
DERMATOLOGICAL DIAGNOSIS EXPLAINABILITY BENCHMARK FOR CONVOLUTIONAL NEURAL NETWORKS

A PREPRINT

Raluca Jalaboi^{1,2}, Ole Winther^{1,3,4}, and Alfiia Galimzianova²

¹Department of Applied Mathematics and Computer Science at the Technical University of Denmark, Richard Petersens Plads, Building 324, DK-2800 Kongens Lyngby, Denmark

²Medable A/S, Havnegade 25, 3., DK-1058 Copenhagen C, Denmark

³Bioinformatics Centre, Department of Biology, University of Copenhagen, Copenhagen, Denmark

⁴Center for Genomic Medicine, Rigshospitalet, Copenhagen University Hospital, Copenhagen, Denmark

ABSTRACT

In recent years, large strides have been taken in developing machine learning methods for various dermatological applications, supported in part by the widespread success of deep learning. To date, diagnosing diseases from images is one of the most explored applications of deep learning within dermatology. Convolutional neural networks (ConvNets) are the most commonly used deep learning method in medical imaging due to their training efficiency and accuracy, although they are often described as black boxes because of their limited explainability. One popular way to obtain insight into a ConvNet’s decision mechanism is gradient class activation maps (Grad-CAM). A quantitative evaluation of the Grad-CAM explainability has been recently made possible by the release of DermXDB, a skin disease diagnosis explainability dataset which enables benchmarking the explainability performance of ConvNet architectures. In this paper, we perform a literature review to identify the most common ConvNet architectures used for this task, and compare their Grad-CAM explainability performance with the explanation maps provided by DermXDB. We identified 11 architectures: DenseNet121, EfficientNet-B0, InceptionV3, InceptionResNetV2, MobileNet, MobileNetV2, NASNetMobile, ResNet50, ResNet50V2, VGG16, and Xception. We pre-trained all architectures on a clinical skin disease dataset, and then fine-tuned them on a subset of DermXDB. Validation results on the DermXDB holdout subset show an explainability F1 score of between 0.35-0.46, with Xception the highest explainability performance, while InceptionResNetV2, ResNet50, and VGG16 displaying the lowest. NASNetMobile reports the highest characteristic-level explainability sensitivity, despite its mediocre diagnosis performance. These results highlight the importance of choosing the right architecture for the desired application and target market, underline need for additional explainability datasets, and further confirm the need for explainability benchmarking that relies on quantitative analyses rather than qualitative assessments.

Keywords deep learning, dermatologists, explainability, benchmark, review

1 Introduction

With an expected shortage of approximately ten million healthcare professionals by 2030 [World Health Organization, 2016], the world is facing a massive healthcare crisis. Automation has been proposed as a solution to the scarcity of medical professionals, with the Food and Drugs Administration in the United States approving medical devices based on artificial intelligence for marketing to the public [U.S. Food and Drug Administration, 2018].

This development is due in part to the advancement in machine learning using unstructured data. Ever since Krizhevsky et al. [2017] won the ImageNet Large Scale Visual Recognition Challenge [Russakovsky et al., 2015] using a convolutional neural network (ConvNet), ConvNets have been at the forefront of machine learning based automation. Employed primarily in healthcare for imaging applications, ConvNets have been used for disease diagnosis [Gao et al., 2019], cell

counting [Falk et al., 2019], disease severity assessment [Gulshan et al., 2016], disease progression estimation [Kijowski et al., 2020], lesion or anatomical region segmentation [Hesamian et al., 2019, Ramesh et al., 2021], etc. Esteva et al. [2017] were the first to demonstrate that ConvNets can achieve expert-level performance in dermatological diagnosis using dermoscopy images. Since then, dermatology has embraced ConvNets as a solution to various diagnosis and segmentation tasks [Esteva et al., 2017, Zhang et al., 2019, Jinnai et al., 2020, Haenssle et al., 2020, Roy et al., 2022].

Despite these considerable advancements in medical imaging, there has not yet been a widespread adoption of machine learning based automation in the clinical workflow. One of the main hurdles that detract from adoption is the lack of ConvNet explainability [Kelly et al., 2019], this issue being enhanced by the recently implemented legislation aimed at ensuring that automated methods can offer an explanation into their decision mechanisms [Goodman and Flaxman, 2017]. Different post-hoc explainability methods have been proposed as a way to explain a ConvNet’s decisions [Bai et al., 2021, Selvaraju et al., 2017, Lundberg and Lee, 2017, Ribeiro et al., 2016]. Gradient class activation maps (Grad-CAM) is currently the most commonly used explainability method within medical imaging, due to its intrinsic ease of interpretation and its low computational requirements. However, validating the resulting explanations is an expensive, time consuming process that requires domain expert intervention, and thus most explainability validations are performed as small, qualitative analyses. With the release of DermXDB [Jalaboi et al., 2022], it became possible to quantitatively analyse the explainability of ConvNets trained for diagnosing six skin conditions: acne, psoriasis, seborrheic dermatitis, viral warts, and vitiligo.

The purpose of this benchmark is to provide the means to quantitatively compare the explainability of the state-of-the-art approaches to dermatological diagnosis using photographic imaging. Our contributions are twofold:

1. We perform a comprehensive systematic review to reveal the usage of the ConvNets for the task of dermatological diagnosis using photographic images,
2. We benchmark the identified ConvNets for diagnostic and explainability performance and compare them with eight expert dermatologists.

2 Background

2.1 Machine learning methods in dermatological diagnosis

After the renewed interest in artificial intelligence and machine learning that started in 2012, practitioners from both academia and the industry began investigating automated methods for dermatological applications [Thomsen et al., 2020, Jeong et al., 2022]. Until 2017, the vast majority of articles applying machine learning methods on dermatological problems were using classical models such as support vector machines [Liu et al., 2012, Sabouri et al., 2014], and linear or logistic regression [Kaur et al., 2015, Kefel et al., 2016]. These models were trained using hand-crafted features or features extracted using classical computer vision methods such as gray-level co-occurrence matrices [Shimizu et al., 2014], Sobel and Hessian filters [Arroyo and Zapirain, 2014], or HOS texture extraction [Shrivastava et al., 2016]. However, the main drawback of classical computer vision approaches is that hand-crafting features is an expensive, time-consuming process, while their automated extraction is too sensitive to the environmental factors of the image acquisition (e.g. lighting, zoom).

Esteva et al. [2017] were the first to propose a ConvNet for diagnosing skin conditions from dermoscopy images. Their ConvNet reached expert-level performance without requiring any hand-crafted features or classical computer vision models, thus paving the way towards the current popularity of ConvNets in dermatological applications.

One key component to the rise of ConvNets was the introduction of large scale dermatological datasets. The International Skin Imaging Collaboration (ISIC) challenge dataset [Codella et al., 2018] is one of the best known open access dermoscopy datasets, containing 25,331 images distributed over nine diagnostic categories. Large clinical image datasets are also available for research purposes, such as SD-260 [Sun et al., 2016] which consists of 20,600 clinical images of 260 different skin diseases, and DermNetNZ [DermNetNZ, 2021] which contains more than 25,000 clinical images.

Aided by the release of increasingly more performant architectures, their publicly available pre-trained weights on the ImageNet [Deng et al., 2009] dataset, and the recently published public dermatological datasets, the vast majority of research contributions in machine learning applications for dermatology rely on ConvNet architectures. ConvNets have been extensively used in lesion diagnosis [Tschandl et al., 2017, Han et al., 2018, Reshma et al., 2022] and lesion segmentation [Yuan et al., 2017, Wu et al., 2022, Baig et al., 2020] on different modalities relevant for the domain. Attempts at explaining the decisions taken by ConvNets were made by several groups [Tschandl et al., 2020, Tanaka et al., 2021], but no quantitative analysis was performed.

Table 1: Search query used on PubMed to identify the list of relevant articles. We searched for articles focused on dermatology, using deep learning methods, written in English. The query was last performed on the 20th of February 2023.

Search term	Search term	Search term
((dermatology[MeSH Terms]) OR (skin disease[MeSH Terms]) OR (skin lesion[MeSH Terms]))	AND ((neural network[MeSH Terms]) OR (machine learning[MeSH Terms]) OR (artificial intelligence[MeSH Terms]) OR (deep learning) OR (deep neural network) OR (convolutional neural network))	AND (English[Language])

2.2 Explainability in convolutional neural networks

ConvNets have, from their very beginning, been notoriously difficult to interpret and explain. Interpretability is generally considered the ability to understand the internal structure and properties of a ConvNet architecture, while explainability is defined as a ConvNet’s capacity to offer plausible arguments in favour of its decision [Roscher et al., 2020]. Within healthcare, explainability is especially important due to its intrinsic ability to interact with domain experts in a common vocabulary [Kelly et al., 2019]. Although some architecture or domain-specific explainability methods exist, most medical imaging research articles employ attribution-based methods due to their ease of use and open source access [Singh et al., 2020, Bai et al., 2021].

There are two main ways of implementing attribution-based methods: through perturbation and by using the ConvNet’s gradients. Perturbation-based methods, such as Shapley values [Lipovetsky and Conklin, 2001], LIME [Ribeiro et al., 2016], or SharpLIME [Graziani et al., 2021], rely on modifying the original image and then evaluating the changes in the ConvNet’s prediction. For example, LIME uses a superpixel algorithm to split the image into sections, and randomly selects a subset of superpixels to occlude. The target ConvNet then performs an inference step on the perturbed image. This procedure is run multiple times to identify the superpixels that lead to the most drastic change in the ConvNet’s prediction. SharpLIME uses hand-crafted segmentations to split the image into relevant sections, and then proceeds with the perturbation process defined in LIME. The main drawback of perturbation based methods is the need to run the prediction algorithm multiple times, which leads to high computational costs and long running times.

Gradient-based methods, such as saliency maps [Simonyan and Zisserman, 2015], guided backpropagation [Springenberg et al., 2014], gradient class-activation maps (Grad-CAM) [Selvaraju et al., 2017], or layer-wise relevance propagation [Bach et al., 2015], use a ConvNet’s backpropagation step to identify the areas in an image that contribute the most to the prediction. In general, gradient-based methods compute the gradient of a given input in relation to the prediction, and apply different post-processing methods to the output. In the case of Grad-CAM, image features are extracted by forward propagating the image until the last convolutional layer. Then, the gradient is set to 0 for all classes except the target class, and the signal is backpropagated to the last convolutional layer. The extracted image features that directly contribute to the backpropagated signal constitute the Grad-CAM for the given class. Since the analysis can be performed at the same time as the inference itself and only requires one iteration, Grad-CAM is often used in research and industrial applications [Pereira et al., 2018, Young et al., 2019, Tschandl et al., 2020, Hepp et al., 2021, Jalaboi et al., 2023]. Due to its popularity, in this paper we will use Grad-CAM to benchmark the explainability of commonly used ConvNet architectures.

3 Material and methods

3.1 Literature review

We performed a systematic literature review on PubMed, following the methodology introduced by Thomsen et al. [2020]. The query, described in Table 1, focused on dermatological applications of deep learning. A total of 3,650 articles were retrieved. We excluded articles that focused on domains other than dermatology, articles that did not include an original contribution in disease classification, articles using modalities other than photographic images, articles using methods other than ConvNets, and articles using proprietary ConvNets.

3.2 Explainability benchmark

3.2.1 Explainability dataset

For explainability benchmarking, we use DermXDB, a skin disease diagnosis explainability dataset published by Jalaboi et al. [2022]. The dataset consists of 524 images sourced from DermNetNZ [DermNetNZ, 2021] and SD-260 [Sun et al., 2016], and labeled with diagnoses and explanations in the form of visual skin lesion characteristics by eight board-certified dermatologists. To match the Grad-CAM output, we focus on the characteristic localization task.

3.2.2 Diagnosis evaluation

For establishing the expert-level diagnosis performance, we compare each dermatologist with the reference standard diagnosis. We follow the same approach for benchmarking the diagnosis performance of the ConvNets. We evaluate the performance using the categorical F1 score, sensitivity, and specificity, defined as:

$$\text{F1 score} = \frac{2TP}{2TP + FP + FN}, \quad (1)$$

$$\text{Sensitivity} = \frac{TP}{TP + FN}, \quad (2)$$

$$\text{Specificity} = \frac{TN}{TN + FP}, \quad (3)$$

where the true positives TP represent correctly classified samples, the false positives FP represent samples incorrectly classified as part of the target class, the false negatives FN represent samples of the target class incorrectly classified as being part of a different class, and the true negatives TN represent samples correctly identified as not being part of the target class.

3.2.3 Explainability evaluation

For establishing expert-level explainability performance, we compare the attention masks of each dermatologist with the aggregated fuzzy union of attention masks created by the other seven dermatologists (explanation maps). More specifically, we define the *image-level explanation maps* as the union of all characteristics segmented by all dermatologists for an image, and the *characteristic-level explanation maps* as the union of all segmentations for each characteristic for an image. Figure 1 illustrates the mask creation process for a psoriasis case. The ConvNet Grad-CAM attention maps are compared with explanations maps derived from all eight dermatologist evaluations.

These two types of explanation maps offer a way to check whether the ConvNets take into account the entire area selected by dermatologists as important to their decision, and whether they focus on specific characteristics when making their decisions. To quantify the similarity between the Grad-CAMs and the explanation maps, we compute the F1 score, sensitivity and specificity following their fuzzy implementation defined in [Crum et al., 2006], described as:

$$\text{F1 score} = \frac{2 \sum_{p \in \text{pixels}} \min(\mathcal{G}_p, \mathcal{E}_p)}{\sum_{p \in \text{pixels}} (\mathcal{G}_p) + \sum_{p \in \text{pixels}} (\mathcal{E}_p)}, \quad (4)$$

$$\text{Sensitivity} = \frac{\sum_{p \in \text{pixels}} \min(\mathcal{G}_p, \mathcal{E}_p)}{\sum_{p \in \text{pixels}} (\mathcal{S}_p)}, \quad (5)$$

$$\text{Specificity} = \frac{\sum_{p \in \text{pixels}} \min(1 - \mathcal{G}_p, 1 - \mathcal{E}_p)}{\sum_{p \in \text{pixels}} (1 - \mathcal{E}_p)}, \quad (6)$$

where \mathcal{G} is the ConvNet-generated Grad-CAM, and \mathcal{E} is the explanation map for a given image.

For characteristics, we report the Grad-CAM sensitivity with regard to the characteristic-level explanation maps. Specificity and F1 score were considered too stringent, as multiple characteristics can be present and essential for a diagnosis, and an explainable ConvNet must detect all of them to plausibly explain the diagnosis.

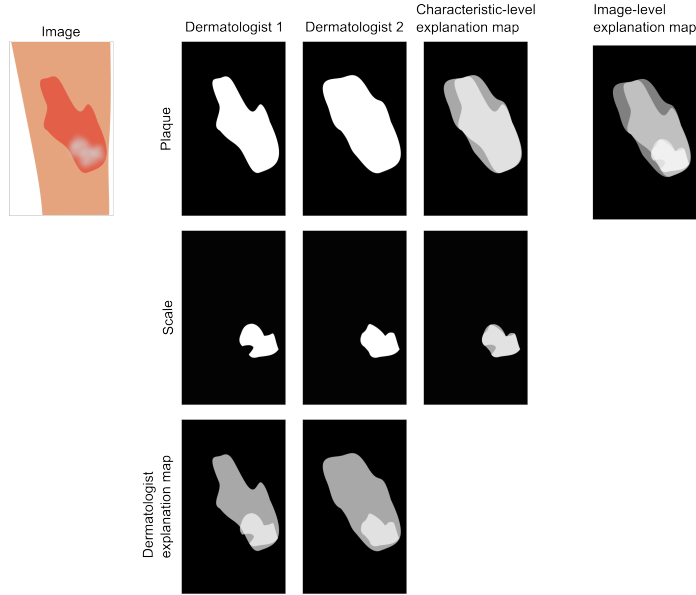


Figure 1: Explanation maps creation example for a psoriasis case evaluated by two dermatologists. Both dermatologists identified plaque and scale as the two characteristics associated with the psoriasis diagnosis, and localized them. By combining the localization maps for each characteristic, we obtain the characteristic-level explanation maps. By combining the localization maps created by each dermatologist, we obtain the individual dermatologist explanation maps. By combining all localization maps, we obtain the image-level explanation map.

3.2.4 Experimental setup

From the 22 articles that fulfilled all inclusion criteria, we selected the set of ConvNets to benchmark based on their reproducibility: we required that all benchmarked ConvNets had been pre-trained on ImageNet due to the limited amount of training data available. Thus, we exclude architectures that do not have publicly available pre-trained ImageNet weights compatible with the deep learning Keras framework [Chollet, 2015], i.e. GoogLeNet [Szegedy et al., 2015], InceptionV4 [Babenko and Lempitsky, 2015], MobileNetV3 [Howard et al., 2019], SENet [Hu et al., 2018], SE-ResNet [Hu et al., 2018], SEResNeXT [Hu et al., 2018], and ShuffleNet [Zhang et al., 2018]. Furthermore, as several articles compare different versions of the same architecture (e.g. EfficientNet-B0 through EfficientNet-B7, see Table 2), we select the smallest version of each architecture for our benchmark to avoid overfitting to the DermXDB dataset.

In the rest of this work, we will focus on the following ConvNets: DenseNet121 [Huang et al., 2017], EfficientNet-B0 [Tan and Le, 2019], InceptionResNetV2 [Szegedy et al., 2017], InceptionV3 [Szegedy et al., 2016], MobileNet [Howard et al., 2017], MobileNetV2 [Sandler et al., 2018], NASNetMobile [Zoph et al., 2018], ResNet50 [He et al., 2016a], ResNet50V2 [He et al., 2016b], VGG16 [Simonyan and Zisserman, 2015], and Xception [Chollet, 2017].

We used the pre-trained weights offered by Keras to initialize the networks in our experiments. Next, all ConvNets were pre-trained on a proprietary clinical photography skin disease dataset collected by a dermatologist between 2004-2018. All images included in the dataset were anonymized, and the patients consented to their data being used for research purposes. More information about the dataset is available in Appendix Table A1. We performed a hyper-parameter search for each ConvNet, with the values used for experimentation and the validation performance being reported in Appendix Table A2 and Appendix Table A3, respectively. We further fine-tuned all ConvNets for 50 epochs with 261 randomly chosen images from the DermXDB dataset. The remaining 263 images were used as the test set. Each ConvNet was trained and tested five times. All results presented in this paper are aggregated over the five test runs. All code used for running the experiments is available at <https://github.com/ralucaj/dermx-benchmark>.

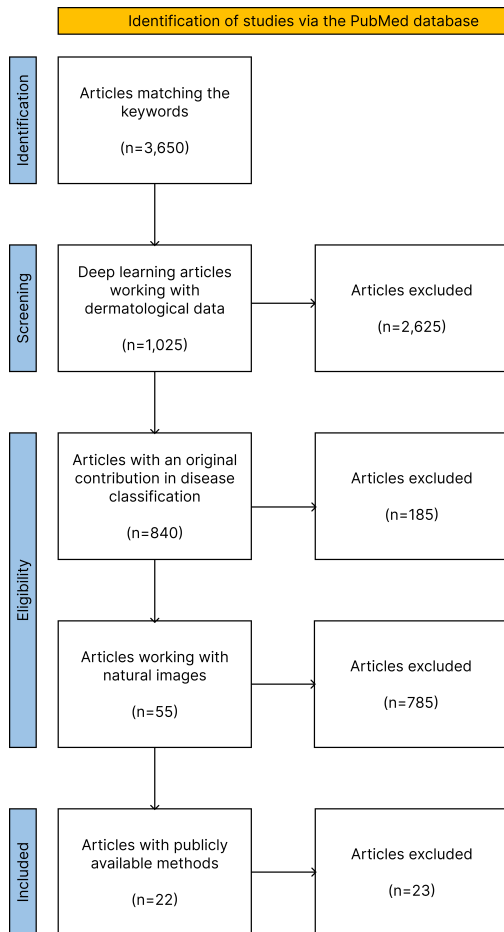


Figure 2: The Preferred Reporting Items for Systematic Reviews and Meta Analyses (PRISMA) statement flowchart of the performed review process for identifying the benchmarked ConvNet architectures. First, we screened articles to ensure that they were using dermatological data and deep learning methods. Afterwards, we excluded review articles and contributions focused on tasks other than classification, and articles that that used non-photographic image data, e.g. dermoscopy, whole slides. Finally, we excluded articles that used proprietary ConvNets, leading to 22 articles serving as the benchmark basis.

4 Results

4.1 Literature review

Figure 2 displays the Preferred Reporting Items for Systematic Review and Meta-Analyses statement flowchart of the performed review, while Figure 3 illustrates the evolution of articles topics over the years. Out of the original 3,650 articles, only 22 fulfilled all the inclusion criteria. Table 2 summarizes the ConvNet architectures, their implementation, and reported performance employed in the final 22 articles selected for benchmarking.

Table 2: Overview of the 22 articles fulfilling all inclusion criteria. All articles use ConvNets for a dermatological classification task using photographic images. Tasks vary between binary or multi-disease diagnosis, disease risk assessment, lesion type classification, and severity assessment.

Publication	ConvNets employed	Task	Data	Performance
Aggarwal [2019]	InceptionV3	Disease diagnosis on five classes	Open source images and images scraped from Google	0.66 F1 score, 0.65 sensitivity, 0.91 specificity, 0.67 precision, 0.91 NPV, 0.57 MCC
Burlina et al. [2019]	ResNet50	Disease diagnosis on four classes	Internet-scraped images	82.79% accuracy, 0.76 kappa score
Zhao et al. [2019]	Xception	Skin cancer risk assessment with three classes	Clinical images	72% accuracy, 0.92-0.96 ROC AUC, 0.85-0.93 sensitivity, 0.85-0.91 specificity
Burlina et al. [2020]	ResNet50, ResNet152, InceptionV3, InceptionResNetV2, DenseNet	Disease diagnosis on eight classes	Clinical and other photographic images scraped using Google and Bing	71.58% accuracy, 0.70 sensitivity, 0.96 specificity, 0.72 precision, 0.96 NPV, 0.67 kappa, 0.72 F1 score, 0.80 average precision, 0.94 AUC
Chin et al. [2020]	DenseNet121, VGG16, ResNet50	Binary skin cancer risk assessment	Smartphone images	0.83-0.86 AUC, 0.72-0.77 sensitivity, 0.85-0.86 specificity
Han et al. [2020]	SENet, SE-ResNet50, VGG19	Disease classification on 134 classes	Clinical images	44.8-56.7% accuracy, 0.94-0.98 AUC
Liu et al. [2020]	InceptionV4	Disease diagnosis on 26 classes	Clinical images	66% accuracy, 0.56 sensitivity
Zhao et al. [2020]	DenseNet121, Xception, InceptionV3, InceptionResNetV2	Binary psoriasis classification	Clinical images	96% accuracy, 0.95-0.98 AUC, 0.96-0.97 specificity, 0.83-0.95 sensitivity
Wu et al. [2021]	SEResNeXt, SE-ResNet, InceptionV3	Disease diagnosis on five classes	Clinical images	0.96-0.97 AUC, 90-91% accuracy, 0.90-0.93 sensitivity, 0.90 specificity
Aggarwal and Papay [2022]	InceptionResNetV2	Disease diagnosis on four classes	Clinical images	0.60-0.82 sensitivity, 0.60-0.82 specificity, 0.33-0.93 precision, 0.33-0.93 NPV, 0.43-0.84 F1 score
Ba et al. [2022]	EfficientNet-B3	Disease diagnosis on 10 classes	Clinical images	78.45% accuracy, 0.73 kappa
Hossain et al. [2022]	VGG16, VGG19, ResNet50, ResNet101, ResNet50V2, ResNet101V2, InceptionV3, InceptionV4, InceptionResNetV2, Xception, DenseNet121, DenseNet169, DenseNet201, MobileNetV2, MobileNetV3Small, MobileNetV3Large, NASNetMobile, EfficientNet-B0 through EfficientNet-B5	Binary Lyme disease classification	Smartphone images	61.42-84.42% accuracy, 0.72-0.90 sensitivity, 0.50-0.81 specificity, 0.61-0.83 precision, 0.63-0.87 NPV, 0.23-0.69 MCC, 0.22-0.69 Cohen’s kappa, 1.46-4.70 positive likelihood ratio, 0.14-0.55 negative likelihood ratio, 0.66-0.85 F1 score, 0.65-0.92 AUC

Hüsers et al. [2022]	MobileNet	Binary wound maceration classification	Clinical images	69% accuracy, 0.69 sensitivity, 0.67 precision
Liu et al. [2022]	InceptionResNetV2	Ulcer characteristic diagnosis on two and three classes	Clinical images	71.2-99.4% accuracy, 0.68-0.99 sensitivity, 0.71-1.00 precision, 0.70-0.94 F1 score
Malihi et al. [2022]	Xception	Binary wound type classification	Clinical images	67-83% accuracy, 0.65-0.94 sensitivity, 0.70-0.75 specificity, 0.65-0.75 precision, 0.70-0.85 F1 score
Munthuli et al. [2022]	DenseNet121	Skin lesion severity classification with five classes	Smartphone images	0.43-0.91 sensitivity, 0.80-0.98 specificity, 0.50-0.87 F1 score
Ni et al. [2022]	DenseNet121, ResNet50	Radiation dermatitis severity classification on four classes	Clinical images	83% accuracy, 0.74-1.00 F1 score
Roy et al. [2022]	ResNet101	Disease diagnosis on 26 classes	Clinical images	62.6% - 75.6% accuracy, 69.3-81.8 AUPR
Sahin et al. [2022]	ResNet18, GoogleNet, EfficientNet-B0, NASNetMobile, ShuffleNet, MobileNetV2	Binary monkeypox classification	Smartphone images	73.33-91.11% accuracy
Xia et al. [2022]	ResNet50	Binary skin cancer classification	Smartphone images	0.77-0.82 AUC, 0.76-0.79 AP
Zhou et al. [2022]	ResNet50	Disease diagnosis on three classes	Clinical images	0.32 error rate, 0.68 sensitivity, 0.69 precision, 0.68 F1 score
Zhang and Ma [2022]	ResNet50	Acne severity classification with three classes	Clinical images	74% accuracy

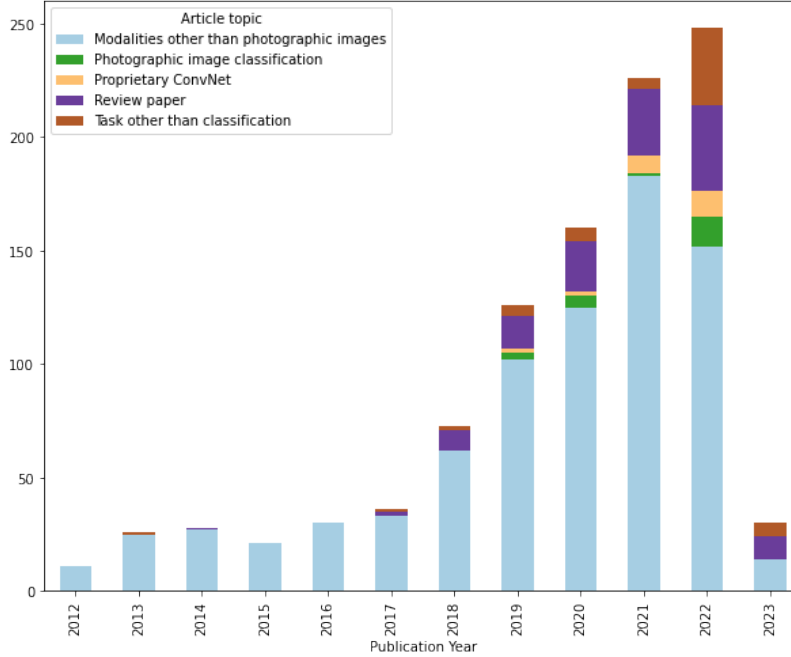


Figure 3: Distribution of retrieved article topics per publication year, based on the search query defined in Table 1 (ran on the 20th of February 2023). 2017 marks an explosion in the number of deep learning applications in dermatology, a fact highlighted by the large increase in articles in the subsequent years, and an increase in review articles. Starting 2019, the industrial involvement in this field became apparent due to the increase in proprietary ConvNets. 2019 also marks the first emergence of dermatological applications using photographic imaging. Finally, although classification is still the most common application, other applications are becoming increasingly more researched.

4.2 Diagnosis results

Table 3 provides an overview of the diagnostic performance of the networks and that of the dermatologists on average in terms of F1 score. As can be seen from the table, although several ConvNets achieve expert-level performance when diagnosing actinic keratosis, seborrheic dermatitis, and viral warts, none of them achieve overall expert-level performance. ConvNets follow the trend also seen in dermatologists of having difficulties correctly diagnosing actinic keratosis and seborrheic dermatitis, while the diagnosis of acne and viral warts displays higher performance. Similar trends can be observed for the sensitivity and specificity performance, as seen in Appendix Table A4 and Appendix Table A5, respectively.

4.3 Explainability results

Table 4 shows the image-level explainability results for each of the benchmarked ConvNets, while Figure 4 shows the relationship between ConvNet diagnosis performance, image-level explainability, and number of parameters. Xception scores the highest on the image-level Grad-CAM F1 score, while InceptionResNetV2, ResNet50, and VGG16 have the lowest performance. DenseNet121 and NASNetMobile report expert-level sensitivity scores, while ResNet50V2 achieves expert-level performance in specificity.

Looking at the characteristic-level sensitivity depicted in Figure 5, NASNetMobile and DenseNet121 achieve the highest overall performance. InceptionResNetV2, ResNet50, ResNet50V2, and VGG16 report the lowest scores. All ConvNets outperform dermatologists in closed comedo, open comedo, and pustule. The opposite is true for dermatoglyph disruption, leukotrichia, patch, plaque, scale, sun damage, and telangiectasia – no ConvNet reaches expert-level.

Figure 6 illustrates the differences in Grad-CAMs between the benchmarked ConvNets. Older ConvNet architectures, such as VGG16, InceptionResNetV2, ResNet50, and ResNet50V2, tend to focus on small areas that contain characteristics relevant for the diagnosis, e.g. focusing on a single plaque in the psoriasis diagnosis example, while more modern ConvNets pay attention to the entire area covered by diagnosis-relevant lesions. Several ConvNets, namely

Table 3: Diagnostic performance of the ConvNets (average \pm standard deviation across five runs) and dermatologists (average \pm standard deviation across eight experts) using F1-score, split by diagnosis. Several ConvNets achieve expert-level per-disease diagnosis performance, in actinic keratosis, seborrheic dermatitis, and viral warts (in **bold**), although none reach the same performance for acne, psoriasis, and vitiligo.

	Acne	Actinic keratosis	Psoriasis	Seborrheic dermatitis	Viral warts	Vitiligo
ConvNets						
DenseNet121	0.80 \pm 0.02	0.63 \pm 0.08	0.66 \pm 0.01	0.69 \pm 0.03	0.88 \pm 0.03	0.74 \pm 0.03
EfficientNet-B0	0.72 \pm 0.03	0.53 \pm 0.10	0.60 \pm 0.06	0.57 \pm 0.08	0.80 \pm 0.07	0.66 \pm 0.02
InceptionV3	0.77 \pm 0.02	0.57 \pm 0.11	0.60 \pm 0.02	0.54 \pm 0.03	0.77 \pm 0.04	0.73 \pm 0.05
InceptionResNetV2	0.73 \pm 0.02	0.52 \pm 0.10	0.53 \pm 0.05	0.56 \pm 0.05	0.69 \pm 0.03	0.53 \pm 0.12
MobileNet	0.72 \pm 0.06	0.55 \pm 0.19	0.51 \pm 0.14	0.57 \pm 0.06	0.68 \pm 0.06	0.56 \pm 0.10
MobileNetV2	0.56 \pm 0.07	0.23 \pm 0.09	0.31 \pm 0.08	0.46 \pm 0.05	0.63 \pm 0.07	0.48 \pm 0.14
NASNetMobile	0.50 \pm 0.05	0.33 \pm 0.12	0.42 \pm 0.07	0.43 \pm 0.05	0.55 \pm 0.11	0.51 \pm 0.05
ResNet50	0.77 \pm 0.04	0.53 \pm 0.17	0.61 \pm 0.03	0.61 \pm 0.19	0.79 \pm 0.02	0.61 \pm 0.07
ResNet50V2	0.76 \pm 0.04	0.62 \pm 0.07	0.59 \pm 0.01	0.57 \pm 0.01	0.76 \pm 0.01	0.75 \pm 0.05
VGG16	0.70 \pm 0.05	0.62 \pm 0.03	0.59 \pm 0.03	0.49 \pm 0.15	0.71 \pm 0.03	0.62 \pm 0.07
Xception	0.80 \pm 0.04	0.64 \pm 0.07	0.70 \pm 0.02	0.60 \pm 0.03	0.81 \pm 0.04	0.81 \pm 0.05
Dermatologists						
Average	0.95 \pm 0.02	0.79 \pm 0.14	0.85 \pm 0.06	0.72 \pm 0.09	0.93 \pm 0.05	0.96 \pm 0.03

EfficientNet-B0, MobileNet, MobileNetV2, and VGG16 seem to have overfit on the training set, focusing on the watermark rather than the image itself when diagnosing the vitiligo case.

5 Discussion

5.1 Literature review

ConvNets have become a default approach when it comes to automated diagnosis using images, aligned with the rise of the deep learning methodology for vision recognition. The continuous breakthroughs in diagnostic performance across a wide variety of medical imaging modalities and disorders have made automated diagnosis as close to integration with practice as ever. In dermatology, the diagnosis performance has achieved that of the expert raters as early as 2017 with a seminal work of Esteva et al. [2017] that disrupted the research field and set the trend that still persists, as can be seen through the trends of the continuous growth outlined in Figure 3. The increased interest of industrial entities that started in 2019, illustrated in Figure 3 by the increase in proprietary methods, is further highlighted by the large number of dermatology-oriented med-tech companies relying on machine learning for their products. Year 2019 also marks the year when research groups began investigating photographic images as a primary modality for diagnosing skin conditions, meaning the rise of machine learning solutions to assist tele dermatology.

The potential of using ConvNets to streamline dermatological tasks is underlined by the diversity of tasks being solved in the retrieved articles. Classification was the first methodology to be approached, with applications in disease diagnosis, risk assessment, lesion type classification, lesion characteristics identification, and disease severity assessment. Segmentation and natural language processing applications are also gaining more traction, as shown by the constant increase in non-classification tasks in Figure 3.

However, this potential has not yet translated into the much-needed transformation of the clinical practice. In part, this is due to regulatory challenges which are often faced due to the limited generalizability and lack of explainability of the methods [Kelly et al., 2019]. By benchmarking the diagnosis and explainability performance of ConvNets, we both

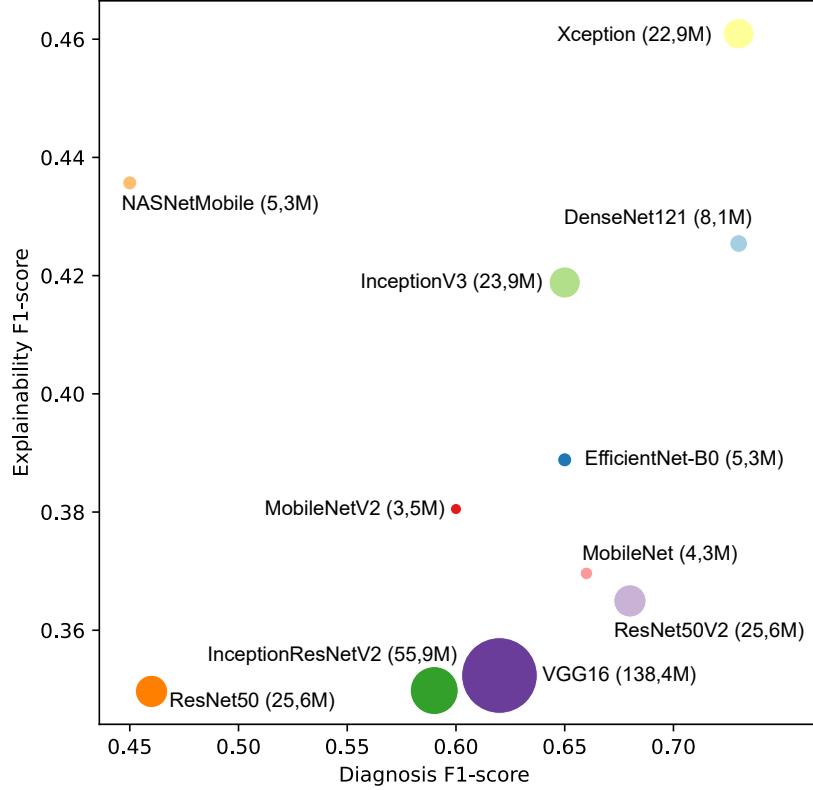


Figure 4: ConvNet explainability as a function of ConvNet performance and their number of parameters. Xception displays both the highest performance and image-level explainability, while ResNet50 performs poorly in both criteria.

enable a comparison among the methods, as well as help the identification gaps between the current state-of-the-art and the clinical practice.

5.2 Diagnosis benchmark

The direct comparison of the diagnostic performance is not possible using reported values from the literature not only due to variability in the choice of the metrics, but more importantly due to the variance in the number of classes and the differences in the datasets used for training and validation (Table 2). By reformulating the task to the diagnosis of six disease classes, utilizing the same initialization, pre-training, and hyperparameter optimization search strategy, and training and validating on the common database, this benchmark minimizes the performance variability related to such implementation details.

We found considerable variability among the diagnostic performance values, with the average F1 scores ranging from 0.50 to 0.80 for acne, from 0.23 to 0.64 for actinic keratosis, from 0.31 to 0.70 for psoriasis, from 0.43 to 0.69 for seborrheic dermatitis, from 0.55 to 0.88 for viral warts, and from 0.51 to 0.81 for vitiligo. These values were aligned with the diagnostic complexity of the diseases as expressed by the performance of the dermatologists, averaging 0.95 for acne, 0.79 for actinic keratosis, 0.85 for psoriasis, 0.72 for seborrheic dermatitis, 0.93 for viral warts, and 0.96 for vitiligo. As such, none of the ConvNets achieved the average dermatologist performance, although there were multiple instances of ConvNets reaching the range of the expert performance for a specific disease (see Table 3). The majority of the benchmarked ConvNets achieved expert level for diagnosis of actinic keratosis and seborrheic dermatitis: seven and six out of 11, respectively. This further confirms the similarity of ConvNet performance with respect to the dermatologists: most ConvNets display a similar difficulty in diagnosing actinic keratosis and seborrheic dermatitis as the eight dermatologists, and a similar ease of diagnosing acne and viral warts.

5.3 Explainability benchmark

While diagnostic performance is recognized as critical for the generalizability of ConvNets, the explainability performance validation has been generally approached as an optional, qualitative, post-hoc analysis. One of the key challenges

Table 4: Explainability performance in terms of the image-level Grad-CAM evaluation for the ConvNets (average \pm standard deviation across five runs), and an explanation map evaluation dermatologists (average \pm standard deviation across eight experts). Older ConvNets, such as ResNet50, ResNet50V2, and VGG16, have lower performance than most other modern ConvNets. Two networks achieve expert-level sensitivity scores, and one achieves expert-level specificity (in **bold**).

	F1 score	Sensitivity	Specificity
ConvNets			
DenseNet121	0.43 \pm 0.01	0.61 \pm 0.01	0.78 \pm 0.00
EfficientNet-B0	0.39 \pm 0.01	0.52 \pm 0.00	0.82 \pm 0.00
InceptionV3	0.42 \pm 0.01	0.56 \pm 0.01	0.82 \pm 0.01
InceptionResNetV2	0.35 \pm 0.01	0.40 \pm 0.01	0.87 \pm 0.01
MobileNet	0.37 \pm 0.02	0.50 \pm 0.01	0.85 \pm 0.01
MobileNetV2	0.38 \pm 0.02	0.49 \pm 0.02	0.87 \pm 0.01
NASNetMobile	0.44 \pm 0.00	0.62 \pm 0.00	0.81 \pm 0.00
ResNet50	0.35 \pm 0.01	0.42 \pm 0.03	0.84 \pm 0.01
ResNet50V2	0.37 \pm 0.01	0.38 \pm 0.01	0.91 \pm 0.00
VGG16	0.35 \pm 0.01	0.40 \pm 0.01	0.86 \pm 0.01
Xception	0.46 \pm 0.01	0.56 \pm 0.00	0.88 \pm 0.01
Dermatologists			
Average	0.66 \pm 0.03	0.67 \pm 0.07	0.93 \pm 0.03

faced by researchers trying to implement a more objective validation of explainability is linking the human-approachable explanations with those feasible for ConvNets. With the use of the labels for dermatological diagnosis explainability available from the recently released DermXDB dataset, our benchmark is quantitative as well as predefined. Thus, we avoid potential biases and limitations stemming from machine learning experts with little domain knowledge performing a visual, qualitative evaluation of Grad-CAMs [Tschandl et al., 2020].

The image-level explainability analysis shows that no ConvNet reaches the same F1 score as the dermatologists, although several ConvNets achieve expert-level sensitivity or specificity. Different ConvNets show different patterns of explanation behaviour (Figure 6): some tend to focus on smaller areas that are highly indicative of the target diagnosis, while others tend to focus on the entire affected area. Extensive user tests with both experts and patients would enable us to learn which of the two options is preferred as an explanation: a single, classical lesion descriptive of the diagnosis, or highlighting the entire affected area.

From a characteristic-level sensitivity perspective, most ConvNets outperform the average dermatologist performance in characteristics smaller than 1cm in diameter [Nast et al., 2016]. For larger characteristics, although NASNetMobile and Xception approach expert-level, no ConvNet exceeds it. The relationship between diseases and their characteristics is visible in the characteristic-level ConvNet explainability: most ConvNets report high sensitivity on characteristics often associated with acne and viral warts (e.g. closed and open comedones, papules, and thrombosed capillaries), while reporting a lower performance on characteristics associated with actinic keratosis and seborrheic dermatitis (e.g. plaque, sun damage, and patch). Characteristic-level explainability may be more relevant for use cases where identifying the differentiating factor between different diseases is the most important component for garnering trust.

These result suggests that while ConvNets have the potential to produce human-approachable explanations, more work is necessary to fully achieve expert-level performance. Part of the necessary work is the creation of additional user-derived explainability datasets that enable quantitative analyses on a ConvNet’s explainability within a domain. A component of this is performing extensive user tests to identify the explainability expectations of an application’s end users. From a machine learning perspective, more research must be devoted to the creation of intrinsically explainable ConvNets, rather than relying solely on post-hoc explanation methods. Such a ConvNet must be aligned with the explainability requirements of its task and its users: a psoriasis diagnosis ConvNet aimed at dermatologists might

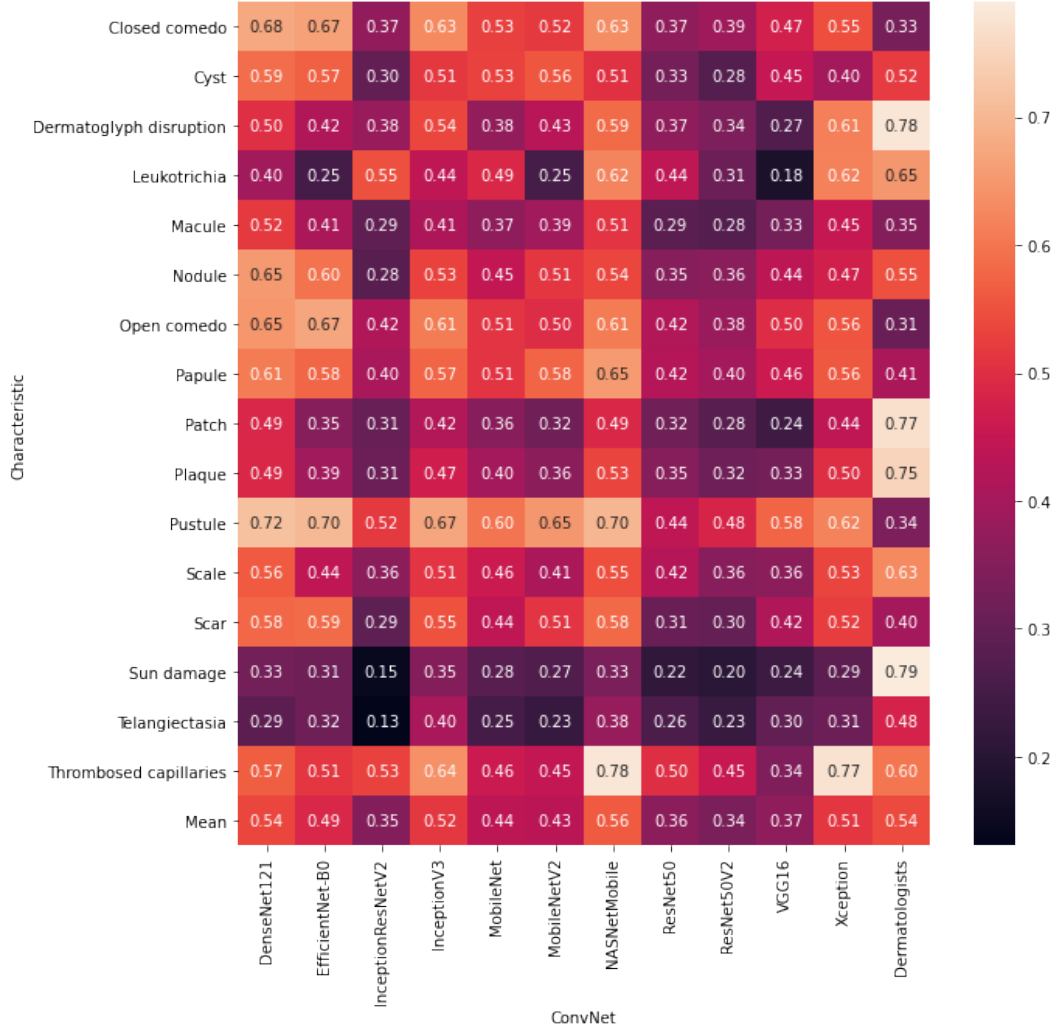


Figure 5: Explainability performance in terms of characteristic-level Grad-CAM sensitivity for the ConvNets (averaged across five runs) and dermatologists (averaged across eight experts). NASNetMobile and Xception outperform expert level in seven characteristics, while no ConvNet achieves expert-level performance in eight characteristics.

require high characteristic-level explainability to offer a constrative explanation against a possible differential diagnosis of atopic dermatitis, while the same ConvNet aimed at patients might require high image-level explainability to reassure the patient that all aspects of their condition are taken into consideration.

5.4 Limitations and future work

Our work has a few limitations. First, the original DermXDB dataset contains little information about the gender, age, and ethnicity of the subjects, leading to difficulties in performing an in-depth bias analysis of our benchmark. Second, the small size of the dataset limits the training capabilities of our benchmark, which may underestimate the performance of the larger ConvNets.

In future work, we plan on expanding this benchmark by using more explainability methods, such as saliency maps and LIME, to also create a benchmark of explainability methods and their performance compared to that of dermatologists. Additionally, with the increased popularity of visual transformers [Khan et al., 2022], an analysis of their Grad-CAM explainability would be of interest to the research world.

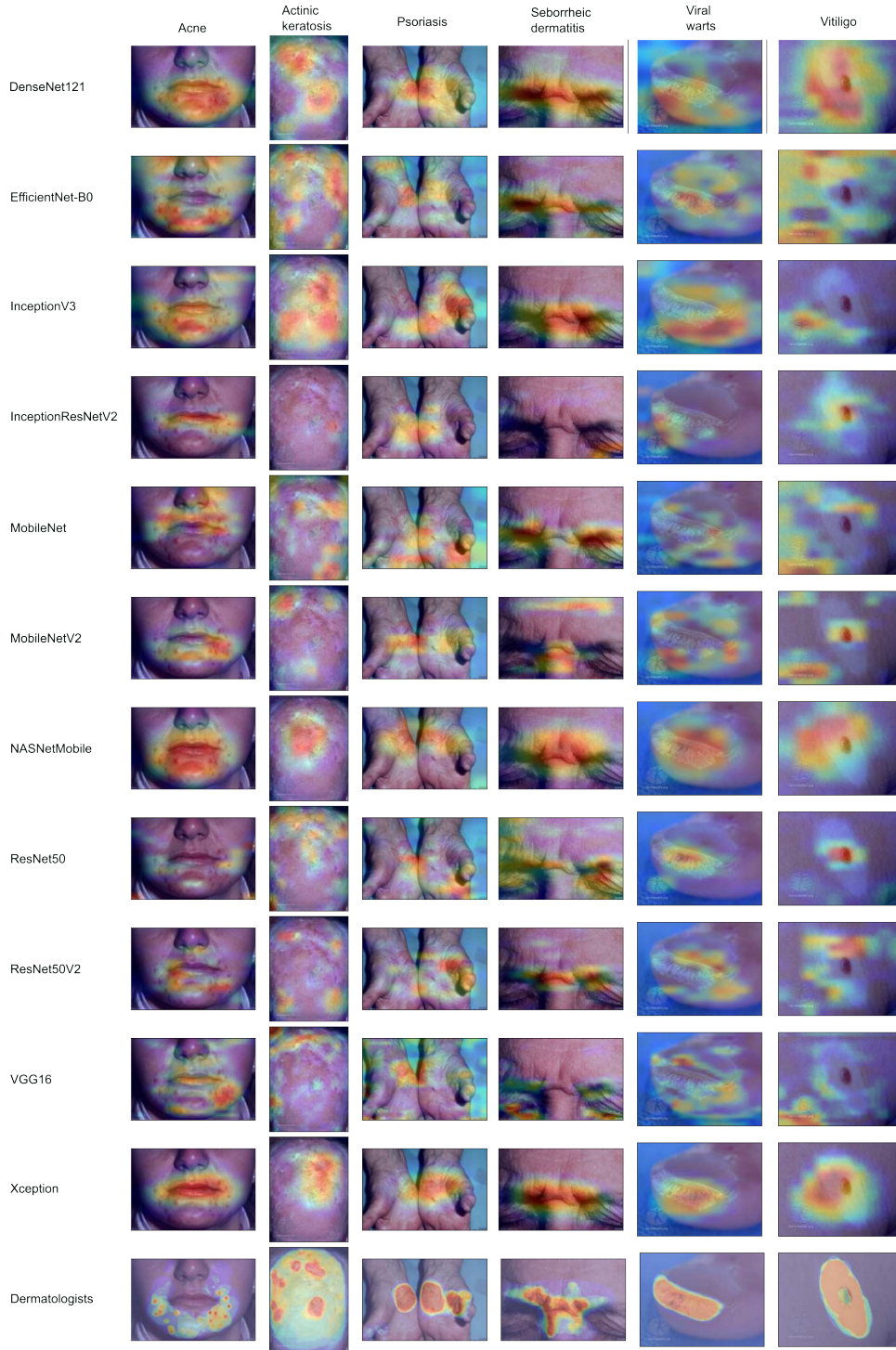


Figure 6: Example of Grad-CAM outputs for six images that were correctly diagnosed by all ConvNets. Older ConvNets, such as VGG16, ResNet50, ResNet50V2, and InceptionResNetV2, tend to focus on a single, highly indicative lesion rather than the whole affected region. More modern ConvNets, such as NASNetMobile, Xception, and EfficientNet, focus on the entire affected area. Some ConvNets overfitted during training, and focus on the watermark when diagnosing vitiligo.

6 Conclusions

In this paper, we performed a systematic literature review to identify the most used ConvNet architectures for the diagnosis of skin diseases from photographic images. We benchmarked the 11 identified ConvNets on DermXDB, a skin disease explainability dataset. Xception stands out as a highly explainable ConvNet, although NASNetMobile outperforms it on characteristic-level sensitivity. Our findings highlight the importance of explainability benchmarking, and will hopefully motivate additional studies within the field of quantitative evaluations for explainability.

Acknowledgments

Funding: RJ’s work was supported in part by the Danish Innovation Fund under Grant 0153-00154A. OW’s work was funded in part by the Novo Nordisk Foundation through the Center for Basic Machine Learning Research in Life Science (NNF20OC0062606). OW acknowledges support from the Pioneer Centre for AI, DNRF grant number P1.

References

- World Health Organization. Working for health and growth: investing in the health workforce. 2016.
- U.S. Food and Drug Administration. FDA permits marketing of artificial intelligence-based device to detect certain diabetes-related eye problems. *News Release, April*, 2018.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- Jun Gao, Qian Jiang, Bo Zhou, and Daozheng Chen. Convolutional neural networks for computer-aided detection or diagnosis in medical image analysis: an overview. *Mathematical Biosciences and Engineering*, 16(6):6536–6561, 2019.
- Thorsten Falk, Dominic Mai, Robert Bensch, Özgün Çiçek, Ahmed Abdulkadir, Yassine Marrakchi, Anton Böhm, Jan Deubner, Zoe Jäckel, Katharina Seiwald, et al. U-net: deep learning for cell counting, detection, and morphometry. *Nature Methods*, 16(1):67–70, 2019.
- Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros, et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *JAMA*, 316(22):2402–2410, 2016.
- Richard Kijowski, Fang Liu, Francesco Caliva, and Valentina Pedoia. Deep learning for lesion detection, progression, and prediction of musculoskeletal disease. *Journal of Magnetic Resonance Imaging*, 52(6):1607–1619, 2020.
- Mohammad Hesam Hesamian, Wenjing Jia, Xiangjian He, and Paul Kennedy. Deep learning techniques for medical image segmentation: achievements and challenges. *Journal of Digital Imaging*, 32:582–596, 2019.
- KKD Ramesh, G Kiran Kumar, K Swapna, Debabrata Datta, and S Suman Rajest. A review of medical image segmentation algorithms. *EAI Endorsed Transactions on Pervasive Health and Technology*, 7(27):e6–e6, 2021.
- Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639):115–118, 2017.
- Jianpeng Zhang, Yutong Xie, Yong Xia, and Chunhua Shen. Attention residual learning for skin lesion classification. *IEEE Transactions on Medical Imaging*, 38(9):2092–2103, 2019.
- Shunichi Jinnai, Naoya Yamazaki, Yuichiro Hirano, Yohei Sugawara, Yuichiro Ohe, and Ryuji Hamamoto. The development of a skin cancer classification system for pigmented skin lesions using deep learning. *Biomolecules*, 10(8):1123, 2020.
- Holger Andreas Haenssle, Christine Fink, Ferdinand Toberer, Julia Winkler, Wilhelm Stolz, Teresa Deinlein, Rainer Hofmann-Wellenhof, Aimilios Lallas, Steffen Emmert, Timo Buhl, et al. Man against machine reloaded: performance of a market-approved convolutional neural network in classifying a broad spectrum of skin lesions in comparison with 96 dermatologists working under less artificial conditions. *Annals of Oncology*, 31(1):137–143, 2020.
- Abhijit Guha Roy, Jie Ren, Shekoofeh Azizi, Aaron Loh, Vivek Natarajan, Basil Mustafa, Nick Pawlowski, Jan Freyberg, Yuan Liu, Zach Beaver, et al. Does your dermatology classifier know what it doesn’t know? detecting the long-tail of unseen conditions. *Medical Image Analysis*, 75:102274, 2022.

- Christopher J Kelly, Alan Karthikesalingam, Mustafa Suleyman, Greg Corrado, and Dominic King. Key challenges for delivering clinical impact with artificial intelligence. *BMC Medicine*, 17(1):1–9, 2019.
- Bryce Goodman and Seth Flaxman. European union regulations on algorithmic decision-making and a “right to explanation”. *AI Magazine*, 38(3):50–57, 2017.
- Xiao Bai, Xiang Wang, Xianglong Liu, Qiang Liu, Jingkuan Song, Nicu Sebe, and Been Kim. Explainable deep learning for efficient and robust pattern recognition: A survey of recent developments. *Pattern Recognition*, 120: 108102, 2021.
- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626, 2017.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 2017.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why should I trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016.
- Raluca Jalaboi, Frederik Faye, Mauricio Orbes-Arteaga, Dan Jørgensen, Ole Winther, and Alfiia Galimzianova. Dermx: an end-to-end framework for explainable automated dermatological diagnosis. *Medical Image Analysis*, page 102647, 2022.
- Kenneth Thomsen, Lars Iversen, Therese Louise Titlestad, and Ole Winther. Systematic review of machine learning for diagnosis and prognosis in dermatology. *Journal of Dermatological Treatment*, 31(5):496–510, 2020. doi:<https://doi.org/10.1080/09546634.2019.1682500>.
- Hyeon Ki Jeong, Christine Park, Ricardo Henao, and Meenal Kheterpal. Deep learning in dermatology: a systematic review of current approaches, outcomes and limitations. *JID Innovations*, page 100150, 2022.
- Zhao Liu, Jiulai Sun, Lyndon Smith, Melvyn Smith, and Robert Warr. Distribution quantification on dermoscopy images for computer-assisted diagnosis of cutaneous melanomas. *Medical & Biological Engineering & Computing*, 50(5): 503–513, 2012.
- Peyman Sabouri, Hamid GholamHosseini, Thomas Larsson, and John Collins. A cascade classifier for diagnosis of melanoma in clinical images. In *2014 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, pages 6748–6751. IEEE, 2014.
- Ravneet Kaur, Peter P Albano, Justin G Cole, Jason Hagerty, Robert W LeAnder, Randy H Moss, and William V Stoecker. Real-time supervised detection of pink areas in dermoscopic images of melanoma: importance of color shades, texture and location. *Skin Research and Technology*, 21(4):466–473, 2015.
- S Kefel, S Pelin Kefel, RW LeAnder, R Kaur, R Kasmi, NK Mishra, RK Rader, JG Cole, ZT Woolsey, and WV Stoecker. Adaptable texture-based segmentation by variance and intensity for automatic detection of semitranslucent and pink blush areas in basal cell carcinoma. *Skin Research and Technology*, 22(4):412–422, 2016.
- Kouhei Shimizu, Hitoshi Iyatomi, M Emre Celebi, Kerri-Ann Norton, and Masaru Tanaka. Four-class classification of skin lesions with task decomposition strategy. *IEEE Transactions on Biomedical Engineering*, 62(1):274–283, 2014.
- Jose Luis García Arroyo and Begoña García Zapirain. Detection of pigment network in dermoscopy images using supervised machine learning and structural analysis. *Computers in Biology and Medicine*, 44:144–157, 2014.
- Vimal K Shrivastava, Narendra D Londhe, Rajendra S Sonawane, and Jasjit S Suri. Computer-aided diagnosis of psoriasis skin images with hos, texture and color features: a first comparative study of its kind. *Computer Methods and Programs in Biomedicine*, 126:98–109, 2016.
- Noel CF Codella, David Gutman, M Emre Celebi, Brian Helba, Michael A Marchetti, Stephen W Dusza, Aadi Kallou, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, et al. Skin lesion analysis toward melanoma detection: A challenge at the 2017 International Symposium on Biomedical Imaging (ISBI), hosted by the International Skin Imaging Collaboration (ISIC). In *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, pages 168–172. IEEE, 2018.
- Xiaoxiao Sun, Jufeng Yang, Ming Sun, and Kai Wang. A benchmark for automatic visual classification of clinical skin disease images. In *European Conference on Computer Vision*, pages 206–222. Springer, 2016.
- DermNetNZ. DermNetNZ. <https://dermnetnz.org/>, 2021. Accessed: 2021-04-01.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. Ieee, 2009.

- P Tschandl, H Kittler, and G Argenziano. A pretrained neural network shows similar diagnostic accuracy to medical students in categorizing dermoscopic images after comparable training conditions. *British Journal of Dermatology*, 177(3):867–869, 2017.
- Seung Seog Han, Gyeong Hun Park, Woohyung Lim, Myoung Shin Kim, Jung Im Na, Ilwoo Park, and Sung Eun Chang. Deep neural networks show an equivalent and often superior performance to dermatologists in onychomycosis diagnosis: Automatic construction of onychomycosis datasets by region-based convolutional deep neural network. *PloS One*, 13(1):e0191493, 2018.
- G Reshma, Chiai Al-Atroshi, V Kumar Nassa, B Geetha, Gurram Sunitha, Mohammad Gouse Galety, and S Neelakanandan. Deep learning-based skin lesion diagnosis model using dermoscopic images. *Intelligent Automation and Soft Computing*, 31(1):621–634, 2022.
- Yading Yuan, Ming Chao, and Yeh-Chi Lo. Automatic skin lesion segmentation using deep fully convolutional networks with jaccard distance. *IEEE Transactions on Medical Imaging*, 36(9):1876–1886, 2017.
- Huisi Wu, Shihuai Chen, Guilian Chen, Wei Wang, Baiying Lei, and Zhenkun Wen. Fat-net: Feature adaptive transformers for automated skin lesion segmentation. *Medical Image Analysis*, 76:102327, 2022.
- Ramsha Baig, Maryam Bibi, Anmol Hamid, Sumaira Kausar, and Shahzad Khalid. Deep learning approaches towards skin lesion segmentation and classification from dermoscopic images-a review. *Current Medical Imaging*, 16(5): 513–533, 2020.
- Philipp Tschandl, Christoph Rinner, Zoe Apalla, Giuseppe Argenziano, Noel Codella, Allan Halpern, Monika Janda, Aimilios Lallas, Caterina Longo, Josep Malvehy, et al. Human–computer collaboration for skin cancer recognition. *Nature Medicine*, 26(8):1229–1234, 2020.
- Masaya Tanaka, Atsushi Saito, Kosuke Shido, Yasuhiro Fujisawa, Kenshi Yamasaki, Manabu Fujimoto, Kohei Murao, Youichirou Ninomiya, Shin’ichi Satoh, and Akinobu Shimizu. Classification of large-scale image database of various skin diseases using deep learning. *International Journal of Computer Assisted Radiology and Surgery*, 16:1875–1887, 2021.
- Ribana Roscher, Bastian Bohn, Marco F Duarte, and Jochen Garcke. Explainable machine learning for scientific insights and discoveries. *IEEE Access*, 8:42200–42216, 2020.
- Amitojdeep Singh, Sourya Sengupta, and Vasudevan Lakshminarayanan. Explainable deep learning models in medical image analysis. *Journal of Imaging*, 6(6):52, 2020.
- Stan Lipovetsky and Michael Conklin. Analysis of regression in game theory approach. *Applied Stochastic Models in Business and Industry*, 17(4):319–330, 2001.
- Mara Graziani, Iam Palatnik de Sousa, Marley MBR Vellasco, Eduardo Costa da Silva, Henning Müller, and Vincent Andrearczyk. Sharpening local interpretable model-agnostic explanations for histopathology: Improved understandability and reliability. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 540–549. Springer, 2021.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS One*, 10(7):e0130140, 2015.
- Sérgio Pereira, Raphael Meier, Victor Alves, Mauricio Reyes, and Carlos A Silva. Automatic brain tumor grading from mri data using convolutional neural networks and quality assessment. In *Understanding and Interpreting Machine Learning in Medical Image Computing Applications*, pages 106–114. Springer, 2018.
- Kyle Young, Gareth Booth, Becks Simpson, Reuben Dutton, and Sally Shrapnel. Deep neural network or dermatologist? In *Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support*, pages 48–55. Springer, 2019.
- Tobias Hepp, Dominik Blum, Karim Armanious, Bernhard Schoelkopf, Darko Stern, Bin Yang, and Sergios Gatidis. Uncertainty estimation and explainability in deep learning-based age estimation of the human brain: Results from the german national cohort mri study. *Computerized Medical Imaging and Graphics*, 92:101967, 2021.
- Raluca Jalaboi, Ole Winther, and Alfia Galimzianova. Explainable image quality assessments in teledermatological photography. *Telemedicine and e-Health*, 2023.

- William R Crum, Oscar Camara, and Derek L G Hill. Generalized overlap measures for evaluation and validation in medical image analysis. *IEEE Transactions on Medical Imaging*, 25(11):1451–1461, 2006.
- François Chollet. keras. <https://github.com/fchollet/keras>, 2015.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- Artem Babenko and Victor Lempitsky. Aggregating deep convolutional features for image retrieval. *arXiv preprint arXiv:1510.07493*, 2015.
- Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1314–1324, 2019.
- Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7132–7141, 2018.
- Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6848–6856, 2018.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4700–4708, 2017.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pages 6105–6114. PMLR, 2019.
- Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-first AAAI Conference on Artificial Intelligence*, 2017.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8697–8710, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016a.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, pages 630–645. Springer, 2016b.
- François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1251–1258, 2017.
- 1st Lt Pushkar Aggarwal. Data augmentation in dermatology image recognition using machine learning. *Skin Research and Technology*, 25(6):815–820, 2019.
- Philippe M Burlina, Neil J Joshi, Elise Ng, Seth D Billings, Alison W Rebman, and John N Aucott. Automated detection of erythema migrans and other confounding skin lesions via deep learning. *Computers in Biology and Medicine*, 105:151–156, 2019.
- Xin-yu Zhao, Xian Wu, Fang-fang Li, Yi Li, Wei-hong Huang, Kai Huang, Xiao-yu He, Wei Fan, Zhe Wu, Ming-liang Chen, et al. The application of deep learning in the risk grading of skin tumors for patients using clinical images. *Journal of Medical Systems*, 43(8):1–7, 2019.
- Philippe M Burlina, Neil J Joshi, Phil A Mathew, William Paul, Alison W Rebman, and John N Aucott. Ai-based detection of erythema migrans and disambiguation against other skin lesions. *Computers in Biology and Medicine*, 125:103977, 2020.

- YPH Chin, ZY Hou, MY Lee, HM Chu, HH Wang, YT Lin, A Gittin, SC Chien, PA Nguyen, LC Li, et al. A patient-oriented, general-practitioner-level, deep-learning-based cutaneous pigmented lesion risk classifier on a smartphone. *The British Journal of Dermatology*, 182(6):1498–1500, 2020.
- Seung Seog Han, Ilwoo Park, Sung Eun Chang, Woohyung Lim, Myoung Shin Kim, Gyeong Hun Park, Je Byeong Chae, Chang Hun Huh, and Jung-Im Na. Augmented intelligence dermatology: deep neural networks empower medical professionals in diagnosing skin cancer and predicting treatment options for 134 skin disorders. *Journal of Investigative Dermatology*, 140(9):1753–1761, 2020.
- Yuan Liu, Ayush Jain, Clara Eng, David H Way, Kang Lee, Peggy Bui, Kimberly Kanada, Guilherme de Oliveira Marinho, Jessica Gallegos, Sara Gabriele, et al. A deep learning system for differential diagnosis of skin diseases. *Nature Medicine*, 26(6):900–908, 2020.
- Shuang Zhao, Bin Xie, Yi Li, X-y Zhao, Yehong Kuang, Juan Su, X-y He, Xian Wu, Wei Fan, Kai Huang, et al. Smart identification of psoriasis by images using convolutional neural networks: a case study in china. *Journal of the European Academy of Dermatology and Venereology*, 34(3):518–524, 2020.
- Haijing Wu, Heng Yin, Haipeng Chen, Moyuan Sun, Xiaoqing Liu, Yizhou Yu, Yang Tang, Hai Long, Bo Zhang, Jing Zhang, et al. A deep learning-based smartphone platform for cutaneous lupus erythematosus classification assistance: Simplifying the diagnosis of complicated diseases. *Journal of the American Academy of Dermatology*, 85(3):792–793, 2021.
- Pushkar Aggarwal and Francis A Papay. Artificial intelligence image recognition of melanoma and basal cell carcinoma in racially diverse populations. *Journal of Dermatological Treatment*, 33(4):2257–2262, 2022.
- Wei Ba, Huan Wu, Wei W Chen, Shu H Wang, Zi Y Zhang, Xuan J Wei, Wen J Wang, Lei Yang, Dong M Zhou, Yi X Zhuang, et al. Convolutional neural network assistance significantly improves dermatologists’ diagnosis of cutaneous tumours using clinical images. *European Journal of Cancer*, 169:156–165, 2022.
- Sk Imran Hossain, Jocelyn de Goër de Herve, Md Shahriar Hassan, Delphine Martineau, Evelina Petrosyan, Violaine Corbin, Jean Beytout, Isabelle Lebert, Jonas Durand, Irene Carravieri, et al. Exploring convolutional neural networks with transfer learning for diagnosing lyme disease from skin lesion images. *Computer Methods and Programs in Biomedicine*, 215:106624, 2022.
- Jens Hüasers, Guido Hafer, Jan Heggemann, Stefan Wiemeyer, Mareike Przysucha, Joachim Dissemond, Maurice Moelleken, Cornelia Erfurt-Berge, and Ursula Hübner. Automatic classification of diabetic foot ulcer images—a transfer-learning approach to detect wound maceration. In *Informatics and Technology in Clinical Care and Public Health*, pages 301–304. IOS Press, 2022.
- Tom J Liu, Mesakh Christian, Yuan-Chia Chu, Yu-Chun Chen, Che-Wei Chang, Feipei Lai, and Hao-Chih Tai. A pressure ulcers assessment system for diagnosis and decision making using convolutional neural networks. *Journal of the Formosan Medical Association*, 121(11):2227–2236, 2022.
- Leila Malihi, Jens Hüasers, Mats L Richter, Maurice Moelleken, Mareike Przysucha, Dorothee Busch, Jan Heggemann, Guido Hafer, Stefan Wiemeyer, Gunther Heidemann, et al. Automatic wound type classification with convolutional neural networks. *Advances in Informatics, Management and Technology in Healthcare*, 295:281, 2022.
- A Munthuli, J Intanai, P Tossanuch, P Pooprasert, P Ingpochai, S Boonyasatian, K Kittithammo, P Thammarach, T Boonmak, S Khaengthanyakan, et al. Extravasation screening and severity prediction from skin lesion image using deep neural networks. In *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 1827–1833. IEEE, 2022.
- Ruiyan Ni, Ta Zhou, Ge Ren, Yuanpeng Zhang, Dongrong Yang, Victor CW Tam, Wan Shun Leung, Hong Ge, Shara WY Lee, and Jing Cai. Deep learning-based automatic assessment of radiation dermatitis in patients with nasopharyngeal carcinoma. *International Journal of Radiation Oncology* Biology* Physics*, 113(3):685–694, 2022.
- Veysel Harun Sahin, Ismail Oztel, and Gozde Yolcu Oztel. Human monkeypox classification from skin lesion images with deep pre-trained network using mobile application. *Journal of Medical Systems*, 46(11):79, 2022.
- Meng Xia, Meenal K Kheterpal, Samantha C Wong, Christine Park, William Ratliff, Lawrence Carin, and Ricardo Henao. Lesion identification and malignancy prediction from clinical dermatological images. *Scientific Reports*, 12(1):15836, 2022.
- Jiancun Zhou, Zheng Wu, Zixi Jiang, Kai Huang, Kehua Guo, and Shuang Zhao. Background selection schema on deep learning-based classification of dermatological disease. *Computers in Biology and Medicine*, 149:105966, 2022.
- Hang Zhang and Tianyi Ma. Acne detection by ensemble neural networks. *Sensors*, 22(18):6828, 2022.
- Alexander Nast, Chris EM Griffiths, Roderick Hay, Wolfram Sterry, and Jean L Bologna. The 2016 international league of dermatological societies’ revised glossary for the description of cutaneous lesions. *British Journal of Dermatology*, 174(6):1351–1358, 2016.

Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM Computing Surveys (CSUR)*, 54(10s):1–41, 2022.

Appendix

Table A1 presents statistics for the proprietary clinical dataset used in the hyper-parameter search and the pre-training step. Table A2 shows the best performing list of parameters identified for each ConvNet. The search space consisted of the following values for each hyperparameter:

- **Rotation:** 10, 20
- **Shear:** 0.00, 0.25, 0.50
- **Zoom:** 0.25, 0.5
- **Brightness ranges:** 0.00-0.50, 0.00-0.25, 0.50-1.00, 0.50-1.50, 0.75-1.25
- **Learning rate:** 0.01, 0.001, 0.0001
- **Last fixed layer:** last convolutional layer, second to last convolutional block
- **Epochs:** 10, 25, 50, 75

Table A1: Dataset statistics for the proprietary pre-training clinical dataset.

Diagnosis	Training	Validation
Acne	832	245
Actinic keratosis	132	33
Psoriasis	771	204
Seborrheic dermatitis	88	25
Viral warts	509	97
Vitiligo	141	37

Table A2: Optimal list of hyperparameters for each ConvNet, as identified after a hyper-parameter search.

ConvNet	Rotation	Shear	Zoom	Brightness	Learning rate	Last fixed layer	Epochs
DenseNet121	20	0.50	0.50	[0.50, 1.50]	0.0001	conv5_block14_concat	75
EfficientNet-B0	20	0.25	0.50	[0.50, 1.50]	0.0001	block6d_add	50
InceptionV3	20	0.50	0.50	[0.50, 1.50]	0.001	activation288	50
InceptionResNetv2	20	0.25	0.50	[0.75, 1.25]	0.0001	block8_9_ac	50
MobileNet	10	0.50	0.50	[0.50, 1.50]	0.0001	conv_pw_12_relu	50
MobileNetV2	10	0.25	0.50	[0.50, 1.50]	0.0001	block_15_add	75
NASNetMobile	20	0.25	0.50	[0.50, 1.00]	0.0001	normal_concat_11	75
ResNet50	20	0.50	0.50	[0.50, 1.50]	0.0001	conv5_block3_out	50
ResNet50V2	20	0.25	0.25	[0.50, 1.00]	0.001	post_relu	75
VGG16	10	0.00	0.25	[0.50, 1.00]	0.01	block5_pool	75
Xception	10	0.25	0.50	[0.50, 1.50]	0.001	block14_sepconv2_act	50

Table A3: Diagnostic performance of ConvNets in terms macro F1-score, sensitivity, and specificity on the validation subset of the proprietary clinical dataset (average \pm standard deviation across five runs).

ConvNet	F1-score	Sensitivity	Specificity
DenseNet121	0.80 \pm 0.01	0.79 \pm 0.01	0.98 \pm 0.00
EfficientNet-B0	0.77 \pm 0.01	0.78 \pm 0.01	0.97 \pm 0.00
InceptionV3	0.76 \pm 0.02	0.74 \pm 0.02	0.96 \pm 0.00
InceptionResNetV2	0.73 \pm 0.02	0.73 \pm 0.02	0.97 \pm 0.00
MobileNet	0.72 \pm 0.02	0.71 \pm 0.02	0.96 \pm 0.00
MobileNetV2	0.72 \pm 0.03	0.73 \pm 0.02	0.96 \pm 0.00
NASNetMobile	0.67 \pm 0.04	0.64 \pm 0.02	0.95 \pm 0.01
ResNet50	0.70 \pm 0.01	0.68 \pm 0.02	0.96 \pm 0.00
ResNet50V2	0.76 \pm 0.01	0.75 \pm 0.02	0.96 \pm 0.00
VGG16	0.66 \pm 0.03	0.67 \pm 0.01	0.95 \pm 0.00
Xception	0.82 \pm 0.03	0.82 \pm 0.02	0.97 \pm 0.00

Table A4: Diagnostic performance of ConvNets (average \pm standard deviation across five runs) and dermatologists (average \pm standard deviation across eight dermatologists) in terms of sensitivity on the DermXDB holdout set, split by diagnosis. Several ConvNets achieve expert-level sensitivity on multiple classes (in **bold**).

	Acne	Actinic keratosis	Psoriasis	Seborrheic dermatitis	Viral warts	Vitiligo
ConvNets						
DenseNet121	0.85 \pm 0.03	0.52 \pm 0.10	0.71 \pm 0.03	0.76 \pm 0.05	0.91 \pm 0.03	0.66 \pm 0.03
EfficientNet-B0	0.71 \pm 0.08	0.43 \pm 0.09	0.63 \pm 0.07	0.64 \pm 0.13	0.66 \pm 0.05	0.84 \pm 0.13
InceptionV3	0.83 \pm 0.06	0.55 \pm 0.15	0.62 \pm 0.08	0.54 \pm 0.13	0.70 \pm 0.10	0.73 \pm 0.12
InceptionResNetV2	0.70 \pm 0.06	0.41 \pm 0.12	0.62 \pm 0.06	0.67 \pm 0.11	0.43 \pm 0.12	0.74 \pm 0.07
MobileNet	0.76 \pm 0.07	0.48 \pm 0.22	0.60 \pm 0.26	0.66 \pm 0.08	0.47 \pm 0.17	0.68 \pm 0.07
MobileNetV2	0.88 \pm 0.08	0.14 \pm 0.06	0.21 \pm 0.10	0.63 \pm 0.22	0.33 \pm 0.14	0.61 \pm 0.23
NASNetMobile	0.79 \pm 0.07	0.24 \pm 0.11	0.33 \pm 0.09	0.48 \pm 0.09	0.42 \pm 0.03	0.48 \pm 0.13
ResNet50	0.79 \pm 0.05	0.40 \pm 0.16	0.69 \pm 0.10	0.74 \pm 0.07	0.54 \pm 0.13	0.80 \pm 0.04
ResNet50V2	0.85 \pm 0.10	0.55 \pm 0.13	0.58 \pm 0.04	0.61 \pm 0.08	0.74 \pm 0.07	0.71 \pm 0.02
VGG16	0.74 \pm 0.10	0.62 \pm 0.09	0.65 \pm 0.06	0.44 \pm 0.18	0.54 \pm 0.11	0.76 \pm 0.07
Xception	0.89 \pm 0.05	0.52 \pm 0.08	0.72 \pm 0.06	0.61 \pm 0.11	0.81 \pm 0.05	0.82 \pm 0.04
Dermatologists						
Average	0.95 \pm 0.03	0.67 \pm 0.18	0.88 \pm 0.06	0.59 \pm 0.11	0.88 \pm 0.09	0.92 \pm 0.05

Table A5: Diagnostic performance of ConvNets (average \pm standard deviation across five runs) and dermatologists (average \pm standard deviation across eight dermatologists) in terms of specificity on the DermXDB holdout set, split by diagnosis. Several ConvNets achieve expert-level specificity (in **bold**).

	Acne	Actinic keratosis	Psoriasis	Seborrheic dermatitis	Viral warts	Vitiligo
ConvNets						
DenseNet121	0.93 \pm 0.01	0.97 \pm 0.01	0.92 \pm 0.01	0.91 \pm 0.02	0.97 \pm 0.01	0.98 \pm 0.01
EfficientNet-B0	0.93 \pm 0.06	0.96 \pm 0.02	0.91 \pm 0.01	0.89 \pm 0.03	0.95 \pm 0.01	0.95 \pm 0.01
InceptionV3	0.92 \pm 0.03	0.92 \pm 0.02	0.92 \pm 0.03	0.91 \pm 0.06	0.96 \pm 0.02	0.97 \pm 0.03
InceptionResNetV2	0.94 \pm 0.02	0.97 \pm 0.01	0.86 \pm 0.03	0.86 \pm 0.02	0.97 \pm 0.01	0.92 \pm 0.04
MobileNet	0.91 \pm 0.05	0.96 \pm 0.02	0.88 \pm 0.07	0.87 \pm 0.08	0.97 \pm 0.02	0.94 \pm 0.02
MobileNetV2	0.66 \pm 0.15	0.99 \pm 0.01	0.98 \pm 0.03	0.79 \pm 0.11	1.00 \pm 0.00	0.94 \pm 0.07
NASNetMobile	0.65 \pm 0.04	0.97 \pm 0.01	0.96 \pm 0.01	0.86 \pm 0.02	0.96 \pm 0.03	0.96 \pm 0.03
ResNet50	0.93 \pm 0.02	0.99 \pm 0.01	0.89 \pm 0.04	0.86 \pm 0.04	0.97 \pm 0.02	0.96 \pm 0.01
ResNet50V2	0.90 \pm 0.06	0.95 \pm 0.03	0.92 \pm 0.01	0.90 \pm 0.04	0.96 \pm 0.02	0.97 \pm 0.01
VGG16	0.90 \pm 0.07	0.92 \pm 0.03	0.90 \pm 0.05	0.95 \pm 0.04	0.97 \pm 0.02	0.92 \pm 0.03
Xception	0.91 \pm 0.04	0.98 \pm 0.01	0.93 \pm 0.02	0.92 \pm 0.02	0.97 \pm 0.02	0.96 \pm 0.03
Dermatologists						
Average	0.99 \pm 0.01	1.00 \pm 0.00	0.96 \pm 0.02	0.99 \pm 0.01	1.00 \pm 0.00	1.00 \pm 0.00