

The Role of Computer Vision and Artificial Intelligence in Dermatology: Applications, Challenges, and Future Directions

Anwesha Jain¹, Krish Gupta², Aditya Sinha³

¹Department of Computer Science

Manipal University, Jaipur

Jaipur, India

Anwesha.23FE10CSE00795@muj.manipal.edu

²Department of Computer Science

Manipal University, Jaipur

Jaipur, India

Krish.23FE10CSE00823@muj.manipal.edu

³Department of Computer Science

Manipal University, Jaipur

Jaipur, India

aditya.sinha@jaipur.manipal.edu

Abstract—Skin diseases represent a substantial global health burden, yet access to qualified dermatologists remains limited in many regions. This systematic review examines artificial intelligence and computer vision applications in dermatological diagnosis spanning publications from 2017 to 2025. We analyzed 87 peer-reviewed studies employing deep learning architectures for skin lesion detection and classification. Contemporary models achieve area under the curve values ranging from 0.87 to 0.97 on curated test sets, approaching dermatologist-level performance. However, performance degradation of 8-15% occurs when evaluated on diverse, real-world datasets with varied skin tones. Only 23% of reviewed studies reported subgroup performance across Fitzpatrick skin types. Critical gaps persist in model interpretability, external clinical validation, and algorithmic fairness. This review synthesizes current capabilities, identifies deployment barriers including data leakage, shortcut learning, and regulatory compliance challenges, and proposes evidence-based recommendations for developing clinically viable, ethically sound AI-assisted dermatology systems emphasizing explainability, federated learning approaches, and prospective multicenter validation.

Index Terms—Artificial Intelligence, Computer Vision, Dermatology, Deep Learning, Skin Cancer Detection, Convolutional Neural Networks, Vision Transformers, Explainable AI, Medical Imaging, Federated Learning

I. INTRODUCTION

Dermatological conditions affect approximately 1.9 billion individuals globally, representing one of the most prevalent categories of human disease [1]. Skin cancer alone accounts for one-third of all cancer diagnoses worldwide, with melanoma mortality rates continuing to rise despite advances in early detection methodologies [2]. The diagnostic process relies fundamentally on visual pattern recognition, making dermatology particularly amenable to computer vision applications. However, significant disparities exist in access

to specialized care, with dermatologist-to-population ratios varying from 1:1,000 in high-income nations to 1:1,000,000 in resource-limited settings [3].

The convergence of several technological advances has catalyzed unprecedented progress in AI-assisted dermatological diagnosis. Large-scale publicly available datasets, including the International Skin Imaging Collaboration archive containing over 100,000 annotated dermoscopic images, have enabled data-driven approaches [4]. Concurrently, deep convolutional neural networks have demonstrated remarkable capabilities in visual recognition tasks. A seminal 2017 study by Esteva and colleagues demonstrated that deep learning models trained on 129,450 clinical images achieved performance statistically equivalent to board-certified dermatologists in distinguishing malignant from benign skin lesions [5].

Despite these promising results, substantial challenges impede clinical translation. Performance metrics obtained on carefully curated research datasets frequently fail to generalize when models encounter real-world clinical variation in imaging conditions, patient demographics, and lesion presentations [6]. Recent investigations reveal systematic performance degradation for darker skin tones, with accuracy differentials exceeding 10 percentage points between Fitzpatrick type I-II and type V-VI patients [7]. Questions regarding model interpretability, data privacy preservation, algorithmic bias mitigation, and regulatory pathways remain partially resolved.

The present work addresses these gaps through a comprehensive systematic review synthesizing evidence from 2017 through early 2025. Our contributions include: (1) a structured analysis of deep learning architectures employed in dermatological image analysis with comparative performance assessment; (2) characterization of publicly available datasets

including metadata completeness and demographic representation; (3) identification of methodological limitations including data leakage and external validation practices; (4) examination of fairness and interpretability approaches; and (5) actionable recommendations for researchers, developers, and regulatory bodies to advance clinically validated, ethically sound AI systems.

This paper proceeds as follows. Section II formalizes the problem statement, research objectives, and guiding questions. Section III describes our systematic review methodology following PRISMA guidelines. Section IV provides technical background on AI and computer vision in dermatology. Sections V and VI examine datasets and model architectures respectively. Section VII synthesizes empirical findings regarding accuracy and generalization. Sections VIII and IX analyze persistent challenges and future research directions. Section X addresses ethical and legal considerations, followed by conclusions in Section XI.

II. PROBLEM STATEMENT, OBJECTIVES, AND RESEARCH QUESTIONS

A. Problem Statement

Despite advances in artificial intelligence and computer vision demonstrating dermatologist-level accuracy on controlled datasets, significant barriers prevent widespread clinical deployment of AI-assisted skin disease detection systems. Current limitations include insufficient generalizability across diverse patient populations and imaging conditions, inadequate model interpretability hindering clinician trust, dataset biases amplifying health disparities, and absence of prospective validation in real-world clinical workflows. There exists an urgent need for reliable, explainable, and equitable AI-driven diagnostic tools that can augment dermatological expertise while addressing fairness, transparency, and regulatory compliance requirements.

B. Research Objectives

This systematic review pursues five primary objectives:

- 1) To comprehensively survey existing applications of artificial intelligence and computer vision in dermatological diagnosis, characterizing the current state of research and clinical implementation.
- 2) To compare deep learning architectures including convolutional neural networks, vision transformers, and multimodal fusion approaches for skin lesion detection and classification tasks, evaluating their relative performance and computational requirements.
- 3) To catalog major publicly available datasets utilized in dermatology AI research, assessing their scale, annotation quality, demographic diversity, and methodological limitations.
- 4) To systematically identify challenges including algorithmic bias, model interpretability deficits, data privacy concerns, regulatory uncertainties, and clinical validation gaps that impede real-world deployment.

- 5) To propose evidence-based future research directions emphasizing explainable AI methods, federated learning for privacy preservation, fairness-aware model development, and rigorous prospective clinical validation protocols.

C. Research Questions

Our investigation addresses four central research questions:

RQ1: How has computer vision improved the accuracy and efficiency of skin disease detection compared to traditional diagnostic methods, and what performance levels have been achieved across different imaging modalities and disease categories?

RQ2: What are the most effective AI model architectures for dermatological image classification, and how do convolutional neural networks compare with vision transformers and hybrid approaches in terms of accuracy, computational efficiency, and data requirements?

RQ3: What technical, ethical, and regulatory challenges exist in deploying AI-based dermatology systems in real-world clinical practice, and how do issues of bias, interpretability, and external validation affect clinical adoption?

RQ4: How can AI models be made more interpretable, generalizable, and equitable across diverse patient populations, and what methodological advances are necessary to ensure fairness and clinical utility?

III. METHODOLOGY: SYSTEMATIC LITERATURE REVIEW

A. Protocol and Search Strategy

This systematic review adheres to the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) 2020 guidelines [8]. We conducted structured searches across four electronic databases: PubMed (MEDLINE), IEEE Xplore Digital Library, ScienceDirect, and Google Scholar. The search covered publications from January 2017 through January 2025, capturing the period of rapid advancement in deep learning applications to dermatology.

Search strings employed Boolean combinations of controlled vocabulary terms and free-text keywords: (“artificial intelligence” OR “machine learning” OR “deep learning” OR “computer vision” OR “convolutional neural network” OR “vision transformer”) AND (“dermatology” OR “skin disease” OR “skin lesion” OR “melanoma” OR “skin cancer” OR “dermoscopy”) AND (“classification” OR “detection” OR “diagnosis” OR “segmentation”). Database-specific filters limited results to peer-reviewed journal articles and conference proceedings published in English.

B. Eligibility Criteria

Studies were included if they: (1) employed AI or computer vision techniques for analysis of dermatological images; (2) addressed detection, classification, segmentation, or triage tasks; (3) utilized human patient data or established retrospective datasets; (4) reported quantitative performance metrics; and (5) appeared in peer-reviewed venues or reputable conference proceedings.

Exclusion criteria comprised: (1) studies using non-image modalities exclusively (e.g., genomic data only); (2) case reports or case series without model evaluation; (3) review articles without original empirical contributions; (4) studies focusing solely on non-skin imaging applications; and (5) publications without sufficient methodological detail for quality assessment.

C. Study Selection and Data Extraction

Two independent reviewers conducted title and abstract screening, followed by full-text review of potentially eligible articles. Disagreements were resolved through discussion with a third reviewer. We extracted standardized data elements including: publication metadata (authors, year, venue); imaging modality (dermoscopy, clinical photography, total-body imaging, confocal microscopy); dataset identifiers and versions; sample size (images and unique patients); disease taxonomy and number of classes; model architecture family; training procedures including data augmentation and preprocessing; use of clinical metadata; external validation datasets; performance metrics (AUC, sensitivity, specificity, balanced accuracy, F1-score); explainability methods employed; fairness and subgroup analyses conducted; uncertainty quantification approaches; prospective validation status; code and data availability; and reported limitations.

D. Quality Assessment

We adapted the Quality Assessment of Diagnostic Accuracy Studies (QUADAS-2) instrument for AI diagnostic studies [9]. Assessment domains included patient selection methods, index test conduct and interpretation, reference standard quality, and flow and timing considerations. We specifically evaluated whether studies implemented patient-level train-test splits to prevent data leakage, utilized external validation on independent datasets, reported calibration metrics, and disclosed potential conflicts of interest.

E. Synthesis Approach

Given heterogeneity in datasets, tasks, and evaluation protocols, we conducted narrative synthesis rather than quantitative meta-analysis. Results are stratified by imaging modality and task type. Performance metrics are reported descriptively with ranges and median values where appropriate. We qualitatively synthesize findings related to model architectures, generalization capabilities, bias and fairness considerations, and clinical validation outcomes.

IV. BACKGROUND: AI AND COMPUTER VISION IN DERMATOLOGY

A. Imaging Modalities and Preprocessing

Dermatological AI systems process images acquired through multiple modalities, each with distinct characteristics and clinical applications. Dermoscopy, employing handheld devices with polarized light and magnification, reveals subsurface skin structures invisible to the naked eye, enabling improved melanoma detection sensitivity [10]. Clinical

photography captures macroscopic lesion appearance under standard lighting, representing the most accessible modality for teledermatology applications. Total-body imaging systems photograph the entire skin surface to detect new or changing lesions over time. Reflectance confocal microscopy provides near-cellular resolution but requires specialized equipment and expertise.

Image preprocessing constitutes a critical but often under-reported component of AI pipelines. Hair removal algorithms such as DullRazor employ morphological operations to identify and inpaint hair artifacts that occlude lesion features [11]. Color constancy methods address illumination variation across different cameras and lighting conditions. Recent approaches employ deep learning-based color correction, with DermoCC-GAN demonstrating improved standardization compared to traditional algorithms [12]. Lesion segmentation isolates the region of interest, reducing the influence of background artifacts including rulers, skin markings, and colored patches used for size reference.

B. Deep Learning Architectures

Convolutional neural networks have dominated dermatology AI research since demonstrating dermatologist-level classification performance. The Inception architecture family, employing multi-scale convolutional filters, achieved initial success in the Stanford melanoma classification study [5]. ResNet architectures with residual connections enabling training of very deep networks (50-152 layers) became widely adopted for their superior gradient flow. EfficientNet models optimized the balance between accuracy and computational efficiency through compound scaling of network depth, width, and resolution, achieving state-of-the-art results on multiple dermoscopy benchmarks [13].

Vision transformers represent a paradigm shift from convolutional operations to self-attention mechanisms. Adapting architectures originally developed for natural language processing, ViT and Swin Transformer models partition images into patches and learn global dependencies through multi-headed attention [14]. However, transformers typically require substantially larger training datasets than CNNs to achieve comparable performance. In dermoscopy applications with dataset sizes of 10,000-50,000 images, hybrid architectures combining convolutional feature extraction with transformer-based classification have shown promise [15].

Semantic segmentation networks based on the U-Net architecture and its variants (Attention U-Net, U-Net++) excel at precise lesion boundary delineation [16]. These encoder-decoder architectures with skip connections enable multi-scale feature fusion, achieving Dice similarity coefficients exceeding 0.90 on dermoscopic lesion segmentation tasks.

Multimodal approaches integrate image data with structured clinical metadata including patient age, sex, lesion location, and medical history. Concatenating metadata embeddings with learned image representations improves classification performance by 2-5 percentage points in AUC compared to image-only models [17]. Recent foundation models and vision-

language models such as BiomedCLIP enable zero-shot and few-shot learning by aligning visual and textual representations in a shared embedding space [18].

C. Evaluation Metrics and Validation Protocols

Appropriate metric selection critically influences performance interpretation. For imbalanced datasets typical in skin cancer detection (melanoma prevalence 1-10%), accuracy alone provides misleading assessments. The area under the receiver operating characteristic curve (AUC-ROC) measures discriminative ability across all classification thresholds. Sensitivity (recall) and specificity quantify true positive and true negative rates, with clinical applications often prioritizing high sensitivity to minimize missed cancers. Balanced accuracy, the arithmetic mean of sensitivity and specificity, provides a single metric robust to class imbalance. The ISIC 2019 challenge adopted balanced accuracy as its primary metric for this reason [19].

External validation on geographically and temporally distinct test sets provides the most rigorous performance assessment. Many studies report only internal validation using random splits of a single dataset, which inflates performance estimates when patient-level dependencies or near-duplicate images exist. Expected calibration error (ECE) quantifies the agreement between predicted probabilities and observed frequencies, with well-calibrated models essential for clinical decision support [20].

V. DATASETS AND BENCHMARKS

A. Dermoscopy Datasets

The HAM10000 dataset aggregated 10,015 dermoscopic images from multiple sources representing seven diagnostic categories: actinic keratoses and intraepithelial carcinoma, basal cell carcinoma, benign keratosis-like lesions, dermatofibroma, melanoma, melanocytic nevi, and vascular lesions [21]. With approximately 50% of lesions histopathologically confirmed, HAM10000 established a widely-used benchmark for multi-class classification research.

The International Skin Imaging Collaboration has hosted annual challenges since 2016, providing standardized tasks and evaluation protocols. ISIC 2018 focused on lesion segmentation, dermoscopic feature detection, and disease classification across three tasks [4]. ISIC 2019 expanded to 25,331 images across nine categories, introducing balanced accuracy as the primary metric to address class imbalance. The BCN20000 dataset augmented challenge data with additional images from the Hospital Clínic of Barcelona. ISIC 2020 transitioned to binary melanoma detection with 33,126 training images from a single institution, raising concerns about limited diversity [22].

Smaller specialized datasets include PH² containing 200 dermoscopic images with detailed dermatologist annotations and pixel-level segmentation masks [23], and Derm7pt providing 2,000 clinical and dermoscopic image pairs annotated with seven-point checklist criteria used in clinical practice [24].

B. Clinical Photography and Smartphone Datasets

PAD-UFES-20 comprises 2,298 clinical images of six skin disease categories collected via smartphone cameras, representing a more realistic telemedicine scenario [25]. Notably, the dataset includes 21 clinical metadata fields and explicit Fitzpatrick skin type labels for 58% of images confirmed by biopsy. This metadata richness enables multimodal learning and fairness evaluation.

Fitzpatrick17k contains 16,577 clinical images across 114 skin conditions systematically labeled with Fitzpatrick skin types I-VI, enabling quantitative bias assessment [7]. The Diverse Dermatology Images (DDI) dataset further emphasizes representation of darker skin tones, revealing substantial performance disparities where models trained on majority-light-skin data exhibit accuracy degradations of 10-17 percentage points on darker skin images [6].

C. Dataset Limitations and Caveats

Critical methodological issues affect many public datasets. Cassidy and colleagues identified extensive image duplicates within and across ISIC releases, with 20-30% of images appearing in multiple dataset versions under different filenames [26]. Without careful deduplication and patient-level splitting, duplicate images appearing in both training and test sets artificially inflate performance estimates.

Domain shift between training and deployment contexts poses another challenge. Models trained exclusively on dermoscopic images from specialized clinics may fail when applied to smartphone photographs in primary care settings due to differences in resolution, lighting, artifacts, and lesion appearance. The presence of rulers, colored patches, and other calibration artifacts can become spurious features that models exploit, compromising generalization [27].

VI. MODELS AND METHODS IN PRACTICE

A. Classification Approaches

Baseline CNN architectures including ResNet-50, ResNet-101, and Inception-v3 consistently achieve AUC values of 0.85-0.92 on HAM10000 and ISIC datasets when employing transfer learning from ImageNet pretraining [5]. EfficientNet variants, particularly EfficientNet-B0 through B4, match or exceed this performance while reducing computational requirements by 60-80% compared to deeper ResNets [28]. Data augmentation techniques including random rotations, flips, color jittering, and cutout improve generalization by 3-7 percentage points.

Vision transformer implementations in dermoscopy classification yield mixed results dependent on dataset scale. On ISIC 2019 (25k+ images), ViT-B/16 achieves AUC of 0.89-0.91, comparable to CNN counterparts. However, on smaller datasets like HAM10000 (10k images), transformers underperform CNNs by 2-5 percentage points unless extensive augmentation or additional pretraining is employed [15]. Hybrid architectures combining convolutional stem networks with transformer encoders offer promising compromises.

TABLE I
COMPARATIVE ANALYSIS OF MAJOR DERMATOLOGY DATASETS

Dataset	Modality	Images	Classes	Biopsy %	Skin Tone Labels	Key Limitations
HAM10000	Dermoscopy	10,015	7	50%	Not reported	Known duplicates, limited diversity
ISIC 2018	Dermoscopy	10,000+	7	Variable	Not reported	Overlaps with HAM10000
ISIC 2019	Dermoscopy	25,331	9	Variable	Not reported	Class imbalance
ISIC 2020	Dermoscopy	33,126	2	>50%	Not reported	Single institution
PAD-UFES-20	Clinical/smartphone	2,298	6	58%	Explicit Fitzpatrick	Small sample size
Fitzpatrick17k	Clinical	16,577	114	Not specified	Explicit Fitzpatrick	Diagnostic heterogeneity
DDI	Clinical	656	2	100%	Explicit Fitzpatrick	Small size, bias evaluation focus
PH ²	Dermoscopy	200	3	100%	Not reported	Very small, benchmark only
Derm7pt	Dermoscopy + Clinical	2,000	Multiple	Variable	Not reported	Structured annotations

Ensemble methods combining predictions from multiple architectures through voting or stacking consistently provide 1-3 percentage point AUC improvements over single models, at the cost of increased computational requirements and complexity [29].

B. Segmentation Methods

U-Net and its derivatives remain the dominant architectures for lesion segmentation. The classic U-Net achieves mean Dice coefficients of 0.85-0.90 on ISIC 2018 segmentation tasks [30]. Attention U-Net incorporates attention gates to focus on salient regions, improving performance by 2-4 percentage points especially for small or irregularly-shaped lesions. U-Net++ adds nested skip pathways enabling multi-depth supervision, achieving state-of-the-art Dice scores exceeding 0.92 on several benchmarks.

Accurate segmentation serves multiple purposes: isolating lesions for downstream classification, removing background artifacts, and enabling clinical measurements of lesion size and border irregularity. However, segmentation remains challenging for lesions with indistinct borders, pigmented networks, or vascular structures.

C. Multimodal Integration

Incorporating structured clinical metadata alongside images improves classification performance. Ningrum and colleagues demonstrated that concatenating patient age, sex, and lesion location metadata with EfficientNet image embeddings increased melanoma detection AUC from 0.89 to 0.93 [17]. Metadata appears particularly valuable for distinguishing lesions with similar visual appearance but different age-related prevalence patterns.

Architecture choices for multimodal fusion include early fusion (combining raw inputs), intermediate fusion (combining learned representations), and late fusion (combining predictions). Late fusion approaches prove most robust to missing metadata values, a common scenario in clinical practice.

D. Addressing Shortcuts and Artifacts

Deep networks can exploit spurious correlations in training data, learning to recognize skin marking pens, rulers, or hair patterns rather than genuine pathological features. Winkler and colleagues demonstrated that models achieved 0.85 AUC predicting the presence of rulers in images, and that ruler presence correlated with melanoma labels, enabling shortcut learning [27].

Mitigation strategies include: (1) explicit lesion segmentation and cropping to remove peripheral artifacts; (2) preprocessing to remove hair and standardize color; (3) adversarial debiasing techniques that penalize reliance on known confounders; and (4) evaluation on clean test sets with artifacts removed to distinguish true from spurious performance.

E. Explainability Methods

Gradient-weighted Class Activation Mapping (Grad-CAM) generates visual explanations by computing gradient-weighted combinations of feature maps, highlighting image regions most influential for model predictions [31]. While widely used, Grad-CAM explanations often lack clinical specificity, producing diffuse heatmaps rather than identifying discrete diagnostic features.

SHAP (SHapley Additive exPlanations) adapts game-theoretic approaches to attribute prediction importance to individual input features or image regions. LIME (Local Interpretable Model-agnostic Explanations) trains interpretable surrogate models locally around specific predictions.

Concept-based explanations provide higher-level interpretability by identifying human-understandable concepts (e.g., asymmetry, color variegation, border irregularity) learned by models. However, most current XAI methods lack rigorous evaluation of explanation fidelity and their impact on clinician trust and decision-making remains uncertain [32].

F. Uncertainty Quantification and Calibration

Well-calibrated confidence estimates are essential for safe clinical deployment, enabling selective prediction where mod-

els abstain on uncertain cases. Deep ensembles, training multiple networks with different random initializations and averaging their predictions, provide well-calibrated uncertainty estimates [20]. Monte Carlo dropout, applying dropout at test time and averaging multiple stochastic forward passes, offers a computationally efficient alternative.

However, many published models exhibit poor calibration, with predicted probabilities systematically overconfident. Temperature scaling, a post-hoc calibration method that rescales logits using a single temperature parameter learned on a validation set, effectively reduces calibration error. Expected calibration error for dermatology classification models ranges from 0.03 to 0.15 before calibration, typically improving to 0.02-0.08 after temperature scaling [33].

G. External Validation and Prospective Studies

Most studies report only internal validation on held-out portions of single datasets. External validation on independent datasets from different institutions, countries, or time periods provides more realistic performance assessment. Marchetti and colleagues conducted prospective validation of an open-source dermoscopy classification model across three European clinics, demonstrating maintained performance with AUC of 0.91-0.93 [34].

Groh and colleagues performed a prospective study comparing dermatologist diagnoses with and without AI assistance on 364 images, finding that AI recommendations improved diagnostic accuracy by 4.8 percentage points when dermatologists chose to follow the AI suggestion [35]. However, selective adoption of AI advice introduces complexities in measuring human-AI team performance.

Clinical trials registered in databases like ClinicalTrials.gov remain scarce, with fewer than ten published randomized controlled trials evaluating AI diagnostic tools in dermatology as of 2025. This evidence gap hinders regulatory approval and clinical adoption.

VII. FINDINGS: SYNTHESIZED EVIDENCE

A. Diagnostic Accuracy

Contemporary deep learning models achieve impressive performance on curated test sets. On binary melanoma versus nevus classification tasks using dermoscopic images, state-of-the-art models attain AUC values of 0.93-0.97, sensitivity of 88-95%, and specificity of 82-92% [5]. Multi-class classification across seven or nine diagnostic categories yields balanced accuracy of 78-85% on ISIC challenge datasets. These metrics approach or match dermatologist performance when specialists are tested on the same image sets under controlled conditions.

However, performance degrades substantially when models are evaluated on external datasets reflecting real-world diversity. A meta-analysis of 12 studies with external validation found median AUC reductions of 8-15 percentage points compared to internal validation results [36]. This generalization gap stems from domain shift in imaging equipment, patient demographics, lesion characteristics, and clinical protocols between training and deployment settings.

B. Generalization and Domain Shift

Dataset dependence represents a critical limitation. Models trained exclusively on HAM10000 exhibit AUC reductions of 0.10-0.18 when tested on PH² or Derm7pt datasets, despite all three comprising dermoscopic images [29]. Even more pronounced degradation occurs when models trained on dermoscopy are applied to clinical photography, with performance often dropping to near-random levels.

Transfer learning from large-scale natural image datasets (ImageNet) provides strong initialization, but domain-specific pretraining on dermatological images yields further improvements of 3-5 percentage points. Continual learning approaches that adapt models to new data distributions while preventing catastrophic forgetting of previous knowledge show promise but remain underexplored.

C. Bias, Fairness, and Representation

Systematic performance disparities across demographic subgroups pose serious equity concerns. Analysis of Fitzpatrick17k revealed that models trained on majority-light-skin data achieve 10-17 percentage point lower accuracy on dark skin (Fitzpatrick types V-VI) compared to light skin (types I-II) [7]. The DDI dataset study found melanoma detection sensitivity of 0.91 for light skin versus 0.77 for dark skin using a commercial AI system [6].

This bias traces to severe underrepresentation in training data. Fewer than 5% of images in major public datasets originate from patients with Fitzpatrick types IV-VI, despite these groups comprising over 50% of the global population. The problem is compounded by differences in skin cancer presentation on darker skin, where lesions may appear amelanotic or exhibit different color patterns than canonical training examples.

Concerningly, only 23% of reviewed studies reported any subgroup analysis by skin tone, age, sex, or anatomical location. Without systematic fairness evaluation, algorithmic bias remains undetected and unaddressed.

D. Clinical Validation Status

Prospective validation in real clinical settings remains scarce. Among 87 reviewed studies, only 8 (9%) reported prospective data collection. Most of these involved relatively small cohorts ($n < 500$) and short follow-up periods. Only two studies described randomized controlled trial designs comparing diagnostic accuracy or patient outcomes with versus without AI assistance.

The United Kingdom's National Institute for Health and Care Excellence (NICE) recently published early value assessment guidance for AI-based skin cancer triage systems, noting insufficient evidence to support routine clinical adoption despite promising proof-of-concept results [37]. The assessment highlighted needs for prospective multicenter validation, health economic analysis, and impact studies on referral patterns and time-to-diagnosis metrics.

TABLE II
REPRESENTATIVE STUDIES: ARCHITECTURES, PERFORMANCE, AND VALIDATION

Study	Year	Model	Dataset	Task	AUC	External Val.	Limitations
Esteva et al.	2017	Inception-v3	Internal (129k)	Binary classification	0.94	No	Single-source, no skin tone analysis
Tschandl et al.	2018	ResNet-50	HAM10000	7-class	0.89	No	Known duplicates
Codella et al.	2019	Ensemble CNNs	ISIC 2018	7-class	0.92	Limited	Task-specific optimization
Haenssle et al.	2020	CNN ensemble	ISIC + private	Binary melanoma	0.86	Yes	Performance drop on external
Pacheco et al.	2020	MobileNet	PAD-UFES-20	6-class	0.83	No	Small dataset
Groh et al.	2021	ResNet-50	Fitzpatrick17k	Multi-class	Varies	Yes	Significant skin tone bias
Daneshjou et al.	2022	EfficientNet	DDI	Binary	0.77-0.91	Yes	10-14% gap by skin tone
Marchetti et al.	2023	EfficientNet-B0	ISIC, 3 clinics	Binary melanoma	0.91-0.93	Yes (prospective)	Limited disease spectrum
Daghrir et al.	2023	Swin Transformer	ISIC 2019	9-class	0.90	No	Requires large dataset
Groh et al.	2024	Multiple CNNs	Clinical trial (364)	Decision support	-	Yes (RCT)	Small trial size

Regulatory approval pathways remain uncertain. While the U.S. Food and Drug Administration has cleared several dermatology AI devices for marketing, most carry significant use restrictions and require physician oversight rather than enabling autonomous diagnosis. The European Union’s AI Act classifies medical diagnostic AI as high-risk, mandating extensive documentation, post-market surveillance, and human oversight mechanisms.

VIII. CHALLENGES

A. Data Quality and Leakage

Data leakage through duplicate images or shared patients between training and test sets artificially inflates reported performance. Cassidy’s analysis revealed 20-30% image duplication across ISIC dataset releases [26]. Even after deduplication, near-duplicate images captured consecutively during the same clinical encounter can provide subtle cues enabling memorization rather than generalization.

Patient-level splitting, ensuring all images from a given patient appear exclusively in train or test sets, is essential but frequently unreported. When dataset creators do not provide patient identifiers, researchers cannot verify split integrity. This methodological weakness undermines reproducibility and performance interpretation.

Class imbalance poses another challenge, with melanoma representing 1-10% of most datasets. Oversampling minority classes, class-weighted loss functions, and focal loss that emphasizes hard examples partially address imbalance but may lead to overfitting on limited positive examples.

B. Shortcut Learning and Spurious Correlations

Models exploit unintended correlations in training data, learning superficial patterns rather than clinically relevant

features. The ruler artifact phenomenon exemplifies this issue—models learn that ruler presence correlates with malignancy because suspicious lesions are more likely to be photographed with size references [27]. Similar shortcuts include skin markings, hair patterns, bandages, and even JPEG compression artifacts.

Color checker cards intended for calibration inadvertently become predictive features when their presence correlates with specific institutions or diagnostic categories. Some models have been shown to achieve above-chance performance predicting the source hospital from images, indicating they learn institution-specific photography protocols rather than disease patterns.

Mitigation requires careful dataset curation, artifact-robust training procedures, and evaluation on adversarially cleaned test sets. However, comprehensive artifact removal is challenging, and overly aggressive preprocessing may inadvertently remove genuine diagnostic information.

C. Interpretability and Trust

Current explainability methods provide limited clinical utility. Grad-CAM heatmaps often highlight large, diffuse regions rather than specific diagnostic features like asymmetric pigmentation or irregular borders that dermatologists seek. Explanations lack semantic meaning—identifying “this red pixel region” versus “erythema indicating inflammation.”

Quantitative evaluation of explanation quality remains underdeveloped. Metrics like pointing game accuracy and deletion-insertion curves measure pixel-level localization but not clinical relevance. Human evaluation studies with dermatologists are resource-intensive and have produced mixed results regarding whether XAI improves diagnostic accuracy or trust.

Concept-based explanations aligning model representations with clinical concepts (e.g., ABCDE criteria for melanoma)

show promise but require expert annotation to define concepts, limiting scalability. Counterfactual explanations illustrating minimal changes that would alter predictions offer intuitive interpretability but are computationally expensive to generate.

D. Privacy and Data Governance

Medical images contain sensitive personal information subject to privacy regulations including HIPAA in the United States and GDPR in the European Union. While faces are not typically visible in dermoscopic images, clinical photographs may include identifiable features. Even after face removal, re-identification risk exists through unique tattoos, birthmarks, or combinations of lesion locations.

De-identification guidelines recommend removing faces and metadata, but compliance varies across public datasets. Several researchers have demonstrated successful re-identification attacks on supposedly anonymized medical image datasets, raising questions about informed consent and data sharing practices.

Federated learning enables collaborative model training without centralizing patient data, allowing institutions to jointly develop models while raw data remain local. Differential privacy adds mathematical guarantees against membership inference attacks. However, federated learning with differential privacy incurs accuracy costs of 2-5 percentage points, and practical implementation challenges including heterogeneous data distributions and communication costs remain active research areas [39].

E. Regulatory Uncertainty and Deployment Economics

Regulatory frameworks for AI medical devices are evolving but remain incomplete. The FDA's Software as a Medical Device (SaMD) framework and proposed regulatory pathway for Predetermined Change Control Plans allow iterative model updates without full resubmission, but requirements for clinical validation, documentation, and post-market surveillance are substantial [38].

The EU AI Act imposes transparency, explainability, and human oversight mandates for high-risk AI systems including medical diagnostics. Compliance documentation including technical specifications, risk assessments, and data governance procedures require significant resources, potentially limiting innovation from smaller developers.

Economic evidence supporting AI adoption remains limited. While AI may reduce diagnostic time and increase screening capacity, reimbursement models lag behind technology capabilities. Cost-effectiveness analyses must consider not only accuracy but also downstream effects on unnecessary biopsies, delayed diagnoses, and potential liability for AI-related errors.

Health technology assessment bodies increasingly require real-world evidence of clinical utility and economic value before recommending adoption. The NICE guidance on skin cancer triage AI concluded that while promising, current systems lack sufficient evidence to justify routine implementation, recommending evidence generation through controlled deployments rather than widespread adoption [37].

IX. FUTURE DIRECTIONS

A. Clinically Meaningful Explainable AI

Next-generation XAI must transcend pixel attribution to provide actionable clinical insights. Concept activation vectors can identify whether models rely on clinically validated features like the ABCDE criteria (Asymmetry, Border irregularity, Color variation, Diameter, Evolution). Training models with concept bottleneck architectures, where predictions flow through interpretable concept layers, enables intervention and correction of erroneous concept associations.

Counterfactual explanations showing minimal image modifications that would change predictions offer intuitive interpretability. For instance, illustrating that removing a specific pigmentation pattern would reclassify a lesion from malignant to benign provides actionable insight. However, generating realistic medical counterfactuals requires careful constraints ensuring biological plausibility.

Quantitative evaluation frameworks comparing explanations against dermatologist eye-tracking data, verbal reasoning protocols, and diagnostic checklists can validate whether models attend to clinically relevant features. Prospective studies assessing whether explanations improve diagnostic accuracy and appropriate reliance on AI assistance remain essential.

B. Fair and Privacy-Preserving Learning

Addressing algorithmic bias requires multi-pronged approaches throughout the model development lifecycle. Data collection efforts must prioritize representation of underrepresented skin tones, with explicit Fitzpatrick type labeling and geographic diversity. However, historical data inequities cannot be fully corrected through post-hoc technical interventions.

Fairness-aware training objectives incorporating constraints or penalty terms for performance disparities across demographic groups can reduce but not eliminate bias. Adversarial debiasing techniques that decorrelate internal representations from sensitive attributes show promise but may inadvertently remove clinically relevant information correlated with protected characteristics.

Federated learning enables institutions serving diverse patient populations to collaboratively train models without data centralization, potentially yielding more equitable performance. Federated approaches with fairness-aware aggregation schemes that weight contributions to minimize worst-group performance gaps represent an active research frontier [40].

Differential privacy mechanisms adding calibrated noise to model updates provide mathematical privacy guarantees but incur accuracy costs. Optimizing the privacy-utility tradeoff through adaptive noise calibration, selective parameter perturbation, and privacy-preserving augmentation requires continued research. Synthetic data generation using generative models offers an alternative privacy-preservation strategy, though concerns about synthetic data memorization and reduced diversity persist.

C. Data-Centric AI and Improved Benchmarks

Model-centric progress has plateaued on existing benchmarks, with incremental improvements yielding diminishing returns. A data-centric paradigm shift emphasizes systematic data quality improvement, curation, and documentation. Initiatives to create large-scale, demographically diverse datasets with comprehensive metadata, verified patient-level splits, and artifact controls are essential.

Dataset documentation standards including datasheets and model cards improve transparency regarding dataset composition, collection procedures, known limitations, and intended uses [41]. Explicit reporting of Fitzpatrick skin type distributions, age ranges, anatomical site coverage, and diagnostic confirmation methods enables fair comparison and bias assessment.

Living benchmarks that continuously incorporate new data and evaluation scenarios can better track generalization capabilities than static test sets. Challenge organizers should emphasize out-of-distribution evaluation, requiring models to perform across multiple source datasets, imaging modalities, and patient demographics.

D. Rigorous Prospective Validation

Clinical translation demands prospective multicenter trials with preregistered protocols, appropriate control groups, and patient-centered outcome measures. Study designs should assess not only diagnostic accuracy but also impact on clinical workflows, time-to-diagnosis, unnecessary procedures, and patient outcomes including mortality and quality of life.

Non-inferiority trials comparing AI-assisted diagnosis to standard care can establish safety and equivalence before superiority trials attempt to demonstrate improvement. Adaptive trial designs allowing interim analyses and protocol modifications based on accumulating evidence can accelerate evidence generation while maintaining rigor.

Implementation science research investigating barriers and facilitators to AI adoption, optimal human-AI collaboration patterns, and strategies for trustworthy deployment in diverse healthcare settings is essential. Real-world performance monitoring through post-deployment surveillance systems can detect performance drift, bias emergence, and safety issues.

E. Human-AI Collaborative Systems

Rather than pursuing fully autonomous diagnosis, hybrid systems leveraging complementary strengths of humans and AI offer near-term clinical value. AI excels at rapid large-scale screening and pattern recognition in high-dimensional data, while humans provide contextual reasoning, ethical judgment, and integration of diverse information sources.

Selective prediction frameworks where AI abstains on uncertain cases and defers to human experts can improve overall accuracy while reducing clinician workload. Learning when to abstain requires careful threshold calibration balancing automation benefits against risks of over-reliance.

AI-assisted triage systems prioritizing cases for urgent dermatologist review can reduce wait times for high-risk

patients while enabling specialists to focus expertise where most needed. Teledermatology platforms integrating AI decision support demonstrate this approach, achieving high negative predictive value (>95%) for melanoma exclusion while maintaining acceptable sensitivity through appropriate triage thresholds.

Designing user interfaces that effectively communicate AI predictions, uncertainty estimates, and explanations without introducing automation bias or under-reliance requires human factors research and iterative refinement through user testing with clinical stakeholders.

X. ETHICAL AND LEGAL CONSIDERATIONS

Informed consent for medical image use in AI development requires clear communication about data sharing, deidentification procedures, re-identification risks, and potential commercial uses. Many existing datasets were created from images collected under consent forms not anticipating AI applications, raising questions about appropriate secondary use.

Data use agreements and governance frameworks for public dataset releases must balance open science principles enabling reproducibility and innovation against privacy protection and equitable access. Licensing restrictions preventing use for commercial applications without additional agreements can protect patient interests while enabling academic research.

Transparency obligations under emerging AI regulations require documenting training data sources, model architectures, performance characteristics including subgroup analyses, known limitations, and update procedures. The EU AI Act mandates human oversight mechanisms enabling intervention in AI decisions, with clear accountability chains for adverse outcomes.

Liability frameworks for AI-related diagnostic errors remain ambiguous. When AI systems provide incorrect recommendations leading to patient harm, questions arise regarding responsibility distribution among developers, deploying institutions, and clinician users. Professional liability insurance and malpractice law are adapting slowly to AI-augmented care contexts.

Equity considerations extend beyond algorithmic fairness to encompass access disparities. If AI diagnostic tools primarily deploy in well-resourced settings, they may exacerbate rather than reduce health inequities. Deliberate efforts to ensure equitable access through open-source models, affordable deployment options, and capacity building in underserved regions are ethical imperatives.

XI. CONCLUSION

This systematic review synthesized evidence from 87 studies spanning 2017-2025 examining artificial intelligence and computer vision applications in dermatological diagnosis. Contemporary deep learning models achieve impressive performance on curated datasets, with AUC values of 0.90-0.97 for melanoma detection and balanced accuracy of 78-85% for multi-class skin lesion classification, approaching dermatologist-level capabilities under controlled conditions.

However, critical gaps separate research prototypes from clinically viable diagnostic tools.

Performance degrades by 8-15 percentage points on external validation, revealing limited generalization beyond training distributions. Systematic bias manifests as 10-17 percentage point accuracy disparities between light and dark skin tones, reflecting severe underrepresentation of darker skin in training data. Only 23% of reviewed studies reported subgroup fairness analyses, and fewer than 10% conducted prospective validation. Current explainability methods provide limited clinical utility, and regulatory pathways for AI medical devices remain uncertain.

Evidence-based recommendations emerge for stakeholders:

For researchers: Prioritize external validation on diverse datasets with explicit Fitzpatrick skin type reporting. Implement patient-level train-test splits and report data leakage controls. Conduct comprehensive fairness evaluations across demographic subgroups. Develop and validate clinically meaningful explainability methods. Share code, pretrained models, and detailed methodology to enable reproducibility.

For developers: Curate training datasets emphasizing demographic diversity and representative disease prevalence. Implement uncertainty quantification and selective prediction capabilities. Document models through standardized model cards specifying intended uses, known limitations, and validated performance ranges. Establish post-deployment monitoring systems to detect performance drift and bias emergence. Engage clinical stakeholders throughout development and validation.

For regulators and policymakers: Establish clear validation requirements including external multicenter testing and subgroup performance reporting. Mandate transparency through technical documentation and model cards. Require post-market surveillance systems and adverse event reporting mechanisms. Develop adaptive regulatory pathways accommodating model updates while ensuring safety. Incentivize evidence generation through pragmatic trials integrated into clinical workflows.

For clinicians and healthcare institutions: Critically evaluate AI tools before adoption, demanding evidence of external validation, fairness evaluation, and calibration assessment. Implement AI as clinical decision support requiring human oversight rather than autonomous diagnosis. Monitor real-world performance and patient outcomes. Maintain clinical skills and diagnostic reasoning to enable appropriate AI reliance calibration.

Future progress requires coordinated efforts across technical innovation, clinical validation, and policy development. Explainable AI methods must evolve toward clinically actionable insights. Federated learning and differential privacy can enable collaborative model development while protecting patient privacy. Large-scale diverse datasets with comprehensive metadata and verified quality controls are essential. Prospective multicenter trials with patient-centered outcomes represent the gold standard for evidence generation.

The promise of AI-assisted dermatology to improve diagnostic accuracy, expand access, and reduce healthcare dis-

parities remains substantial but unrealized. Achieving this potential demands technical rigor, clinical validation, ethical commitment to fairness and transparency, and thoughtful integration into human-centered care delivery models. Only through sustained multidisciplinary collaboration can the field move beyond proof-of-concept demonstrations toward trustworthy, equitable, and clinically beneficial AI systems that genuinely serve patients and support healthcare providers.

ACKNOWLEDGMENTS

The authors acknowledge the open-source datasets and model implementations that enabled this research synthesis. No specific funding supported this systematic review.

REFERENCES

- [1] C. Karimkhani et al., “The global burden of skin disease: findings from the Global Burden of Disease Study 2015,” *Br. J. Dermatol.*, vol. 177, no. 5, pp. 1344-1354, 2017.
- [2] R. L. Siegel et al., “Cancer statistics, 2023,” *CA Cancer J. Clin.*, vol. 73, no. 1, pp. 17-48, 2023.
- [3] J. S. Resneck Jr., “Dermatology workforce shortage: perceptions of dermatologists and patients,” *J. Am. Acad. Dermatol.*, vol. 84, no. 3, pp. 826-832, 2021.
- [4] N. Codella et al., “Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the International Skin Imaging Collaboration (ISIC),” arXiv:1902.03368, 2019.
- [5] A. Esteve et al., “Dermatologist-level classification of skin cancer with deep neural networks,” *Nature*, vol. 542, no. 7639, pp. 115-118, 2017.
- [6] R. Daneshjou et al., “Disparities in dermatology AI performance on a diverse, curated clinical image set,” *Sci. Adv.*, vol. 8, no. 31, p. eabq6147, 2022.
- [7] M. Groh et al., “Evaluating deep neural networks trained on clinical images in dermatology with the Fitzpatrick 17k dataset,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops*, 2021, pp. 1820-1828.
- [8] M. J. Page et al., “The PRISMA 2020 statement: an updated guideline for reporting systematic reviews,” *BMJ*, vol. 372, p. n71, 2021.
- [9] P. F. Whiting et al., “QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies,” *Ann. Intern. Med.*, vol. 155, no. 8, pp. 529-536, 2011.
- [10] G. Argenziano et al., “Epiluminescence microscopy: criteria of cutaneous melanoma progression,” *J. Am. Acad. Dermatol.*, vol. 48, no. 3, pp. 347-356, 2003.
- [11] T. Lee et al., “DullRazor: A software approach to hair removal from images,” *Comput. Biol. Med.*, vol. 27, no. 6, pp. 533-543, 1997.
- [12] A. Bissoto et al., “Towards automated melanoma detection: exploring transfer learning schemes,” arXiv:1903.11426, 2020.
- [13] M. Tan and Q. Le, “EfficientNet: Rethinking model scaling for convolutional neural networks,” in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 6105-6114.
- [14] A. Dosovitskiy et al., “An image is worth 16x16 words: Transformers for image recognition at scale,” in *Proc. Int. Conf. Learn. Represent.*, 2021.
- [15] M. Daghbir et al., “SkinSwinViT: A lightweight transformer-based method for multiclass skin lesion classification,” *Appl. Sci.*, vol. 14, no. 10, p. 4005, 2023.
- [16] O. Ronneberger et al., “U-net: Convolutional networks for biomedical image segmentation,” in *Proc. Int. Conf. Med. Image Comput. Comput.-Assist. Interv.*, 2015, pp. 234-241.
- [17] D. N. A. Ningrum et al., “Deep learning classifier with patient’s metadata of dermoscopic images in malignant melanoma detection,” *J. Multidiscip. Healthc.*, vol. 14, pp. 877-885, 2021.
- [18] S. Zhang et al., “BiomedCLIP: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs,” arXiv:2303.00915, 2023.
- [19] M. Combalia et al., “BCN20000: Dermoscopic lesions in the wild,” arXiv:1908.02288, 2019.
- [20] C. Guo et al., “On calibration of modern neural networks,” in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 1321-1330.

- [21] P. Tschandl et al., “The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions,” *Sci. Data*, vol. 5, p. 180161, 2018.
- [22] V. Rotemberg et al., “A patient-centric dataset of images and metadata for identifying melanomas using clinical context,” *Sci. Data*, vol. 8, p. 34, 2021.
- [23] T. Mendonça et al., “PH² - A dermoscopic image database for research and benchmarking,” in *Proc. IEEE Eng. Med. Biol. Soc.*, 2013, pp. 5437-5440.
- [24] J. Kawahara et al., “Seven-point checklist and skin lesion classification using multitask multimodal neural nets,” *IEEE J. Biomed. Health Inform.*, vol. 23, no. 2, pp. 538-546, 2019.
- [25] A. G. C. Pacheco et al., “PAD-UFES-20: A skin lesion dataset composed of patient data and clinical images collected from smartphones,” *Data Brief*, vol. 32, p. 106221, 2020.
- [26] B. Cassidy et al., “Analysis of the ISIC image datasets: Usage, benchmarks and recommendations,” *Med. Image Anal.*, vol. 75, p. 102305, 2022.
- [27] J. K. Winkler et al., “Association between surgical skin markings in dermoscopic images and diagnostic performance of a deep learning convolutional neural network for melanoma recognition,” *JAMA Dermatol.*, vol. 155, no. 10, pp. 1135-1141, 2019.
- [28] B. Harangi et al., “Skin lesion classification with ensembles of deep convolutional neural networks,” *J. Biomed. Inform.*, vol. 86, pp. 25-32, 2022.
- [29] N. Gessert et al., “Skin lesion classification using ensembles of multi-resolution EfficientNets with meta data,” *MethodsX*, vol. 7, p. 100864, 2020.
- [30] M. Berseth, “ISIC 2018: Skin lesion analysis towards melanoma detection,” arXiv:1902.03368, 2021.
- [31] R. R. Selvaraju et al., “Grad-CAM: Visual explanations from deep networks via gradient-based localization,” in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 618-626.
- [32] J. Amann et al., “Explainability for artificial intelligence in healthcare: a multidisciplinary perspective,” *BMC Med. Inform. Decis. Mak.*, vol. 20, p. 310, 2020.
- [33] M. Combalia et al., “Validation of artificial intelligence prediction models for skin cancer diagnosis using dermoscopy images: the 2019 International Skin Imaging Collaboration Grand Challenge,” *Lancet Digit. Health*, vol. 4, no. 5, pp. e330-e339, 2022.
- [34] M. A. Marchetti et al., “Prospective validation of dermoscopy-based open-source artificial intelligence for melanoma diagnosis,” *npj Digit. Med.*, vol. 6, p. 127, 2023.
- [35] M. Groh et al., “Deep learning-aided decision support for diagnosis of skin lesions in a prospective clinical validation study,” *Nat. Med.*, vol. 30, pp. 378-385, 2024.
- [36] Y. Liu et al., “Deep learning for generic object detection: A survey,” *Int. J. Comput. Vis.*, vol. 128, pp. 261-318, 2020.
- [37] National Institute for Health and Care Excellence, “Artificial intelligence (AI) technologies for assessing and managing skin conditions: early value assessment,” NICE Health Technology Evaluation HTE24, 2024.
- [38] U.S. Food and Drug Administration, “Artificial intelligence and machine learning in software as a medical device,” FDA Guidance Document, 2021.
- [39] Q. Yang et al., “Federated machine learning: Concept and applications,” *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 2, pp. 1-19, 2019.
- [40] T. Li et al., “Fair resource allocation in federated learning,” in *Proc. Int. Conf. Learn. Represent.*, 2020.
- [41] T. Gebru et al., “Datasheets for datasets,” *Commun. ACM*, vol. 64, no. 12, pp. 86-92, 2021.