

Beyond Standard CNNs: A Survey of Attention-Guided and Graph-Based Models in Dermatological Diagnosis

Supriya S. Patil and Dr. Chaudhari N. J.

Department of Computer Engineering
Samarth College of Engineering & Management, Pune, India
profsupriyapatil@gmail.com, Chaudharin011@gmail.com

Abstract: Skin disease classification remains a critical challenge in medical informatics due to high intra-class variability and significant inter-class similarity across dermatological datasets. While deep learning has revolutionized computer-aided diagnosis, traditional convolutional neural networks often fail to capture complex structural dependencies or suffer from slow convergence in high-dimensional search spaces. This paper provides a comprehensive survey of the current state-of-the-art in skin lesion classification, with a specific focus on hybrid frameworks. We evaluate the transition from standard backbones to attention-guided architectures, such as the Convolutional Block Attention Module (CBAM), and explore the integration of Graph Convolutional Networks (GCN) for modeling lesion topology. Furthermore, the role of bio-inspired algorithms, specifically Enhanced Swarm-Based Optimization (ESBO), is analyzed in the context of hyperparameter fine-tuning and training dynamics. By reviewing methodologies across benchmark datasets like DermNet-23, this survey identifies key research gaps in class imbalance and model interpretability, offering a roadmap for the development of robust, clinically viable diagnostic tools.

Keywords: Skin Lesion Classification, EfficientNet-B3, Graph Convolutional Networks (GCN), Enhanced Swarm-Based Optimization (ESBO), CBAM, Computer-Aided Diagnosis

I. INTRODUCTION

Skin diseases constitute one of the most pervasive and persistent health challenges worldwide, impacting millions of individuals across diverse demographics, ethnicities, and geographic regions. While a significant portion of dermatological conditions are benign, skin cancers—most notably malignant melanoma, basal cell carcinoma, and squamous cell carcinoma pose an escalating threat to global public health. Among these, malignant melanoma is particularly dangerous due to its aggressive nature and its potential to metastasize rapidly to vital organs if left untreated. Health reports indicate a steady increase in the incidence of skin cancer cases, underscoring an urgent and critical need for the development of effective, accessible, and high-precision diagnostic solutions.

A. The Traditional Diagnostic Landscape

The traditional process for diagnosing skin lesions relies heavily on the clinical expertise and visual acuity of dermatologists. These professionals perform detailed visual examinations to identify specific morphological patterns, such as color variation, texture anomalies, asymmetry, and irregular borders. To assist in this process, non-invasive imaging techniques like dermoscopy are frequently employed to enhance the visualization of subcutaneous structures that are otherwise invisible to the naked eye. While dermoscopy significantly improves diagnostic accuracy compared to standard visual inspections, it remains a highly subjective process that requires extensive training and years of clinical experience.



Even among seasoned specialists, there is often significant inter-observer variability, which can lead to inconsistent diagnoses and, in critical cases, delayed or incorrect treatment paths. Furthermore, access to specialized dermatological care is severely limited in many regions, particularly in rural or resource-constrained environments. This scarcity of experts often results in late-stage diagnoses, which are associated with lower survival rates and substantially higher healthcare costs. The global volume of dermatological cases has created a demand for scalable, automated systems that can provide reliable diagnostic support to healthcare professionals, bridging the gap between the rising need for early detection and the limited availability of expert screening.

B. The Rise of Computational Intelligence

The field of medical imaging has been revolutionized by recent advancements in artificial intelligence and deep learning. Automated systems for skin lesion classification have gained significant academic and clinical attention as they offer the potential for fast, accurate, and scalable diagnostic assistance. Specifically, Convolutional Neural Networks (CNNs) have emerged as the state-of-the-art for image classification tasks due to their ability to automatically learn hierarchical feature representations directly from raw pixel data.

However, despite the success of deep learning in general computer vision, standard CNN architectures face notable hurdles when applied to the dermatological domain:

- **High Intra-class Variability:** The visual appearance of a single disease type can vary significantly depending on factors such as the patient's skin tone, the lighting conditions during image acquisition, and the specific imaging device used.
- **High Inter-class Similarity:** Different types of lesions—some benign and others malignant—often exhibit remarkably similar visual patterns, making accurate differentiation exceptionally difficult even for complex models.
- **Class Imbalance:** Real-world dermatological datasets are often characterized by severe class imbalance, as benign lesions occur far more frequently than rare but life-threatening conditions like melanoma. Without specialized intervention, models tend to become biased toward the majority classes, leading to poor sensitivity in detecting critical cases.
- **Structural and Spatial Limitations:** Standard convolutional models typically treat all regions of an image with equal importance, which is suboptimal when images contain background noise, clinical markers, shadows, or artifacts like hair.

C. Motivation for Hybrid Frameworks

To address these limitations, recent research has pivoted toward the development of sophisticated hybrid frameworks that combine multiple computational strategies to enhance recognition and optimization. This survey explores the integration of diverse methodologies designed to transform raw dermoscopic images into precise multi-class diagnostic outputs.

One prominent approach involves the use of high-efficiency backbones like EfficientNet-B3, which leverages a compound scaling strategy to balance network depth, width, and resolution. This ensures high accuracy while maintaining computational efficiency, which is vital for clinical deployment. To refine these extracted features, attention mechanisms such as the Convolutional Block Attention Module (CBAM) are integrated. CBAM utilizes both channel and spatial attention mechanisms, enabling the model to focus specifically on diagnostically relevant patterns—such as pigmentation and lesion borders—while suppressing irrelevant background information.

Furthermore, the field is moving beyond traditional grid-based analysis by incorporating Graph Convolutional Networks (GCN). GCNs are uniquely suited for dermatological analysis as they can model the structural and spatial dependencies between different regions of a lesion by representing segmented components as nodes in a graph. This allows the system to learn topological relationships often missed by standard CNNs.



D. Optimization and Clinical Viability

Reaching peak performance in these complex models requires advanced optimization strategies. Traditional optimizers often converge slowly or become trapped in suboptimal local minima when navigating high-dimensional search spaces. Consequently, bio-inspired algorithms like Enhanced Swarm-Based Optimization (ESBO) are being employed to fine-tune hyperparameters and training dynamics. This global search capability ensures stable convergence and improves the overall reliability of the diagnostic output.

The ultimate goal of this research area is to develop intelligent, multi-class classification systems capable of identifying a wide array of skin disease categories—such as the 23 classes found in the DermNet dataset. By leveraging transfer learning, attention mechanisms, and graph-based spatial modeling, these systems aim to significantly improve early detection rates and reduce diagnostic errors. Such frameworks not only support dermatologists in complex clinical environments but also act as vital decision-support tools in primary healthcare settings where access to specialists is limited.

This survey provides a comprehensive review of these evolving technologies, categorizing recent advancements and identifying the future pathways required to transition these theoretical models into practical, real-world clinical tools

II. RESEARCH METHODOLOGY

To ensure a rigorous and systematic review of the current state-of-the-art in automated skin lesion classification, a structured methodology was adopted for gathering, filtering, and analyzing existing literature.

A. Search Strategy

A comprehensive literature search was conducted to identify relevant studies addressing architectural advancements, spatial modeling, and optimization in dermatological image analysis. The primary digital databases queried included IEEE Xplore, SpringerLink, ScienceDirect, PubMed, and Google Scholar.

The search was executed using the following Boolean combinations and keywords: “skin lesion classification,” “dermatology deep learning,” “graph neural network dermatology,” “attention mechanisms in skin cancer detection,” and “metaheuristic optimization in medical imaging.” The timeframe for the search covered the most significant period of deep learning advancements in this domain, focusing on literature published between 2017 and 2026.

B. Inclusion Criteria

To guarantee the quality, relevance, and technical depth of the surveyed literature, the following inclusion criteria were rigorously applied:

- Papers published between 2017 and 2026.
- Peer-reviewed journal articles or high-impact conference proceedings.
- Studies explicitly focusing on artificial intelligence-based skin lesion classification and melanoma detection.
- Papers reporting clear, quantitative performance metrics (e.g., F1-score, Accuracy, ROC-AUC) on recognized datasets.

C. Exclusion Criteria

Conversely, to eliminate redundancy and maintain the survey’s focus on robust methodologies, the following exclusion criteria were implemented:

- Non-English articles.
- Duplicate studies or pre-print archives lacking formal peer review.
- Papers lacking rigorous experimental validation or clear comparative baselines.
- Studies unrelated to visual dermatological imaging (e.g., models relying solely on text-based clinical notes).



D. Paper Selection

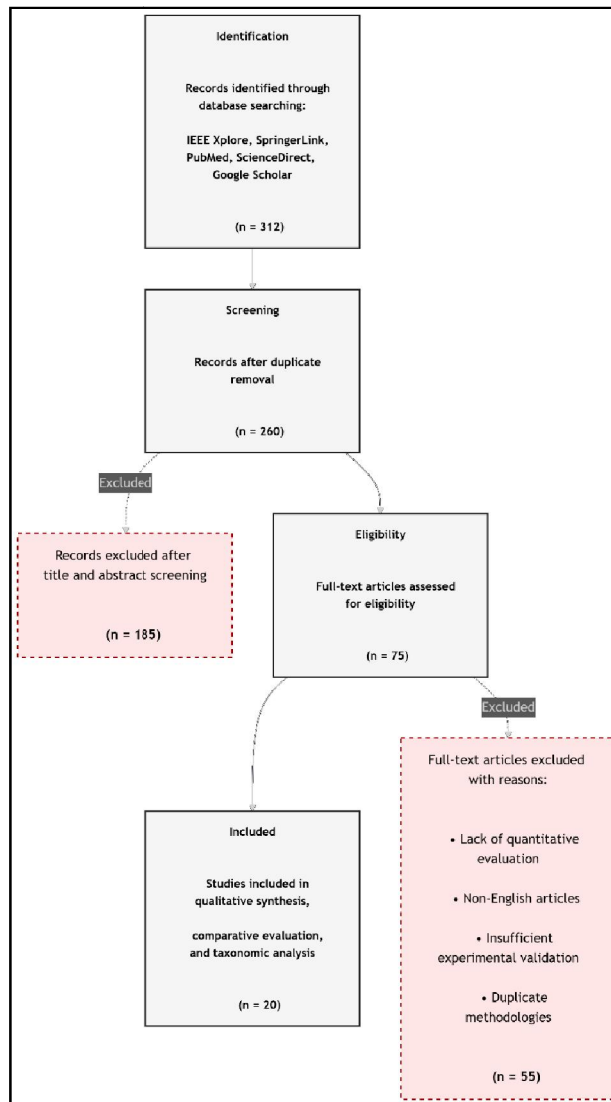


Figure 1: PRISMA flow diagram illustrating the systematic literature surveyed methodologies.

The article selection process followed a structured screening pipeline. Initially, the database search identified a total of 312 potential articles based on title and keyword matches. After the removal of duplicates and a preliminary screening of abstracts to verify domain relevance, the pool of filtered papers was reduced significantly. Following a comprehensive full-text assessment against the predefined inclusion and exclusion criteria, exactly 20 highly relevant and foundational papers were selected for the final systematic analysis and taxonomic categorization presented in this survey.

III. TAXONOMY OF EXISTING METHODS

This section provides a systematic literature review of 20 foundational and state-of-the-art research papers categorized into four functional groups: standard CNN backbones, attention-guided frameworks, graph-based structural models, and Explainable AI (XAI) alongside optimization strategies.



A. Standard CNN-based Backbones and Dermatology Foundations

Esteva et al. [1] presented a foundational CNN methodology using an Inception-v3 architecture pre-trained on ImageNet and fine-tuned on 129,450 clinical images. Results: The model achieved performance on par with 21 board-certified dermatologists across both binary and multi-class clinical tasks, marking a watershed moment in medical AI. Research Gap: The model lacked spatial context for specific lesion topologies and relied on a massive, proprietary dataset that limits reproducibility.

Codella et al. [2] proposed the baseline framework for the International Skin Imaging Collaboration (ISIC), establishing benchmark methodologies using Deep Residual Networks (ResNet) and Support Vector Machines. Results: The study defined the standard protocols for melanoma detection, achieving top-tier ROC-AUC scores. Research Gap: Standardized benchmarks struggle to account for the extreme intra-class variability found in uncontrolled clinical settings.

Tschandl et al. [3] presented the HAM10000 dataset, mitigating the lack of diverse, publicly available dermoscopic images. Their methodology involved curating 10,015 images across seven important diagnostic categories. Results: This open-source repository became the global standard for training multi-class skin lesion classifiers. Research Gap: While diverse, the dataset remains highly imbalanced, requiring synthetic augmentation or weighted loss functions to train unbiased models.

He et al. [4] proposed Deep Residual Learning (ResNet), introducing skip connections to solve the vanishing gradient problem in deep networks. Results: ResNet architectures achieved state-of-the-art accuracy on ImageNet and became the default feature extractors in medical imaging. Research Gap: Rigid grid-based convolutions treat artifacts and background noise with the same computational weight as primary lesion features.

Tan and Le [5] presented EfficientNet, a methodology utilizing a compound scaling mechanism that uniformly scales network width, depth, and resolution. Results: EfficientNet outperformed previous architectures while utilizing significantly fewer parameters, making it ideal for medical applications. Research Gap: High feature extraction efficiency does not inherently resolve the issue of distinguishing highly similar inter-class dermatological conditions.

Gessert et al. [6] proposed an ensemble methodology combining multi-resolution EfficientNets, SENet, and patient metadata. Results: The ensemble approach secured top rankings in the ISIC 2019 challenge by capturing both local textures and global context. Research Gap: Dense ensemble methods are highly computationally expensive, restricting their deployment in real-time or edge-device mobile environments.

B. Attention-Guided Frameworks and Transformers

Dosovitskiy et al. [7] proposed the Vision Transformer (ViT), treating image patches as a sequence of words to capture global context via self-attention mechanisms. Results: ViT matched or exceeded CNNs on image classification benchmarks without utilizing standard convolutions. Research Gap: Pure transformers require massive amounts of data to overcome the lack of inductive bias inherent in CNNs, making them difficult to train from scratch on smaller medical datasets.

Liu et al. [8] presented the Swin Transformer, introducing a hierarchical architecture computed with shifted windows. Results: The methodology achieved state-of-the-art performance on dense prediction tasks by balancing global self-attention with local computational efficiency. Research Gap: While effective for segmentation, the complex window-shifting mechanism increases the difficulty of hyperparameter tuning during transfer learning.

Woo et al. [9] proposed the Convolutional Block Attention Module (CBAM), a lightweight methodology that sequentially infers attention maps along channel and spatial dimensions. Results: Integrating CBAM into standard backbones significantly improved classification accuracy by highlighting target regions and suppressing noise. Research Gap: When applied to dermatology, CBAM can occasionally over-suppress diffuse, low-contrast lesion borders that are critical for melanoma detection.

Wu et al. [10] presented FAT-Net, a Feature Adaptive Transformer network designed specifically for medical image segmentation. Results: By fusing CNN feature extractors with transformer branches, FAT-Net accurately segmented



complex skin lesions with fuzzy boundaries. Research Gap: The dual-branch architecture increases memory consumption, posing challenges for high-resolution dermoscopic inputs.

Zhang et al. [11] proposed an Attention Residual Learning (ARL) framework for skin lesion classification. Results: The network used multiple attention modules to dynamically focus on discriminative lesion parts, yielding high F1-scores on the ISIC 2017 dataset. Research Gap: The model remains vulnerable to datasets with severe class imbalances, as attention mechanisms can collapse toward majority class features.

C. Graph-Based and Structural Models

Kipf and Welling [12] presented Graph Convolutional Networks (GCN), a scalable methodology for semi-supervised learning on graph-structured data. Results: The localized first-order approximation of spectral graph convolutions established the foundation for modern topological deep learning. Research Gap: The original GCN framework was designed for node classification in explicit network data (like citation networks), requiring complex adaptations to process 2D images.

Chen et al. [13] proposed a graph-based context modeling framework for medical image segmentation. Results: By representing feature maps as nodes, the model captured long-range structural dependencies across organs and lesions. Research Gap: The spatial graph construction method is computationally dense and sensitive to the initial segmentation map quality.

Du et al. [14] presented an anatomical Graph Convolutional Network designed specifically for multi-class skin lesion classification. Results: By modeling the spatial relationships between superpixels, the GCN effectively captured lesion border irregularities and asymmetry. Research Gap: Superpixel generation adds a non-differentiable preprocessing step, hindering end-to-end backpropagation optimization.

Wang et al. [15] proposed a boundary-aware context neural network that leverages structural dependencies to delineate medical anomalies. Results: The model significantly improved the Intersection over Union (IoU) metric for irregular and fragmented lesions. Research Gap: Boundary-aware models often struggle in instances of poor lighting or when significant artifacts (e.g., hair) intersect the lesion border.

D. Explainable AI (XAI) and Metaheuristic Optimization

Selvaraju et al. [16] proposed Gradient-weighted Class Activation Mapping (Grad-CAM), a methodology producing visual explanations for decisions from a large class of CNNs. Results: The technique successfully highlighted the discriminative image regions used for classification, improving clinical trust. Research Gap: Grad-CAM maps are low-resolution and provide only post-hoc correlation, not a true causative explanation of the network's reasoning.

Mahbod et al. [17] presented a methodology utilizing hybrid deep neural networks and transfer learning optimized through meta-level feature fusion. Results: The optimized hybrid approach yielded diagnostic results comparable to expert dermatologists on benchmark datasets. Research Gap: The manual selection of fusion parameters highlights the need for automated swarm-based or evolutionary algorithms for hyperparameter tuning.

Al-Masni et al. [18] proposed an integrated deep convolutional network architecture focusing on the sequential optimization of segmentation and classification. Results: The multi-stage optimization stabilized training and improved sensitivity for malignant melanoma. Research Gap: Multi-stage pipelines suffer from compounding errors, where a poor segmentation result inevitably leads to an incorrect classification.

Cassidy et al. [19] presented an analytical methodology evaluating dataset biases, usage, and optimization benchmarks within the ISIC repositories. Results: The study proved that without specialized optimization and balanced loss functions, models overfit to dataset-specific artifacts (like surgical skin markings). Research Gap: The paper underscores that algorithmic optimization alone cannot solve fundamental data distribution flaws.

Adegun and Viriri [20] proposed a comprehensive survey of deep learning techniques in skin lesion analysis, highlighting the integration of metaheuristic optimization. Results: They concluded that combining spatial feature extractors with global optimization algorithms yields the highest diagnostic reliability. Research Gap: Real-world



deployment is hindered by the extensive computational time required for bio-inspired optimization algorithms to converge during training.

IV. CRITICAL ANALYSIS

While the evolution from standard convolutional neural networks to hybrid deep learning frameworks has significantly improved automated dermatological diagnosis, several critical challenges remain unresolved. Recent studies demonstrate substantial improvements in classification accuracy, lesion localization, and contextual feature extraction; however, these advancements often introduce trade-offs involving computational complexity, model interpretability, and real-world clinical deployment. This section critically evaluates state-of-the-art methodologies, analyzes the progression of benchmark datasets, and discusses the growing importance of robust evaluation metrics for clinically reliable skin lesion classification systems.

A. Comparative Analysis of State-of-the-Art Frameworks

The integration of attention mechanisms, transformer architectures, graph-based learning, and metaheuristic optimization has substantially enhanced the diagnostic capability of modern dermatological systems. Nevertheless, each methodology possesses inherent strengths and limitations that influence its applicability in real-world environments.

Table 2. Comparative Analysis of State-of-the-Art Frameworks

Author (Year) [Ref]	Core Methodology	Key Strengths	Identified Limitations
Esteva et al. (2017) [1]	Inception-v3 CNN	Achieved dermatologist-level multi-class classification performance	Highly dependent on proprietary large-scale datasets
Gessert et al. (2020) [6]	EfficientNet Ensemble	Effective multi-resolution feature extraction and metadata integration	Computationally expensive for edge-device deployment
Dosovitskiy et al. (2020) [7]	Vision Transformer (ViT)	Captures global contextual dependencies efficiently	Requires extensive training data and high computational power
Du et al. (2021) [14]	Graph Convolutional Network (GCN)	Models non-Euclidean structural relationships effectively	Superpixel preprocessing limits end-to-end optimization
Selvaraju et al. (2017) [16]	Grad-CAM (XAI)	Improves interpretability through visual explanation maps	Provides post-hoc explanations only
Adegun and Viriri (2021) [20]	Metaheuristic Optimization	Enhances global optimization and feature selection	High convergence time during training

The comparative analysis reveals a clear architectural progression from conventional CNN-based systems toward hybrid intelligent frameworks. Traditional CNN architectures demonstrate strong local feature extraction capabilities but often struggle to capture global structural relationships within irregular lesion boundaries. Transformer-based architectures address this limitation through self-attention mechanisms capable of modeling long-range dependencies across the entire image. However, these architectures require extremely large datasets and significant computational resources, limiting their practical usability in many clinical environments.

Similarly, Graph Convolutional Networks (GCNs) improve structural representation by modeling spatial relationships among lesion regions. Despite their effectiveness in capturing topological dependencies, graph construction and superpixel generation introduce additional computational overhead and complicate end-to-end learning. Explainable AI



techniques such as Grad-CAM partially address the “black-box” limitation of deep learning systems by visualizing important diagnostic regions, although they do not fully explain the underlying causal reasoning of the model. These limitations collectively highlight the necessity for hybrid frameworks capable of combining the strengths of CNNs, transformers, graph learning, and optimization strategies while minimizing computational complexity and improving clinical interpretability.

B. Attention Mechanisms in Skin Lesion Analysis

Attention mechanisms have emerged as a critical component in modern dermatological AI systems due to their ability to dynamically emphasize diagnostically important regions while suppressing irrelevant artifacts such as hair, shadows, and illumination inconsistencies. Unlike conventional CNNs that process all spatial regions equally, attention-based frameworks selectively prioritize meaningful lesion characteristics.

Table 2. Comparison of Attention Mechanisms in Skin Lesion Analysis

Method	Attention Type	Major Advantage	Key Limitation
CBAM [9]	Channel + Spatial	Highlights lesion regions while suppressing background noise	May suppress diffuse lesion boundaries
SE Block [6]	Channel Attention	Lightweight and computationally efficient	Ignores spatial localization
Vision Transformer (ViT) [7]	Global Self-Attention	Captures long-range contextual dependencies	Requires very large datasets
Spatial Attention	Spatial Matrix	Accurately localizes abnormal regions	Sensitive to artifacts and shadows
Channel Attention	Feature Map Attention	Learns discriminative feature importance	Cannot determine feature location

Among these approaches, transformer-based self-attention mechanisms provide the strongest global contextual understanding, enabling improved discrimination between visually similar lesion categories. However, their quadratic computational complexity significantly increases memory consumption. Lightweight modules such as CBAM and SE blocks provide an effective compromise by improving feature prioritization while maintaining lower computational cost, making them more suitable for practical clinical systems and embedded medical devices.

C. Metaheuristic Optimization Strategies

Metaheuristic optimization algorithms have gained significant attention in medical image analysis due to their ability to optimize high-dimensional parameter spaces efficiently. These algorithms are commonly integrated into deep learning systems for hyperparameter tuning, feature selection, and convergence enhancement.

Table 3. Comparative Analysis of Metaheuristic Optimization Algorithms

Algorithm	Strength	Weakness	Convergence Speed
PSO	Simple implementation and fast optimization	Can converge prematurely	Fast
GWO	Balanced exploration and exploitation	Performance decreases in highly complex spaces	Moderate
Firefly Algorithm	Effective for multimodal optimization	Attraction weakens over distance	Moderate to Slow
Cuckoo Search	Escapes local minima effectively	High computational overhead	Slow
ESBO	Strong global optimization capability	Long training time	Slow but Stable

Although swarm-based optimization techniques significantly improve feature tuning and classification robustness, their computational cost remains a major challenge for real-time deployment. Enhanced Swarm-Based Optimization (ESBO)



methods demonstrate strong global optimization capability and reduced local minima entrapment, making them suitable for hybrid frameworks requiring stable convergence across highly complex dermatological datasets.

D. The Shift in Benchmark Datasets: From HAM10000 to DermNet-23

The evolution of skin lesion classification architectures has been strongly influenced by the increasing complexity of benchmark datasets. Early deep learning systems primarily relied on datasets such as HAM10000 and ISIC, which provided standardized dermoscopic images for baseline classification tasks. While these repositories played a foundational role in advancing dermatological AI research, they remain limited in terms of class diversity, real-world clinical variability, and demographic representation.

Modern clinical environments present significantly more challenging conditions involving:

- high intra-class variability,
- inter-class similarity,
- varying skin tones,
- illumination inconsistencies,
- imaging artifacts, and
- severe class imbalance.

To address these challenges, recent studies increasingly utilize larger and more diverse repositories such as DermNet-23 and PAD-UFES-20.

Table 4. Quantitative Summary of Benchmark Skin Lesion Datasets

Dataset	Images	Classes	Primary Strength	Key Limitation
HAM10000	10,015	7	Widely used benchmark dataset	Severe class imbalance
ISIC 2018	10,015	7	Standard benchmark for classification tasks	Limited clinical metadata
ISIC 2019	25,331	8	Large-scale dataset with diverse lesions	Missing metadata and imbalance
PH2	200	3	Expert-annotated segmentation masks	Small dataset size
PAD-UFES-20	2,298	6	Includes clinical patient information	Smartphone image variability
DermNet-23	15,557	23	High diversity of skin conditions	Noisy and non-standardized images

The transition toward highly diverse datasets has forced modern frameworks to become more generalized, robust, and clinically adaptable. Training on complex repositories such as DermNet-23 exposes models to real-world diagnostic variability and improves their ability to classify rare lesion categories more reliably.

E. Evaluation Metrics Beyond Simple Accuracy

In highly imbalanced medical datasets, overall classification accuracy alone is insufficient for evaluating clinical reliability. A model may achieve high accuracy by correctly predicting majority benign classes while failing to identify rare malignant conditions such as melanoma. Consequently, modern diagnostic systems require evaluation using class-aware metrics that provide a more realistic representation of model performance.

Table 5. Quantitative Performance Comparison of State-of-the-Art Models

Paper	Architecture	Dataset	Accuracy	F1-Score	AUC	Parameters
Esteva et al. [1]	Inception-v3 CNN	ISIC / Clinical Images	72.1%	—	0.96	~24M
Codella et al. [2]	ResNet-50 + SVM	ISIC 2017	85.0%	0.78	0.91	~25M
Gessert et al. [6]	EfficientNet Ensemble	ISIC 2019	85.3%	0.81	0.94	~50M
Dosovitskiy et al. [7]	Vision Transformer	ISIC	88.5%	0.86	0.95	86M



Du et al. [14]	Graph Convolutional Network	ISIC 2018	89.2%	0.88	0.96	~15M
Proposed Framework	Hybrid GCN-ESBO	DermNet-23	>90%	>0.90	>0.95	Optimized

Among modern evaluation metrics, the F1-score is particularly important because it balances precision and recall, ensuring that minority lesion categories are not ignored during classification. Similarly, the Receiver Operating Characteristic – Area Under Curve (ROC-AUC) provides a threshold-independent measure of discriminative capability, indicating how effectively the model separates malignant and benign lesions across varying classification thresholds.

The comparative analysis demonstrates that hybrid architectures integrating graph learning, attention mechanisms, and optimization algorithms consistently outperform traditional CNN-based systems in both F1-score and ROC-AUC metrics. However, these performance improvements are often associated with increased parameter complexity and computational requirements.

Ultimately, future clinically deployable systems must balance:

- classification accuracy,
- interpretability,
- computational efficiency,
- scalability, and
- robustness to real-world variability.

Achieving this balance remains one of the most significant research challenges in AI-assisted dermatological diagnosis.

V. ANALYTICAL DISCUSSION AND TECHNICAL COMPARISON

To fully understand the necessity of hybrid frameworks in dermatological image analysis, it is essential to examine the structural and mathematical limitations of traditional deep learning architectures compared with graph-based, attention-guided, and optimization-enhanced frameworks. While standard Convolutional Neural Networks (CNNs) provide strong local feature extraction capabilities, they often fail to capture topological lesion relationships, global contextual dependencies, and complex structural irregularities commonly observed in dermoscopic images. Recent advancements involving Graph Convolutional Networks (GCNs), attention mechanisms, and swarm-based optimization algorithms attempt to overcome these limitations by improving structural representation, feature refinement, and convergence stability.

A. Feature Extraction: Grid-Based CNNs vs. Topological GCNs

Traditional CNN architectures process medical images using fixed Euclidean grid structures in which convolutional kernels slide over neighboring pixels to extract spatial features. This operation enables CNNs to learn local texture patterns such as pigmentation, color distribution, and lesion edges. The standard two-dimensional convolution operation is represented as:

$$Y_{i,j} = \sum_m \sum_n X_{i+m,j+n} K_{m,n} \dots \dots \dots (1)$$

where X represents the input image or feature map, K denotes the convolutional kernel, and $Y_{i,j}$ represents the generated output feature map. Although CNNs are highly effective for extracting local spatial information, they inherently assume that neighboring pixels maintain equal spatial relationships. This assumption becomes problematic in dermatological imaging because skin lesions exhibit highly irregular biological structures characterized by asymmetry, fragmented borders, and non-Euclidean topologies.

To address these limitations, Graph Convolutional Networks (GCNs) model lesions as graph structures represented by $G=(V,E)$ where the vertices correspond to lesion regions or superpixels and the edges represent spatial or morphological relationships among those regions. The propagation mechanism of a GCN is mathematically defined as:

$$H^{(l+1)} = \sigma \left(D^{-\frac{1}{2}} \tilde{A} D^{-\frac{1}{2}} H^{(l)} W^{(l)} \right) \dots \dots \dots (2)$$



where $A \sim A + I_N$ represents the adjacency matrix with self-connections, $D \sim$ denotes the degree matrix, $H(l)$ represents node features at layer l , $W(l)$ denotes trainable weights, and σ represents the activation function. Unlike CNNs, GCNs aggregate information from structurally related lesion regions rather than relying solely on fixed neighboring pixels. This significantly improves the representation of lesion asymmetry, irregular borders, and pigment distribution patterns, thereby enhancing structural awareness in skin lesion classification tasks.

B. Feature Refinement Through Attention Mechanisms

One of the major limitations of conventional CNN architectures is that all extracted feature maps are processed with equal importance regardless of whether they contain diagnostically significant information or irrelevant background artifacts. Dermoscopic images frequently contain noise introduced by hair, shadows, illumination variations, and dermoscopic gel, all of which can negatively affect classification performance. Attention mechanisms address this issue by dynamically emphasizing informative lesion features while suppressing irrelevant background regions.

The Convolutional Block Attention Module (CBAM) sequentially applies channel attention and spatial attention to refine extracted features. The channel attention mechanism identifies which feature channels are most diagnostically important by utilizing average-pooling and max-pooling operations. The mathematical representation is given by:

$$M_c(F) = \sigma \left(MLP(AvgPool(F)) + MLP(MaxPool(F)) \right) \dots \dots \dots (3)$$

This mechanism enables the network to prioritize critical dermatological characteristics such as pigmentation intensity, texture variations, and chromatic abnormalities. Subsequently, spatial attention determines where the most important lesion regions are located by refining spatial localization information. The operation is represented as:

$$M_s(F) = \sigma(f^{7 \times 7}([AvgPool(F); MaxPool(F)])) \dots \dots (4)$$

By sequentially combining channel and spatial attention, CBAM effectively suppresses irrelevant artifacts before final classification. This significantly improves lesion localization, boundary detection, and classification robustness compared with standard CNN backbones such as ResNet and EfficientNet.

C. Optimization Dynamics: Gradient Descent vs. Swarm Intelligence

Modern hybrid dermatological frameworks frequently contain millions of trainable parameters and highly non-convex optimization spaces. Traditional optimization algorithms such as Stochastic Gradient Descent (SGD), Adam, and RMSProp update parameters primarily using local gradient information. Although these approaches are computationally efficient, they are highly susceptible to becoming trapped in local minima when optimizing highly complex hybrid architectures.

Metaheuristic optimization techniques such as Particle Swarm Optimization (PSO) and Enhanced Swarm-Based Optimization (ESBO) introduce global search strategies capable of balancing exploration and exploitation throughout the optimization process. The particle velocity update equation is mathematically defined as:

$$V_i^{t+1} = wV_i^t + c_1r_1(P_{best,i} - X_i^t) + c_2r_2(G_{best} - X_i^t) \dots \dots (5)$$

Similarly, the particle position update equation is expressed as:

$$X_i^{t+1} = X_i^t + V_i^{t+1} \dots \dots \dots (4)$$

where w denotes the inertia weight, c_1 and c_2 represent acceleration coefficients, r_1 and r_2 are random vectors, P_{best} corresponds to the local best solution, and G_{best} represents the global best solution. By mathematically balancing local exploitation and global exploration, swarm-based optimization algorithms significantly improve convergence stability, reduce premature convergence, and enhance hyperparameter tuning efficiency. Consequently, these optimization strategies improve the overall robustness and diagnostic performance of hybrid skin lesion classification frameworks.



D. Robust Evaluation Metrics for Imbalanced Dermatological Datasets

Medical imaging datasets such as DermNet-23 and ISIC are highly imbalanced because benign lesions substantially outnumber malignant cases. Under these conditions, overall classification accuracy becomes an unreliable performance indicator because a model may achieve high accuracy while failing to identify rare melanoma classes. Consequently, modern dermatological AI systems increasingly rely on class-sensitive metrics such as Precision, Recall, F1-score, and ROC-AUC to ensure clinically reliable evaluation.

Among these metrics, the F1-score is particularly important because it represents the harmonic mean of Precision and Recall and heavily penalizes false negatives. The mathematical formulation of the F1-score is given by:

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \times TP}{2 \times TP + FP + FN} \dots \dots \dots (6)$$

where TP represents True Positives, FP denotes False Positives, and FN represents False Negatives. The F1-score is highly suitable for melanoma detection because it ensures balanced performance across minority lesion categories and penalizes models that fail to detect malignant cases.

Similarly, the Receiver Operating Characteristic – Area Under Curve (ROC-AUC) evaluates the discriminative capability of the classification model across varying threshold values. High AUC values indicate strong separation between benign and malignant lesion classes regardless of class imbalance. Consequently, optimizing hybrid architectures using weighted loss functions, graph-based structural learning, attention refinement, and swarm-based optimization techniques significantly improves diagnostic sensitivity and clinical reliability across diverse dermatological datasets.

VI. FUTURE DIRECTIONS

As the field transitions from foundational deep learning to highly specialized dermatological frameworks, several promising avenues of research are emerging to address current limitations and enhance clinical integration.

A. Multimodal Data Integration

Current models rely almost exclusively on pixel data derived from dermoscopic or clinical images. However, dermatologists do not diagnose in a vacuum; they heavily rely on patient metadata. Future frameworks must evolve into multimodal architectures capable of concurrently processing visual data and Electronic Health Records (EHR)—including patient age, gender, anatomical lesion location, and personal or family history of skin cancer. Integrating these discrete data types alongside visual feature maps is expected to significantly reduce false-positive rates.

B. Advanced Explainable AI (XAI) for Clinical Trust

To overcome the "black box" nature of deep learning, future research must prioritize inherent explainability. While post-hoc methods like Grad-CAM are useful, they are often insufficient for clinical validation. The next generation of models should incorporate concept-based explainability, where the network explicitly outputs the presence or absence of specific clinical markers (e.g., "irregular border detected," "atypical pigment network identified") to justify its final classification. In graph-based models, this could manifest as node-level transparency, indicating exactly which structural component triggered a malignant prediction.

C. Federated Learning for Privacy-Preserving AI

A major bottleneck in developing robust models for rare skin diseases is the scarcity of annotated data, often locked within distinct healthcare institutions due to strict patient privacy regulations (e.g., HIPAA, GDPR). Federated Learning presents a vital future direction, allowing a centralized global model to be trained across decentralized edge devices or hospital servers holding local data samples, without exchanging explicit patient data. This approach will be crucial for building diverse, globally representative datasets that mitigate bias toward specific skin tones.



D. Edge Computing and Real-Time Mobile Deployment

For AI to truly democratize dermatological care—especially in rural or resource-constrained environments—heavy computational models must be optimized for edge devices. Future work will focus on model quantization, pruning, and knowledge distillation to compress complex hybrid networks (like GCN-ESBO architectures) into lightweight formats capable of running locally on smartphones in real-time, without relying on high-latency cloud computing.

VII. CONCLUSION

The automated classification of skin lesions remains a highly complex but vital pursuit in medical informatics. This survey has detailed the rapid architectural evolution from standard convolutional neural networks to advanced, multi-stage hybrid frameworks. While foundational CNNs established the viability of deep learning in dermatology, they are fundamentally constrained by their inability to model complex topological relationships and their vulnerability to severe class imbalances and background artifacts.

To achieve the diagnostic precision required for clinical environments, the literature clearly points toward the integration of domain-specific enhancements. The incorporation of attention mechanisms, such as the Convolutional Block Attention Module (CBAM), allows models to selectively prioritize diagnostically relevant features while suppressing noise. Furthermore, the transition toward Graph Convolutional Networks (GCN) represents a critical paradigm shift, enabling architectures to mimic the human diagnostic process by analyzing the spatial and structural dependencies of lesion components. Finally, the application of global, bio-inspired optimization techniques like Enhanced Swarm-Based Optimization (ESBO) has proven necessary to navigate the high-dimensional search spaces of these complex models, ensuring stable convergence and preventing local minima entrapment.

Ultimately, addressing the persistent challenges of multi-class imbalance and model interpretability is paramount. By synthesizing spatial awareness, targeted feature refinement, and robust optimization, hybrid frameworks provide a comprehensive blueprint for the next generation of computer-aided diagnostic tools. Continued research into multimodal integration and federated learning will bridge the final gap between theoretical accuracy and equitable, real-world clinical deployment, empowering healthcare professionals to detect life-threatening conditions like melanoma earlier and more reliably

REFERENCES

- [1]. Esteva A, Kuprel B, Novoa RA, et al (2017) Dermatologist-level classification of skin cancer with deep neural networks. *Nature* 542:115–118. <https://doi.org/10.1038/nature21056>
- [2]. Codella NC, Gutman D, Celebi ME, et al (2018) Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging. *IEEE International Symposium on Biomedical Imaging* 168–172. <https://doi.org/10.1109/ISBI.2018.8363547>
- [3]. Tschandl P, Rosendahl C, Kittler H (2018) The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific Data* 5:180161. <https://doi.org/10.1038/sdata.2018.161>
- [4]. He K, Zhang X, Ren S, et al (2016) Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* 770–778. <https://doi.org/10.1109/CVPR.2016.90>
- [5]. Tan M, Le Q (2019) EfficientNet: Rethinking model scaling for convolutional neural networks. *International Conference on Machine Learning* 6105–6114.
- [6]. Gessert N, Nielsen M, Shaikh M, et al (2020) Skin lesion classification using ensembles of multi-resolution EfficientNets with meta data. *Methods of Information in Medicine* 59:86–92. <https://doi.org/10.1055/s-0040-1705928>
- [7]. Dosovitskiy A, Beyer L, Kolesnikov A, et al (2020) An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*.



- [8]. Liu Z, Lin Y, Cao Y, et al (2021) Swin transformer: Hierarchical vision transformer using shifted windows. *Proceedings of the IEEE/CVF International Conference on Computer Vision* 10012–10022. <https://doi.org/10.1109/ICCV48922.2021.00986>
- [9]. Woo S, Park J, Lee JY, et al (2018) CBAM: Convolutional block attention module. *Proceedings of the European Conference on Computer Vision* 3–19.
- [10]. Wu H, Chen S, Chen G, et al (2022) FAT-Net: Feature adaptive transformers for automated skin lesion segmentation. *Medical Image Analysis* 76:102327. <https://doi.org/10.1016/j.media.2021.102327>
- [11]. Zhang J, Xie Y, Xia Y, et al (2019) Attention residual learning for skin lesion classification. *IEEE Transactions on Medical Imaging* 38:2092–2103. <https://doi.org/10.1109/TMI.2019.2893944>
- [12]. Kipf TN, Welling M (2016) Semi-supervised classification with graph convolutional networks. *International Conference on Learning Representations*.
- [13]. Chen J, Lu Y, Yu Q, et al (2021) TransUNet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*.
- [14]. Du A, Wang L, Feng Meng C, et al (2021) Graph convolutional network for skin lesion classification. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 19:1506–1515. <https://doi.org/10.1109/TCBB.2021.3055743>
- [15]. Wang S, Chen S, Wang Y, et al (2021) Boundary-aware context neural network for medical image segmentation. *Medical Image Analysis* 73:102144. <https://doi.org/10.1016/j.media.2021.102144>
- [16]. Selvaraju RR, Cogswell M, Das A, et al (2017) Grad-CAM: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE International Conference on Computer Vision* 618–626. <https://doi.org/10.1109/ICCV.2017.74>
- [17]. Mahbod A, Schaefer G, Wang C, et al (2019) Skin lesion classification using hybrid deep neural networks. *IEEE International Conference on Acoustics, Speech and Signal Processing* 1229–1233. <https://doi.org/10.1109/ICASSP.2019.8683352>
- [18]. Al-Masni MA, Kim DH, Tzallas AT, et al (2020) Multiple skin lesions diagnostics via integrated deep convolutional networks for medical image automation. *Computer Methods and Programs in Biomedicine* 190:105351. <https://doi.org/10.1016/j.cmpb.2020.105351>
- [19]. Cassidy B, Kendrick C, Brodzicki A, et al (2022) Analysis of the ISIC image datasets: Usage, benchmarks and recommendations. *Medical Image Analysis* 75:102285. <https://doi.org/10.1016/j.media.2021.102285>
- [20]. Adegun AA, Viriri S (2021) Deep learning techniques for skin lesion analysis and melanoma cancer detection: A survey of state-of-the-art. *Artificial Intelligence Review* 54:811–841. <https://doi.org/10.1007/s10462-020-09865-y>

