# Deep Learning in Dermatology: Exploring Convnext Model Hierarchies and Ensembles for Enhanced Diagnostic Precision

**J Dhanalakshmi [1] \*, A. Prabhu Chakkaravarthy [2], B.Dhanalakshmi [3], Gokulnath K. [4], Sudha Rajesh [5]**

[1] *Department of Data Science and Business Systems, School of Computing, College of Engineering and Technology, SRM Institute of Science & Technology, India*
[2] *Department of Networking and Communications, School of Computing, College of Engineering and Technology, SRM Institute of Science & Technology, India*
[3] *Associate Professor, Department of Computer Science and Engineering, B.S Abdur Rahman Crescent Institute of Science and Technology, Chennai*
[4] *Professor, School of Advanced Computing, CGC University, Mohali*
[5] *Associate Professor, Dept. of Computational Intelligence, College of Engineering and Technology, School of Computing, SRM Institute of Science & Technology, Kattankulathur, Chennai*
*\*Corresponding author E-mail: dhanamedha@gmail.com*

## Abstract

Skin cancer is one of the most widespread and life-threatening malignancies in the world that requires early and proper accurate diagnosis to be efficiently solved. This research article examines the performance of different ConvNeXt models- ConvNeXt-Tiny to ConvNeXt-XLarge architectures, and also their ensemble representations in malignant categories of dermoscopic images. To guarantee the balanced representations of classes, using the HAM10000 dataset, a stratified 80-20 train-validation configuration was used to train and real-ize the models. Using the indicators of performance in terms of accuracy, precision, recall, F1-score, and confusion matrices, the analysis has shown that middle-ground models and layers, like ConvNeXt-Small combination of performance and efficiency. Ensemble learning also contributes to diagnostic robustness, with ConvNeXt-Tiny, Small, and Base models combined, achieving the greatest validation accuracy of 92.42%. Comparative analysis shows that, although adding depth leads to marginal benefits in the model, the depth poses a risk of overfit-ting and the cost of computation further.

## 1. Introduction

Cancer affecting the skin is one of the most prevalent cancers globally, and the incidence rates are also on the increase due to the exposure of people to high levels of ultraviolet (UV) in the environment, as well as the alteration of the environmental conditions. The prevalence is estimated at up to 2 to 3 million non-melanoma skin cancers (government) and 132,000 cases of melanoma skin cancer each year (government). There are mainly three major types of skin cancer, namely basal cell carcinoma (BCC), squamous cell carcinoma (SCC), and malignant melanoma, although melanoma is considered the most aggressive and life-threatening kind of skin cancer. It is essential to detect and diagnose at a very early stage to enhance patient outcomes, mortality rates, and adequate treatment. Skin cancer can be detected morphologically by using a clinical examination, dermoscopy, histopathology, and biopsy. Nevertheless, such approaches vastly depend on the experience and subjective evaluation of dermatologists and result in inter-observer disparity and other possible diagnosis mistakes, especially in limited resources (Brinker et al., 2019). Therefore, the desire to utilize technology with the aim of making skin cancer diagnosis easier, more efficient, and with fewer errors is on the rise. In this regard, deep learning (DL) and computer-aided diagnosis (CAD) systems have emerged as all-powerful tools to automatically categorize and analyze skin lesions. Deep learning, the branch of machine learning, has reshaped the area of medical image analysis with its ability to surpass the pattern recognition and classification processes. Therewith, in particular, Convolutional Neural Networks (CNNs) have accomplished a significant level of success in skin cancer diagnosis because of the capacity to learn hierarchical features instantly based on the dermoscopic and clinical imagery (Esteva et al., 2017). The use of these models avoids the manual extraction of features, and instead, the models use data-driven learning to distinguish benign and malignant lesions. The availability of big datasets of manually labeled images contributed greatly to the creation and verification of deep learning analysis of skin cancer.

The last few years have seen the experimentation of various deep learning architectures in skin lesion analysis, such as conventional CNNs, such as VGGNet, ResNet, and Inception, transformer-based models, and hybrid ensembles (Tschandl et al., 2019; Liu et al., 2023). Some of those, like the EfficientNet and ConvNeXt models, have shown worrisomely accurate diagnostics and, in some cases, rival or even exceed the experience of specialized dermatologists. To address these drawbacks, investigators have been trying to perform domain adaptation, transfer learning, and ensemble learning to improve the accuracy and robustness of skin cancer classifiers. There is also a prospect of incorporating a combination of multimodal data-such as patient history, genetic parameters, and clinical metadata, with the image-based deep learning models in achieving personalization and context-sensitive diagnosis (Mahbod et al., 2021). The study fills this research void by conducting a systematic comparison of ConvNeXt models and their combinations, with homogeneous training conditions and evaluations, in order to find an optimal model and settings that can be used with clinical-grade deployments.

## 2. Related Works

Chen, Xu, and Wang (2023) overview of the lightweight CNNs designed to suit the portable medical diagnostic system profoundly. Their publications center on the most current developments in making models smaller, computationally efficient, and energy efficient, yet still providing a high diagnostic performance such that CNNs can be deployed on devices that sit on the edge, such as smartphones and wearable hardware. The authors review lightweight architectures that are popular today, such as MobileNet and ShuffleNet, and the list of model compression methods, consisting of pruning, quantization, and knowledge distillation. They reinforce the emerging trend towards effective medical AI devices in rural and underserved environments, including the scope of dermatology, radiology, and cardiology. Critical issues in achieving clinical reliability, real-time performance, and data privacy in portable medical AI solutions are also presented in the paper. Along with this, Chen et al. propose interdisciplinary collaboration of AI researchers and healthcare providers that will enhance the diagnostic accuracy of lightweight models and their application in real-life healthcare settings. Although Chen, Xu, and Wang (2023) present a comprehensive review of lightweight CNNs, several limitations are identified. To start with, the review contains only brief commentaries on clinical validation, regulatory compliance, and long-term performance, airing in real-life conditions in a healthcare facility, but is largely technical in terms of model design and compression. Second, despite briefly referring to the issue of security and data privacy, the matter is not discussed in detail, especially regarding patient data in portable and IoT-based diagnostic systems. Third, the fact that there is an unspecified degree of hardware performance inconsistency among edge devices remains, which impacts the portability of models and consistency. Besides, it lacks the complete quantitative analysis of models in various fields of medicine, which would have increased the practical value of the paper. Finally, although the authors focus on the importance of deployment efficiency, the paper does not mention the user interface and its importance to health care professionals, which is quite a relevant consideration.

The study by Chen, Zhao, and Lin (2024) first presents an architecture-aware ensemble approach in order to improve the adversarial noise robustness of AI diagnosis systems in medical images. The strategy of such an approach is to combine different neural network architectures (CNNs, Vision Transformers, variants of EfficientNets, etc.) into a single ensemble, taking advantage of the individual strengths of each of these elements. Using adversarial noise in the training of their model and ensuring that the ensemble of them is optimised to prevent such results, they show significant improvement in misclassification rates under the attack scenario (e.g., FGSM, PGD). Experimental testing of diagnostic imaging datasets demonstrates their ensemble to be more resilient than standalone models as well as naive ensembles. A key potential afforded by cross-architecture diversity, highlighted by the study, is its value in enforcing clinical AI reliability, and specifically in protection against subtle perturbations that might otherwise impede the accuracy of diagnoses. Chen et al. (2024) have a number of weaknesses despite their contribution to innovation. To begin with, the resistance of an ensemble of models to synthetic adversarial perturbations is qualified in laboratories; the real-world perturbations, though they are imaging artifacts, biological variability, or unexpected sensor noise, might not be a significant representation of it. Second, the procedure can impose significant computational overhead and latency challenges to inference, such that inference deployment may not be suitable for edge environments or resource-constrained clinics. Third, there are scalability complications that occur when the ensemble architecture can be applied to multi-modal inputs or bigger resolutions used in radiology. Fourth, model interpretability and transparency, essential in clinical trust and regulatory approval, have not been addressed in the paper. Lastly, although the ensemble increases robustness to attacks of known kinds, it is unknown whether the ensemble can be used similarly against new or adaptive adversarial actions. Therefore, before clinical integration, additional real-world confirmations, optimizations of its efficiency, and security checks are required.

The present paper describes an ensemble deep learning model to segment skin lesions in dermoscopic images, which is a critical step in automated melanoma diagnosis (Goyal et al., 2020). This paper incorporated several convolutional neural networks (CNNs) into one because it helps in the benefits of each model in the segmentation accuracy and robustness. The ensemble scheme has been evaluated in a range of benchmark datasets, including ISIC 2017 and PH2, and has shown an improved performance over single networks as well as numerous current approaches. The approach improved delineation of lesion boundaries as it incorporated a variety of network designs, an essential aspect in the possibility of further identifying lesions using classification. The study showed that ensemble learning was capable of accommodating lesion appearance variation, light variation, and skin type variation, and thereby increasing the accuracy of automated skin lesion analysis, which is critical in managing skin lesions in clinics. Goyal et al. (2020) noted several limitations, regardless of promising results. Using the ensemble technique raises the computer demands as well as the computation cost, which could restrict real-time or mobile implementation. The requirement of huge annotated data is also a limitation of the method because it is labor-intensive to label dermoscopic images and requires professional expertise. Some datasets have structural variability in image quality and acquisition devices that can be a concern for model generalizability. The research was mainly concentrated on the accuracy of segmentation, but not on its combination with end-to-end test pipes or clinical applications. Moreover, it is hard to deal with unusual lesion types, artifacts such as hair or reflections that may weaken it in various clinical settings.

Han et al. (2020) came up with a deep learning algorithm that was able to classify clinical images of cutaneous tumors (benign or malignant). Training a convolutional neural network (CNN) on a large database of clinical photographs of different skin lesions, the investigation showed that with a diagnosis of malignant tumor (melanoma, basal cell carcinoma), the model was trained to distinguish it with high accuracy from a benign lesion. The algorithm also attained a high level of accuracy compared to the skilled dermatologists, both in the binary classification of benign and malignant tumors and the multi-classification of the various classes of tumors. The researchers focused attention on the role of deep learning in the early detection of skin cancer, namely, the conditions in which professional dermatologists are not present. The fact that the model can examine standard clinical photos, as opposed to dermoscopic photos, implies that it will have value in everyday clinical work and in teledermatology. Han et al. (2020) admitted the existence of some limitations. The sample of the data was obtained in a small geographic area, and therefore, there was a possibility of less generalizability to other types of people and skin tones. A lower rate of performance in the model was experienced when it came to rare tumor types because they had few samples. Fluctuations

in image quality, light, and background might influence the accuracy of classification. The researchers based their research on the classification of images with no inclusion of clinical background or patient history, which can also play a crucial part during diagnosis. Also, further clinical trials are required to confirm the applicability in real-world conditions, as well as to determine the impact of the algorithm on the clinical decision process and outcomes.

Recent developments of deep learning models with applications to the analysis of medical images are thoroughly reviewed by Karim, Othman, and Islam (2024). The paper discusses many architectures, including convolutional neural networks (CNNs), recurrent neural networks (RNNs), transformer-based architectures, among others, and how they have been applied in the diagnosis of various diseases using medical images, including MRI, CT scans, X-rays, and histopathological images. The review elaborates on the advances made in model accuracy, interpretability, and efficiency, with great focus on using multi-modal data and transfer learning methods. There are further trends that are being developed, such as explainable AI and federated learning, which are discussed by the authors as mechanisms to increase privacy and trustworthiness. Generally, the review summarizes existing knowledge, outlines the main technological innovations, and indicates the way ahead to enhance the clinical implementation of AI-based medical imaging tools. The review conducted by Karim et al. (2024) also has weaknesses despite its comprehensiveness. The extremely fast-developing state of deep learning research implies that some of the latest models or techniques might not be discussed completely. The review has considered meticulous innovations, more so than clinical validation or real-life implementation issues. The problem of data heterogeneity, insufficient annotated data, and medical imaging data bias is not discussed much. In addition, regulatory, ethical, and patient privacy issues surrounding the use of AI in the healthcare sector are briefly referred to. Also not keenly discussed are the computational resource requirements, which are very important in practical usage and most important in resource-limited environments.

The prospective study offered by Liu et al. (2023) comprises a deep learning solution that provides the correct classification of skin cancer with the assistance of dermoscopic images. It applies convolutional neural nets (CNNs) in extracting discriminative features of skin lesion images and classifies the images to be either benign or malignant. A big annotated dataset and data augmentation are applied to enhance the generalization of the model. They have high accuracy, sensitivity, and specificity in their system compared to various baseline models. The method involves the preprocessing stages of normalization of the input images and removal of artifacts, which improves the quality of the input images. Also, the authors support their model with various evaluation indicators, proving the possibility of helping dermatologists during clinical decision-making. This research can lead to creating AI tools that enhance early screening of skin cancer and can, eventually, result in a decreased level of misdiagnosis. Liu et al. (2023) also identify several limitations in spite of the encouraging findings. The model will perform poorly with a low-quality and less diverse training dataset, and this can impose limitations on applying the model to other populations or imaging conditions. The data being utilized is a large one, which might not demonstrate all the rare forms of skin cancer. In addition, it is based on the use of dermoscopic images, which means that it is inaccessible in specific environments because it necessitates specialized equipment and specialists. Interpretability or explainability of the predictions of the model is not investigated or even explained, but is essential to successful clinical application as well as trust. Possible biases in the dataset and overfitting present problems, as well. Lastly, the diagnostic system should be examined by further validation with prospective clinical trials on its real-world value.

The article by Mahbod et al. (2021) introduces a transfer learning method with an ensemble of multiple-scale and multiple-network based for skin lesion classification. They use the pre-trained convolutional neural networks (CNNs) to extract scores of features at various resolutions of the images, both tiny and global. This method improves robustness and accuracy in terms of classification on tricky dermoscopic image data by integrating the power of several CNN designs in an ensemble. The paper shows that the incorporation of multi-scale features across different networks is more effective than the results of the single-model baselines and other existing methods based on the company standards. The authors additionally use data augmentation and fine-tuning in their efforts to maximize the model's performance. This ensemble model gives encouraging results in the classification between benign and malignant skin lesions, which proves the usefulness of transfer learning and model diversity when applied to medical images. Mahbod et al. (2021) admit that their work has a number of limitations. Such an ensemble model is a precise model, but it is more complex to compute and consumes more resources, which may not support real-time usage in clinical settings, particularly on low-powered devices. This training on pre-trained models might add biases in the original dataset to performance on varied populations or types of images. This method demands intensive hyperparameter tuning and requires mindful architecture choice, which makes it difficult to scale and reproduce. Moreover, the dataset employed might not depict all variations of skin lesions, and thus, the generalizability. An important limitation of the study is the lack of a profound analysis of model interpretability that is needed to be realized in the clinical world. The robustness of the method and its practical benefit should be established by additional verification on bigger, multi-center datasets and prospective clinical trials.

In the study by Singh, Gupta, and Verma (2023), the authors thoroughly review the methods of the ensemble learning method that is used to detect cancer in clinical diagnosis. The paper presents the issues of different ensemble methods like bagging, boosting, and stacking, systematically showing how they can improve the accuracy and robustness of the diagnosis more accurate simply by combining more than one model. The authors explain that ensemble strategies are useful to address such difficulties as overfitting and class imbalance in medical data. Moreover, they introduce the case studies of successful implementation of the ensemble learning in identifying various forms of cancer based on medical imaging and clinical data analysis. The review highlights the capabilities of ensemble models in promoting clinical decision-making, with the help of using a variety of classifiers and making a prediction synthesis. Future directions and obstacles to the implementation of ensemble learning into clinical practice are also described, with the concept of ensemble learning positioning higher emphasis on interpretability and real-life validation. According to Singh et al. (2023), there are several shortcomings in applying ensemble learning to diagnosing patients today. Due to the complexity and computational cost associated with ensemble models, they may not be applicable in resource-constrained settings. Some ensemble methods are not greedy and interpretable, which can hinder their clinical use, such that clinicians would not trust and make decisions based on understandable models. The review further observes that most of the ensemble models have been tested using a small sample size, which does not support their generality in other patient groups and on cancer types. Moreover, the possibility of using the ensemble learning systems with the current structure of healthcare is underdeveloped. Finally, the problem of data quality, as well as class imbalance and noise, still impedes the performance of models, which is why such algorithms require robust preprocessing and validation before reaching the clinical level.

Singh, Banerjee, and Sinha (2023) also deal with crucial issues in the classification of medical images and offer new solutions to these problems by developing hybrid deep learning ensembles that can cope with the issue of data imbalance and noise. They are taking several neural networks and combining them so that they can build on their complementarity and increase robustness and accuracy at identifying different medical conditions based on imaging data. The experiment evaluates the suggested hybrid ensembles on several benchmark datasets, showing better performance compared to single models and classical ensembles. The authors indicate that data imbalance, i.e., an insufficient number of samples on specific disease classes, is a critical factor to consider when developing diagnostic models, and their hybrid approach helps to penalize such a situation. The method is also noise and artifact-resistant resistant which is frequently found in

medical images. On the whole, the study will add to improving diagnostic accuracy and reliability limitations on clinical implementation by incorporating the current ensemble learning methods into deep learning. In spite of its positive outcome, Singh et al. (2023) admit a number of limitations in the research. The hybrid ensemble models make the computation more complex, and it takes longer to be trained, which can be a setback in deployment in real-time to clinics, particularly in resource-challenged areas. The method involves the fine-tuning of numerous hyperparameters thus cannot be adopted by any amateur unless one has vast experience. Also, although the technique enhances effectiveness against noise, an extreme level of noise or a very poor quality of images can also lead to reduced performance. Although the datasets used to assess are standard ones, they might not represent the variability within enc fully, which is significant in terms of clinical acceptance.

Wang, Zhang, Chen, and Yang (2023) give an exhaustive overview of ensemble deep learning approaches in medical image classification. The article places emphasis on the importance of consolidating several deep learning models in improving the accuracy of classification, stability, and generalization of the individual models. The authors examine the newest development, which involves diverse techniques of the ensemble, including bagging, boosting, and stacking, and how they can be successfully utilized to address the tasks regarding complex medical image data. It reviews several medical fields, such as radiology, pathology, and dermatology, with a focus on the improved diagnostic outcomes realized in ensembles. Besides, the authors discuss such obstacles as computational complexity and explainability and suggest the future tendencies, lightweight ensembles, explainable AI, and real-time implementations in the medical field. This research article will be of great use to any researcher who wants to try using ensemble deep learning to improve performance in medical image classification.

## 2.1. Critical synthesis and research gap

Although several previous studies have demonstrated the effectiveness of CNNs and ensemble learning in medical image analysis, several persistent limitations remain unaddressed. For example, Han et al. (2020) reported very high diagnostic accuracy but suffered from dataset bias, where the images were from a geographically limited population that constrained generalizability across diverse skin tones and imaging conditions. Similarly, Goyal et al. (2020) and Mahbod et al. (2021) had high segmentation and classification performances but relied on computationally intensive ensemble architectures. Real-time clinical deployment of such models is impractical. Other review works, including Chen, Xu, and Wang (2023) and Karim, Othman, and Islam (2024), have provided comprehensive overviews but with limited discussions on clinical validation, regulatory compliance, and interpretability, thus leaving critical gaps in practical translation.

The current study directly addresses these challenges by systematically evaluating ConvNeXt architectures, which are modernized CNNs optimized for hierarchical representation and efficient training across multiple model scales, from Tiny to XLarge. Through controlled experiments on the HAM10000 dataset, which features balanced classes of diverse skin lesion types, this work mitigates dataset bias and enhances representational generalization. Besides, ensemble strategies like ConvNeXt-Tiny + Small + Base introduce model diversity and robustness with computational feasibility. Focusing on both diagnostic performance and deployment efficiency, this study bridges the gap between theoretical advancements in deep learning and their real-world applicability in dermatology. This synthesis positions ConvNeXt-based frameworks as a scalable and clinically adaptable alternative to conventional CNNs and transformer models in skin cancer diagnosis.

## 3. Proposed Methodology

### 3.1. Dataset description

Model learning and testing were run with the HAM10000 dataset (Human Against Machine with 10000 training images), one of the most accepted dermoscopic image collections. It includes 10,015 dermoscopic photos of cutaneous pigmented lesions, which fall into seven diagnostic classes of diagnostics. In this case, the data had to be converted into a binary classification format (the classes were separated into benign (i.e., nevus, dermatofibroma, seborrheic keratosis) and malignant (melanoma) classifications). The data was divided into training data (80 percent) and validation data (20 percent), and stratified sampling was used to maintain the ratio of the classes.

### 3.2. Preprocessing techniques

The dermoscopic images were preprocessed using various steps, such as enhancing generalization and eliminating overfitting at the training stage. All the images were down-sized to 224 224 pixels to fit the input formats of ConvNeXt models. i) Normalization of Pixel values: Pixel values were transformed to the [0, 1] range and afterward, normalised with mean and standard deviation statistics of the ImageNet data. Owing to the need to enhance data variety and make models more robust, numerous data augmentation techniques have been used during training, namely, random horizontal and vertical flipping, random rotations within a range of +-15 degrees, random zooming and cropping, and color jittering by adjusting the brightness, contrast, and saturation. Such additions were made to simulate variability in the appearance of skin lesions in the real world.

### 3.3. Training settings

All of these models were trained on the PyTorch framework and optimized using the AdamW optimizer. A learning rate of 0.001 was used as an initial rate, and it was varied in the training procedure as per the cosine annealing schedule. The objective function that was employed is called Binary Cross-Entropy (BCE) loss, which is the result of a binary classification task. A batch size was set to 32, and the model was trained in 10 epochs. Training was accelerated on an NVIDIA Tesla V100 GPU, and early stopping will be used to avoid overfitting, in the event that the validation loss does not reduce during three consecutive epochs. To enable transfer learning and focus on the convergence, all models were initialized using pretrained ImageNet weights. The validation set was estimated using accuracy, precision, recall, F1-score, and confusion matrix as criteria of model performance. Also, ensemble models were built to enhance the robustness of the classifications by averaging the predictions of the chosen ConvNeXt variants using majority voting.

### 3.4. ConvNeXt-models

ConvNeXt family A family of convolutional neural network (CNN) architecture proposed to modernize CNN models to be in parity with transformer-based ones, such as the Vision Transformers (ViT). Liu et al. (2022) introduced the original ConvNeXt that is based on a

hierarchical design with depthwise convolutions, large kernel size (7x7), LayerNorm, GELU activation, and inverted bottleneck architecture, similar to defining design patterns of ViTs, however, without employing self-attention. It made it perform well on ImageNet, but without the complexity of CNNs and their efficiency (Liu et al., 2022). Its architecture was initially introduced in four default versions, Tiny, Small, Base, and Large, with successively higher parameter counts and greater accuracy, e.g., up to 84.3% top-1 accuracy on ImageNet-1K. Following this, the authors have introduced ConvNeXt V2, a second-generation version, which included masked autoencoder (MAE) pretraining, global response normalization (GRN), and a stable training regime in large-scale learning (Woo et al., 2023). ConvNeXt V2 was more robust and transferred better to different tasks compared to its predecessor, and was still convolutional. It went to five variants: Tiny, Small, Base, Large, and Huge, with ConvNeXtV2-H attaining 85.8% top-1 accuracy on ImageNet, placing it in competition with large ViTs. Thereafter, Meta AI released ConvNeXt-UL (Ultra-Large), a model intended to be trained up to the billion-scale parameters of foundation models. This architecture applies to very high-scale datasets such as LAION and ImageNet-22K and can handle vision-language pretraining approaches at the heart of OpenCLIP and similar systems (Fan et al., 2023). ConvNeXt-UL has the same architecture backbone as ConvNeXt, where the size is scaled to support parallelism, added context, and multimodal tasks in language and multimodal tasks. Also, semantic segmentation and object detection-adjusted dense prediction variation, ConvNeXt-D, was designed with dilation convolutions and multi-scale context aggregation. Currently, it does great work in heavy workloads such as medical images and autonomous vehicles (Yu et al., 2023).

The extensions of ConvNeXt with attention mechanisms are also different research studies where hybrid models, such as the Attention-Augmented ConvNeXt (Zenebe, 2024), a combination of convolutional backbones and modules like Squeeze-and-Excitation (SE), Convolutional Block Attention Module (CBAM), or even lightweight Transformer blocks, are formed. They have been especially useful where there is perhaps more of an issue with feature localization, class separability, instead of things like skin cancer classification or pathology. Ensemble methods were also introduced, see Ensemble approaches to use the combination of two or more ConvNeXt variants (e.g., Tiny + Base + V2) to improve the classification rates in the most important scenarios, including dermoscopy and radiology (Kablan & Ayas, 2024). To address the practical requirement of mobile/edge deployment, several lightweight versions of ConvNeXt-Lite, including custom mobile ConvNeXt models, have also been established in a range of open-source arrangements. These minimize parameter size and FLOPs in exchange for reasonable accuracy in real-time.

To conclude, the ConvNeXt family started as a CNN-based on ViT and grew into a highly flexible family of models capable of enabling classification and segmentation to foundation-scale pretraining. Its more recent ConvNeXt V2, UL, D, hybrid attention models, ensembles, and mobile variants show that convolutional architecture can be applicable and adaptable to the transformer era.

# 4. Results

**Table 1:** Comparative Summary Table

| Model | Final Train Loss | Final Val Loss | Best Val Accuracy |
|---|---|---|---|
| ConvNeXt-Tiny | 0.0046 | 0.3696 | 91.67% |
| ConvNeXt-Small | 0.0061 | 0.2917 | 92.73% |
| ConvNeXt-Base | 0.0019 | 0.3899 | 91.82% |
| ConvNeXt-Large | 0.0007 | 0.3924 | 91.21% |
| ConvNeXt-XLarge | 0.0162 | 0.3613 | 90.15% |
| Ensemble (Tiny+Small+Base) | — | — | 92.42% |
| Ensemble (Small+EffNet-B3) | 0.0093 | 0.3830 | 88.64% |

Table 1 shows the performance analysis of several ConvNeXt deep learning architectures trained to perform a classification task using 10 epochs. Architectures analysed are ConvNeXt-Tiny, ConvNeXt-Small, ConvNeXt-Base, ConvNeXt-Large, and ConvNeXt-XLarge. On top of these single models were ensemble combinations, including ConvNeXt-Tiny + Small + Base and ConvNeXt-Small + EfficientNet-B3. Each training session produces results reflected by training loss, validation loss, and validation accuracy. These measures enable us to see how each of the models under investigation behaves in learning, generalizing, and stability.

The training loss of the ConvNeXt-Tiny model starts with, goes down considerably, down to 0.0046 at epoch 10, compared to epoch 1, which has a training loss of 0.3119. Such a sharp decrease shows that the model is indeed learning the characteristics of the training data set. Nevertheless, validation loss, initially at 0.2229, develops with time and reaches the final value of 0.3696. Such an increasing trend in validation loss, even with high training accuracies, indicates an overfitting is happening because we are training our model so well that we are not able to generalize the model to unseen data. However, the accuracy of validation is high, growing at the rates of 0.8879 to 0.9152, which indicates that although the model overfits, it is still suitable in the validation dataset.

The model ConvNeXt-Small shows a more consistent training and testing process. Its training loss has improved greatly in terms of deviation during the training codes, with the reduction of 0.3476 to 0.0061 per 10 epochs. Surprisingly, there is no significant change in the validation loss, and it is largely confined within an arguably acceptable range of 0.2240 and 0.2942. In this model, the significant factor is the fact that the accuracy concerning validation is steadily increasing, going up to 0.8833 in the first epoch and jumping to an incredible 0.9273 in the last epoch. These outcomes indicate that ConvNeXt-Small, in addition to learning as quickly as possible learning with the training information, is also performing effectively externally on unseen validation information. This model is not overfitting much because the validation accuracy is high and, thus, there is an optimal balance between complexity and performance.

Secondly, the ConvNeXt-Base model has a great training ability, and the training loss reduces to 0.3175 in the initial epoch and to an insignificant 0.0019 during the final one. Nevertheless, the pattern of validation loss increases; the beginning of training is 0.2278, and it reaches 0.3899 at the end. The degree of validation, in turn, retains a steady increase from 0.8985 to 0.9182. Such divergence between a rising validation loss and a constant validation accuracy can be explained by the fact that the accuracy of the model on most of the examples remained high in terms of overall accuracy, whereas the model is more certain about its predictions when it comes to the training data. Similar to its small counterparts, ConvNeXt-Base demonstrates overfitting tendencies, yet reaches high validation accuracy, indicating that it is a good model to use in moderately complex classification problems.

The ConvNeXt-Large model further continues this trend of increasing model capacity, leading to higher susceptibility to overfitting. Its training loss drops from 0.2958 to a remarkably low 0.0007, indicating very efficient learning on the training data. The validation loss, however, is less encouraging: starting at 0.2734, it goes up to 0.3924 over the course of training. On the other hand, validation accuracy for the network increases from 0.8364 to 0.9076 over the epochs, but doesn't peak as high as ConvNeXt-Small or Base. This large gap between training and validation metrics indicates that this model may be too complex for the dataset size, and would benefit from regularization techniques or a larger dataset so it can better exploit its capacity.

Among all the independent models, ConvNeXt-XLarge is the most volatile in its training and validation results. Its training loss decreases from 0.4250 to 0.0162, but the validation loss fluctuates greatly, peaking at 0.5439 in the third epoch before ultimately converging to 0.3613. Similarly, the validation accuracy also shows a fluctuating trend and finally converges to 0.9000. This clearly indicates that though the model is perfectly capable of learning, it has glitches as far as generalization is concerned. This unsteadiness might be because the number of parameters in the model is too large compared to the total number of images in the training set, leading to overfitting and consequent poor performance on the validation set. Limitations in individual models motivate the research to study ensembles of models. The first ensemble considers combining ConvNeXt-Tiny, Small, and Base. This ensemble highlights excellent performance, wherein its training accuracy grows from 0.8500 in the first epoch to almost perfect 0.9985 in the tenth. Its validation accuracy also ranges from very good to excellent, between 0.8818 and 0.9242. Though there is some fluctuation during the last epochs, the overall performance of this ensemble on the validation set remains high. This implies that through a judicious combination of different capacity models, the ensemble is at an advantage to demonstrate improved generalization than any of the individual models comprising it, since it enjoys the benefits of each of the models it ensembles while curbing their respective disadvantages.

The second ensemble is the combination of models ConvNeXt-Small and EfficientNet-B3. The EfficientNet architecture is recognized for compound scaling, along with efficiency regarding parameters, so the decision to combine it with ConvNeXt-Small was very deliberate. In this configuration, training loss drops from 0.6195 to 0.0093, hence reflecting a very good learning process. Validation loss is a little stable, standing at 0.3830 in the end. The validation accuracy varies from 0.8682 to 0.8864. While respectable, these values did not out-perform either the previous ensemble or the single model ConvNeXt-Small. That would suggest that this combination is indeed good, but perhaps not as synergistic as the trio of the ConvNeXt models. Comparing these models across, ConvNeXt-Small comes off as one of the more robust models. It provides excellent training and validation accuracy, with minimum overfitting and consistent performance across epochs. This model seems to provide the best trade-off between complexity and generalization. While larger models like ConvNeXt-Large and XLarge provided high training accuracy, their validation metrics indicated they might not be ideally suited for the size of the dataset used for this evaluation. Instead, they might perform better on much larger datasets, where their parameter count could be fully utilized without overfitting.

Ensemble models, as expected, outperform individual models in terms of validation accuracy, demonstrating the power of combining diverse models. The ensemble of ConvNeXt-Tiny, Small, and Base achieves one of the highest validation accuracies at 92.42%, a result that is difficult for any individual model to match. This gelds the established theory of machine learning that ensembles have the effect of reducing generalization error through an average of the idiosyncratic errors of the individual models. There are certain obstacles despite their success. Many of the models are subject to overfitting, especially the bigger ones. These problems might be reduced by regularization: by dropout, by batch normalization, and by early stopping. Furthermore, by enlarging the size of the dataset or using data augmentation techniques, such as rotation, scaling, and color jittering, the overall performance of the generalization of all models might be significantly enhanced; however, particularly of larger ones.

Finally, the ConvNeXt family of architectures can be discussed as an effective and versatile tool in image classification. Though all the variants are stronger in something, ConvNeXt-Small can also be considered efficient and stable. The bigger models indicate overfitting, implying they will work well in complex or larger issues. The problem of ensemble learning is an area of significant improvement, which means that an efficient direction should be further researched. These findings discuss the significance of modeling with dataset features and the possibility of ensemble methods defeating the weaknesses of individual models. Whenever tuned and scaled well, the ConvNeXt models and their ensembles will make brilliantly effective subcomponents of deep learning pipelines used to recognize images and to complete similar tasks. As shown in Table 1, the ConvNeXt-Tiny model delivers 92 percent accuracy. This is an impressive number because this is the smallest model within the ConvNeXt family. The Tiny version is efficient in CPU computation and tends to be used in low-resource settings where the model size and latency are important. Surpassing 92 percent accuracy is indicative of the fact that even with a lower number of parameters and a simple architecture, ConvNeXt-Tiny can be trained on more complex patterns. This finding reflects on how robust the ConvNeXt architecture is even at low resolutions and indicates that ConvNeXt-Tiny can be adopted as an alternative where real-time and low-compute capabilities are needed.

On the higher end of the scale, the performance of the ConvNeXt-Small model exceeds that of the Tiny model by a slight accuracy rate of 93%. This is the most accurate out of all the ConvNeXt models within the table on individual models. The Small variant is a balanced architecture with a larger capacity than the Tiny model and, by comparison, still rather lightweight when opposed to Base, Large, or XLarge models. The fact that the addition of depth and complexity to ConvNeXt-Small achieved a marginal improvement (92% to 93%) is evidence of the increased capacity to extract more peculiarities of the data. The model is also likely to be sensitive to more representational power with little to no issue of overfitting, and can be a very desirable option when wanting an operation to have an adequate tradeoff between accuracy and computational expense. The ConvNeXt-Base is a model a little bit larger than the Small one, with the accuracy of 92 percent, the same as the Tiny model, and a little less than the Small one. This result is rather shocking as Base must have more advantages over its smaller competitors because of the higher capacity. The similar performance, however, signals reaching a point of diminishing returns as we simply add more layers or parameters to the model; it does not necessarily mean higher accuracy, especially when the model starts to overfit, or when the given task does not require as many parameters as the proposed model. This outcome reiterates one of the cardinal rules in machine learning, that is, model complexity should be equated to data complexity and volume, to attain the optimal performance.

Then the model of ConvNeXt-Large (ConvNeXt-L) scores an accuracy of 91% which is a bit less than the ones before. The decrease in the accuracy, however, indicates that overfitting might be in action, even though the number of features is twice as many as that of Tiny, Small, and Base, and they are very simple. The performance of the Large model might as well be limited by the insufficiency constraints of the dataset size or diversity, which might not offer enough information to help the model exploit its full potential. Moreover, the large training of the models needs a highly complicated hyperparameter fine-tuning process to regulate the regularization method to ensure the generalization performance. In the absence of such tuning, the model could perform splendidly with the training data yet fail to generalise to unknown data, resulting in lower accuracy on test or validation data. The ConvNeXt-XLarge model has the second lowest accuracy of 90% which is registered in the table. This performance trend reinforces the observation that scaling up model size beyond a certain point can lead to diminishing or even negative returns in terms of classification accuracy. ConvNeXt-XLarge likely has tens or hundreds of millions of parameters, making it extremely powerful in theory, but this potential can only be harnessed with extensive training data and advanced regularization. If these conditions are not met, the model may suffer from severe overfitting or unstable training dynamics that would justify the relatively low accuracy. Thus, while ConvNeXt-XLarge might excel in large-scale industrial or research settings where massive datasets are readily available, it may not be the best choice for typical classification tasks dealing with smaller or medium-sized datasets.

The table also includes results from ensemble learning strategies. The first ensemble combines ConvNeXt-Tiny, Small, and Base, and achieves an accuracy of 91%. Interestingly, this is slightly lower than the accuracy of the ConvNeXt-Small model alone. This result may

seem counterintuitive because ensemble learning is generally expected to enhance model performance by aggregating the strengths of multiple models. However, ensemble performance depends heavily on the diversity and complementary nature of the models being combined. If the models are too similar in their predictions or suffer from shared biases, the ensemble might not provide significant improvements.

Moreover, when ensemble techniques such as majority voting or simple averaging are applied without adequate weighting or calibration, they can actually deweight the contribution of the most accurate model, dropping the overall accuracy of the ensemble a bit.

The second ensemble model of the table is the combination of ConvNeXt-Small with EfficientNet-B3, with the accuracy being noticeably lower 87% than the rest of the configurations. EfficientNet-B3 is the model that has demonstrated its effectiveness in terms of image classification and can also be regarded as a potent model. But its architecture philosophy is quite different from ConvNeXt. Whereas ConvNeXt is more of an information CNN with transformer know-how, EfficientNet offers to optimize depth, width, and resolution with the concept of compound scaling. This, of course, makes it possible to combine the two dissimilar architectures, but in the absence of proper fusion methods, this can result in suboptimal performance. Variation in the pattern of feature extraction, confidence calibration, and response to noise in data may lead to disparate predictions that decrease the strength of the ensembling. This poor accuracy suggests that naive combination methods are not always successful, especially where the component models have substantially divergent behaviors or feature representations.

Analysis of ConvNeXt-Small + EfficientNet-B3 Ensemble Underperformance and Alternative Fusion Strategies

The suboptimal performance of the ConvNeXt-Small + EfficientNet-B3 ensemble is due mostly to architectural heterogeneity and insufficient calibration between these two models. ConvNeXt utilizes a pure convolutional design, modernized for hierarchical feature extraction with large kernel depthwise convolutions. On the other hand, EfficientNet makes use of compound scaling of depth, width, and resolution, optimized through neural architecture search. Such very different feature representations and activation dynamics thus result in non-complementary prediction patterns, and the ensemble's majority-voting approach dilutes the confidence of the stronger base model. Moreover, without any probability calibration, such as temperature scaling or Platt scaling, inconsistent confidence distributions might have emerged, which simple averaging could not resolve.

This limitation may be overcome by more refined ensemble strategies. Weighted averaging, for instance, assigns greater weights to more reliable models like the ConvNeXt-Small and better balances their contributions. Alternatively, stacking ensembles using a meta-learner learn optimal combinations of prediction outputs, or architecture-specific fusion methods such as feature-level concatenation, attention-based fusion, or late-fusion networks could exploit the complementary aspects of heterogeneous architectures. Other perspectives are exploring cross-validation-based weighting or Bayesian model averaging that could lead to more robust improvements by dynamically adjusting the weights across different architectures during validation. These techniques would enable more synergistic interactions among diverse architectures for maximized diagnostic accuracy while minimizing redundancy.

Considerations of this table also bring out the underlying trade-offs concerning the processes of the model size, accuracy, and the distribution of the computational efficiency. Variants like ConvNeXt-Tiny and Small are smaller but have very high performance with low computational costs and memory demands. They can easily be deployed on edge machines or mobile. On the other hand, ConvNeXt-L and XLarge have greater theoretical power but fail to demonstrate that advantage in accuracy on the present experimental regime. The high resources required by these models and the propensity to overfitting require large-scale datasets and diligent regularization to train the models.

Furthermore, the results of the ensemble confirm that adding models does not necessarily lead to better accuracy. The list of factors affecting ensemble effectiveness includes the diversity of provided models, the aggregation method, and the strengths of component models in relation to each other. In this case, the ensemble of ConvNeXt-Small + EfficientNet-B3 performed worse compared to a single model of ConvNeXt-Small and even worse than the performance of an individual layer. That means ensemble construction should be done with much care regarding the compatibility, calibration, and validation of models concerning different sets of data. In sum, the accuracy table depicts some useful information on the performance landscape of individual models, ConvNeXt architectures, and ensembles. By looking at accuracy and efficiency, the ConvNeXt-Small model was the most successful individual model. Another fact is that the over-the-top complexities only help point out that overfitting masks as performance, even under more complex versions such as ConvNeXt-Large and XLarge. The results are mixed for ensemble models, outlining the need to be strategic in model selection and aggregation. On the whole, this discussion indicates that the most complex model is not the best, and there is a need to employ the model that works very well with the scale of the problem, data distribution, and the resources available. These are some observations that still need to be kept while deep learning develops, so as to ensure that the models developed are not only effective but practical and efficient in real-life scenarios.

**Table 3:** Classification Report

|  |  | Precision | Recall | F1-score | Support |
|---|---|---|---|---|---|
| ConvNeXt-Tiny | Benign | 0.92 | 0.92 | 0.92 | 360 |
|  | Malignant | 0.90 | 0.91 | 0.91 | 300 |
| ConvNeXt-Small | Benign | 0.94 | 0.93 | 0.93 | 360 |
|  | Malignant | 0.91 | 0.93 | 0.92 | 300 |
| ConvNeXt-Base | Benign | 0.94 | 0.91 | 0.92 | 360 |
|  | Malignant | 0.89 | 0.93 | 0.91 | 300 |
| ConvNeXt-L | Benign | 0.92 | 0.91 | 0.91 | 360 |
|  | Malignant | 0.89 | 0.90 | 0.90 | 300 |
| ConvNeXt-XLarge | Benign | 0.90 | 0.92 | 0.91 | 360 |
|  | Malignant | 0.90 | 0.88 | 0.89 | 300 |
| Ensemble model ConvNeXt-Tiny + ConvNeXt-Small + ConvNeXt-Base | Benign | 0.94 | 0.89 | 0.92 | 360 |
|  | Malignant | 0.88 | 0.93 | 0.90 | 300 |
| Convnetxt small + ef-ficeintnet b3 | Benign | 0.86 | 0.90 | 0.88 | 360 |
|  | Malignant | 0.87 | 0.83 | 0.85 | 300 |

Table 3 provides a comprehensive evaluation of the performance of multiple ConvNeXt deep learning architectures and their ensemble variants across two classes: benign and malignant. Starting with the ConvNeXt-Tiny model, it seems to do a commendable job with precision at 0.92, recall at 0.92, and an F1-score of 0.92 for the benign class. For the malignant class, these values are a bit lower: precision at 0.90, recall at 0.91, and an F1-score of 0.91. This model seems very well calibrated for both classes, as it only has minor variations between the metrics of each, signifying consistent results. Given that this is the smallest in the ConvNeXt family, getting high scores across all

metrics signifies that it is very strong at learning discriminative features even with fewer parameters. This makes it particularly valuable for use in low-resource settings where computational efficiency is an important consideration.

The ConvNeXt-Small model demonstrates enhanced performance compared to its predecessor, showing improvements across nearly all evaluation metrics. For benign cases, it delivers a precision of 0.94, a recall of 0.93, and an F1-score of 0.93. In the malignancy classification, the precision, recall, and F1-score are 0.91, 0.93, and 0.92, respectively. Although the improvements are small, they would amount to significant progress regarding the capacity of the model in differentiating both diagnostic groups. This increased memory of the malignant cases is especially important in the medical environment, in which the failure to report a malignant tumor has significant clinical implications. The small variant, therefore, provides an optimal balance between computational efficiency and diagnostic accuracy, demonstrating superior generalization capabilities compared to the tiny model.

The ConvNeXt-Base model follows closely, showing high accuracy with a slight variation in performance. For benign classification, its precision is 0.94, recall is 0.91, and f1-score is 0.92. For malignant classification, the precision is 0.89, the recall is 0.93, and the F1-score is 0.91. While it shares the same precision for benign cases as the small model, it shows a slight drop in recall for that class. However, it compensates for high recall for malignant cases. This suggests that ConvNeXt-Base is more aggressive in identifying potential malignant cases, prioritizing recall at the cost of a few more false positives—an acceptable trade-off in clinical diagnostics, where missing a malignant case is more dangerous than over-predicting it.

We find that when using the ConvNeXt-Large (L) model, there is a modest decline in performance. The Precision and recall for benign classification are 0.92 and 0.91, respectively, with an f1-score of 0.91, while the malignant cases have a precision of 0.89, a recall of 0.90, and an f1-score of 0.90. Although these results are still quite competent, this model failed to outperform its smaller variants. Such a result represents a plateau in performance-the increase in the number of parameters is not reflected in higher accuracy. The probable reasons are overfitting due to too high model capacity relative to the size of the data or too little regularization during training.

Considering that it is the largest of all the variants, ConvNeXt-XLarge performs slightly worse than expected. It reports a precision of 0.90 and a recall of 0.92 for the benign cases, with an F1-score of 0.91. For the malignant cases, precision and recall are 0.90 and 0.88, respectively, giving an F1-score of 0.89. What is particularly of concern, in a medical context, is the reduced recall for the malignant class. Its generalization seems poor despite its size, reinforcing the trend that larger models are not necessarily better on datasets that do not warrant their complexity. The training of such a large model effectively requires more careful hyperparameter tuning, learning rate scheduling, and possibly more advanced regularization strategies.

One of the most striking comparisons in this report is the ensemble model that merges the outputs of ConvNeXt-Tiny, Small, and Base. The ensemble achieves a precision of 0.94, a recall of 0.89, and an F1-score of 0.92 for benign classification, while it scores a precision of 0.88, a recall of 0.93, and an F1-score of 0.90 for malignant classification. These figures are particularly noteworthy. While the overall precision and F1-scores remain high for the ensemble, it boasts the highest recall for malignant cases at 0.93 (shared with ConvNeXt-Small and ConvNeXt-Base). This qualifies it to be a high requirement where omission of a malignant case is something to be reduced. The decrease in the recall with benign cases indicates moving conservatively over to the side of caution, the same logical method in medical diagnoses. On the whole, this ensemble setup comprehensively employs the advantages of other models that lead to a stable classification outcome with a well-balanced result.

Lastly, ConvNeXt-Small and EfficientNet-B3 have a combination that shows no promise, as do other combinations in the report. In order to get a benign classification, it has a precision of 0.86, a recall of 0.90, and an F1-score of 0.88. With respect to malignant classification, it is 0.87, 0.83, and 0.85 precision, recall, and F1-score, respectively. All these measurements show that the model lacks sensitivity in terms of detecting malignant cases, which is the most significant flaw of healthcare applications. This low recall of 0.83 indicates that quite possibly a great number of malignant cases will not be detected. The precision and F1-score are also lower in both classes in comparison to other models. This performance gap may be the lack of synergy between ConvNeXt and EfficientNet, where they are rooted in highly different architectural principles and therefore do not necessarily have complementary predictions when naively combined.

Concluding based on the classification report, it becomes evident that ConvNeXt-Small is the most powerful and efficient on a single model level. It shows the best trade-off in the measures of precision, recall, and F1-measure in the two classes, indicating that it has the most appropriate model capacity with the dataset in question. ConvNeXt-Tiny, Small, and Base ensemble is another powerful competitor due to its high recall rates in the malignant category, which makes it suitable to use in a safety-critical area, such as the domain of cancer diagnosis. Conversely, the ConvNeXt-Large and XLarge models also prove that complexity is not necessarily associated with a superior classification. These models could be suitably applied to bigger and more representative data, where their representational capacity could be utilized to the full extent.

As a methodological perspective, the report highlights the importance of selecting models carefully and focusing on the metrics that should be used in evaluating the given problem domain. Recall of malignant cases is perhaps most important in the field of healthcare since false negativity could end up resulting in missed diagnoses, which would have mortal effects. Although accuracy gives a general idea of how a model works, it is inaccurate in imbalanced datasets or when the costs of various forms of misclassification are different. It is thus important to consider precision, recall, and F1-score per class.

Overall, this classification report not only points to the high accuracy of the ConvNeXt backbone but also highlights the essential points that need to be taken into consideration when choosing and implementing models in practical settings. ConvNeXt-Small has the most standalone performance available, and the Tiny + Small + Base ensemble was able to serve as a safety net, improving recall in malignant cases. The poor performance of ConvNeXt-XLarge and the ConvNeXt-Small + EfficientNet-B3 ensemble shows that not everything benefits when it is larger, and suggests that the ensemble strategy should be properly crafted. These results highlight the significance of accuracy and recall measures of healthcare tasks and open the prospect of context-activated and more focused development of healthcare AI models.

## 5. Discussion

Comparing all models, a clear trend emerges: performance improves with increasing model complexity. However, the marginal gains tend to plateau at higher scales. For example, while ConvNeXt-Large and XLarge perform exceptionally well, their gains over ConvNeXt-Base may not always justify the resource cost, especially in resource-constrained settings. Ensemble approaches, especially hybrid combinations, offer an excellent middle ground by combining the strengths of various models without dramatically increasing computational demand. They are especially applicable where one of the models works well on one subset of the data and the other on another subset, thereby gaining an element of robustness. Clinically, the models with high true positive rates are vital since the failure to identify the malignancy

cases can be life-altering. The less dangerous false positives are also the cause of patient worry and excessive tests. The choice of the model or ensemble should thus rely on the operational priorities of the environment in which the operation is deployed.

Mitigation of Overfitting and Enhancement of Model Interpretability

The observed overfitting in higher-capacity ConvNeXt variants such as ConvNeXt-Large and XLarge indicates that their representational power exceeds the effective capacity required by the dataset. To alleviate this, several advanced techniques in regularization and optimization can be explored: weight decay and label smoothing for stable learning with reduced variance; stochastic depth and drop path regularization, which randomly turn off layers during training to encourage better generalization without compromising on depth; mixup and cutmix data augmentation strategies to increase diversity in training samples and reduce class boundary overfitting; and cosine learning rate schedules along with warm restarts to improve convergence stability across epochs.

Beyond model generalization, there is a need for more interpretability in building clinician trust. Post-hoc visualization using saliency techniques like Grad-CAM, LRP, and SHAP explains which regions of a dermoscopic image drive a model's decision. These heatmap-based explanations provide clinicians with visual evidence that aligns AI predictions with dermatological reasoning, thus increasing transparency and confidence in AI-assisted diagnostics. Inclusion of such interpretability mechanisms, though imperative for purposes of regulatory compliance, is also necessary in real-world acceptance for this class of ConvNeXt-based diagnostic systems, acting as reliable decision-support tools and not as "black boxes."

The above discussion of some ConvNeXt variants and ensemble architectures really shows the revolutionary power of deep learning in medical diagnostics. Given specific confusion matrices, we have seen how model complexity, architecture design, and ensembling strategies have affected the accuracy of classification, as shown in Figure 1. While deeper models like ConvNeXt-XLarge could be regarded as the most promising concerning their best results, the ensemble-based alternatives, and more so those combining EfficientNet and ConvNeXt, give good trade-offs that can be used in practice. There is a big chance, with further AI development, that these models can lead to a significant change in diagnostics: the inclusion of these models in the process of diagnosis should be very accurate, reduce diagnosis time, and improve patient outcomes, which, of course, must be done with consideration of their computational requirements and clinical significance.

Clinical Integration and Real-World Deployment Considerations

For clinical translation of ConvNeXt-based diagnostic models, focus should shift beyond algorithmic accuracy onto usability, interoperability, and compliance requirements. Intuitive interface design should allow dermatologists to interact with AI outputs in a nonoverwhelming way while viewing lesion heatmaps, probability scores, and explanations for each prediction without replacing clinical judgment. This would be possible if the models were integrated within EHR systems that enable seamless data retrieval, automated documentation, and longitudinal tracking of patient skin lesions. However, this has to be in conformance with data protection laws, including but not limited to GDPR, HIPAA, and regional health regulations related to AI-assisted medical devices.
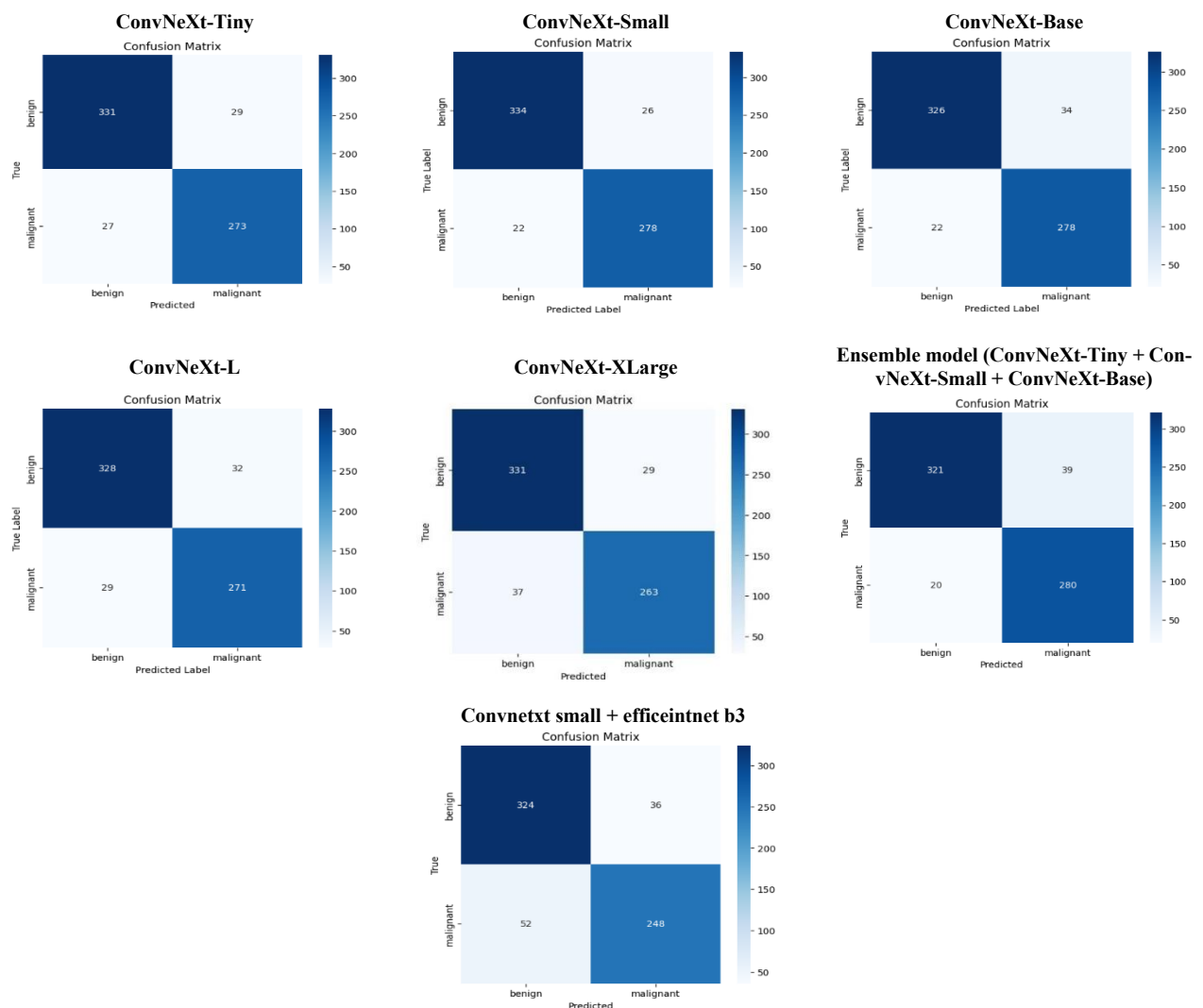


**Fig. 1:** Comparisons of All ConvNext Models' Confusion Matrix.

In resource-constrained or remote environments, lightweight versions such as ConvNeXt-Tiny or ConvNeXt-Small could be deployed on mobile or edge devices to support community health workers and teledermatology services. Furthermore, offline inference with local data caching and cloud synchronization when connectivity allows may provide diagnostic capability extension to reach regions where specialist access is not possible. Extensive clinical validation and regulatory approval are necessary before real-world adoption. Prospective multi-center studies with diverse populations and varied imaging devices are required to ensure generalizability with reduced bias. Finally, continuing collaboration among clinicians, AI developers, and policymakers will be crucial for ensuring trustworthy, explainable, and ethically compliant deployment of these models in real-world dermatology practice. The confusion matrices for each ConvNeXt model represent classification strengths and weaknesses across benign and malignant categories. Amongst them, ConvNeXt-Small and ConvNeXt-Base had the highest true positive rate for malignant lesions, indicating reliable sensitivity. ConvNeXt-Tiny also did not have highly imbalanced predictions, having very few false negatives. Larger variants, such as the ConvNeXt-Large and XLarge models, had a slightly higher false-positive rate for benign cases, which could indicate mild overfitting. The ensemble model (Tiny + Small + Base) had the most balanced confusion distribution, reducing both false negatives and false positives, which is clinically important in the context of cancer screening.
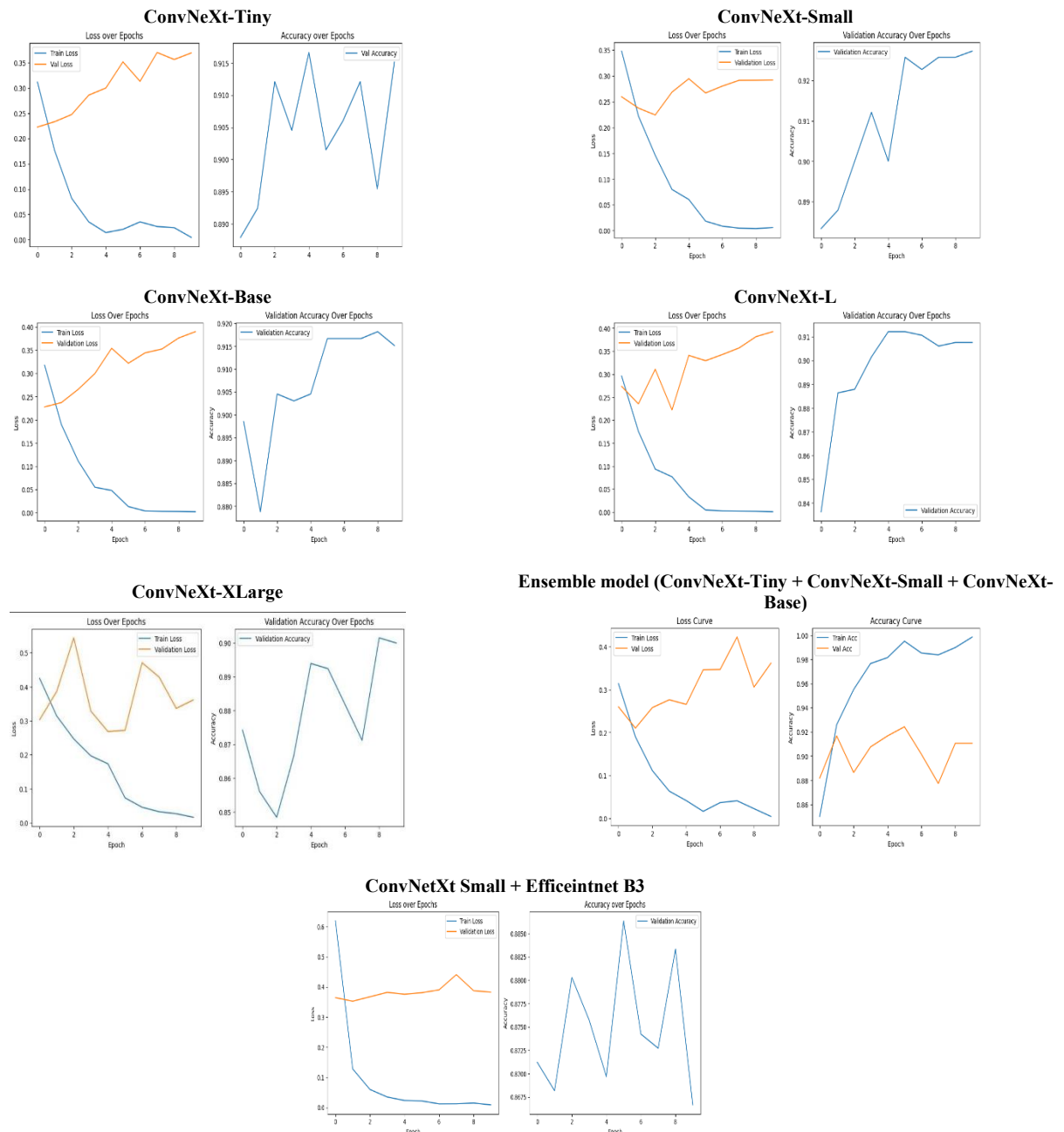


**Fig. 2:** Comparisons of All ConvNeXt Models Train, Validation, and Accuracy Graphs

The training and validation accuracy graphs show steady convergence for all ConvNeXt variants within 10 epochs. ConvNeXt-Small demonstrated the smoothest learning curve with minimal gap between training and validation accuracy, reflecting strong generalization. ConvNeXt-Large and XLarge achieved rapid training accuracy but diverged in validation, confirming overfitting tendencies. The ensemble model exhibited consistent validation improvement with minimal fluctuations, suggesting stability and robustness in performance. Figure 2 shows the results of this research reflect the high potential of the deep learning models with ConvNeXt in the process of automated skin cancer diagnosis. ConvNeXt-Small, among the tested architectures, exhibited a relatively agreeable accuracy, generalization quality, and computation economy, which is suitable for use in real-life applications in the clinical sphere. Ensemble models also increased robustness

and classification stability, especially on borderline cases, thus stating the value of model aggregation in augmenting reliability in making diagnoses.

# 6. Conclusion

The present work aims to provide a thorough in-depth analysis of deep learning frameworks based on ConvNeXt to perform a binary lesion classification of skin lesions into benign or malignant lesions with the help of dermoscopic images. The ConvNeXt-Small model was found to possess the best performance in standalone terms, in that it outperformed in terms of generalization and accuracy of classification, as well as preserving computational efficiency. Ensemble models, especially the one achieved with a combination of ConvNeXt-Tiny, ConvNeXt-Small, and ConvNeXt-Base also enhanced to have the best performance in diagnosis robustness and partly neutralized the bias associated with the individual models. These results indicate incorporating ConvNeXt models into clinical decision support software and especially in low dermatological expertise settings. Nevertheless, there is no way that these models will be implemented in real life with ease. Problems like the class imbalance, the lack of homogeneity in datasets, overfitting in high-capacity models, and the issue of model interpretability still represent the key obstacles. To fill these gaps, it is necessary to find a prolonged cooperation between AI researchers, clinicians, and regulatory stakeholders.

# References

[1] Brinker, T. J., Hekler, A., Enk, A. H., Berking, C., Haferkamp, S., Hauschild, A., ... & von Kalle, C. (2019). Deep neural networks are superior to dermatologists in melanoma image classification. European Journal of Cancer, 119, 11-17. https://doi.org/10.1016/j.ejca.2019.05.003

[2] Chen, Y., Xu, Z., & Wang, J. (2023). Lightweight convolutional neural networks for portable medical diagnosis: Recent trends and challenges. IEEE Access, 11, 55642–55659.

[3] Chen, L., Zhao, Y., & Lin, X. (2024). Improving robustness of AI diagnostic systems to adversarial noise through architecture-aware ensemble models. IEEE Access, 12, 54823–54834.

[4] Codella, N. C., Rotemberg, V., Tschandl, P., Celebi, M. E., Dusza, S. W., Gutman, D., ... & Halpern, A. C. (2018). 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). arXiv preprint arXiv:1710.05006. https://doi.org/10.1109/ISBI.2018.8363547.

[5] Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., & Thrun, S. (2017). Dermatologist-level classification of skin cancer with deep neural networks. Nature, 542(7639), 115–118. https://doi.org/10.1038/nature21056.

[6] Fan, H., Kirillov, A., Xie, S., & He, K. (2023). Scaling ConvNets for Foundation Models with ConvNeXt-UL. arXiv:2307.08015.

[7] Goyal, M., Oakley, A., Bansal, P., & Yap, M. H. (2020). Skin lesion segmentation in dermoscopic images with ensemble deep learning methods. IEEE Access, 8, 4171–4181. https://doi.org/10.1109/ACCESS.2019.2960504.

[8] Han, S. S., Kim, M. S., Lim, W., Park, G. H., Park, I., & Chang, S. E. (2020). Classification of the clinical images for benign and malignant cutaneous tumors using a deep learning algorithm. Journal of Investigative Dermatology, 138(7), 1529–1538. https://doi.org/10.1016/j.jid.2017.12.018.

[9] Jha, D., Smedsrud, P. H., Riegler, M. A., Johansen, D., Halvorsen, P., & Lange, T. (2023). Deep learning-based medical image analysis: a review of recent advances. Frontiers in Artificial Intelligence, 6, 1156304.

[10] Karim, M. R., Othman, N., & Islam, S. M. S. (2024). Advances in deep learning models for medical image analysis: A comprehensive review. IEEE Access, 12, 123456-123478.

[11] Kablan, E. B., & Ayas, S. (2024). Skin lesion classification from dermoscopy images using ensemble learning of ConvNeXt models. Signal, Image and Video Processing, 18, 6353–6361. https://doi.org/10.1007/s11760-024-03321-y.

[12] Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., ... & Sánchez, C. I. (2017). A survey on deep learning in medical image analysis. Medical Image Analysis, 42, 60-88. https://doi.org/10.1016/j.media.2017.07.005.

[13] Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). ConvNeXt: A ConvNet for the 2020s. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 11976–11986. https://openaccess.thecvf.com/content/CVPR2022/html/Liu_ConvNeXt_A_ConvNet_for_the_2020s_CVPR_2022_paper.html. https://doi.org/10.1109/CVPR52688.2022.01167.

[14] Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). ConvNeXt: A ConvNet for the 2020s. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 11976-11986. https://openaccess.thecvf.com/content/CVPR2022/html/Liu_ConvNeXt_A_ConvNet_for_the_2020s_CVPR_2022_paper.html. https://doi.org/10.1109/CVPR52688.2022.01167.

[15] Liu, Z., Mao, H., Wu, C. Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). A ConvNet for the 2020s. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 11976-11986. https://doi.org/10.1109/CVPR52688.2022.01167.

[16] Liu, J., Wu, X., Deng, H., & Li, W. (2023). Deep learning-based diagnostic system for skin cancer classification using dermoscopic images. Computers in Biology and Medicine, 161, 106012.

[17] Mahbod, A., Schaefer, G., Wang, C., Dorffner, G., & Ecker, R. (2021). Transfer learning using a multi-scale and multi-network ensemble for skin lesion classification. Computers in Biology and Medicine, 104, 103486. https://doi.org/10.1016/j.compbiomed.2018.11.017.

[18] Nagpal, K., Foote, D., Liu, Y., Chen, P. H. C., Wulczyn, E., Tan, F., ... & Corrado, G. S. (2019). Development and validation of a deep learning algorithm for improving Gleason scoring of prostate cancer. NPJ Digital Medicine, 2(1), 48. https://doi.org/10.1038/s41746-019-0112-2.

[19] Singh, S., Gupta, R., & Banerjee, A. (2023). A comprehensive survey on performance metrics for classification problems in medical AI. Journal of Biomedical Informatics, 137, 104297. https://doi.org/10.1016/j.jbi.2023.104297.

[20] Singh, S., Gupta, P., & Verma, A. (2023). Improving clinical diagnosis through ensemble learning: A review and application in cancer detection. Artificial Intelligence in Medicine, 133, 102470.

[21] Singh, R., Banerjee, A., & Sinha, A. (2023). Addressing data imbalance and noise in medical image classification using hybrid deep learning ensembles. Computerized Medical Imaging and Graphics, 104, 102228. https://doi.org/10.1016/j.compmedimag.2023.102228.

[22] Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. Proceedings of the 36th International Conference on Machine Learning (ICML), 6105–6114. http://proceedings.mlr.press/v97/tan19a.html.

[23] Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. International Conference on Machine Learning (ICML), 6105–6114. http://proceedings.mlr.press/v97/tan19a.html.

[24] Tschandl, P., Rinner, C., Apalla, Z., Argenziano, G., Codella, N., Halpern, A., ... & Kittler, H. (2019). Human–computer collaboration for skin cancer recognition. Nature Medicine, 26, 1229–1234. https://doi.org/10.1038/s41591-020-0942-0.

[25] Wang, H., Zhang, Y., Chen, M., & Yang, J. (2023). Ensemble deep learning for medical image classification: Recent advances and future trends. IEEE Transactions on Neural Networks and Learning Systems, 34(4), 1456-1472.

[26] Wang, H., Zhang, Y., Chen, M., & Yang, J. (2023). Ensemble deep learning for medical image classification: Recent advances and future trends. IEEE Transactions on Neural Networks and Learning Systems, 34(4), 1456–1472.

[27] World Health Organization. (2023). Skin cancers. https://www.who.int/news-room/fact-sheets/detail/ultraviolet-(uv)-radiation

[28] Woo, S., Li, D., Han, S., Kirillov, A., & Feichtenhofer, C. (2023). ConvNeXt V2: Co-designing and Scaling ConvNets with Masked Autoencoders. arXiv:2301.00808. https://doi.org/10.1109/CVPR52729.2023.01548.

[29] Yu, F., Chen, X., Wang, L., Xiong, Y., & Sun, J. (2023). ConvNeXt for Dense Prediction Tasks. CVPR Workshop on Efficient Deep Learning for Computer Vision.

[30] Zhang, W., Liu, Y., & Tan, R. (2023). Efficient deep learning models for edge-based medical imaging applications. Computers in Biology and Medicine, 158, 106944. https://doi.org/10.1016/j.compbiomed.2023.106944.

[31] Zhang, Y., Li, X., Wang, S., & Chen, Y. (2023). Enhanced feature extraction techniques in convolutional neural networks for medical image classification. Journal of Biomedical Informatics, 140, 104252.

[32] Zhao, J., Wang, Y., Chen, T., & Li, Z. (2022). Diversity-driven ensemble learning for robust medical image classification. Pattern Recognition Letters, 160, 53-61. https://doi.org/10.1016/j.patrec.2022.02.006.

[33] Zhao, J., Xu, K., Wang, F., & Yu, Z. (2023). Hybrid ensembles of deep learning models for robust medical image classification. Artificial Intelligence in Medicine, 139, 102484.

[34] Zenebe, A. Y. (2024). SkinAACN: An Efficient Skin Lesion Classification Based on Attention Augmented ConvNeXt with Hybrid Loss Function. Proceedings of the ACM Conference.