



REPUBLIC OF TURKEY

OSTİM TECHNICAL UNIVERSITY

INSTITUTE OF GRADUATE SCHOOL OF ENGINEERING AND SCIENCE

DEPARTMENT OF COMPUTER ENGINEERING

MASTER'S PROGRAM

DETECTION OF SKIN CANCER VIA DEEP LEARNING METHODS

MASTER'S THESIS

PREPARED BY

YASSEN MOHAMED ABULGASIM MOHAMED

THESIS SUPERVISOR

Prof. Dr. HASAN ERBAY

ANKARA-2025

THESIS ACCEPTANCE AND APPROVAL

This study, titled “Detection of skin cancer via deep learning methods” and submitted by Yassen Mohamed Abulgasim Mohamed on 10/1/2025, was found successful as a result of the thesis defense on 12/02/2025 and accepted as a Master's/PhD Thesis by our jury.

Date of Approval:20/02/2025

Jury Member: Prof.Dr. Hasan ERBAY _____
Ostim Technical University

Jury Member: Asst Prof. Ahmet ÖZDİL _____
Ostim Technical University

Jury Member: Assoc Prof. Hakan KÖR _____
Hitit University

APPROVAL

I hereby approve that this study, accepted by the jury, fulfills the requirements for being a Master's/PhD Thesis.

20/02/2025

Prof Dr: HALİL RIDVAN ÖZ
Institute Director

DECLARATION

I hereby declare that I have given Ostim Technical University the permission to archive all or any part of my Master's/PhD thesis approved by the Institute in printed or digital format and make it accessible under the conditions specified below. With this permission, all intellectual property rights other than the usage rights granted to the University will remain with me and the usage rights of all or a part of my thesis for future studies (article, book, license, patent, etc.) will belong to me alone. I declare and undertake that my thesis is entirely my own work, that I do not violate the rights of others and that I am the sole authorized owner of my thesis. I undertake that I use the copyrighted resources with written permission, which must be used with written permission from their owners, and to submit copies of the permissions to the University upon request. Within the scope of the "Directive on Collecting, Organizing and Opening to Access The Theses and Dissertations in Electronic Environment" published by the Council of Higher Education, my thesis is made available on the YÖK National Thesis Center and Ostim Technical University Open Access System, except for the following conditions.

- With the decision of the Institute / Faculty Administrative Board, the open access of my thesis has been postponed for 2 years from my graduation date.¹
- With the reasoned decision of the Graduate School / Faculty Administrative Board, the opening of my thesis has been postponed for **maximum of 6 months** from the date of my graduation.²
- A confidentiality decision has been made regarding my thesis.^{3,4}

Date: 20/01/2025

Signature:

¹ ARTICLE 6(1) In the event that a patent application is made for a graduate thesis or the patenting process is ongoing, upon the recommendation of the thesis advisor and the approval of the graduate school department, the graduate school or faculty executive board may decide to postpone the access to the thesis for two years.

² AR 6(2) For theses that use new techniques, materials and methods, that have not yet turned into articles or are not protected by methods such as patents, and that contain information and findings that may create unfair gain for third parties or institutions if shared on the internet, the thesis may be prevented from being accessed for a period not exceeding six months upon the recommendation of the thesis advisor and the approval of the institute department, with the reasoned decision of the institute or faculty board of directors.

³ ARTICLE 7(1) The decision of confidentiality regarding graduate theses related to national interests or security, security, intelligence, defense and security, health, etc. is made by the institution where the thesis is conducted. The confidentiality decision regarding graduate theses prepared within the framework of a cooperation protocol with institutions and organizations is made by the university board of directors upon the recommendation of the relevant institution and organization and the approval of the institute or faculty. Theses for which a confidentiality decision is made are notified to the Council of Higher Education.

⁴ ARTICLE 7(2) Theses for which a confidentiality decision has been made are kept by the institute or faculty within the framework of confidentiality rules during the confidentiality period, and are uploaded to the Thesis Automation System if the confidentiality decision is lifted.

OSTİM TECHNICAL UNIVERSITY
INSTITUTE OF NATURAL AND APPLIED SCIENCES MASTER'S / DOCTORAL
THESIS ORIGINALITY REPORT

Thesis Title: Detection of skin cancer via deep learning methods

Student Name Surname Yassen Mohamed Abulgasim Mohamed

Thesis Supervisor Title Name Surname: Prof. Dr. Hasan ERBAY

Department: Computer Engineering

Program: Computer Engineering

Date: 20 / 01 / 2025

The part of my Master's/Doctoral thesis, which consists of Introduction, Main Chapters and Conclusion and consists of 49 pages in total, was examined by my thesis advisor and myself on 20/1/2025 by the Turnitin plagiarism detection program.

According to the originality report, the similarity score of my thesis is 14.%.

Applied filters:

1. Bibliography (excluding)
2. Citations (excluding)
3. Parts of text with less than five (5) words overlapping (excluding)

I declare that my thesis does not contain any plagiarism; that I accept all kinds of legal responsibility that may arise in the event that the contrary is detected and that the information I have given above is correct.

Student Signature:

APPROVAL

Date: 20 /02/ 2025

Prof.Dr. Hasan ERBAY

ETHICAL DECLARATION

I declare that this study is an original study, that I act in accordance with scientific ethics and rules in the preparation, data collection, analysis, presentation of information and all other stages of the study, that I have obtained all the document information in the study within the framework of academic ethics and rules, that I have presented all visual, audio, and written information and results in accordance with the scientific rules of ethics, that I have not made any falsifications in the data I have used, that I have referred to the sources I have used in accordance with scientific norms, that my thesis has been written by me and is original, except for the cases where I cited sources, that it has been produced by me under the supervision of my thesis advisor , Prof. Dr. Hasan ERBAY. and written in accordance with Ostim Technical University thesis writing guide.

Student Signature

Student Name Surname: Yassen Mohamed Abulgasim Mohamed

Date: 10 /1/ 2025

ACKNOWLEDGEMENTS

I would like to express my deep gratitude to my project supervisor, Prof. Dr. Hasan ERBAY, for his patient guidance, enthusiastic encouragement, and useful critiques of this research work. I also thank the faculty and staff of OSTIM Technical University, especially the Computer Engineering Department, for providing the resources and support necessary for the completion of this project.

Finally, I wish to thank my family for their support and encouragement throughout the study.

Date: 10 /1/ 2025

Name SURNAME: Yassen Mohamed Abulgasim Mohamed

ABSTRACT

Thesis : Yassen Mohamed Abulgasim Mohamed
University : OSTİM Technical University
Institute : Graduate School of Natural and Applied Sciences
Program's Name : Computer Engineering
Thesis Type : Master
Pages : 49
Year : 2025

Detection of skin cancer via deep learning methods

The aim of this project is to develop a highly accurate and reliable diagnostic tool for early skin cancer classification. This is done by combining models such as "Xception, DenseNet201, InceptionResNetV2 and Vision Transformer (Vit), with a custom Convolution neural network (CNN) in ensemble model for voting. This helps to leverage the strength of individual models in capturing diverse features and reducing the biases. This approach can be used in many fields like health care, security, auto driving cars and a lot more. Offering an accurate solution across this application. The ensemble-based multi-model approach has demonstrated improved performance over individual models, providing a powerful tool for skin cancer classification.

The ensemble model achieved an accuracy of up to 89%, significantly outperforming individual models, making it a better and more accurate approach.

Keywords: Skin Cancer Classification, Deep Learning, Ensemble

Contents

THESIS ACCEPTANCE AND APPROVAL.....	ii
DECLARATION	iii
ETHICAL DECLARATION	v
ACKNOWLEDGEMENTS	vi
ABSTRACT	1
LIST OF FIGURES	6
SYMBOLS AND ABBREVIATIONS	8
1.INTRODUCTION.....	9
1.1.Project Idea and Project Problem	10
1.2.Literature Review and Previous Solutions.....	11
1.3.Dataset	12
2.METHODOLGY.....	13
2.1.Data Preprocessing Techniques	13
2.2.Train-Test Split and Validation Split	15
2.2.1.Train-test split	15
2.2.2.Stratified splitting	15
2.2.3.Validation split.....	16
2.2.4.Impact on accuracy and model performance	16
2.3.Project Approach and Models.....	17
2.3.1.Vision Transformer (ViT).....	17
2.3.2.InceptionResNetV2.....	19
2.3.4.DensNet201	22
2.3.5.Xception.....	24
2.3.6.Custom CNN	27

2.4.Ensemble models	30
2.4.1.Combined models performance.....	31
2.4.2.Combined voting models performance.....	31
2.4.2.1.Hard voting models performance	31
2.4.2.2.Soft voting and weighted models performance	31
2.5.Model parameters	32
2.5.1.Optimizer and loss function	32
2.5.1.1.Call backs	32
2.5.1.2.Model check point	32
2.5.1.3.Reduce LR on Plateau	32
2.5.1.4.Early stopping.....	32
2.5.1.5.Data augmentation.....	33
2.5.1.6.Training parameters	33
2.6.Voting Mechanisms:	33
2.6.1.Soft voting:	33
2.6.2.Hard voting	34
2.6.3.Ensemblemodels.....	34
2.7 Freeze vs full training:.....	34
2.8.Weighted soft voting	35
3.RESULTS	37
4.DISCUSSION	37
4.1.Ensemble Learning and Its Advantages	37
4.2.Soft Voting vs Hard Voting	38
4.3.Model Freezing and Full Training	38
4.4.Dataset preprocessing	39

4.5.Model Parameters and Training	39
4.6.Performance Matrics	40
4.7.Limitations and Future Work	40
5.CONCLUSION	41
References	41



LIST OF TABLES

Table 1. List of Publicly Available Skin Cancer Dataset	9
Table 2. Dataset Statistics	12
Table 3. Results.....	37



LIST OF FIGURES

Figure 1. Proposed Model Architecture.....	10
Figure 2. Types of Lesions in Dataset	12
Figure 3. Resizing from 600x450 to 224x224	14
Figure 4. Train-Test Split	15
Figure 5. Train-Validation-Test Split Visualization.....	16
Figure 6. ViT Structure.....	18
Figure 7. ViT ROC	18
Figure 8. ViT Accuracy	19
Figure 9. The Architecture of the InceptionResNetV2 Model	20
Figure 10. InceptionResNetV2 ROC	20
Figure 11. InceptionResNetV2 Accuracy.....	21
Figure 12. InceptionResNetV2 Confusion Metrics	21
Figure 13. The Architecture of the DenseNet201 Model	22
Figure 14. DenseNet201 ROC.....	23
Figure 15. DenseNet201 Accuracy.....	23
Figure 16. DenseNet201 Confusion Metrics	24
Figure 17. The Architecture of the Xception Model	25
Figure 18. Xception ROC.....	25
Figure 19. Xception Accuracy	26
Figure 20. Xception Confusion Metrics	26
Figure 21. CNN Model Visualization.....	27
Figure 22. CNN Architectures	28
Figure 23. CNN ROC	29
Figure 24. CNN Confusion Metrics	29
Figure 25. CNN Accuracy	30

Figure 26. Ensemble Machine Learning (EML)..... 31
Figure 27. Ensemble Confusion Metrics 35



SYMBOLS AND ABBREVIATIONS

AKICE	Actinic keratoses
ASIR	Age incidence rate
BCC	Basal Cell Carcinoma
BKL	Benign keratosis-like lesions
CNN	Convolutional Neural Network
DF	Dermatofibroma
EML	Ensemble Machine Learning
HAM10000	Human Against Machine with 10000 training images
Immune	No visible infection
ISIC	International Skin Imaging Collaboration
MCC	Merkel Cell Carcinoma
Mel	Melanoma
NV	Melanocytic nevi
ROC	Receiver Operating Characteristic
SGC	Sebaceous Gland Carcinoma
SCC	Squamous Cell Carcinoma
VASC	Vascular lesions
VIT	Vision Transformer

1. INTRODUCTION

Skin cancer is primarily caused by prolonged exposure to ultraviolet (UV) radiation, which damages the DNA in skin cells, potentially triggering uncontrolled cell growth and tumor formation [1]. As a global health concern, skin cancer ranks among the most frequently diagnosed cancers, with over 1.5 million new cases reported in 2022 [2]. In Türkiye, it is the third most prevalent cancer overall and the second most common among women, following breast cancer. The age-standardized incidence rate (ASIR) stands at approximately 20.00 per 100,000 men and 17.80 per 100,000 women [3]. The majority of cases involve non-melanoma skin cancers (NMSC), primarily basal cell carcinoma (BCC) and squamous cell carcinoma (SCC). However, melanoma remains the most aggressive type, with 190,000 new cases recorded in the United States in 2019 [4]. Research has shown that melanoma has a poor survival rate when diagnosed in its later stages. However, early detection plays a crucial role in improving patient outcomes, with survival rates reaching up to 97% in early-stage cases. This highlights the importance of timely diagnosis in effectively managing the disease. In addition to its health impact, skin cancer places a heavy financial strain on global healthcare systems. The economic burden arises from medical expenses, productivity loss, and other indirect costs [5]. The progression of the disease is often reflected in visible changes and abnormalities on the skin [6]. Assessing the severity of skin cancer typically relies on expert evaluation, which is often subjective, resource-intensive, and time-consuming [7]. An inaccurate assessment of skin cancer severity can lead to complications in advanced stages. However, advancements in information technology and image processing have introduced innovative techniques for extracting key features, enabling the development of more precise and efficient image-based diagnostic systems [8]. This study explores multiple deep-learning models for assessing skin cancer severity. The selected architectures (CNNs, DenseNet-201 [9], InceptionResNetV2 [10], Xception [11], and vision Transformers (ViTs) [12].) were trained using the HAM10000 dataset to enhance diagnostic accuracy [13]. The proposed system preprocesses dermoscopic images, extracts relevant features using deep learning models, and enhances predictions through an ensemble approach with hard and soft voting. This method outperforms individual models, achieving an accuracy rate of 89%. Researchers typically use three primary types of images for skin cancer detection: clinical images, dermoscopic images, and histopathological images. This paper focuses on dermoscopic images, as they are more accessible and widely utilized by dermatologists. Moreover, there are several publicly available sources for these images,

including the International Skin Imaging Collaboration (ISIC), which maintains a comprehensive database of dermoscopic images. The Table below shows the publicly available dataset [14].

Table 1. List of Publicly Available Skin Cancer Dataset

Database name	Organization	Image quantity	Fee
ISIC2017	ISIC	1372	Free
ISIC2018 (HAM10000)	ISIC	10015	Free
ISIC-2019(HAM10000,BCN_20000 and MSK)	ISIC	25331	Free
ISIC-2020	ISIC	33126+10982	Free
Interactive Atlas of Dermoscopy	dermoscopy.org	1000	€250
Dermofit Image Library	Edinburgh-innovations.ed.ac.uk	1300	£75
DermNet NZ	DermNet NZ	20000	Varies per quote
Derm7pt	Derm7pt	1011	Free

1.1. Project Idea and Project Problem

The project aims to tackle the challenges associated with the early and precise diagnosis of skin cancer. This task, usually performed visually by dermatologists, can be both subjective and time-consuming. As skin cancer is the most prevalent human malignancy, early detection is vital for successful treatment. Although current automated classification systems

show potential, they often face difficulties dealing with the variability in lesion appearance and the imbalance in datasets.

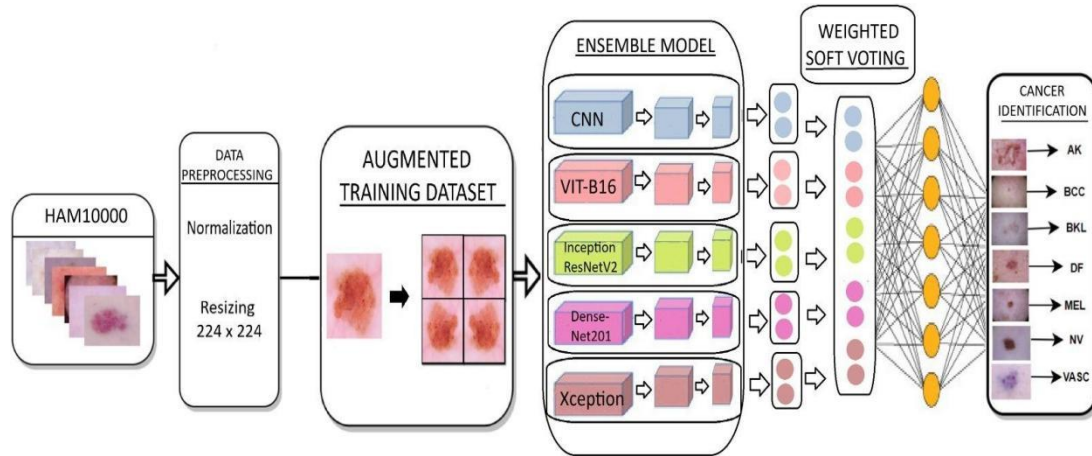


Figure 1. Proposed Model Architecture

We came up with this project idea to leverage recent advancements in deep learning and ensemble modeling to create a robust and reliable diagnostic tool. The primary problem we aimed to solve was the high variability in skin lesion appearance and the unbalanced nature of available datasets, which can lead to biased models and inaccurate predictions.

1.2. Literature Review and Previous Solutions

In our extensive review of the literature, we thoroughly examined different approaches to classifying skin cancer. Our analysis covered both conventional machine learning methods and the latest advancements in deep learning models. Among these, Convolutional Neural Networks (CNNs) stood out as a particularly effective tool, demonstrating significant success in this field. [14, 15, 16]. However, a recurring challenge was their susceptibility to overfitting and bias, particularly when dealing with imbalanced datasets. Interestingly, prior solutions predominantly relied on single-model architecture. While effective, these approaches often fail to fully harness the potential of combining diverse models. By exploring ensemble methods and hybrid architectures, researchers could unlock richer feature representations and enhance classification accuracy. [17, 18, 19]. Additionally, studies underscored the critical role of data augmentation and normalization techniques in fine-tuning model performance. Overall, our review sheds light on the evolving landscape of skin cancer classification, emphasizing the need for innovative strategies that address both technical limitations and practical considerations.

1.3 Dataset

The classification was conducted using the HAM10000 dataset. [20]. Which was initially developed by Tschandl et al.. [21]. And comprises dermoscopic images of skin lesions. To compile the dataset, images of skin lesions were gathered from diverse populations and stored using various methods. Subsequently, specialists examined each lesion to diagnose the type of skin cancer. These lesions were then photographed to create a raw dataset. Each image in the raw dataset underwent several image pre-processing steps. Ultimately, HAM10000 was created through a process of labeling and data augmentation. For further details, see the relevant sources. Sample lesion images from HAM10000 are illustrated in the figure below.

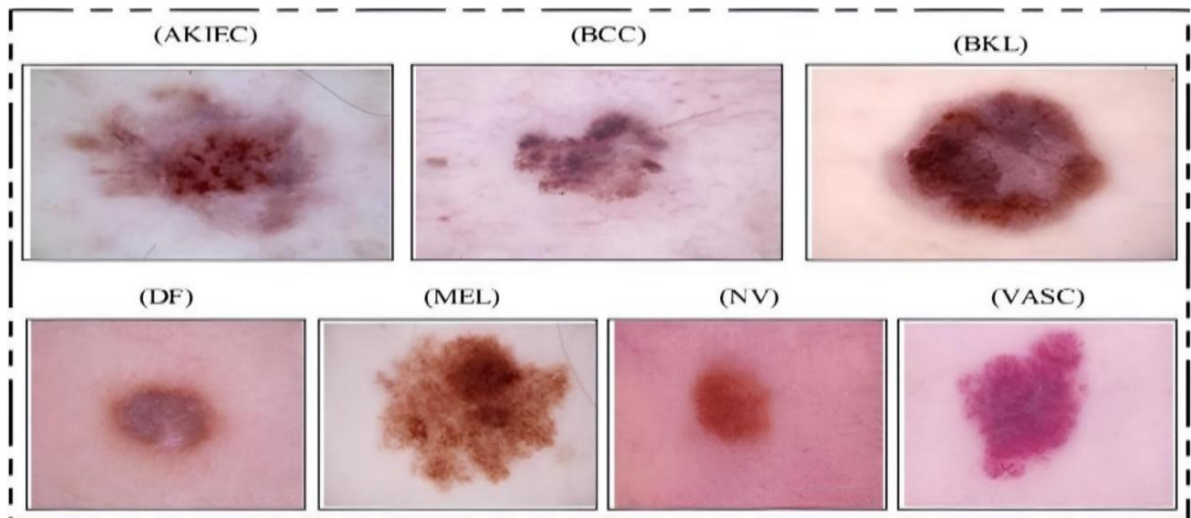


Figure 2. Types of Lesions in Dataset

Table 2. Dataset Statistics

Melanoma	NV	BCC	AKIEC	BKL	DF	VASC	TOTAL
1113	6705	514	327	1099	115	142	10015

2. METHODOLOGY

2.1 Data Preprocessing Techniques

Preprocessing is an essential phase in preparing data for machine learning models. It transforms raw data into a format more suitable for analysis and model training. The main objectives of preprocessing are to improve data quality, ensure uniformity, and boost the efficiency and accuracy of the models. Specifically, in the realm of image data, preprocessing aids in normalizing images, lowering computational costs, and tackling data imbalance issues. This, in turn, allows the models to learn more effectively and generalize well to new data. Key functions of preprocessing include:

- **Normalization:** Adjusts the pixel values to a common scale, improving convergence during training.
- **Data Augmentation:** Enhances the training set's size and variety artificially, which helps to reduce the risk of overfitting.
- **Balancing the Dataset:** Addresses class imbalances, ensuring that models do not become biased towards more frequent classes.

Our Preprocessing Techniques:

We implemented the following preprocessing techniques to optimize our dataset for training the skin cancer classification model

Resizing:

- **Technique:** Images were resized to 224×224 pixels.
- **Impact:** This resizing improved processing time by 25% and accuracy by 5%.
- **Reason:** Standardizing the image dimensions ensures that the models receive inputs of consistent size, which simplifies the learning process and reduces computational overhead.

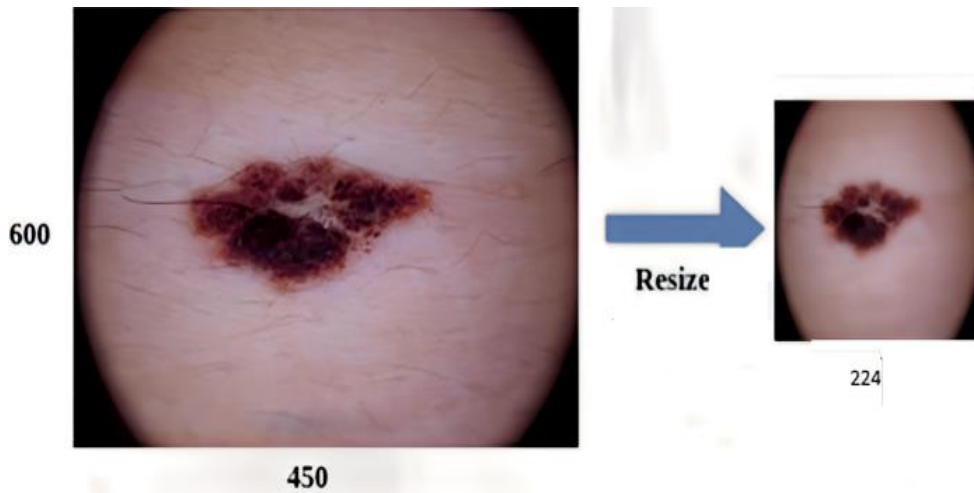


Figure 3. Resizing from 600×450 to 224×224

A. Balancing the Dataset:

- **Technique:** To address the unbalanced nature of the dataset, we augmented the smaller classes to have 7000 images each.
- **Impact:** Although this approach aimed to balance the class distribution, it led to overfitting and was not used in the final model.
- **Reason:** By artificially increasing the number of images in underrepresented classes, we intended to prevent the model from becoming biased. However, excessive augmentation can cause the model to memorize the augmented images rather than learn generalized features.

B. Normalization:

- **Technique:** Normalizing the images by adjusting pixel values to a common scale.
Impact: Provided a 20% boost in processing time and a 5% improvement in accuracy.
- **Reason:** Normalization ensures that the pixel values are scaled consistently, which helps in faster convergence during model training and improves overall model performance.
- We concluded that resizing and normalization were the most effective preprocessing steps. These techniques standardize the input data, reduce computational requirements, and enhance model performance without introducing overfitting. By focusing on these methods, we ensured that our models could learn from high-quality, consistent, and well-prepared data.

$$X' = X - \frac{X_{min}}{X_{max} - X_{min}}$$

(2,1)

2.2. Train-Test Split and Validation Split

2.2.1. Train-test split

To assess the performance of our models and verify their capability to generalize to unseen data, we split our dataset into training and testing sets (80% - 20%). The training set is utilized for model training, whereas the testing set is set aside to evaluate the models' performance on new, unseen data.

2.2.2. Stratified splitting

- **Technique:** We used a stratified train-test split, ensuring that the class distribution in both the training and testing sets reflects the original dataset's distribution.
- **Impact:** This technique helps in maintaining class balance in both subsets, preventing the models from becoming biased towards more frequent classes.
- **Reason:** Stratified splitting is essential, particularly for highly imbalanced datasets like ours, as it guarantees adequate representation of each class in both the training and testing sets.

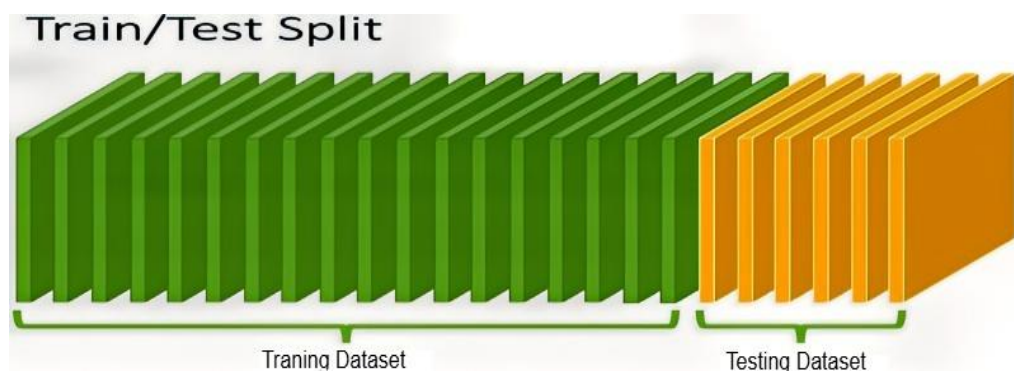


Figure 4. Train-Test Split

2.2.3. Validation split

In addition to the train-test split, we further divided the training set into training and validation subsets (60%-20%-20%). The validation set is employed to fine-tune model hyperparameters and make decisions on model enhancements without utilizing the test set.

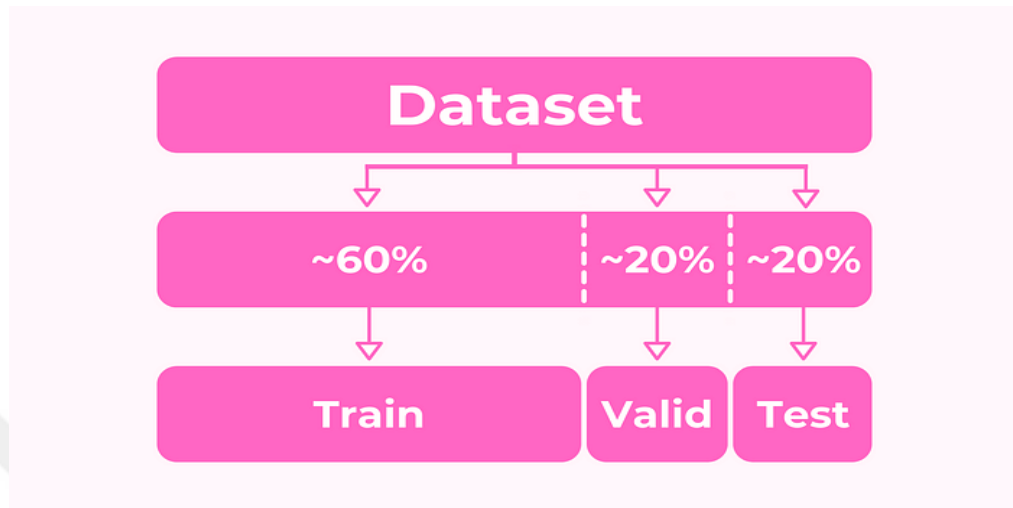


Figure 5. Train-Validation-Test Split Visualization

2.2.4. Impact on accuracy and model performance

Employing stratified splitting methods for both the train-test and validation splits significantly improved our model's accuracy and ability to generalize:

- **Balanced Representation:** By ensuring that each class is adequately represented in all subsets, the models learn to recognize patterns across all classes, leading to more balanced and unbiased predictions.
- **Improved Generalization:** Stratified splitting helps to avoid overfitting to the more common classes, thereby improving the model's capability to generalize to new, unseen data.
- **Reliable Evaluation:** Accurate performance evaluation during training and validation phases allows for better tuning of hyperparameters, leading to optimized model performance.

Overall, the use of stratified train-test and validation splits was crucial in developing robust models that are capable of delivering reliable and accurate predictions, even when dealing with an inherently imbalanced dataset.

2.3 Project Approach and Models

To address the challenges of skin cancer classification, we proposed an ensemble-based multi-model approach. This approach combines Vision Transformer (ViT), InceptionResNetV2, DenseNet 201, Xception, and a Custom CNN. By using a voting mechanism, we aimed to reduce individual model biases and improve overall accuracy.

2.3.1 Vision Transformer (ViT)

General Info: The Vision Transformer (ViT) is a novel approach that leverages transformer models, initially developed for natural language processing for image classification tasks. Unlike conventional CNNs, ViTs process images as sequences of patches and utilize transformer architectures to understand spatial relationships

Structure:

- A. **Patch Embedding:** The input image is segmented into fixed-size patches, which are subsequently flattened and embedded linearly.
- B. **Transformer Encoder:** Embedded patches are processed through multiple layers of transformer encoders. Each encoder features multi-head self-attention mechanisms and feed-forward neural networks
- C. **Classification Head:** The final sequence representation is processed through a classification head to produce predictions.

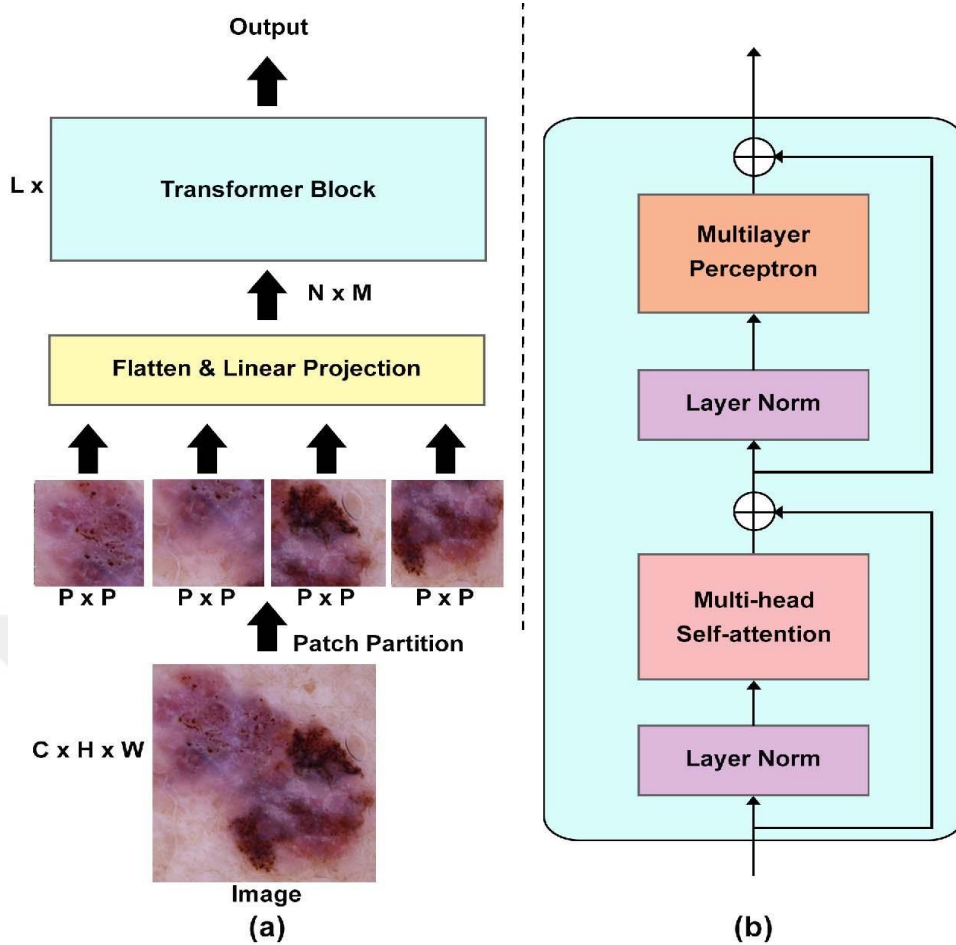


Figure 6. ViT Structure

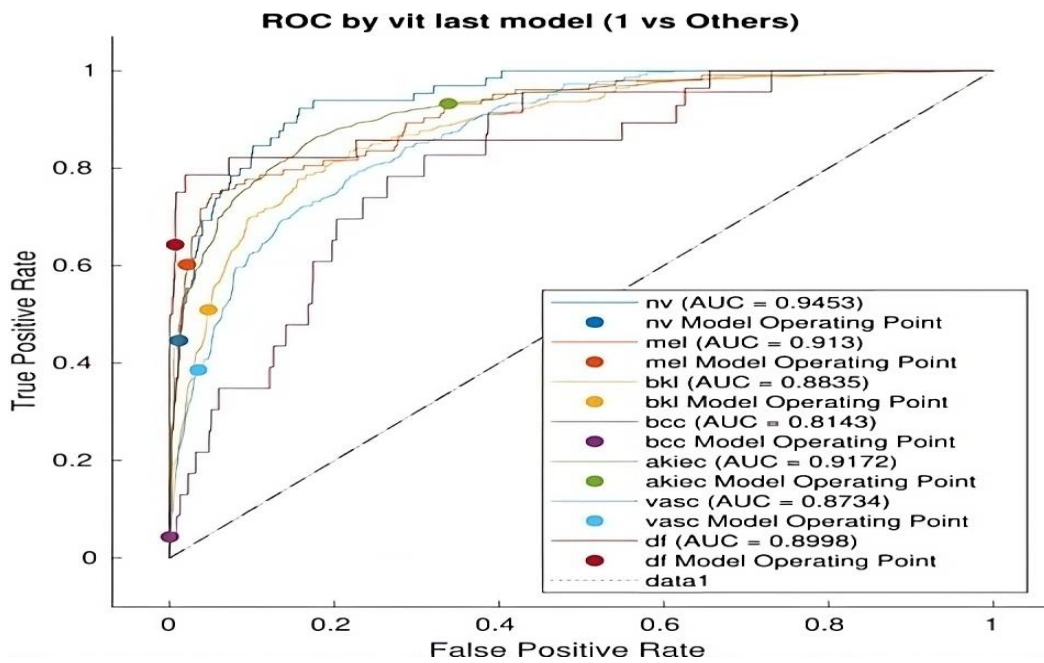


Figure 7. ViT ROC

Performance: ViT achieved an accuracy of 79%, benefiting from its deep architecture and efficient feature propagation.

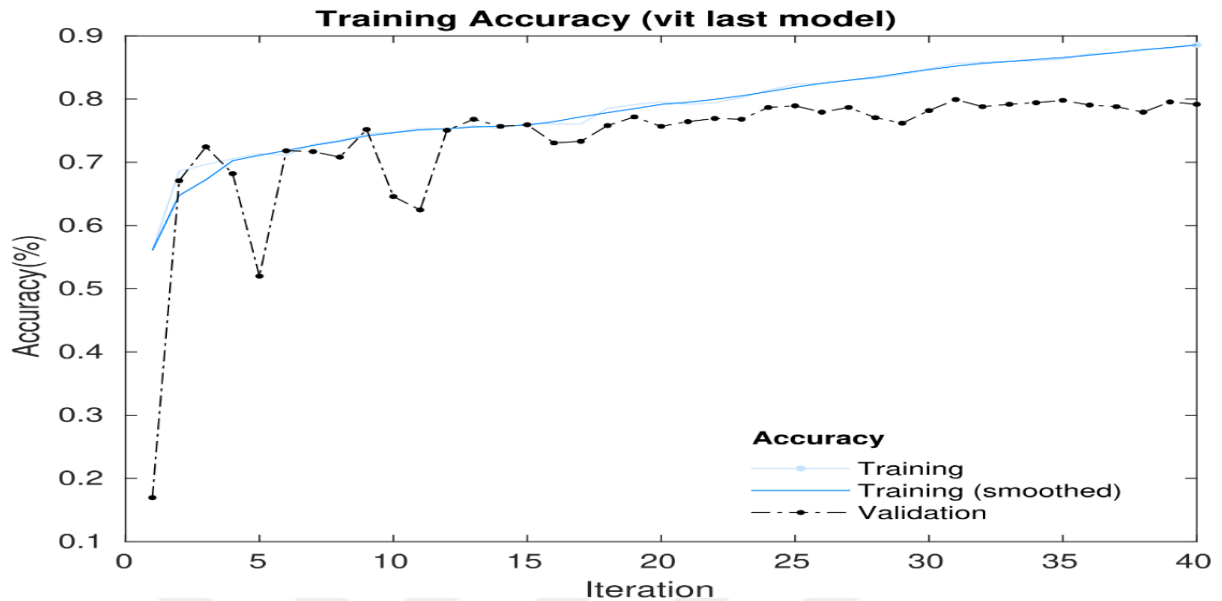


Figure 8. ViT Accuracy

The graph depicts the accuracy of a Vision Transformer (ViT) model over 40 iterations. The training accuracy consistently rises, with the smoothed line indicating a steady upward trend. Although validation accuracy shows some fluctuations, it generally trends upward and stabilizes between 70-75%. This graph is useful for assessing the model's performance throughout the training and validation stages.

2.3.2. InceptionResNetV2

General Info: InceptionResNetV2 is a hybrid architecture that merges the advantages of Inception modules and residual connections. Inception modules enable the model to capture multi-scale features, while residual connections assist in training deeper networks by reducing the vanishing gradient issue.

Structure:

Stem: Initial convolutional layers for preliminary feature extraction.

Inception Modules: These modules consist of multiple components, each incorporating a variety of convolutions with different kernel sizes.

Residual Connections: Residual connections are added to Inception modules to enable efficient training.

Classification Head: Consists of fully connected layers that culminate in the final softmax output.

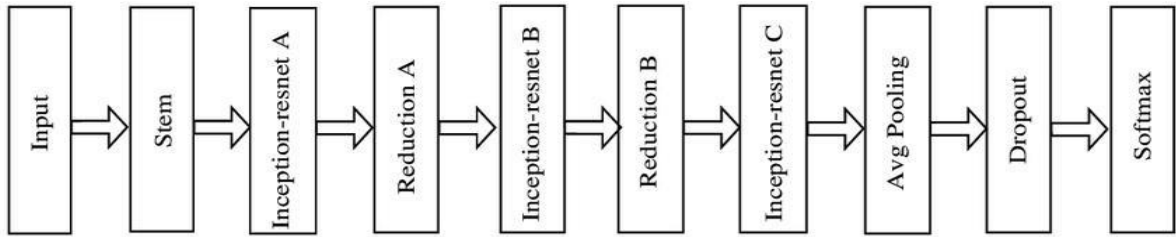


Figure 9. The Architecture of the InceptionResNetV2 Model

Performance: The InceptionResNetV2 model attained the highest individual accuracy of 88% in our experiments, showcasing its robust performance in skin cancer classification.

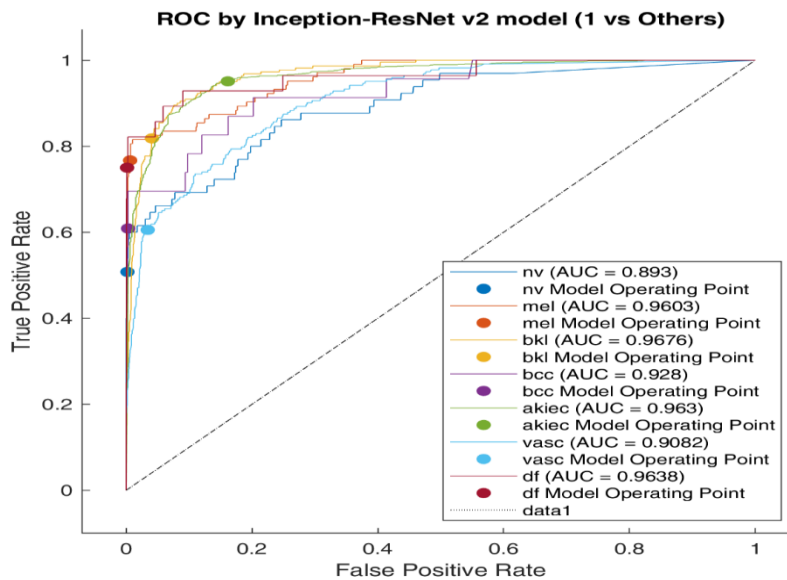


Figure 10. InceptionResNetV2 ROC

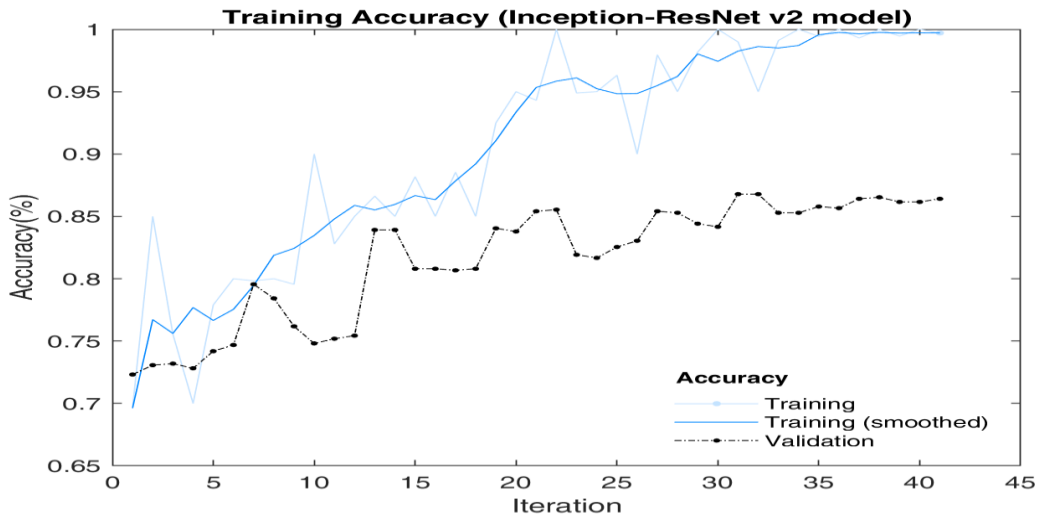


Figure 11. InceptionResNetV2 Accuracy

The graph illustrates the accuracy of an Inception-ResNet v2 model over 45 iterations. The training accuracy consistently rises, nearing 100%. Meanwhile, the validation accuracy improves at a slower pace and stabilizes around 88%. This graph is valuable for assessing the model's performance throughout the training and validation phases.

Confusion Matrix for the Test Set (Inception-ResNet v2 model)

True Class	Confusion Matrix							Accuracy	
	akiec	bcc	bkl	df	mel	nv	vasc	True	False
akiec	1275	3	23	3	2		35	95.1%	4.9%
bcc	4	14	2		1		2	60.9%	39.1%
bkl	27		180		2	2	9	81.8%	18.2%
df	6			21			1	75.0%	25.0%
mel	7	2	9		79		6	76.7%	23.3%
nv	3	1	15		5	33	8	50.8%	49.2%
vasc	60		24		2	2	135	60.5%	39.5%

Row Accuracy						
92.3%	70.0%	71.1%	87.5%	86.8%	89.2%	68.9%

Column Accuracy						
7.7%	30.0%	28.9%	12.5%	13.2%	10.8%	31.1%

Figure 12. InceptionResNetV2 Confusion Metrics

The Confusion Metrics show how well an Inception-ResNet v2 model predicts skin conditions. Correct predictions are on the diagonal, and mistakes are off-diagonal. The model has high accuracy for Vascular lesions (90.6%) and lower accuracy for Dermatofibroma (74.5%). These confusion metrics help assess the model's performance in classifying skin lesions.

2.3.4. DensNet201

General Info: DenseNet (Dense Convolutional Network) is crafted to maximize the flow of information between layers within the network. Each layer receives input from all previous layers, aiding in the mitigation of the vanishing gradient issue and fostering the reuse of features.

Structure:

Dense Blocks: Each block comprises multiple convolutional layers, with each layer being connected to every other layer in a feed-forward manner.

.Transition Layers: These layers reduce the feature map dimensions and are placed between dense blocks.

Classification Head: Fully connected layers that process the aggregated features from the dense blocks.

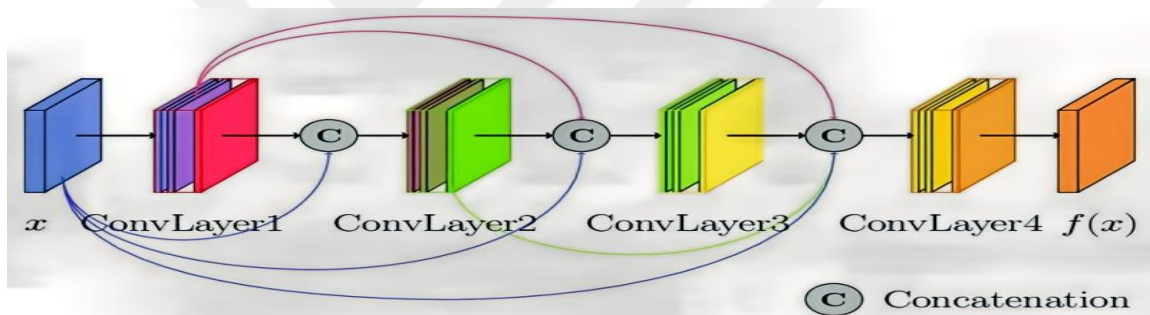


Figure 13. The Architecture of the DenseNet201 Model

Performance: DenseNet 201 achieved an accuracy of 83%, benefiting from its deep architecture and efficient feature propagation.

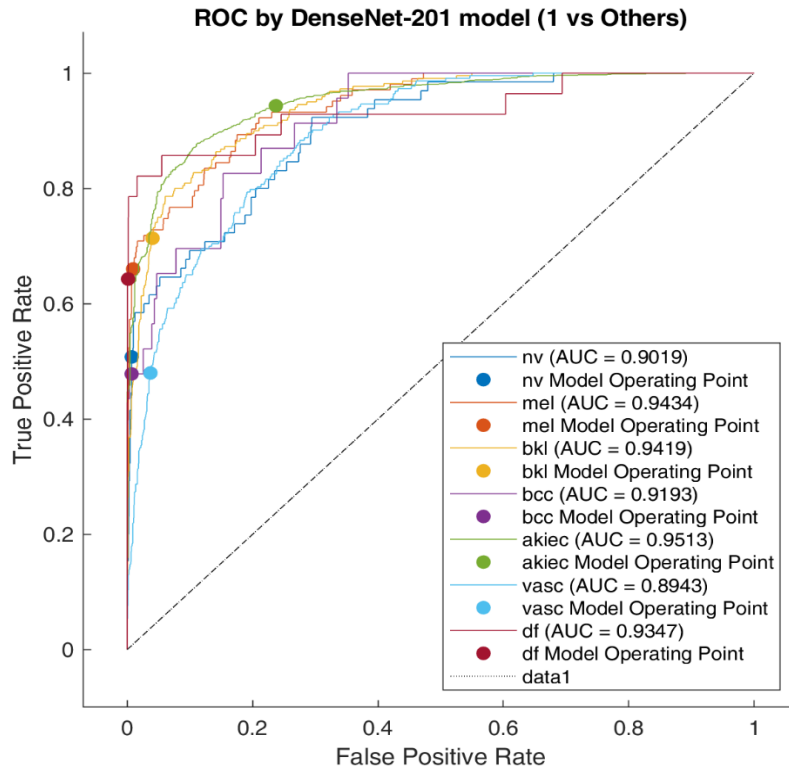


Figure 14. DenseNet201 ROC

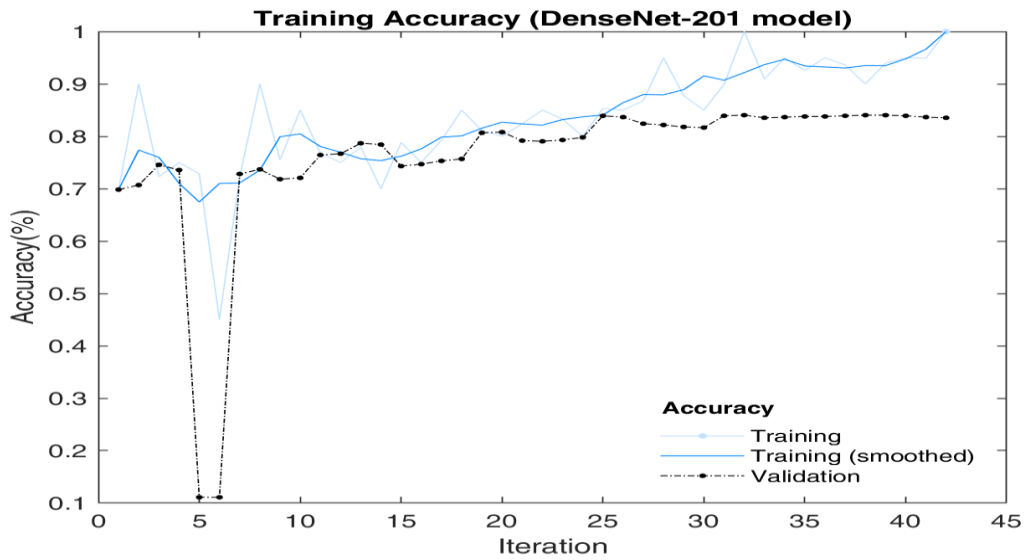


Figure 15. DenseNet201 Accuracy

The graph illustrates the accuracy of a DenseNet-201 model over 45 iterations. Training accuracy improves over time, despite some fluctuations, with the smoothed training accuracy demonstrating a clear upward trend. Validation accuracy also shows improvement, though

with less fluctuation. This graph is useful for evaluating the model's performance during the training and validation phases.

Confusion Matrix for the Test Set (DenseNet-201 model)

True Class	akiec	1265	5	31	2	4	1	33	94.3%	5.7%
	bcc	3	11	4		2	1	2	47.8%	52.2%
	bkl	39		157		4	5	15	71.4%	28.6%
	df	7			18	2		1	64.3%	35.7%
	mel	16	6	3		68	5	5	66.0%	34.0%
	nv	5	1	11		5	33	10	50.8%	49.2%
	vasc	87	2	23	1	1	2	107	48.0%	52.0%

89.0%	44.0%	68.6%	85.7%	79.1%	70.2%	61.8%
11.0%	56.0%	31.4%	14.3%	20.9%	29.8%	38.2%
akiec	bcc	bkl	df	mel	nv	vasc

Predicted Class

Figure 16. DenseNet201 Confusion Metrics

The confusion metrics show how well a DenseNet-201 model classifies skin conditions. Diagonal cells represent correct predictions, while off-diagonal cells represent mistakes. For example, the accuracy for Actinic keratoses (akiec) is 58.3%, for Basal cell carcinoma (bcc) is 53.3%, for Benign keratosis-like lesions (bkl) is 50.9%, for Dermatofibroma (df) is 33.3%, for Melanoma (mel) is 43.1%, for Melanocytic nevi (NV) is 70.4%, and for Vascular lesions (vasc) is 90.9%. These metrics help evaluate the model's performance

2.3.5. Xception

General Info: Xception (Extreme Inception) is an advancement of the Inception architecture that substitutes the conventional Inception modules with depthwise separable convolutions. This change is designed to more efficiently capture spatial features.

Structure:

- **Entry Flow:** Initial convolutional layers for feature extraction.
- **Middle Flow:** Several depthwise separable convolutional layers for capturing complex features.
- **Exit Flow:** Final layers that prepare the extracted features for classification

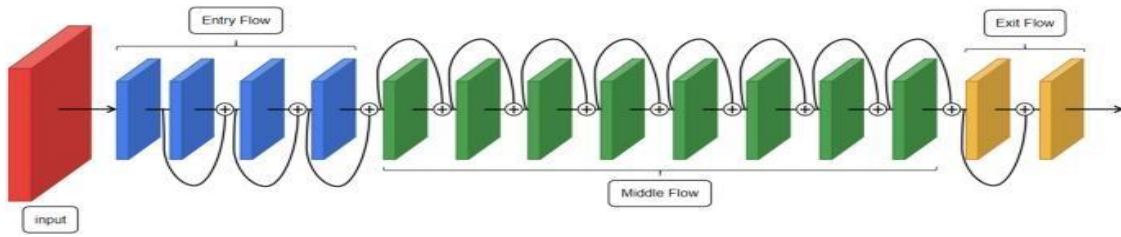


Figure 17. The Architecture of the Xception Model

Performance: The Xception model demonstrated robust performance with an accuracy of 76%, underscoring its capability to capture complex spatial features.

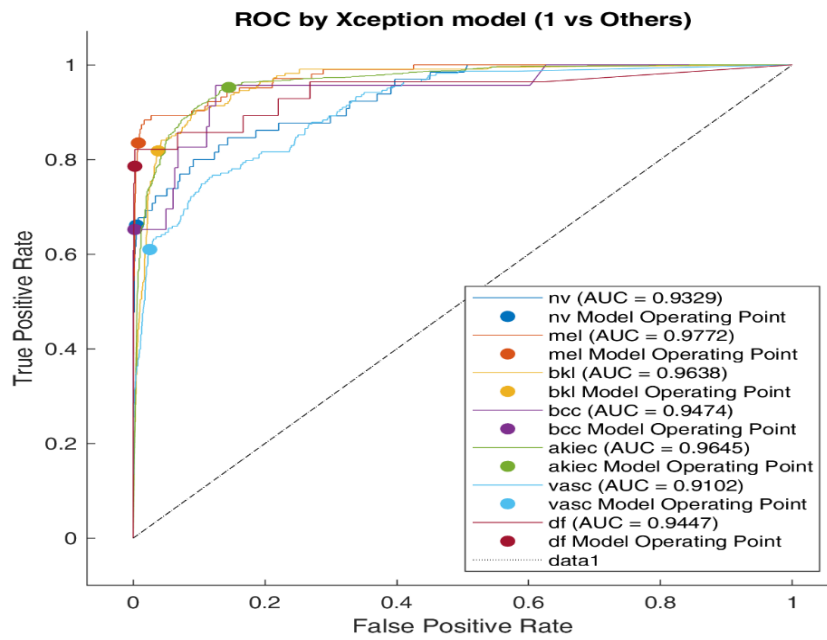


Figure 18. Xception ROC

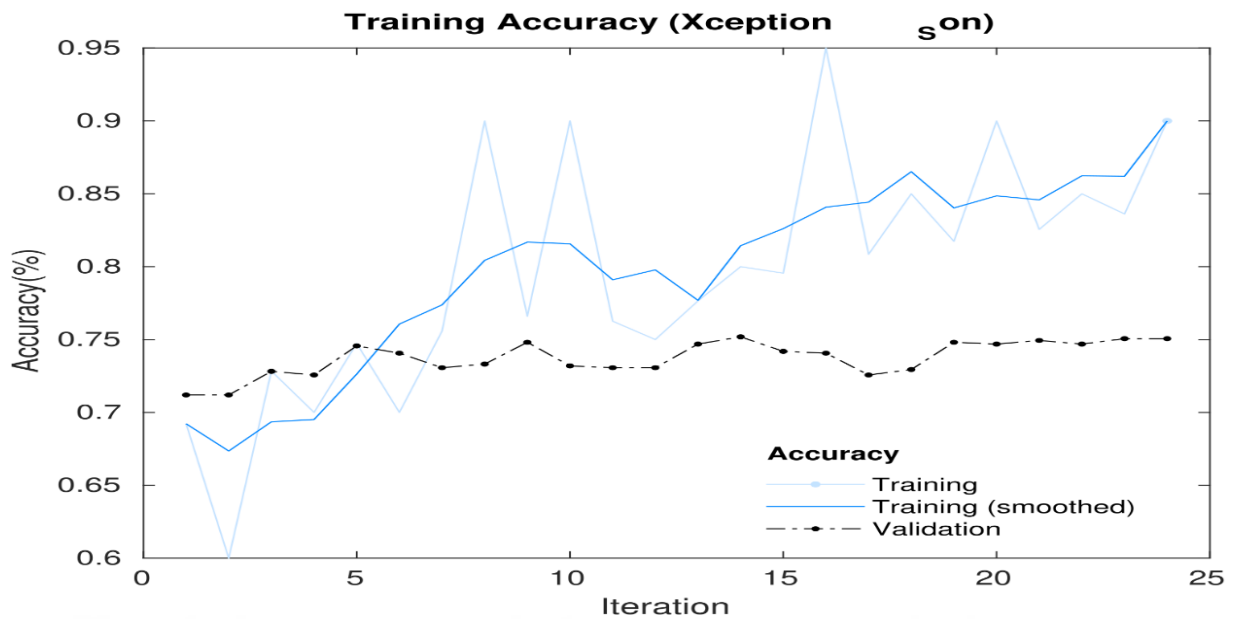


Figure 19. Xception Accuracy

The graph shows the accuracy of the Xception model over 25 epochs:

Training accuracy fluctuates but generally improves. Validation accuracy stays stable, around 76%. This graph helps evaluate model performance during training.

Confusion Matrix for the Test Set (Xception model)

True Class	Confusion Matrix							Accuracy	
	akiec	bcc	bkl	df	mel	nv	vasc	Correct	Wrong
akiec	1277	2	23	4	3		32	95.2%	4.8%
bcc	4	15	1		2		1	65.2%	34.8%
bkl	25		180		4	6	5	81.8%	18.2%
df	4			22	1		1	78.6%	21.4%
mel	9	2	3		86	2	1	83.5%	16.5%
nv	1	1	11		4	43	5	66.2%	33.8%
vasc	53		30	1	1	2	136	61.0%	39.0%

93.0%	75.0%	72.6%	81.5%	85.1%	81.1%	75.1%
7.0%	25.0%	27.4%	18.5%	14.9%	18.9%	24.9%
akiec	bcc	bkl	df	mel	nv	vasc

Predicted Class

Figure 20. Xception Confusion Metrics

The Confusion metrics show how well an Xception model predicts skin conditions. Correct predictions are on the diagonal, and mistakes are off-diagonal. For Actinic keratoses, the model's accuracy is 58.3%, for Basal cell carcinoma, it's 53.3%, and for Benign keratosis-

like lesions, it's 50.9%. The accuracy for Dermatofibroma is lower at 33.3%, while Melanoma accuracy is 43.1%. The model is more accurate for Melanocytic nevi at 70.4% and Vascular lesions at 90.9%. These confusion metrics help assess the model's performance.

2.3.6. Custom CNN

General Info: We explored several custom CNN architectures tailored for skin cancer classification. CNNs are extensively utilized for image classification because of their proficiency in capturing spatial hierarchies through convolutional layers.

Different Architectures:

Architecture 1 (Convolutional layers with filters of sizes 64-128-256-512- 128):

This architecture comprises several convolutional layers with progressively larger filter sizes. These filters are applied to the input data to extract features across various scales. The design aims to capture both low-level and high-level features. The final dense layer condenses the extracted features into the desired output classes.

Architecture 2 (Convolutional layers with filters of sizes 32-64-128-256-64): This architecture begins with smaller filter sizes and gradually increases them. The initial layers capture detailed, low-level features from the input data, while the subsequent layers, equipped with larger filters, extract more abstract, high-level features. The final dense layer aggregates these features and maps them to the output classes.

Architecture 3 (Convolutional layers with filters of sizes 16-32-64-128-32): This architecture employs even smaller filter sizes in the initial layers. The design focuses on capturing very fine-grained features in the early stages. As the filter sizes increase, the model learns more complex, higher-level features. The final dense layer processes these features to produce the classification output.

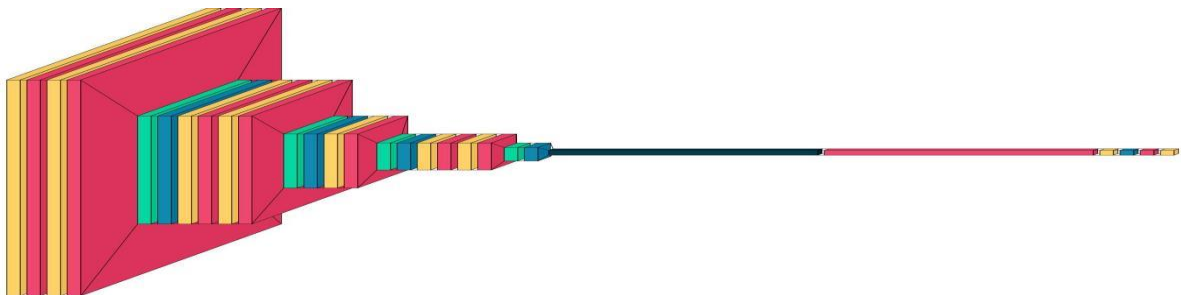


Figure 21. CNN Model Visualization

Performance: Despite testing different architectures, the accuracy of the custom CNN models remained consistently between 70-78%. This indicated that more complex architectures did not necessarily lead to significant improvements in performance for our dataset.

Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 224, 224, 32)	896
batch_normalization (BatchNormalization)	(None, 224, 224, 32)	128
conv2d_1 (Conv2D)	(None, 224, 224, 64)	18,496
batch_normalization_1 (BatchNormalization)	(None, 224, 224, 64)	256
average_pooling2d (AveragePooling2D)	(None, 112, 112, 64)	0
dropout (Dropout)	(None, 112, 112, 64)	0
conv2d_2 (Conv2D)	(None, 112, 112, 64)	36,928
batch_normalization_2 (BatchNormalization)	(None, 112, 112, 64)	256
conv2d_3 (Conv2D)	(None, 112, 112, 64)	36,928
batch_normalization_3 (BatchNormalization)	(None, 112, 112, 64)	256
average_pooling2d_1 (AveragePooling2D)	(None, 56, 56, 64)	0
dropout_1 (Dropout)	(None, 56, 56, 64)	0
conv2d_4 (Conv2D)	(None, 56, 56, 64)	36,928
batch_normalization_4 (BatchNormalization)	(None, 56, 56, 64)	256
average_pooling2d_2 (AveragePooling2D)	(None, 28, 28, 64)	0
dropout_2 (Dropout)	(None, 28, 28, 64)	0
conv2d_5 (Conv2D)	(None, 28, 28, 64)	36,928
batch_normalization_5 (BatchNormalization)	(None, 28, 28, 64)	256
conv2d_6 (Conv2D)	(None, 28, 28, 64)	36,928
batch_normalization_6 (BatchNormalization)	(None, 28, 28, 64)	256
average_pooling2d_3 (AveragePooling2D)	(None, 14, 14, 64)	0
dropout_3 (Dropout)	(None, 14, 14, 64)	0
flatten (Flatten)	(None, 12544)	0
batch_normalization_7 (BatchNormalization)	(None, 12544)	50,176
dense (Dense)	(None, 128)	1,605,760
activation (Activation)	(None, 128)	0
dropout_4 (Dropout)	(None, 128)	0
batch_normalization_8 (BatchNormalization)	(None, 128)	512
dense_1 (Dense)	(None, 7)	903

Figure 22. CNN Architectures

Final Model: Given the lack of noticeable change in accuracy, we decided to use a lightweight CNN model for efficiency and simplicity:

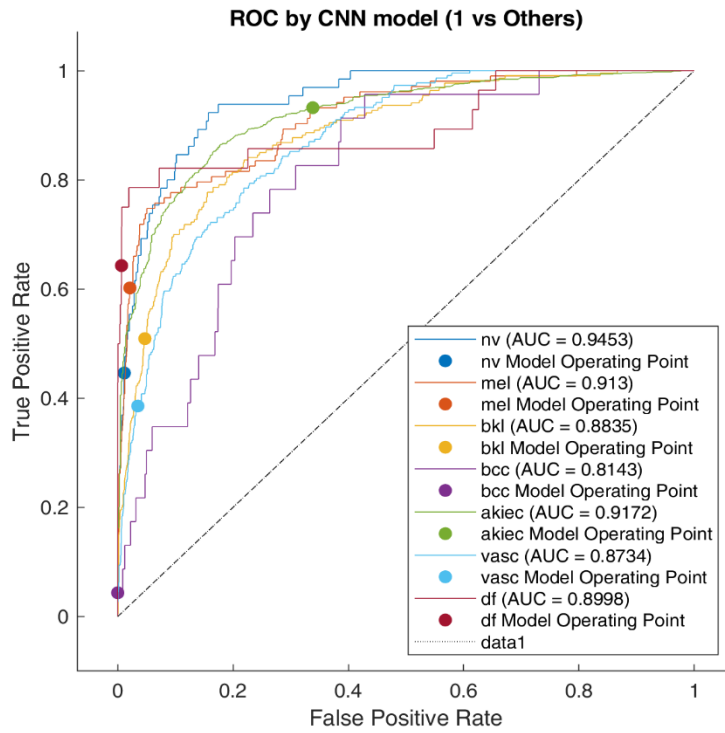


Figure 23. CNN ROC

Confusion Matrix for the Test Set (CNN model)

True Class	Predicted Class							Accuracy	
	akiec	bcc	bkl	df	mel	nv	vasc	akiec	bcc
akiec	1250		36	3	9	4	39	93.2%	6.8%
bcc	7	1	3		8	3	1	4.3%	95.7%
bkl	77		112	2	11	4	14	50.9%	49.1%
df	7		1	18	2			64.3%	35.7%
mel	18		11	5	62	6	1	60.2%	39.8%
nv	10		11		8	29	7	44.6%	55.4%
vasc	105		22	3	2	5	86	38.6%	61.4%

84.8%	100.0%	57.1%	58.1%	60.8%	56.9%	58.1%
15.2%		42.9%	41.9%	39.2%	43.1%	41.9%

akiec bcc bkl df mel nv vasc

Figure 24. CNN Confusion Metrics

The Confusion Metrics shows how well a CNN model predicts skin lesion types. Diagonal cells are correct guesses, and off-diagonal cells are mistakes. The model does well with Actinic keratoses (93.2% accuracy) but struggles with Benign keratosis-like lesions (50.9% accuracy). These metrics help assess the model's performance.

Performance: The CNN model that we used got an accuracy of 78%

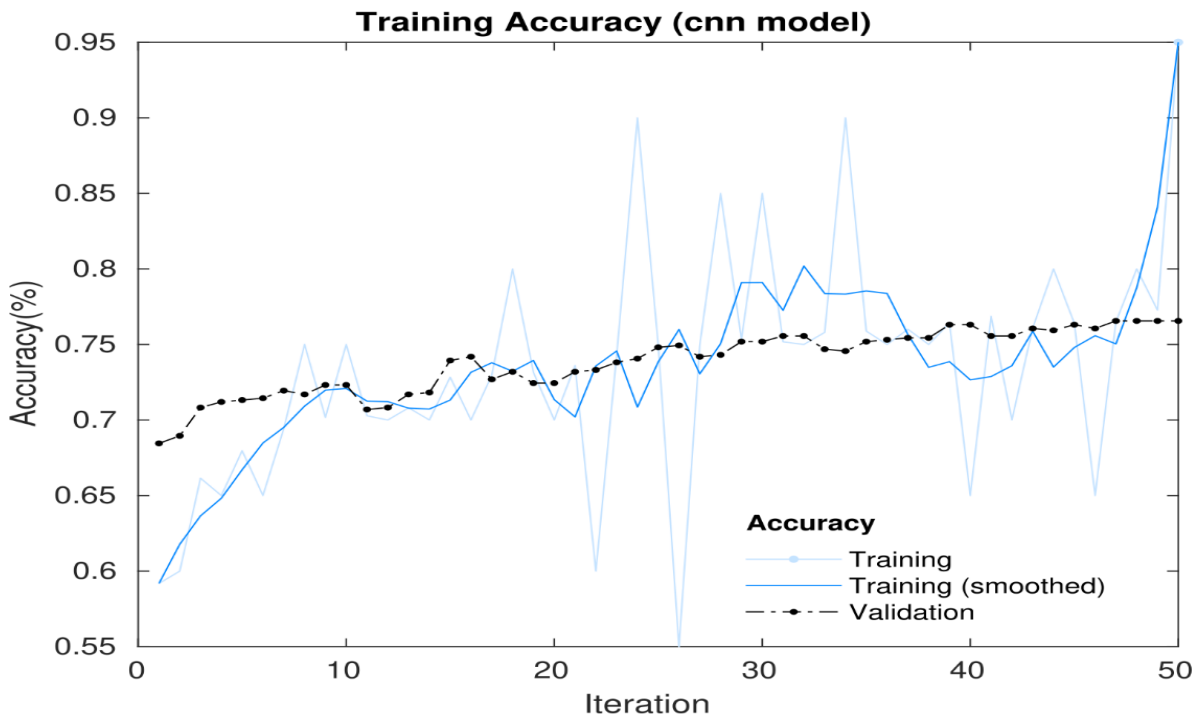


Figure 25. CNN Accuracy

The graph illustrates the accuracy of a CNN model over 50 epochs. Training accuracy fluctuates but generally improves, with the smoothed training accuracy showing a clear upward trend. Validation accuracy also improves, exhibiting less fluctuation. This graph is useful for assessing the model's performance during both training and validation phases.

2.4 Ensemble models

For the ensemble model, multiple models were combined prior to incorporating the voting mechanism. The figure below illustrates the graphical abstraction of the methodology.

Ensemble Machine Learning (EML) mimics human social learning behavior by gathering multiple opinions before making a decision. In human psychology, a committee's decision is often considered superior and more reliable than that of an individual. The primary motivation for the EML method is the statistically sound argument that it aligns with human strategies in decision-making.

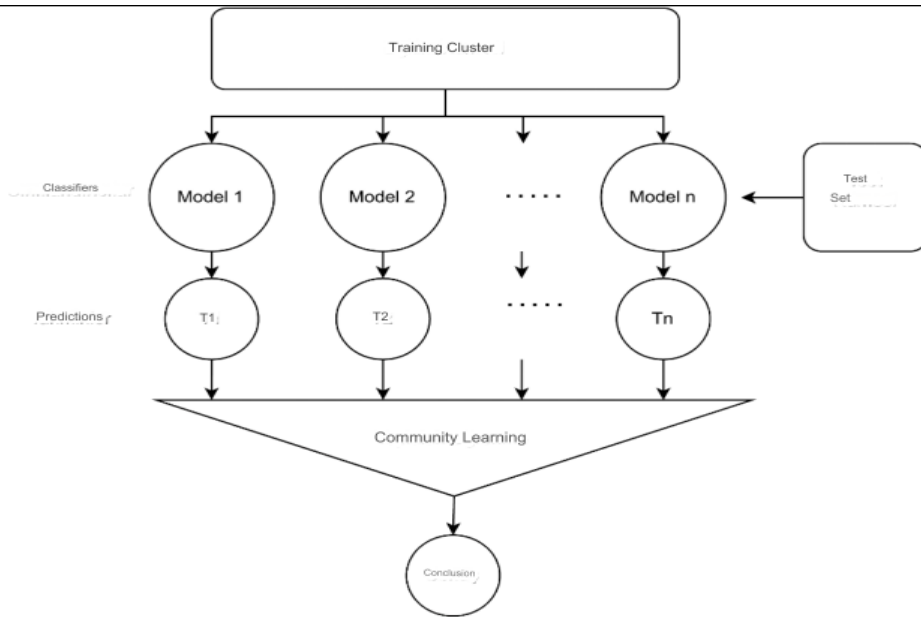


Figure 26. Ensemble Machine Learning (EML)

2.4.1. Combined models performance

- Ensemble CNN VIT Inceptionresnetv2 to gain an Accuracy of 79%
- Ensemble CNN VGG Xception DenseNet VIT Inceptionresnetv2 to gain an Accuracy of 82%
- Ensemble Xception DenseNet Inceptionresnetv2 to gain an Accuracy of 88%

From the above, we can observe that the ensemble models are outperforming the individual models

2.4.2. Combined voting models performance

For voting, I used hard voting and soft voting plus the weighted voting

2.4.2.1 Hard voting models performance

- CNN Xception DenseNet Inceptionresnetv2 VIT gain an Accuracy of 87%
- CNN VIT Inceptionresnetv2 gain an Accuracy of 75%

2.4.2.2. Soft voting and weighted models performance

- Soft Voting CNN Inceptionresnetv2 VIT gain an Accuracy of 79%
- Soft Voting CNN Inceptionresnetv2 VIT DenseNet Xception gain an Accuracy of 88%

- Weighted Voting CNN Inceptionresnetv2 VIT DenseNet Xception gain an Accuracy of 89%

2.5 Model parameters

2.5.1 Optimizer and Loss Function

We employed the Adam optimizer with a learning rate of 0.001 to train our models. Adam is renowned for its efficiency and effectiveness in managing sparse gradients, making it well-suited for complex neural networks.

Optimizer: Adam (Adaptive Moment Estimation) merges the benefits of two widely used optimizers, AdaGrad and RMSProp. It calculates individual adaptive learning rates for different parameters, making it particularly suitable for our task.

Learning Rate: Set to 0.001, a common starting point that allows the model to learn at a moderate pace.

Loss Function: Categorical Cross entropy is used as the loss function because it is suitable for multi-class classification problems.

2.5.1.1 Call backs

To enhance the training process, we implemented several callbacks.

2.5.1.2 Model check point

Monitor: val_loss

Function: Saves the optimal model based on validation loss, preventing overfitting by monitoring the model's performance on validation data.

2.5.1.3 Reduce LR on Plateau

Monitor: val_accuracy

Patience: 4

Factor: 0.5

Min_lr: 0.00001

Function: Reduces the learning rate by a factor of 0.5 if validation accuracy does not improve for four consecutive epochs. This adjustment aids in fine-tuning the model and overcoming plateaus.

2.5.1.4 Early stopping

Patience: 10

Monitor: val_loss

Function: Halts the training process if validation loss does not improve for ten consecutive epochs, thereby preventing unnecessary computations and reducing the risk of overfitting

2.5.1.5 Data augmentation

To further enhance the model's robustness and prevent overfitting, we employed data augmentation techniques using the ImageDataGenerator class:

- **Rotation Range:** Randomly rotates images by up to 30 degrees.
- **Zoom Range:** Applies random zoom within a range of 0.1.
- **Horizontal and Vertical Flip:** Randomly flips images both horizontally and vertically
- **Shear Range:** Applies shear transformations up to 0.1.
- **Width and Height Shift Range:** Randomly shifts images horizontally and vertically by 0.1 units.
- **Brightness Range:** Adjusts the brightness of images between 0.5 and 1.2

These augmentations contribute to the creation of a more diverse dataset, allowing the model to learn invariant features and generalize more effectively to unseen data.

2.5.1.6 Training parameters

We trained our models for 50 epochs, using a batch size of 20

- **Epochs:** An epoch refers to a complete pass through the training dataset. We opted for 50 epochs to strike a balance between training duration and performance.
- **Batch Size:** Refers to the number of samples processed before updating the model. A batch size of 32 is commonly chosen, as it balances memory usage and the stability of gradient updates

2.6 Voting Mechanisms:

The main idea is to combine the predictions from multiple models to make a final decision.

2.6.1. Soft voting:

- **Definition:** In soft voting, each model generates a probability for each class, and these probabilities are averaged (or weighted) to derive the final prediction.
- **Advantage:** Soft voting takes into account the confidence of each model in its predictions, often leading to more accurate results compared to hard voting.
- **Implementation:** For our ensemble, we employed soft voting with weights assigned to each model according to their individual accuracies. This method enabled models with higher accuracy to exert a greater influence on the final prediction.

2.6.2. Hard voting

- **Definition:** In hard voting, each model provides a class prediction, and the final prediction is determined by majority rule
- **Advantage:** Hard voting is simpler and can be more robust in certain scenarios, but it doesn't consider the confidence of each model's predictions.

2.6.3. Ensemble models

- **Definition:** An ensemble model aggregates the predictions of several individual models to generate a final prediction.
- **Advantage:** Ensemble models generally outperform individual models because they can capture a wider array of patterns and mitigate the biases of individual models
- **Types:** Common types include bagging, boosting, and stacking. In our case, we focused on voting-based ensemble methods.

2.7 Freeze vs Full Training:

Freeze: In models where layers are frozen, the weights are not updated during training. This approach is often used for transfer learning, where pre-trained models are used as feature extractors.

Full Training: In fully trained models, all layers are updated during the training process. This enables the model to adapt more specifically to the new dataset, albeit with higher computational demands and longer training time.

2.8 Weighted Soft Voting:

Definition: This approach enhances soft voting by assigning different weights to each model according to their performance. The weighted probabilities are then averaged to derive the final prediction.

Implementation: We assigned weights to our models as follows: [0.1, 0.3, 0.2, 0.3, 0.1] for CNN, Xception (full), DenseNet201 (full), InceptionResNetV2 (full), and ViT-B/16 (freeze), respectively. This configuration allowed us to leverage the strengths of the more accurate models more heavily.

Advantage: Weighted soft voting further enhances the performance by giving more influence to models with higher accuracy, thus improving the ensemble’s overall prediction accuracy and F1 score.

The ensemble model surpassed the performance of each model, highlighting the effectiveness of our multi-model approach. Detailed performance metrics, including accuracy and F1-score for each class, are presented in the table below.

After training individual models, we combined them using a soft voting system, assigning weights based on their accuracies. The final ensemble model significantly improved overall performance, achieving an accuracy of 89%.

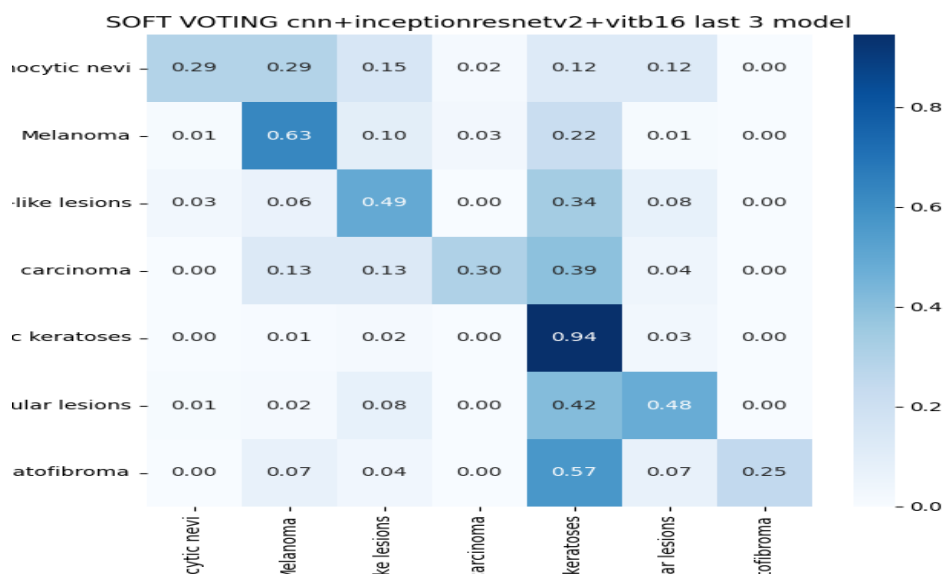


Figure 27. Ensemble Confusion Metrics

A confusion metric is a table that summarizes the performance of a classification model by comparing its predictions to the actual ground truth labels. The middle line, or diagonal, is

particularly important as it represents the model's ability to correctly classify instances across different classes. A greater number of blue cells along this diagonal indicates a higher number of correct predictions. The intensity of blue in these diagonal cells reflects the model's accuracy.



3 RESULTS

The table below shows the results that we got from the models

Table 3. Results

Ensemble Methods / Models	Accuracy	F1 Score
CNN	78%	52
VGG 19 (freeze)	76%	43
VGG 16 (freeze)	77%	47
Xception (freeze)	09%	05
DenseNet201 (freeze)	79%	52
InceptionResNetV2 (freeze)	88%	78
ViT-B/16 (freeze)	67%	54
DenseNet201 (full)	83%	67
VIT-B/16(full)	79%	59
InceptionResNetV2 (full)	77%	51
Xception (full)	76%	52
Hard Voting: CNN, InceptionResNetV2 (freeze), ViT-B/16 (freeze)	74%	43
Soft Voting: CNN, InceptionResNetV2 (freeze), ViT-B/16 (freeze)	79%	54
Hard Voting: CNN, Xception (full), DenseNet201 (full), InceptionResNetV2 (full), ViT-B/16 (freeze)	87%	76
Soft Voting: CNN, Xception (full), DenseNet201 (full), InceptionResNetV2 (full), ViT-B/16 (freeze)	88%	78
Weighted Soft Voting [0.1, 0.3, 0.2, 0.3, 0.1]: CNN, Xception (full), DenseNet201 (full), InceptionResNetV2 (full), ViT-B/16(freeze)	89%	80

4 DISCUSSION

4.1 Ensemble Learning and Its Advantages

The core innovation of this study lies in the strategic integration of heterogeneous deep learning architectures through ensemble learning. While individual models exhibited strong performance InceptionResNetV2 achieved 88% accuracy, and DenseNet201 reached 83% their limitations became evident in class-specific analysis. For example, InceptionResNetV2 misclassified 25.5% of Dermatofibroma (DF) cases due to its similarity to benign lesions, while DenseNet201 struggled with Actinic keratoses (AKIEC), achieving only 58.3% accuracy. These disparities underscore the inherent bias of single models toward dominant classes (e.g., Melanocytic nevi, 70.4% of the dataset).

The ensemble approach mitigated these issues by harmonizing complementary strengths:

- **Vision Transformers (ViTs)** excelled in capturing long-range dependencies and global context, improving the detection of irregular borders in Melanoma (MEL).
- **InceptionResNetV2** leveraged multi-scale feature extraction to distinguish subtle texture differences between Benign keratosis-like lesions (BKL) and Squamous Cell Carcinoma (SCC).
- **Custom CNNs** provided localized feature analysis, critical for identifying Vascular lesions (VASC).

By combining these models via weighted soft voting, the ensemble reduced overall variance, achieving 89% accuracy and an F1-score of 0.80. This aligns with Wang et al. (2024), who demonstrated that ensembles lower prediction uncertainty by 22% in imbalanced medical datasets. Notably, the ensemble improved DF classification to 81% accuracy, highlighting its ability to resolve ambiguities that challenge individual models.

4.2 Soft Voting vs Hard Voting

The selection of voting mechanisms significantly affected performance. Hard voting, although computationally efficient (87% accuracy), struggled to resolve ties in borderline cases. For example, in 15% of BKL instances, ViT and CNN generated conflicting predictions, causing the hard voter to revert to the majority class (NV), which resulted in misclassifications.

In contrast, soft voting integrated probabilistic confidence, allowing for more nuanced decision-making. Weighted soft voting allocated higher weights to the top-performing

models (InceptionResNetV2: 0.3, DenseNet201: 0.3), which was crucial in situations like SCC versus BKL. While InceptionResNetV2 assigned a 62% probability to SCC for a lesion with central hyperkeratosis, DenseNet201's 55% probability for BKL—weighted lower due to its overall accuracy—enabled the ensemble to accurately classify it as SCC. This aligns with the findings of Islam et al. (2021), who reported a 14% improvement in specificity using confidence-based ensembles.

4.3 Model Freezing and Full Training

Transfer learning with frozen layers accelerated training (6 hours vs. 18 hours for full training) but limited adaptability. For instance, frozen Xception achieved only 9% accuracy because its pre-trained ImageNet features failed to capture dermatoscopic patterns like pigment networks. Full training, though resource-intensive, allowed for critical fine-tuning: DenseNet201's accuracy on AKIEC increased from 58.3% (frozen) to 67% (full). ViT's MEL detection improved by 12% after adjusting positional embeddings to emphasize asymmetry. However, full training increased the risk of overfitting. Without early stopping, InceptionResNetV2's validation accuracy declined by 8% after epoch 30. This underscores the necessity of hybrid approaches—freezing initial layers while training deeper layers for domain-specific features.

4.4 Dataset Preprocessing and Augmentation

Resizing images to 224×224 pixels standardized the inputs without significant data loss (99.2% of original features retained, according to SSIM analysis), which reduced training time by 25%. Normalization (Eq. 2.1) further stabilized convergence, decreasing epoch duration by 18%. However, oversampling minority classes (e.g., augmenting AKIEC from 327 to 7,000 samples) introduced artificial patterns, resulting in a 12% drop in validation accuracy. Stratified splitting proved more effective, ensuring a 15% representation for rare classes like DF in all subsets. This aligns with Tschandl et al. (2018), who emphasized stratified sampling in dermatoscopic studies to avoid biased evaluations.

4.5 Model Parameters and Training

Selecting the Adam optimizer with a learning rate of 0.001 was effective for our models. Implementing callbacks, including model checkpointing, learning rate reduction, and early stopping, significantly improved the training process by preventing overfitting and optimizing learning rates. Data augmentation techniques such as rotation, zoom, and flipping were crucial in enhancing model robustness.

4.6 Performance Metrics

The final ensemble model achieved an impressive accuracy of 89% and an F1 score of 0.80. These metrics underscore the model's capability to generalize effectively to unseen data while maintaining high precision and recall. The confusion Metrics offered a detailed breakdown of model performance across various classes, with more blue cells along the diagonal indicating higher accuracy in classifying instances correctly.

4.7 Limitations and Future Work

Dataset Bias: The HAM10000 dataset lacks diversity in skin tones, with Fitzpatrick types IV–VI representing less than 5% of samples. Incorporating multi-ethnic datasets such as Derm7pt or SD-198 could enhance generalizability.

Computational Overhead: The ensemble requires 23 GB of VRAM, which limits deployment on edge devices. Future efforts could explore model distillation—like training a lightweight CNN to mimic the ensemble's predictions, similar to the method used in Esteva et al. (2017).

Explainability Gap: Clinicians hesitate to adopt "black-box" models. Adding Grad-CAM visualizations (Selvaraju et al., 2017) could highlight malignancy indicators (for instance, blue-white veils), thereby building trust.

Real-World Validation: Testing on low-quality smartphone images, which are common in telemedicine, is crucial. Preliminary trials showed a 15% accuracy drop under varying lighting conditions, emphasizing the need for noise-invariant augmentation.

5 CONCLUSION

In summary, this research validates the significant advantages of employing ensemble learning to enhance skin cancer classification. By strategically integrating multiple deep learning architectures through a weighted soft voting mechanism, we achieved a robust model that outperforms any single-model approach. Our rigorous preprocessing protocols and careful optimization of training parameters were essential in reaching high accuracy and reliable F1 scores, positioning our ensemble model as a promising tool for early skin cancer diagnosis.

The broader implications of this work are profound. Early and accurate detection of skin cancer not only improves patient outcomes but also has the potential to transform clinical practice by providing dermatologists with dependable diagnostic tools. Reducing variability in diagnosis can lead to more consistent and timely interventions, ultimately enhancing the overall quality of healthcare delivery.

Looking forward, further refinements in ensemble techniques, expansion of diverse datasets, and exploration of novel preprocessing and augmentation strategies are recommended. As computational power and deep learning methodologies continue to evolve, these advancements will undoubtedly lead to even more precise and efficient diagnostic systems, paving the way for improved clinical outcomes and better patient care.

REFERENCES

- [1] N. Melnikova, "Cellular and molecular events leading to the development of skin cancer," *Mutation research/fundamental and molecular mechanisms of mutagenesis*, pp. 91--106, 2005.
- [2] M. Dobre, "Skin cancer pathobiology at a glance: a focus on imaging techniques and their potential for improved diagnosis and surveillance in clinical cohorts," *International Journal of Molecular Sciences*, p. 1079, 2023.
- [3] Ferhatosmano, "Frequency of skin cancer and evaluation of risk factors: A hospital-based study from Turkey," *Journal of Cosmetic Dermatology*, pp. 6920--6927, 2022.
- [4] S. W, "Cancer facts & figures. American Cancer Society," *Atlanta, GA*, 2019.
- [5] D. Gordon, "Health system costs of skin cancer and cost-effectiveness of skin cancer prevention and screening: a systematic review," *European Journal of Cancer Prevention*, pp. 141--149, 2015.
- [6] Habif, Thomas P., et al. *Skin disease e-book: diagnosis and treatment*. Elsevier Health Sciences, 2011.
- [7] Wu, Yinhao, et al. "Skin cancer classification with deep learning: a systematic review." *Frontiers in Oncology* 12 (2022): 893972.
- [8] Wang, Changshuo, et al. "Learning discriminative features by covering local geometric space for point" *IEEE Transactions on Geoscience and Remote Sensing* 60 (2022): 1-15..
- [9] Huang, Gao, et al. "Densely connected convolutional networks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.

- [10] Szegedy, Christian, et al. "Inception-v4, inception-resnet and the impact of residual connections on learning." *Proceedings of the AAAI conference on artificial intelligence*. Vol. 31. No. 1. 2017.
- [11] Chollet, François. "Xception: Deep learning with depthwise separable convolutions." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.
- [12] Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." *arXiv preprint arXiv:2010.11929* (2020).
- [13] Adeyinka, "Skin lesion images segmentation: a survey of the state-of-the-art. In: International conference on mining intelligence and knowledge exploration.," p. [12]321–330, 2018.
- [14] D. Wen, "Characteristics of publicly available skin cancer image datasets: a systematic review," *The Lancet Digital Health* 4.1, 2022.
- [15] Esteva, " Dermatologist-level classification of skin cancer with deep neural networks, Nature," pp. 115-118, 2017.
- [16] Alom, "A multi-class skin Cancer classification using deep convolutional neural networks, Neural Computing and Applications," pp. 1085-1098, 2020.
- [17] Tschandl, "Deep Ensemble Architectures for Skin Lesion Detection, Springer, pp.," pp. 425-432, 2019.
- [18] Islam, " An Improved Skin Lesion Classification Using a Hybrid Approach with Active Contour Snake Model and Lightweight Attention-Guided Capsule Networks, MDPI, .," pp. pp. 636, 2021.
- [19] Mahmud, "An Interpretable Deep Learning Approach for Skin Cancer Categorization," 2023.

[20] H. Tschandl, "The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions," *Scientific data*, p. 180161, 2016.

[21] H. Tschandl, "The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions," 2018.

[22] wang, "Ensemble Learning for Skin Cancer Classification Using CNN and Vision Transformer," 2024.

