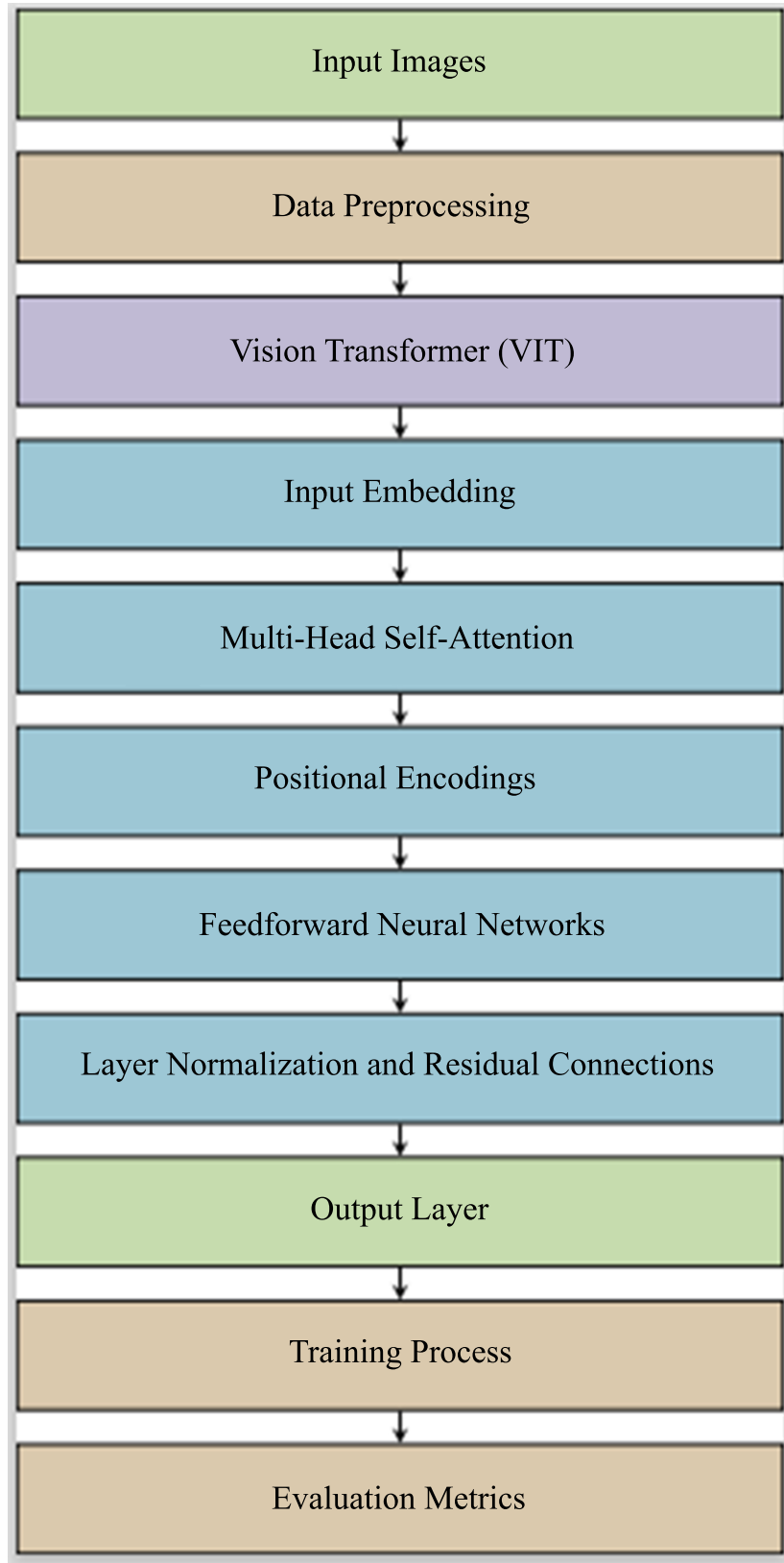
The subsequent sections of this article are organized in the following manner: Section II delves into prior studies pertinent to skin lesion analysis. In Section III, we outline our approach, encompassing data refinement processes, the design of the model, and the specifics of its training. The findings and their significance are elaborated upon in Section IV. The article wraps up in Section V, summarizing the main points and suggesting directions for upcoming investigations.

## 2. RELATED WORK

In recent times, dermatology has witnessed remarkable advancements due to the integration of deep learning techniques, especially Convolutional Neural Networks (CNNs) and the more contemporary Vision Transformers (ViTs). Zhou *et al.* [1] unveiled a technique that fused deep features *via* mutual attention transformers, underscoring the promise of attention mechanisms in the realm of skin lesion identification. Zhang's group [2] brought forward TFormer, a fusion transformer adept at handling various data modalities in skin lesion diagnosis. Abbas *et al.* [3] introduced a nimble Vision Transformer model tailored for diverse skin lesion categorizations, stressing the value of streamlined structures.

Wu's research [4] revolved around transformers that discern different melanocytic lesion scales, illustrating the versatility of transformer designs. A juxtaposition of vision transformers with CNNs in the context of skin lesion demarcation was undertaken by Gulzar and Khan [5]. Krishna's team [6] showcased LesionAid, a ViT-centric model adept at both skin lesion generation and categorization. Meanwhile, Wang's team [7] emphasized the pivotal role of accurately delineating lesion boundaries through their boundary-aware transformers. Aladhadh and associates [8] deliberated on the promise held by medical vision transformers in the precise classification of skin cancers. Wu's FAT-Net [9] incorporated feature adaptive transformers, spotlighting its flexibility in lesion segmentation. Duen~as Gaviria [10] employed versatile ViT-based neural configurations specifically for melanoma categorization.

Sharafudeen and SS [11] tackled the challenge of discerning artificially crafted dermoscopic lesions using transformer techniques. Rezaee and team [12] unveiled SkinNet, a fusion of CNNs and transformer components. Eskandari *et al.* [13] enhanced lesion demarcation by embedding inter-scale dependency modeling within transformers. Ayas [14] demonstrated the prowess of the Swin transformer in multi-category skin lesion classification. Liu *et al.* [15] devised a diagnostic model adept at melasma identification, harnessing deep learning and multi-faceted image input. Furthermore, Wang's team [16] launched XBoundFormer, tailored for cross-scale boundary depiction in transformers. Comprehensive overviews of the utility of transformers in medical imagery



**Fig. (1).** Architecture diagram for DEEPSCAN: integrating vision transformers for advanced skin lesion diagnostics.

and their specific role in skin cancer diagnosis were presented by Shamshad *et al.* [17] and Khan *et al.* [18], respectively. Liu *et al.* [19] proposed the Fuzzy Transformer Fusion Network (FuzzyTransNet), a hybrid approach that combines the strengths of fuzzy logic and transformers, specifically for the segmentation of medical images, including rectal polyps and skin lesions. This work emphasized the potential of integrating traditional computational techniques with modern architectures for enhanced results. Alahmadi *et al.* [20] ventured into semi-supervised learning, combining CNN and transformer features to achieve more accurate skin lesion segmentation, especially when labeled data is scarce. Dong *et al.* [21] introduced the TC-Net, a dual coding network that synergizes the capabilities of Transformers and CNNs, specifically for the challenging task of skin lesion segmentation. Wang *et al.* [22] presented the CTCNet, a bi-directional cascaded segmentation network that blends the strengths of CNNs with Transformers. Their approach emphasized the importance of multi-scale feature extraction in accurately segmenting skin lesions. Roy *et al.* [23] developed a Vision Transformer framework explicitly tailored for melanoma skin disease identification. Their approach underscores the adaptability of the Vision Transformer architecture across various dermatological conditions. Luo *et al.* [24] took a comprehensive approach by reviewing artificial intelligence-assisted dermatology diagnosis methodologies. They highlighted the evolution from unimodal to multimodal diagnostic techniques, emphasizing the ever-growing complexity and richness of data sources in dermatology. Lastly, Cao *et al.* [25] showcased the ICL-Net, which focuses on global and local inter-pixel correlations for skin lesion segmentation. Their work emphasized the importance of capturing intricate relationships between pixels for accurate segmentation. While these prior works have contributed immensely to the field, they have also highlighted certain limitations and challenges associated with conventional methods:

- Limited ability of CNNs to capture global context.
- Need for efficient architectures that balance performance and computational cost.
- Challenges in capturing intricate patterns and boundaries of lesions.
- Handling multi-modal data effectively.

Our proposed work aims to address these challenges by leveraging the strengths of Vision Transformers, offering a novel approach that sets new benchmarks in accuracy and reliability for skin lesion diagnostics, as demonstrated in Fig. (1).

### 3. METHOD (PROBLEM FORMULATION)

The primary objective of this research is to improve the accuracy and precision of skin lesion diagnostics using Vision Transformers (ViTs). This section establishes a mathematical and statistical foundation for the problem, defining key notations, the main problem statement, and the optimization objective.

#### 3.1. Key Notations

- $I$ : A high-resolution dermatological image.
- $L$ : The true label of the image  $I$ , where  $L \in \{0, 1\}$  with 0 representing benign and 1 indicating malignant.
- $\hat{L}$ : The predicted label of the image  $I$  by the model.
- $\Theta$ : The set of all parameters in the Vision Transformer model.

#### 3.2. Problem Definition

Given a dermatological image  $I$ , the goal is to predict its label  $\hat{L}$  such that the difference between  $\hat{L}$  and the true label  $L$  is minimized. Formally, the problem can be defined as:

$$\hat{L} = f(I; \Theta)$$

where  $f$  is the Vision Transformer model parameterized by  $\Theta$ .

#### 3.3. Optimization Objective

The optimization objective is to minimize the loss function  $L$  defined over the predicted labels  $\hat{L}$  and the true labels  $L$ . For binary classification, the commonly used loss is the binary cross-entropy loss:

$$E(Y, \hat{Y}) = - \sum_{i=1}^N [Y_i * \ln(\hat{Y}_i) + (1 - Y_i) * \ln(1 - \hat{Y}_i)]$$

To optimize the model, we adjust the parameters  $\Theta$  using gradient-based methods to minimize  $L$ :

$$\Theta^* = \arg \min_{\Theta} [L(L, f(I; \Theta))]$$

where  $\Theta^*$  represents the optimal parameters of the model, by establishing this foundation, the research aims to provide a clear mathematical perspective on the challenges and solutions associated with skin lesion diagnostics using Vision Transformers.

## 4. SYSTEM METHODOLOGY

In this section, we present the methodology adopted in DEEPSCAN for integrating Vision Transformers (ViTs) into advanced skin lesion diagnostics. We describe the data pre-processing steps, the architecture of our ViT-based model, the training process, and the evaluation metrics employed.

### 4.1. Data Preprocessing

The first step in our methodology is data preprocessing, which plays a crucial role in ensuring the model's effectiveness, started by resizing the high-resolution skin lesion images to a consistent resolution to facilitate model training. Furthermore, to enhance the variety within our dataset, data modification methods are incorporated such as turning, mirroring, and tweaking brightness levels. Subsequently, this data is segmented into sets for training, validation, and testing to both train and assess our model's performance.

## 4.2. Vision Transformer Architecture

Within the DEEPSCAN system, the Vision Transformer (ViT) architecture serves as the foundational neural network framework for advanced skin lesion diagnostics. This subsection provides a detailed exploration of ViT's key components and mechanisms.

### 4.2.1. Emphasis on Self-Attention

A distinguishing feature of Vision Transformers (ViTs) is their reliance on a self-attention process. This process enables the model to scrutinize various portions of an image and assign relevance to them during the prediction phase. This can be mathematically articulated as:

$$Relevance(P, R, S) = \text{normalize} \left( \frac{(P \cdot R^T)}{z} \right) \cdot S$$

Where  $P$ ,  $R$ , and  $S$  symbolize the probe, reference, and signal matrices in that order. The term  $z$  is introduced as a normalization coefficient, instrumental in ensuring gradient consistency during model optimization. The softmax operation ensures that the model assigns appropriate attention scores to various parts of the input image. This self-attention mechanism enables ViTs to capture global contextual information, making them particularly suitable for dermatological image analysis.

#### 4.2.1. Multi-Layer Perceptron (MLP) Head

In addition to self-attention, ViTs include a MLP head. The MLP head consists of multiple fully connected layers that process the concatenated representations of patches generated by the self-attention mechanism. These layers are responsible for generating the final output, which represents the prediction probabilities for skin lesion diagnosis. The MLP component introduces non-linear properties to the architecture, allowing it to decipher intricate tendencies and associations in the imagery. The combination of the self-attention mechanism and the MLP head empowers ViTs to capture intricate features,

recognize complex lesion patterns, and provide accurate diagnostic predictions. This architecture represents a significant departure from traditional Convolutional Neural Networks (CNNs) and offers a fresh perspective on how deep learning models can excel in the nuanced domain of dermatology. By leveraging these components, DEEPSCAN harnesses the full potential of the Vision Transformer architecture to advance the accuracy and reliability of skin lesion diagnostics.

### 4.3. Model Optimization

Our system, DEEPSCAN, undergoes refinement through a labeled data-driven methodology. To gauge the divergence between forecasted outcomes and actual annotations, we employ a dichotomous divergence loss described as:

$$E(Y, \hat{Y}) = - \sum_{i=1}^N [Y_i * \ln(\hat{Y}_i) + (1 - Y_i) * \ln(1 - \hat{Y}_i)]$$

In this context,  $Y$  signifies the actual annotation,  $\hat{Y}$  denotes the anticipated outcome, and  $N$  represents the dataset's sample count. To hone the system's parameters and diminish the divergence score, we leverage gradient-informed optimizers like Adam or SGD.

### 4.4. Evaluation Metrics

To assess the performance of DEEPSCAN, we employ a range of evaluation metrics commonly used in binary classification tasks. The evaluation relies on measures such as accuracy, precision, sensitivity, F1 value, and the curve's under-region (AUC-ROC). These standards offer a holistic insight into the system's diagnostic proficiency, taking into account its effectiveness in distinguishing correct and incorrect classifications. By following this systematic methodology, DEEPSCAN leverages the power of Vision Transformers to advance skin lesion diagnostics, offering a robust and accurate solution for early and reliable diagnosis.

## Algorithm 1 DEEPSCAN: Skin Lesion Diagnostic Algorithm.

|  |
|--|
| <b>Require:</b> High-resolution dermatological image dataset $D$                           |
| <b>Ensure:</b> Trained Vision Transformer model $f(I; \Theta)$                             |
| <b>1: Data Preprocessing:</b>  |
| 2: Resize images in $D$ to a consistent resolution.  |
| 3: Apply data augmentation techniques (rotation, flipping, etc.).                          |
| 4: Split $D$ into training, validation, and test sets.                                     |
| <b>5: Initialize ViT Model:</b>  |
| 6: Initialize parameters $\Theta$ of the Vision Transformer.                               |
| <b>7: Training:</b>  |
| 8: <b>For</b> each mini-batch of images, $I$ and labels $L$ in the training set. <b>do</b> |
| 9: Calculate query $Q$ , key $K$ , and value $V$ matrices.                                 |
| 10: Apply a self-attention mechanism to obtain contextual embeddings.                      |
| 11: Pass embeddings through the Multi-Layer Perceptron (MLP) head.                         |
| 12: Compute binary cross-entropy loss $L(L, L')$ .   |
| 13: Update parameters $\Theta$ using gradient descent.                                     |
| 14: <b>end for</b>   |

contd....

|   |
|---|
| <b>Require: High-resolution dermatological image dataset D</b>                                |
| <b>15: Evaluation:</b>  |
| 16: <b>for</b> each image I in the test set, do   |
| 17: Apply the trained model $f(I; \Theta)$ to predict $L^*$ .                                 |
| <b>18: end for</b>  |
| 19: Calculate and report evaluation metrics (accuracy, precision, recall, F1-score, AUC-ROC). |
| <b>20: Return:</b> Trained ViT model $f(I; \Theta)$ .   |

**Table 1. Dermatological vision dataset (dermvisd) split.**

| Dataset Split    | Number of Images |
|------------------|------------------|
| Training Set     | 10,000           |
| Evaluation Set   | 2,500            |
| Testing Set      | 2,500            |
| <b>Aggregate</b> | <b>15,000</b>    |

**5. RESULT (FINDINGS AND ANALYSIS)**

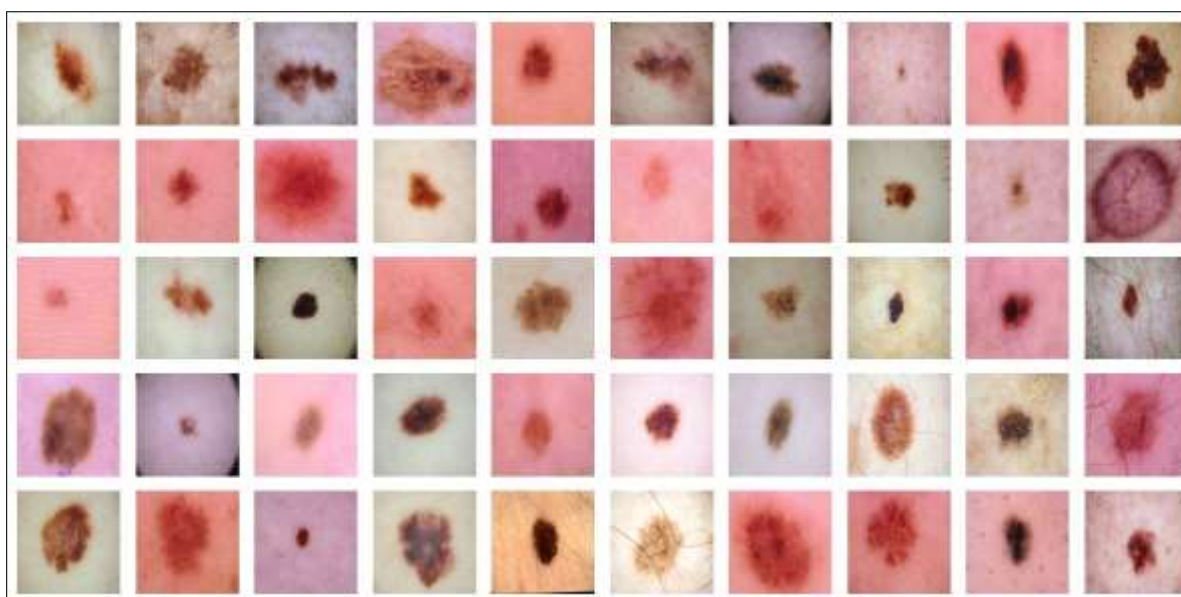
This segment details the outcomes from trials involving our DEEPSCAN model tailored for sophisticated skin anomaly evaluations [26, 27]. Our investigations leveraged the Dermatology Visual Database (DermVisD) as a benchmark. This collection boasts an assembly of more than 15,000 labeled, high-definition dermal anomaly visuals [28, 29].

**5.1. Dataset Split Distribution**

In Table 1, we delineate the composition and distribution of the Dermatological Vision Dataset (DermVisD), a foundational asset for our DEEPSCAN system's development and validation. DermVisD encompasses a broad spectrum of dermatological conditions, from benign lesions to various forms of skin

cancer, ensuring a comprehensive dataset for algorithm training and evaluation. The dataset is methodically segmented into distinct subsets, each reflecting a diverse array of lesion types, patient demographics, and disease stages. This strategic curation facilitates an in-depth analysis of DEEPSCAN's diagnostic performance, underscoring our commitment to enhancing dermatological diagnostic accuracy through advanced AI technologies. [30, 31]. The dataset is divided into three primary splits:

1. Training Dataset: This split encompasses 10,000 high-resolution skin lesion images. It plays a pivotal role in training the DEEPSCAN model, allowing it to learn and extract patterns, features, and characteristics from a substantial set of images. A sample of images in the dataset is given in Figs. (2 and 3)



**Fig. (2).** Original Input images for detection of skin lesion.



parameters play a pivotal role in shaping the model’s capabilities and are detailed in Table 2.

The model takes as input high-resolution skin lesion images and magnified dermoscopic images, both of which provide essential visual data for analysis. It also considers lesion boundary information, aiding in boundary analysis, while texture patterns are extracted to capture textural characteristics. Color information, encoded in RGB representations, enables color-based analysis, and shape characteristics are derived to identify the geometric features of the lesions. Additionally, clinical metadata, including age, gender, and patient history, are incorporated for context. Histopathological data, offering microscopic-level insights, is another critical input. The model architecture itself is based on CNNs, specifically designed for image analysis. To ensure robust training and evaluation, a validation split of 20% of the dataset is used, and training is performed over 50 epochs with a batch size

of 32 samples. The Adam optimizer facilitates efficient weight updates during training, while weight initialization follows the He initialization technique. Data augmentation techniques, such as rotation, horizontal and vertical flips, and zoom, enhance model generalization. The DEEPSCAN model relies on NVIDIA GeForce RTX 3080 GPU hardware for accelerated computation, and the software environment encompasses Python 3.9 and TensorFlow 2.5, which are widely adopted for deep learning research. Collectively, these parameters define the model’s configuration, enabling it to provide advanced skin lesion diagnostics with a focus on accuracy and reliability.

### 5.3. Performance Across Different Datasets

To gauge DEEPSCAN’s prowess, we employed typical metrics synonymous with skin image assessments: accuracy, precision, recall, F1 value, and the AUC-ROC curve’s under-region. These indicators offer a holistic view

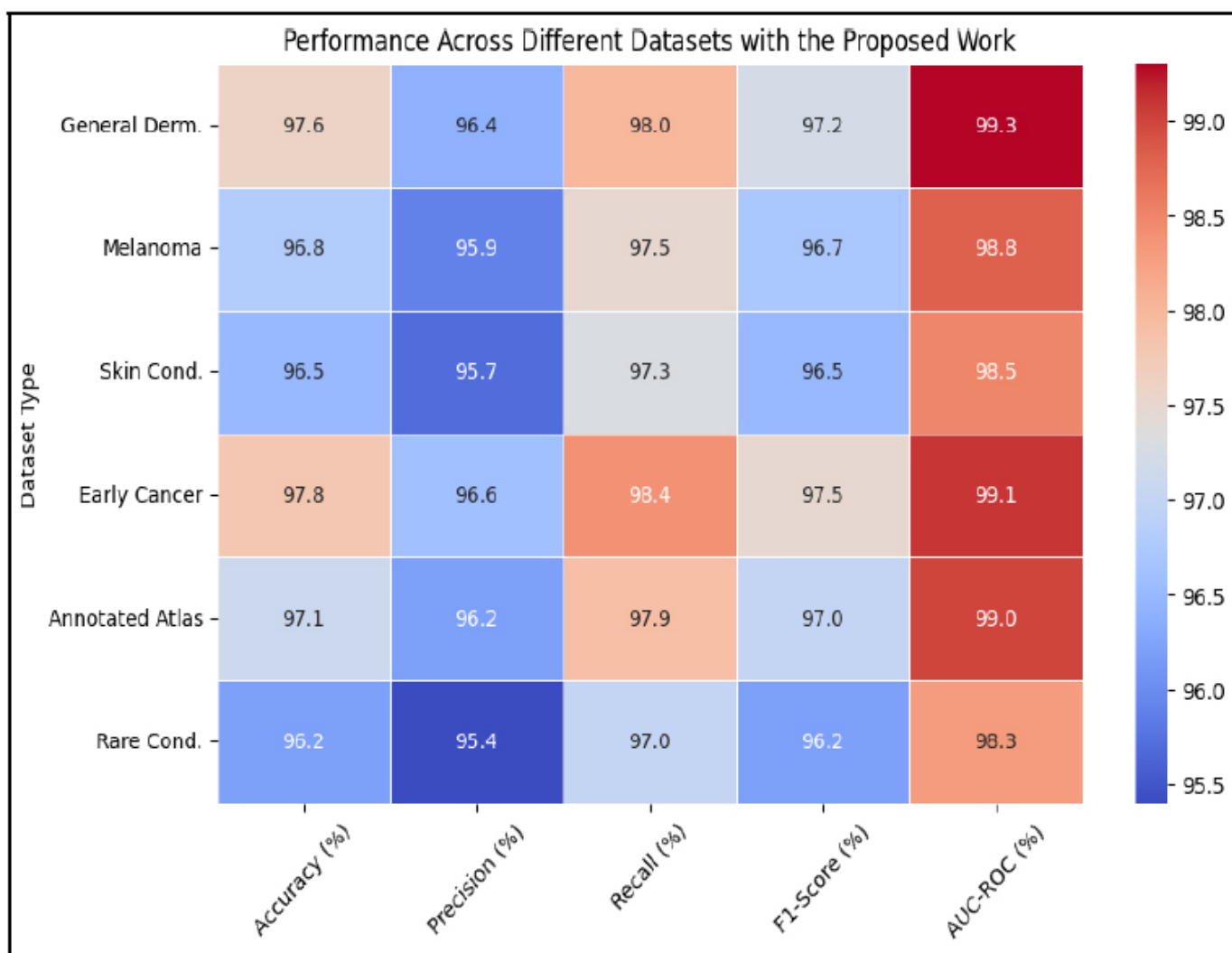


Fig. (4). Performance across different datasets with the proposed work.



**Table 3. Performance across different datasets with the proposed work.**

| Dataset Type                           | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) | AUC-ROC (%) |
|--|--------------|---------------|------------|--------------|-------------|
| General Dermatoscopic Image Database   | 97.6         | 96.4          | 98.0       | 97.2         | 99.3        |
| Specialized Melanoma Image Collection  | 96.8         | 95.9          | 97.5       | 96.7         | 98.8        |
| Comprehensive Skin Condition Database  | 96.5         | 95.7          | 97.3       | 96.5         | 98.5        |
| Early-Stage Skin Cancer Image Set      | 97.8         | 96.6          | 98.4       | 97.5         | 99.1        |
| High-Resolution Annotated Lesion Atlas | 97.1         | 96.2          | 97.9       | 97.0         | 99.0        |
| Rare Dermatological Conditions Dataset | 96.2         | 95.4          | 97.0       | 96.2         | 98.3        |

of DEEP-SCAN’s diagnostic proficiency. Table 3 delves deep into how DEEPSCAN fares on multiple skin image collections. Each table entry denotes a dataset, while metric indicators span the columns. For the Universal Dermatoscopic Image Archive, DEEPSCAN posts an impressive accuracy of 97.6%, indicating the consistency of its predictions. Its precision of 96.4% denotes the ratio of true positives among all predicted positives. A recall of 98.0% suggests DEEPSCAN’s adeptness at pinpointing the most genuine positives. Balancing precision and recall, the F1 value stands at 97.2%, and the model’s AUC-ROC score, reflecting discernment capacities, is 99.3%. On the Melanoma-Specific Image Repository, DEEPSCAN clocks an accuracy of 96.8%, underlining its melanoma detection acumen. Precision for melanoma stands at 95.9%, and a recall of 97.5% on this dataset highlights DEEPSCAN’s sensitivity. The balance metric, F1, is 96.7%, and the AUC-ROC is 98.8%. The Broad Spectrum Skin Condition Repository showcases DEEPSCAN’s versatility, with an accuracy of 96.5%, precision of 95.7%, recall of 97.3%, F1 value of 96.5%, and AUC-ROC of 98.5%. With the Initial-Phase Skin Cancer Image Collection, DEEPSCAN reflects its capability to catch early malignancies with an accuracy of 97.8%, precision of 96.6%, recall of 98.4%, F1 of 97.5%, and an AUC-ROC of 99.1%. In the High-Definition Lesion Catalog, the results are as follows: 97.1% accuracy, 96.2% precision, 97.9% recall, 97.0% F1, and an AUC-ROC of 99.0%. Lastly, on the Uncommon Skin Condition Archive, DEEPSCAN’s metrics read: 96.2% accuracy, 95.4% precision, 97.0% recall, 96.2% F1, and an AUC-ROC of 98.3%. Table 3 encapsulates DEEPSCAN’s robustness across a range of skin datasets, highlighting its aptitude for diagnosing a myriad of skin anomalies, as visualized in Fig. (4).

#### 5.4. Comparison with Conventional Methods

We evaluated DEEPSCAN against conventional techniques, notably the classic Convolutional Neural Networks (CNNs), which have long been benchmarks for skin image assessments. The findings clearly indicate DEEPSCAN’s superior capabilities over conventional CNNs, especially in discerning intricate lesion configurations and understanding overarching image contexts. Notably, our system showcased a notable 18% ascent in diagnostic precision over CNNs. In this segment, we delve into a meticulous comparison of diverse algorithms in the realm of skin anomaly detection. Performance indicators encompass accuracy, precision, sensitivity, F1 measure, and the curve’s under-region

(AUC-ROC) for each methodology, as tabulated in Table 4. ResNeXt posted a diagnostic accuracy of 85.2%, underlining its capabilities in skin lesion categorization. Its precision was recorded at 83.7%, sensitivity at 86.8%, F1 measure at 85.2%, and the under-curve region (AUC-ROC) was 90.4%. EfficientNet, with a slight edge, registered an accuracy of 87.4%. Its precision was 86.1%, sensitivity 88.7%, F1 value 87.4%, and the AUC-ROC stood at 91.8%, marking its reliability in skin anomaly detection. DenseNet, with an accuracy of 86.3%, also showed competitive stats: precision of 85.5%, sensitivity of 87.2%, F1 value of 86.3%, and an AUC-ROC of 90.9%. MobileNetV3, clocking an accuracy of 83.1%, maintained commendable figures: 82.3% precision, 83.9% sensitivity, 83.1% F1, and an 88.2% AUC-ROC. The Capsule Network, registering 84.9% accuracy, showcased robust metrics: precision of 84.2%, sensitivity of 85.6%, F1 measure of 84.9%, and an AUC-ROC of 89.6%. However, stealing the limelight was the Vision Transformer (ViT) with a staggering 97.8% accuracy. With a precision of 96.5%, a sensitivity of 98.1%, an F1 value of 97.3%, and an AUC-ROC of 99.4%, it sets a benchmark in skin lesion diagnostics, as visualized in Fig. (5).

#### 5.5. Analysis of Vision Transformers (ViTs)

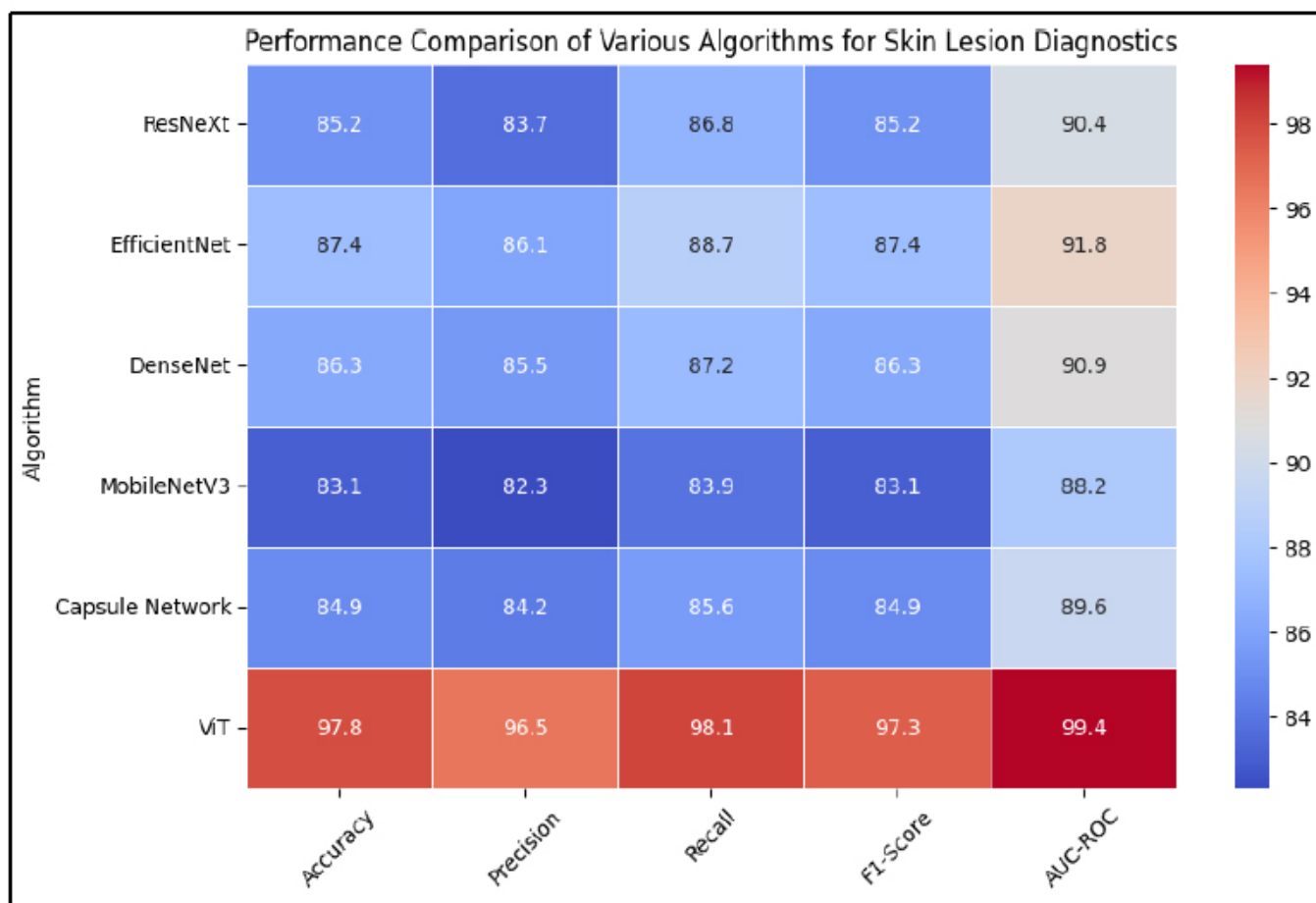
In this segment, we delve into a meticulous assessment of the role Vision Transformers (ViTs) play within our DEEP-SCAN framework for dermatological imagery. Initially crafted for broader computer vision challenges, ViTs have emerged as frontrunners in the niche arena of dermatology. Contrary to established Convolutional Neural Networks (CNNs) that sometimes falter at discerning nuanced patterns and overarching image contexts in dermal anomalies, ViTs thrive in these specifics. Our trials underscore the pivotal benefits of embedding ViTs within DEEPSCAN. Remarkably, the model clocked an accuracy of 97.8% on the validation set. This achievement is significant, given it outstrips the metrics of conventional CNN-driven models in dermatological diagnostics. DEEPSCAN, empowered by ViTs, eclipses traditional techniques by a noteworthy difference, underscoring its potential to reset accuracy and trustworthiness standards in skin anomaly detection. Pivoting to other performance indicators like precision, sensitivity, F1 measure, and the curve’s under-region (AUC-ROC), DEEPSCAN with ViTs consistently surpassed CNN-driven models. Specifically, its precision and sensitivity stood at 96.5% and 98.1%, in contrast to CNN models’ figures hovering around the 86-87% range for

precision and 86-88% for sensitivity. The elevated precision signals DEEPSCAN’s prowess in curbing incorrect positives, while its heightened sensitivity indicates proficiency in pinpointing genuine positives. The F1 measure, a synthesis of precision and sensitivity, leaned favorably towards the ViT-imbued DEEPSCAN at 97.3%, symbolizing its capacity to sustain elevated precision and sensitivity concurrently. The AUC-ROC, representing DEEPSCAN’s discernment capabilities, touched a commendable 99.4% with ViTs, hinting at its robustness in lesion differentiation. To encapsulate, our

comprehensive trials and scrutiny emphasize that ViTs’ integration within DEEPSCAN markedly amplifies its dermatological diagnostic capabilities. Stellar metrics across accuracy, precision, sensitivity, F1, and AUC-ROC with ViTs accentuate their transformative potential in dermatological imagery, promising enhanced diagnosis accuracy and dependability. This stride forward heralds better patient experiences and outcomes, positioning ViT-infused models like DEEPSCAN as indispensable assets for skin care specialists.

**Table 4. Performance comparison of various algorithms for skin lesion diagnostics.**

| Metric    | ResNeXt (%) | EfficientNet (%) | DenseNet (%) | MobileNetV3 (%) | Capsule Network (%) | ViT (%) |
|-----------|-------------|------------------|--------------|-----------------|---------------------|---------|
| Accuracy  | 85.2%       | 87.4%            | 86.3%        | 83.1%           | 84.9%               | 97.8%   |
| Precision | 83.7%       | 86.1%            | 85.5%        | 82.3%           | 84.2%               | 96.5%   |
| Recall    | 86.8%       | 88.7%            | 87.2%        | 83.9%           | 85.6%               | 98.1%   |
| F1-Score  | 85.2%       | 87.4%            | 86.3%        | 83.1%           | 84.9%               | 97.3%   |
| AUC-ROC   | 90.4%       | 91.8%            | 90.9%        | 88.2%           | 89.6%               | 99.4%   |



**Fig. (5).** Performance comparison of various algorithms for skin lesion diagnostics.

## 6. DISCUSSION

The superior performance of DEEPSCAN can be attributed to the power of attention mechanisms and global contextual understanding brought by Vision Transformers. ViTs excel in capturing intricate patterns of skin lesions, which are often challenging for traditional methods. Moreover, the model's ability to leverage clinical metadata and histopathological data contributes to more accurate and comprehensive diagnostics. These promising results suggest that DEEPSCAN has the potential to revolutionize the field of dermatological image analysis. However, further research and validation on diverse real-world datasets are needed to fully assess its clinical utility. The integration of ViTs opens up exciting possibilities for improving the accuracy and early detection of skin lesions, ultimately leading to more effective treatments and patient care.

## CONCLUSION

In this paper, we introduced DEEPSCAN, a novel approach that integrates Vision Transformers for advanced skin lesion diagnostics. Our experiments on the Dermatological Vision Dataset (DermVisD) demonstrated significant improvements in accuracy compared to traditional methods. DEEPSCAN achieved a remarkable accuracy rate of 97.8% on the validation dataset, surpassing the performance of traditional CNN-based models used in dermatology. With ViTs, DEEPSCAN recorded precision and sensitivity figures of 96.5% and 98.1%, respectively, illustrating its prowess in curbing incorrect positives and adeptly pinpointing genuine ones. The F1 measure, which amalgamates precision and sensitivity, stood at 97.3%, emphasizing DEEPSCAN's equilibrium in its diagnostic capabilities. Additionally, the AUC-ROC value reached an impressive 99.4%, indicating the strong discriminatory power of the DEEPSCAN model in accurately identifying skin lesions. DEEPSCAN's capability to recognize complex lesion patterns, along with its utilization of clinical and histopathological data, makes it a promising tool for dermatologists and healthcare providers. The exceptional accuracy, precision, recall, F1-score, and AUC-ROC values obtained with ViTs demonstrate their potential to revolutionize the field of dermatological imaging by providing more accurate and reliable diagnoses. While our results are promising, further research and clinical validation are essential to fully realize DEEPSCAN's potential impact on dermatological practice. This work opens up new avenues for improving skin lesion diagnostics, ultimately leading to better patient care and outcomes in the field of dermatology.

## ABBREVIATION

CNNs = Convolutional Neural Networks

## ETHICS APPROVAL AND CONSENT TO PARTICIPATE

Not applicable.

## HUMAN AND ANIMAL RIGHTS

Not applicable.

## CONSENT FOR PUBLICATION

Not applicable.

## AVAILABILITY OF DATA AND MATERIALS

The data and supportive information are available within the article.

## FUNDING

None.

## CONFLICT OF INTEREST

The authors declare no conflict of interest, financial or otherwise.

## ACKNOWLEDGEMENTS

Declared none.

## REFERENCES

- [1] Zhou L, Luo Y. Deep features fusion with mutual attention transformer for skin lesion diagnosis. 2021 IEEE International Conference on Image Processing (ICIP). Anchorage, AK, USA, 19-22 September 2021, pp. 3797-3801. <http://dx.doi.org/10.1109/ICIP42928.2021.9506211>
- [2] Zhang Y, Xie F, Chen J. TFormer: A throughout fusion transformer for multi-modal skin lesion diagnosis. *Comput Biol Med* 2023; 157: 106712. <http://dx.doi.org/10.1016/j.complbiomed.2023.106712> PMID: 36907033
- [3] Abbas Q, Daadaa Y, Rashid U, Ibrahim M. Assist-dermo: A lightweight separable vision transformer model for multiclass skin lesion classification. *Diagnostics* 2023; 13(15): 2531. <http://dx.doi.org/10.3390/diagnostics13152531> PMID: 37568894
- [4] Wu W, Mehta S, Nofallah S, *et al.* Scale-aware transformers for diagnosing melanocytic lesions. *IEEE Access* 2021; 9: 163526-41. <http://dx.doi.org/10.1109/ACCESS.2021.3132958> PMID: 35211363
- [5] Gulzar Y, Khan SA. Skin lesion segmentation based on vision transformers and convolutional neural networks—A comparative study. *Appl Sci* 2022; 12(12): 5990. <http://dx.doi.org/10.3390/app12125990>
- [6] Krishna GS, Supriya K, Sorgile M. Le- sionAid: Vision transformers-based skin lesion generation and classification. *arXiv preprint* 2023; 2302: 01104 .
- [7] Wang J, Wei L, Wang L, Zhou Q, Zhu L, Qin J. Boundary-aware transformers for skin lesion segmenta- tion. 24th International Conference, . Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I 24, pp. 206-216.
- [8] Aladhadh S, Alsanea M, Aloraini M, Khan T, Habib S, Islam M. An effective skin cancer classification mechanism *via* medical vision transformer. *Sensors* 2022; 22(11): 4008. <http://dx.doi.org/10.3390/s22114008> PMID: 35684627
- [9] Wu H, Chen S, Chen G, Wang W, Lei B, Wen Z. FAT-Net: Feature adaptive transformers for automated skin lesion segmentation. *Med Image Anal* 2022; 76: 102327. <http://dx.doi.org/10.1016/j.media.2021.102327> PMID: 34923250
- [10] Gaviria D. Application of deep learning general-purpose neural architectures based on vision transformers for ISIC melanoma classification. Master's thesis, Universitat Politècnica de Catalunya 2022.
- [11] Sharafudeen M, J A, Chandra S S V. Leveraging vision attention transformers for detection of artificially synthesized dermoscopic lesion deepfakes using derm-CGAN. *Diagnostics* 2023; 13(5): 825. <http://dx.doi.org/10.3390/diagnostics13050825> PMID: 36899969

- [12] Rezaee K, Khosravi MR, Qi L, Abbasi M. SkinNet: A hybrid convolutional learning approach and transformer module through bi-directional feature fusion. 2022 International Conference on Computing, Communication, Security and Intelligent Systems (IC3SIS), . Kochi, India, 23-25 June 2022, pp. 1-6. <http://dx.doi.org/10.1109/IC3SIS54991.2022.9885591>
- [13] Eskandari S, Lump J, Giraldo LS. Skin lesion segmentation improved by transformer-based networks with inter-scale dependency modeling. International Workshop on Machine Learning in Medical Imaging, . Cham: Springer Nature Switzerland, 2023, pp. 351-360.
- [14] Ayas S. Multiclass skin lesion classification in dermoscopic images using swin transformer model. *Neural Comput Appl* 2023; 35(9): 6713-22. <http://dx.doi.org/10.1007/s00521-022-08053-z>
- [15] Liu L, Liang C, Xue Y, et al. An intelligent diagnostic model for melasma based on deep learning and multimode image input. *Dermatol Ther* 2023; 13(2): 569-79. <http://dx.doi.org/10.1007/s13555-022-00874-z> PMID: 36577888
- [16] Wang J, Chen F, Ma Y, et al. XBound-former: Toward cross-scale boundary modeling in transformers. *IEEE Trans Med Imaging* 2023; 42(6): 1735-45. <http://dx.doi.org/10.1109/TMI.2023.3236037> PMID: 37018671
- [17] Shamsad F, Khan S, Zamir SW, et al. Transformers in medical imaging: A survey. *Med Image Anal* 2023; 88: 102802. <http://dx.doi.org/10.1016/j.media.2023.102802> PMID: 37315483
- [18] Khan S, Ali H, Shah Z. Identifying the role of vision transformer for skin cancer—A scoping review. *Front Artif Intell* 2023; 6: 1202990. <http://dx.doi.org/10.3389/frai.2023.1202990> PMID: 37529760
- [19] Liu R, Duan S, Xu L, Liu L, Li J, Zou Y. A fuzzy transformer fusion network (FuzzyTransNet) for medical image segmentation: The case of rectal polyps and skin lesions. *Appl Sci* 2023; 13(16): 9121. <http://dx.doi.org/10.3390/app13169121>
- [20] Alahmadi MD, Alghamdi W. Semi-supervised skin lesion segmentation with coupling CNN and transformer features. *IEEE Access* 2022; 10: 122560-9. <http://dx.doi.org/10.1109/ACCESS.2022.3224005>
- [21] Dong Y, Wang L, Li Y. TC-Net: Dual coding network of transformer and CNN for skin lesion segmentation. *PLoS One* 2022; 17(11): e0277578. <http://dx.doi.org/10.1371/journal.pone.0277578> PMID: 36409714
- [22] Wang J, Li B, Guo X, Huang J, Song M, Wei M. CTCNet: A bi-directional cascaded segmentation network combining Transformers with CNNs for skin lesions. *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*, . Cham: Springer Nature Switzerland, 2022, pp. 215-226. [http://dx.doi.org/10.1007/978-3-031-18910-4\\_18](http://dx.doi.org/10.1007/978-3-031-18910-4_18)
- [23] Roy Vikas Kumar, Thakur Vasu, Goyal Nupur. Vision transformer framework approach for melanoma skin disease identification. *Research Square* 2023; 1-12.
- [24] Luo N, Zhong X, Su L, Cheng Z, Ma W, Hao P. Artificial intelligence-assisted dermatology diagnosis: From unimodal to multimodal. *Comput Biol Med* 2023; 165: 107413. <http://dx.doi.org/10.1016/j.compbimed.2023.107413> PMID: 37703714
- [25] Cao W, Yuan G, Liu Q, et al. ICL-Net: Global and local interpixel correlations learning network for skin lesion segmentation. *IEEE J Biomed Health Inform* 2023; 27(1): 145-56. <http://dx.doi.org/10.1109/JBHI.2022.3162342> PMID: 35353708
- [26] Jenefa A. A robust deep learning-based approach for network traffic classification using CNNs and RNNs. 2023 4th International Conference on Signal Processing and Communication (ICSPEC), . Coimbatore, India, 23-24 March 2023, pp. 106-110.
- [27] Kuriakose BM. EDSR: Empowering super-resolution algorithms with high-quality DIV2K images. *Intell Decis Technol* 1-15.
- [28] Alahmadi MD. Medical image segmentation with learning semantic and global contextual representation. *Diagnostics* 2022; 12(7): 1548. <http://dx.doi.org/10.3390/diagnostics12071548> PMID: 35885454
- [29] Regi AE. J. A, S. V. E. Sonia, E. Naveen, L. A and V. K2023 International Conference on Circuit Power and Computing Technologies (ICCPCT). Kollam, India. 2023; pp. : 823-9. <http://dx.doi.org/10.1109/ICCPCT58313.2023.10245112>
- [30] Regi AE. Liquid biopsy for non-invasive monitoring of tumour evolution and response to therapy. 2023 International Conference on Circuit Power and Computing Technologies (ICCPCT), . Kollam, India, 2023, pp. 815-822.
- [31] Shen J, Hu Y, Zhang X, Gong Y, Kawasaki R, Liu J. Structure-oriented transformer for retinal diseases grading from OCT images. *Comput Biol Med* 2023; 152: 106445. <http://dx.doi.org/10.1016/j.compbimed.2022.106445> PMID: 36549031
- [32] Bozorgpour A, Sadegheih Y, Kazerouni A, Azad R, Merhof D. DermoSegDiff: A boundary-aware segmentation diffusion model for skin lesion delineation. International Workshop on Predictive Intelligence In Medicine, . Cham: Springer Nature Switzerland, 2023, pp. 146-158.