

# Image Driven Skin Cancer Diagnosis

---

- [Proposal](#)
- [Mid-Term](#)
- [Final](#)

## Final

---

### Introduction and Background

---

One in five Americans will be diagnosed with skin cancer by the time they turn 70. More than two people die of skin cancer every hour. Unlike other aggressive cancers, skin cancer has a five-year survival rate greater than 95% provided it is detected and treated early. This makes timely and accurate diagnosis the need of the hour.

Traditionally, skin lesions are diagnosed manually by visual dermatologist examinations followed by biopsies. Given the growing shortage of dermatologists, automating the diagnostic process would reduce the burden on our medical systems in terms of time, cost and effort, and improve timely accessibility of diagnostic mechanisms to wider segments of society.

### Dataset

---

[ISIC 2024 - Skin Cancer Detection with 3D-TBP](#) is a collection of around 400,000 images of isolated skin lesions taken from 3D full body images. Each lesion image is cropped and labeled to indicate if the lesion is benign or malignant.

The data also includes significant metadata, like the patient's age, gender, which can improve training of the machine learning models. The images resemble regular, close-up smartphone pictures in order to ensure the quality is representative of the kind of pictures regularly submitted in telehealth settings.

### Literature Review

---

Skin cancer detection has significantly benefited from the advancements in deep learning methodologies, offering improved accuracy and efficiency compared to traditional diagnostic approaches. [1] demonstrated the effectiveness of CNNs in classifying skin cancer with accuracy comparable to dermatologists. Their model was trained on over 129,000 clinical images, which shows

the scalability of deep learning models. Studies cited in [2] demonstrated that transfer learning significantly enhances the performance of models when applied to limited datasets. The authors also noted that integrating deep learning models into clinical workflows can aid dermatologists in diagnosing skin cancer.

[3] employed ResNet-152 to achieve high accuracy in skin cancer detection. Their approach demonstrated the efficiency of using pre-trained models on large datasets and fine-tuning them on medical images, showcasing how transfer learning can compensate for the limited availability of annotated medical data. Researchers at Google Health presented DermAssist [4], a teledermatology AI system aimed at diagnosing around 26 different skin conditions through deep learning models trained on around 16K data points collected from telehealth services. They found their system to be non-inferior in performance in comparison to six dermatologists and clearly superior to six primary care physicians and six nurse practitioners.

[5] presents 10 clear ethical risks involved in using AI within health-care settings including explainability, reliability, privacy, responsibility and dehumanization.

## Problem Definition

---

### Problem

Of the three major types of skin cancer (Basal Cell, Squamous Cell and Melanoma), Melanoma, while less common, is the deadliest and is estimated to be diagnosed 200,000 times in the US in 2024 with 9000 projected to succumb to it. Skin cancer has a nearly 95% cure rate provided the patient has an early diagnosis and treatment. This makes diagnosis of skin cancer a time-sensitive problem. Additionally, studies report a decline in the number of dermatologists available and in the post-covid era, our health care systems are operating under heavy burdens.

Accessible diagnostic methods are the need of the hour. Telehealth allows patients who otherwise lack access to timely diagnosis receive crucial care. The challenge is to develop an automated system for malignancy prediction in skin lesions using their images. This system can be used in telehealth settings where users can upload pictures of lesions to get a preliminary understanding of their skin condition.

### Motivation

By building an ML model that can classify skin lesion images as cancerous or non-cancerous, this project aims to:

- Improve detection rates at early stages
- Provide a second opinion for dermatologists to rely on
- Provide access to regions lacking professional diagnostic services

# Methods

---

## Data Preprocessing

### Down-sampling Benign Tumors

One of the primary concerns with the dataset is the large imbalance between the number of benign tumors (~400k) and the number of malignant tumors (~400). To address this, we downsampled the number of benign tumors by randomly sampling using the reservoir sampling algorithm, ensuring uniform random selection of elements without replacement. We retained 10k images of benign tumors using this process.

### Focused Augmentation of Malignant Tumors

In order to improve data variability and address class imbalance by increasing representation of malignant tumors we did focused augmentation of malignant tumors. We performed 16 augmentations on each image. The augmentations performed include geometric transformations like flipping, rotating, distorting as well as color and lighting modifications including modifying brightness, contrast, saturation, hues, CLAHE. Adding noise and blurring effects including motion and gaussian blur and gauss noise and finally simulating occlusions and missing data by cutting out random rectangular sections from the image.

### Image Normalizations

We further normalized values of all pixels to be between 0 and 1 in order to improve faster convergence, aid optimal weight initialization and promote model stability and generalization.

### Image Resizing

Finally all images were resized to a standard dimension (128x128) to ensure consistency, reduce computational loads, improve training time and memory efficiency.

## Exploratory Data Analysis

Post pre-processing the data, exploratory data analysis was conducted to better understand the data. This involved calculating image means, deviations, color channels and comparing their distributions across the two target classes.

## Unsupervised Learning

Ground Truth data is hard to come by in medical applications due to privacy concerns. Being able to solve the diagnosis problem with an unsupervised approach would enable significant progress in the area. Thus, despite having ground truth labels, we attempt to find coherent clusters and map available ground truth with the clusters and evaluate the performance of the clustering algorithms in solving the lesion identification problem. We clustered images based on image mean, standard deviations, skewness, RGB channel means and other essential features highlighted in Figure 1. We determined through reading and experimentation that the number of clusters in our dataset is around 4. We used this as our k parameter and conducted clustering through K-means and Gaussian Mixture Models. In order to conduct clustering through GMM, we had to ensure feature normality. Two features - image standard deviation and image skew needed further normalization before clustering.

## Supervised Learning

Currently, five transfer learning based approaches have been explored and their performance compared. The goal is to continue to train and compare more state of the art computer vision models before finalizing on our primary model.

ImageNet is a large image collection (~14M images and 20K+ classes) which contains diverse data for general object recognition. While not a disease focused dataset, it has shown success in various medical image classification tasks. In this approach, models pre-trained on ImageNet are used as a starting point. Low-level features like edge and texture based analysis could be transferable to lesion identification.

Particularly at this stage - NASNet, EfficientNetB0, Xception, ResNet50 and Densenet were trained with ImageNet weights as starting points. All five are powerful CNN models in the computer vision (CV) space, with architectures optimized for efficient, high-performance image classification.

**NASNet:** NASNet is designed using Neural Architecture Search (NAS), enabling the automated discovery of optimal network architectures. This project uses NASNetMobile, a lightweight version ideal for mobile and embedded applications. This CNN uses a combination of normal and reduction cells obtained from NAS, and it has approximately 5.3 million parameters.

**EfficientNetB0:** EfficientNetB0, developed by Google, is the baseline model in the EfficientNet family. It employs compound scaling, focusing on the uniform scaling of network width, depth, and resolution to balance accuracy and efficiency. This CNN is based on MobileNetV2, utilizing inverted residual blocks with added squeeze-and-excitation (SE) blocks for enhanced performance. It also has about 5.3 million parameters.

**Xception:** Xception is inspired by the Inception architecture but replaces Inception modules with depthwise separable convolutions, resulting in significantly deeper models with fewer connections and improved efficiency for specific tasks. Xception has approximately 22.9 million parameters.

**ResNet50:** contains exactly 50 layers in total that are divided into convolutional layers, pooling layers and fully connected layers. Resnet uses residual blocks which include skip/shortcut connections which

aid the network in learning residual functions for layer inputs instead of learning unreferenced functions. It allows for easier optimization and improved accuracy in comparison to plainer networks of the same depth. At ~25.6 M parameters, it is one of the bigger models.

**DenseNet:** Densely Connected Convolutional Networks is a CNN model that has dense connectivity between layers. Each layer here in DenseNet is directly connected to every other layer in a feed-forward fashion and feature reuse is promoted by allowing each layer to access feature maps from all previous layers. It has a compact structure as it requires fewer parameters compared to traditional CNNs. Its most popular variant - DenseNet 121 has 8M parameters making it relatively lightweight.

**InceptionResNetV2:** This architecture is very powerful when it comes to hierarchical learning of complex features. We hypothesized that this ability to capture such intricate patterns across scales would probably suit our dataset, especially where the minority class may show subtle or distinctive features. Also, the efficient design endowed with automatic robustness to overfitting makes this a strong candidate for the test, particularly in situations with small samples for the minority class.

**VGG19:** It's known for its strength in tasks requiring fine-grained categorization, so it matched well with our intention of bringing reliability even in the most unbalanced classes. Additionally, due to its comparably larger size (144 million parameters), it would be able to model an even more complex data distribution, possibly further differentiating between the classes: that of the majority and that of the minority. VGG19's simple architecture made it a good baseline against which to compare more complicated models like InceptionResNetV2 and those previously explored.

Unlike NASNet and EfficientNet, which come in a variety of sizes to meet different computational environments, Xception is not designed to scale flexibly based on computational needs. NASNet is generally known for achieving high accuracy, though it is computationally expensive. In contrast, EfficientNetB0 is valued for maintaining a good balance between accuracy and efficiency. Efficiency may be a key concern when deploying this application to mobile devices to enhance accessibility for patients. Densenet introduces a hyperparameter that provides the trainer more fine-grained control over the number of new features added to each layer. DenseNet also uses features at all complexity levels leading to smoother decision boundaries.

For each of the models, after freezing the ImageNet weights for our initial base layer, custom base layers were added. A GlobalAveragePooling2D layer was added to reduce spatial dimensions and to prioritize important features. A dense layer with ReLU activation was added to aid additional feature extraction and a layer with sigmoid activation was added to aid binary classification. The models were compiled using Adam optimizer (LR: 0.0001), binary cross entropy was used as the loss function and accuracy was our main evaluation metric.

The challenge of imbalance class in the dataset is tried to be managed by applying Synthetic Minority Oversampling Technique. This is a commonly used data-augmenting technique where interpolation is done among the data in the minority class to create synthetic examples within the class it enhances the balance in the dataset by increasing the amounts of samples in the underrepresented classes of the dataset. However, SMOTE does not make any betterment on our image data. This can be attributed to

the fact that it is mostly developed for tabular data, and thus may not be easily related to the spatial multidimensional presentations of image data.

Since SMOTE failed to cope with the class imbalance of our image data, we have proceeded to study further models, InceptionResNetV2 and VGG19. The reason is that these architectures have unique feature extraction capabilities which may offer better handling of imbalanced datasets.

Thus, we aim to differentiate these architectures to assess if they can capture class-specific features well in a balanced way without having to depend on SMOTE. This also renders a good scope to compare their performance with the five initial models (NASNet, EfficientNetB0, Xception, ResNet50, DenseNet) and identifying the most promising one for our classification task.

## Results and Discussion

---

### Exploratory Data Analysis

From figure 1 we notice that both classes (Benign = 0 and Malignant = 1) have a similar distribution of image means, deviations and skewness as a whole and in each quartile. However, benign tumors generally have slightly higher mean pixel intensities and greater variation in pixel intensity (higher standard deviation), suggesting that they are brighter on average. The most noticeable difference is in the color channels (red, green, blue), where benign images exhibit slightly higher tones across all channels compared to malignant ones..

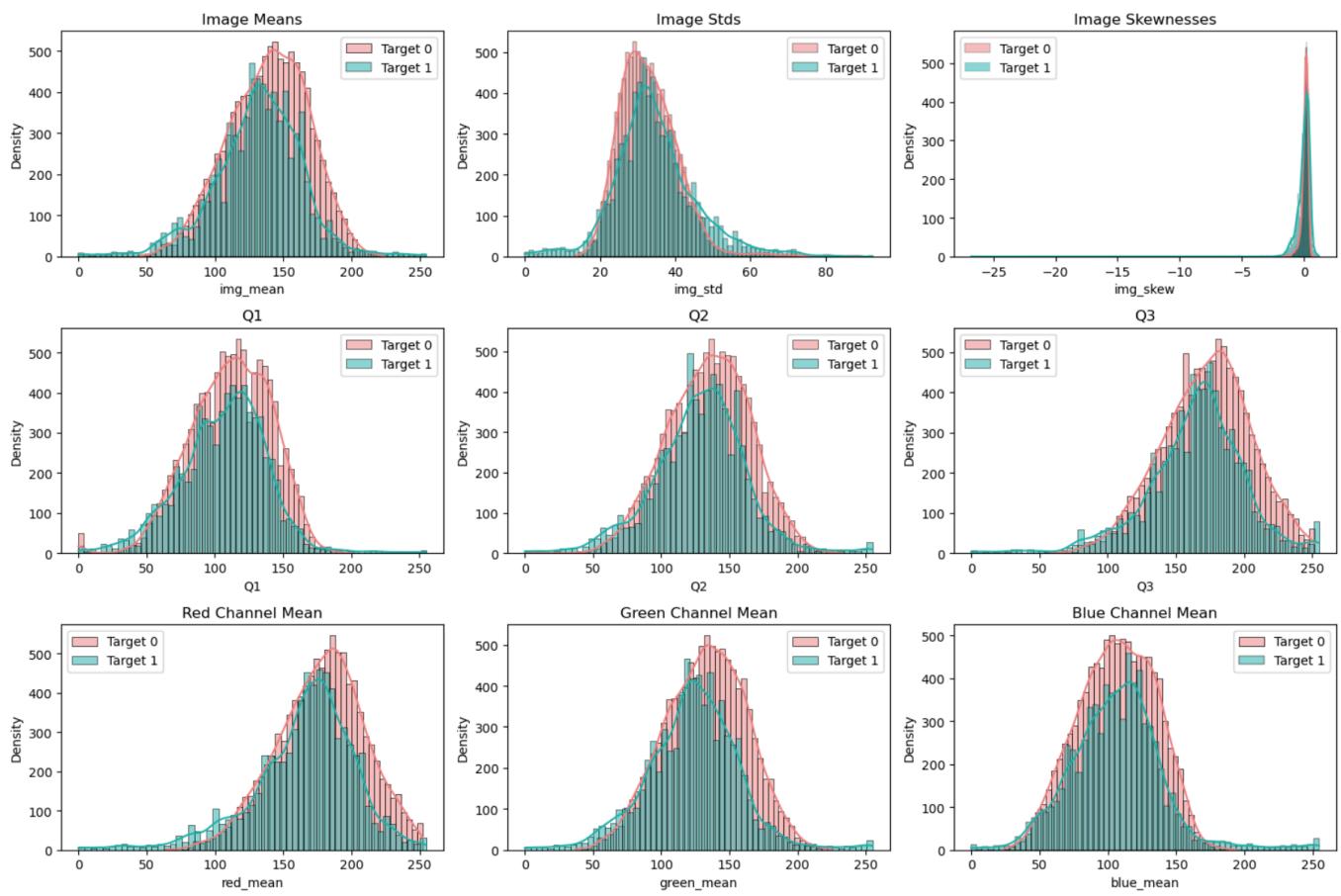
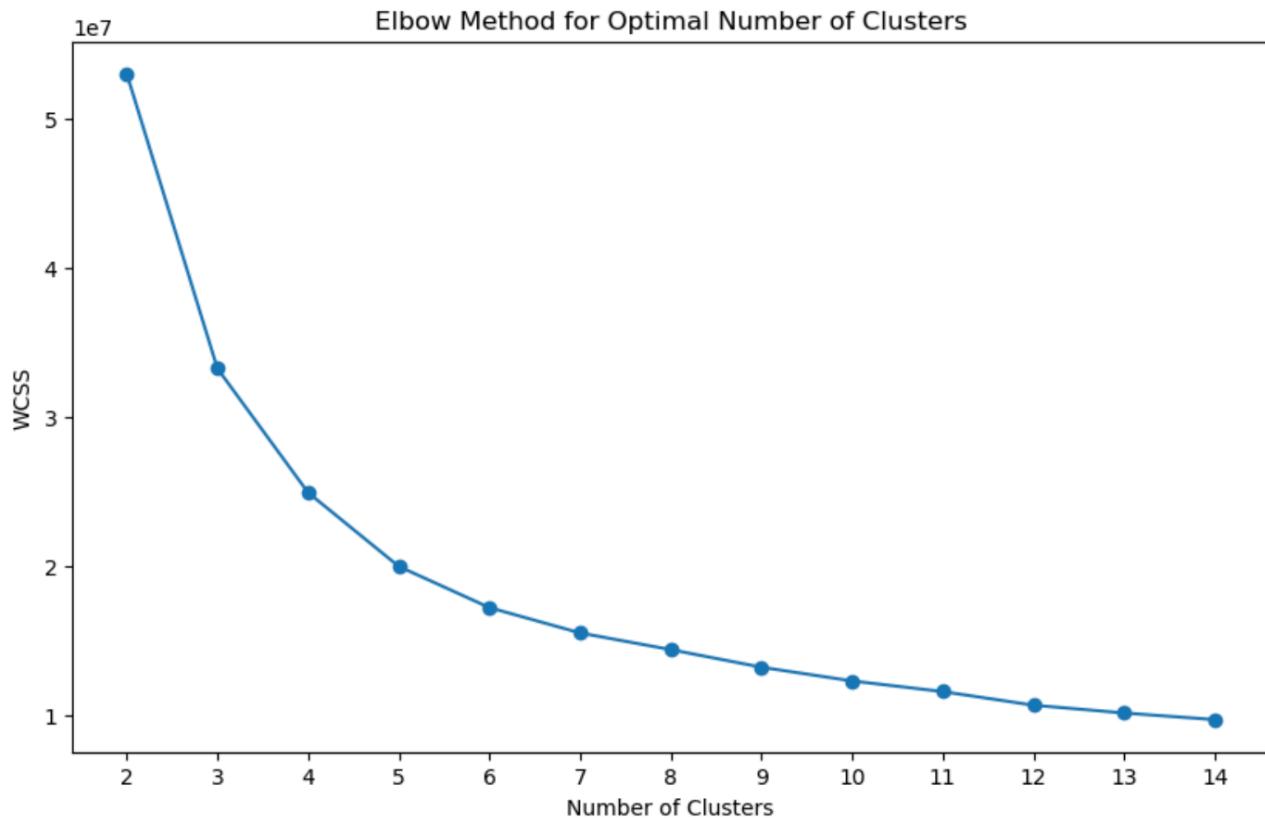


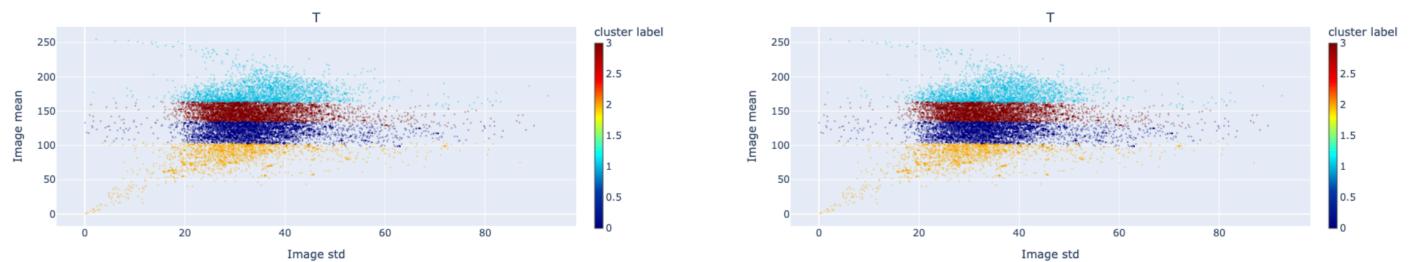
Figure 1. Comparing Distributions of Benign (Target 0) and Malignant (Target 1) Features

## Unsupervised Learning

### Clustering on Image Data

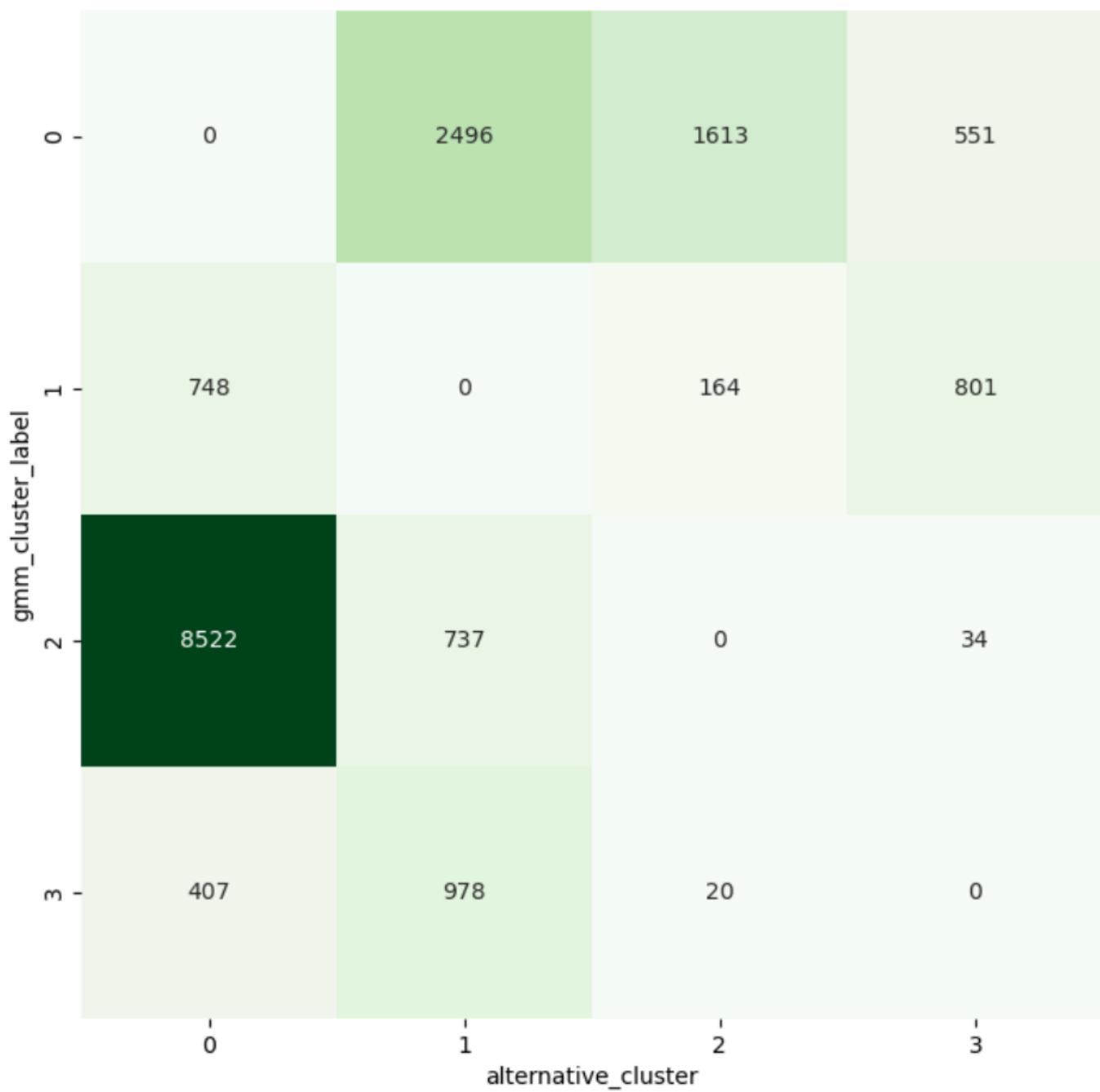


**Figure 2. Elbow Method to determine number of clusters**



**Figure 3A. KMeans Clusters**

**Figure 3B. GMM Clusters**



**Figure 4. GMM top 2 cluster comparisons**

The elbow method was used to identify the number of clusters as 4 (Figure 2). We see four clearly demarcated clusters using both K-mean and GMM (Figure 3). The soft clustering properties of GMM were leveraged to get a better understanding of the dataset. The counts of all images that had the same first and second most likely clusters were plotted in order to understand feature similarity in the dataset. From Figure 4 we conclude that the images closer geometrically on the graph are also easier to miscategorize since they are very close to other groups.

## Clustering on the Metadata

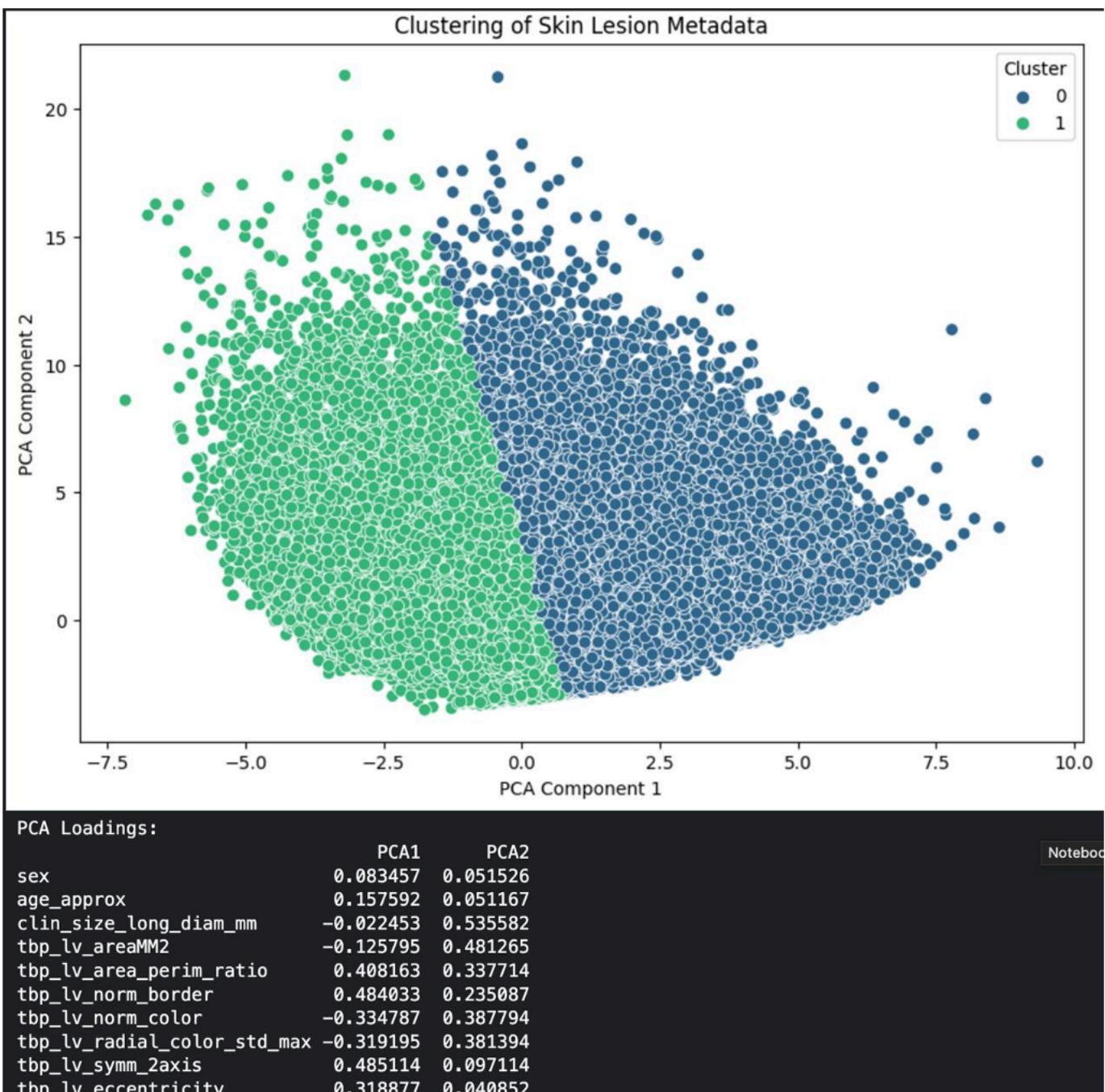


Figure 5: Dimensionality Reduction and K-Means Clustering on the Metadata

Accuracy: 0.6285  
Precision: 0.0010  
Recall: 0.3648  
F1 Score: 0.0019

Classification Report:

	precision	recall	f1-score	support
0	1.00	0.63	0.77	387274
1	0.00	0.36	0.00	381
accuracy			0.63	387655
macro avg	0.50	0.50	0.39	387655
weighted avg	1.00	0.63	0.77	387655

Confusion Matrix

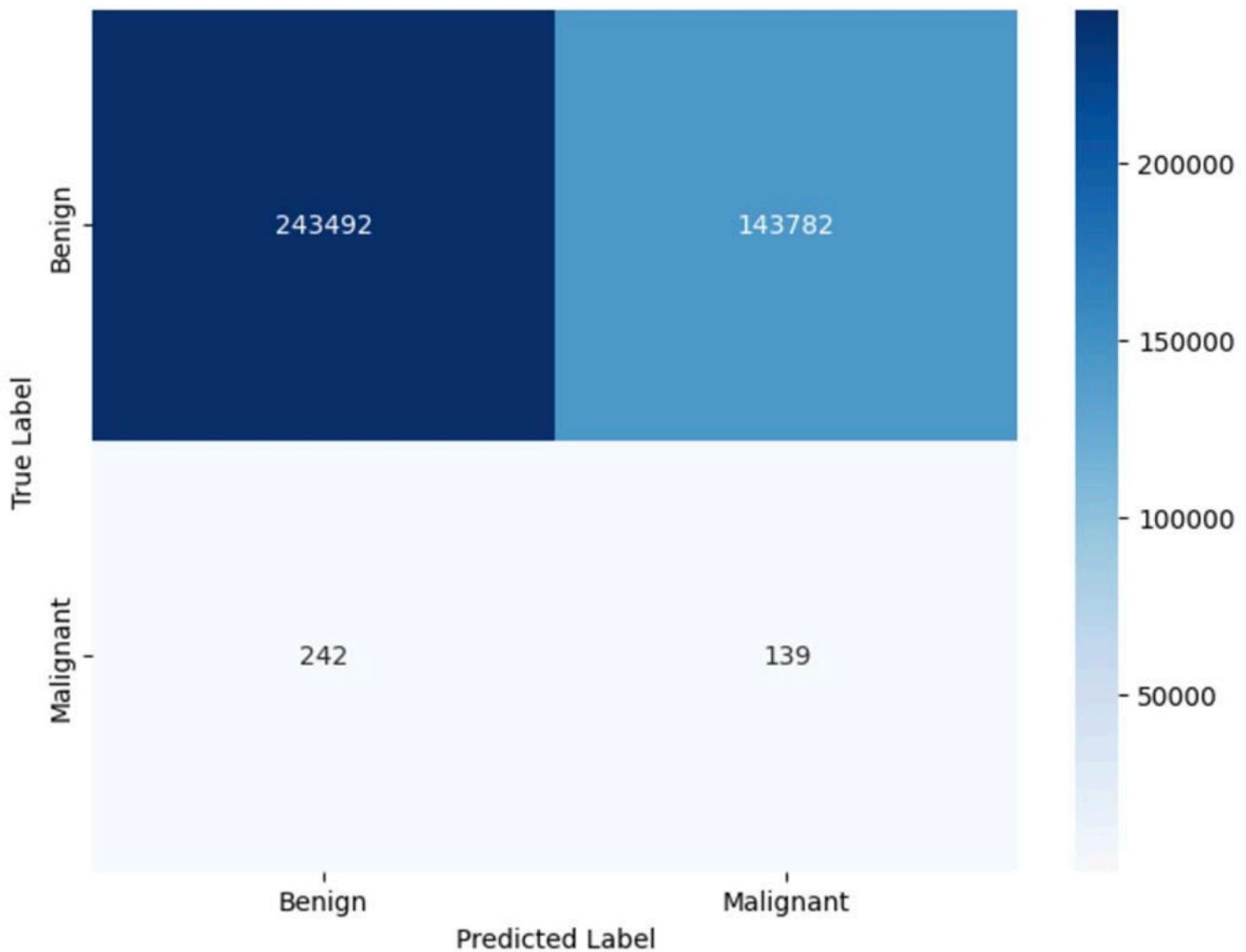
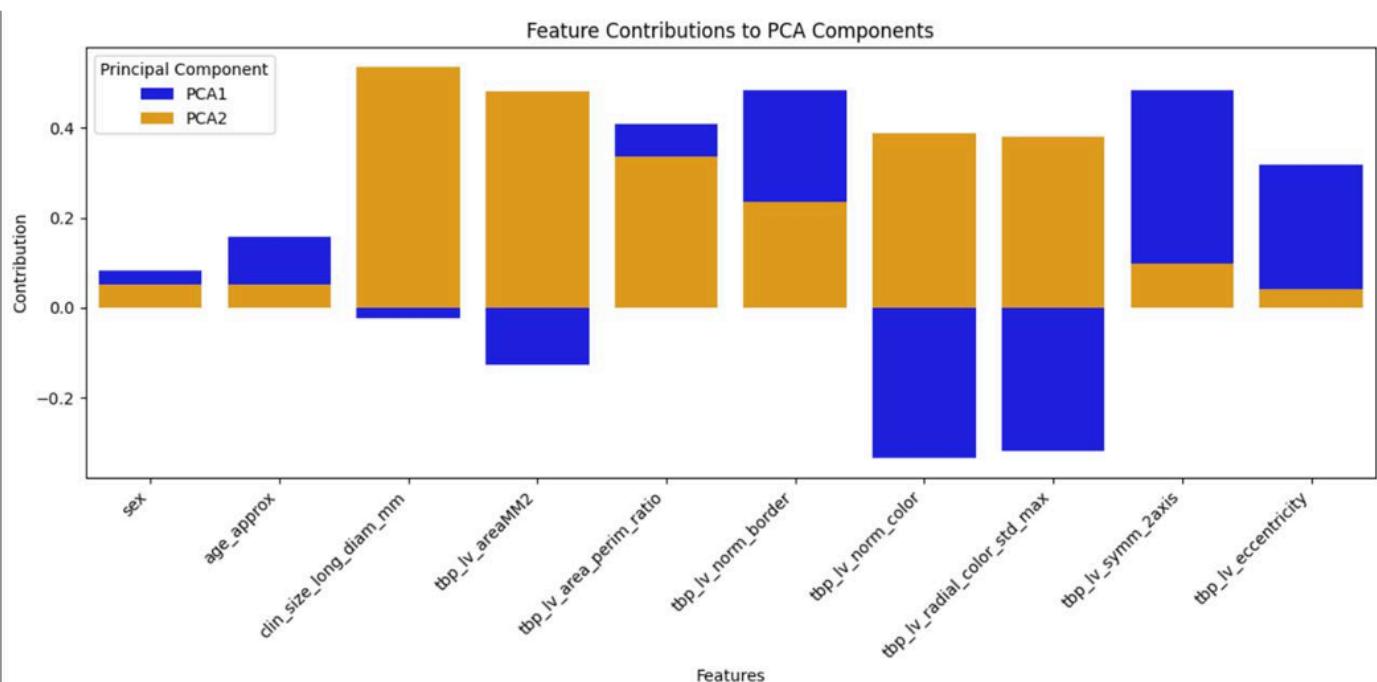


Figure 6: Classification Report and Confusion Matrix



**Figure 7: Feature Contributions to PCA Components**

The objective of this analysis was to evaluate the significance of metadata. We applied a K-Means to cluster the dataset based on metadata. We performed dimensionality reduction to project the data into two dimensions. This helped us visualize how the samples were distributed (Figure 5). However, despite using K-Means for clustering, no clear clusters emerged based on the metadata features alone. This lack of clear separation suggested that the metadata were not strongly indicative of distinct groupings in the data. As a result, clustering on the metadata did not provide any meaningful insights into the data's structure.

Figure 6 shows the classification report where the model achieved an accuracy of 62.8%. While this accuracy is moderate, a closer look at the precision and recall metrics revealed a significant issue. The model's performance was heavily skewed towards benign images, with much higher precision and recall for benign samples compared to malignant ones. This suggests that the model may be biased or facing challenges due to an imbalance in the data.

Furthermore, figure 7 shows the feature contributions table which revealed that patient-specific and non-image related data were not important features for explaining the variance in the data. These metadata features showed minimal impact on clustering or classification. In contrast, image-based features accounted for a much larger portion of the variance, suggesting that image data was far more relevant for distinguishing meaningful patterns in the dataset.

Overall, the analysis emphasized that while metadata features might provide limited insight, the focus should remain on optimizing and refining the image-based features to improve classification performance.

## Supervised Learning

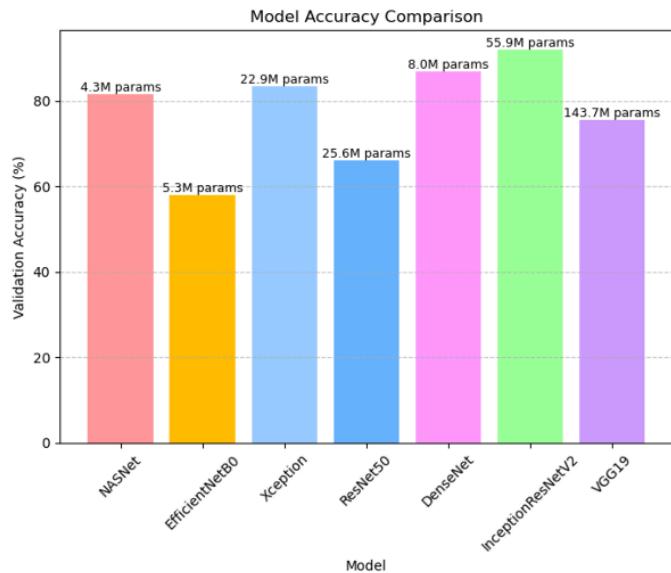


Figure 8: Model Accuracy Comparison

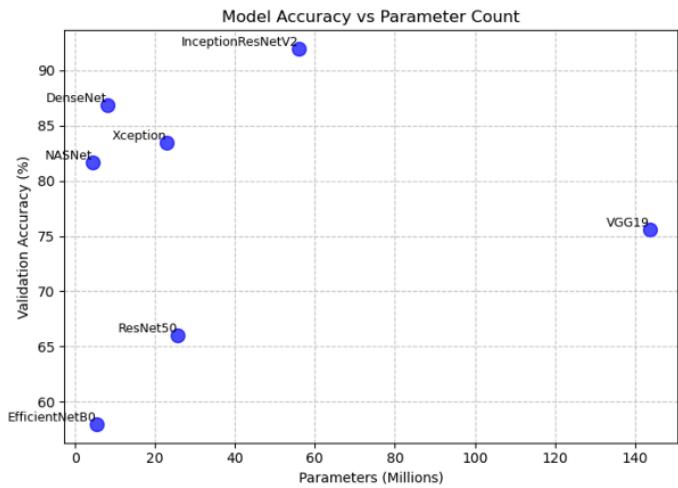


Figure 9: Model Accuracy vs Parameter Count

On comparing model accuracy (Figure 8), InceptionResNetV2 achieved the highest accuracy among the tested models at 91.93%. However, the tradeoff here is that InceptionResNetV2 has 55.9M parameters thereby making it computationally intensive and less suited for situations of limited memory and processing power. DenseNet and Xception perform well with accuracies of 86.87% and 83.47% respectively. However, there is a large difference in the number of parameters between the two models with DenseNet having only 8M parameters compared to Xception's 22.9M parameters and InceptionResNetV2's 55.9M parameters. VGG19 on the other hand achieved an accuracy of 75.57%. This score is significantly lower than DenseNet and InceptionResNetV2, despite its large parameter size of 143.7M. Hence, while optimizing for efficiency (Figure 9) in parameter size might not always be a good idea - as demonstrated by EfficientNet which has an accuracy of ~57% with 5.3M parameters, DenseNet provides the best accuracy-efficiency trade-off by having the highest accuracy for a relatively smaller number of parameters. However, InceptionResNetV2's highest accuracy makes it the most promising model for tasks where the performance is the top priority.

<b>Model</b>	<b>Parameters</b>	<b>Accuracy</b>
NasNet	4.3M	81.63%
EfficientNetB0	5.3M	58%
Xception	22.9M	83.48%
ResNet50	25.6M	66.02%
DenseNet	8M	86.88%
InceptionResNetV2	55.9M	91.93%
VGG 19	143.7M	75.57%

Model Descriptive Metrics

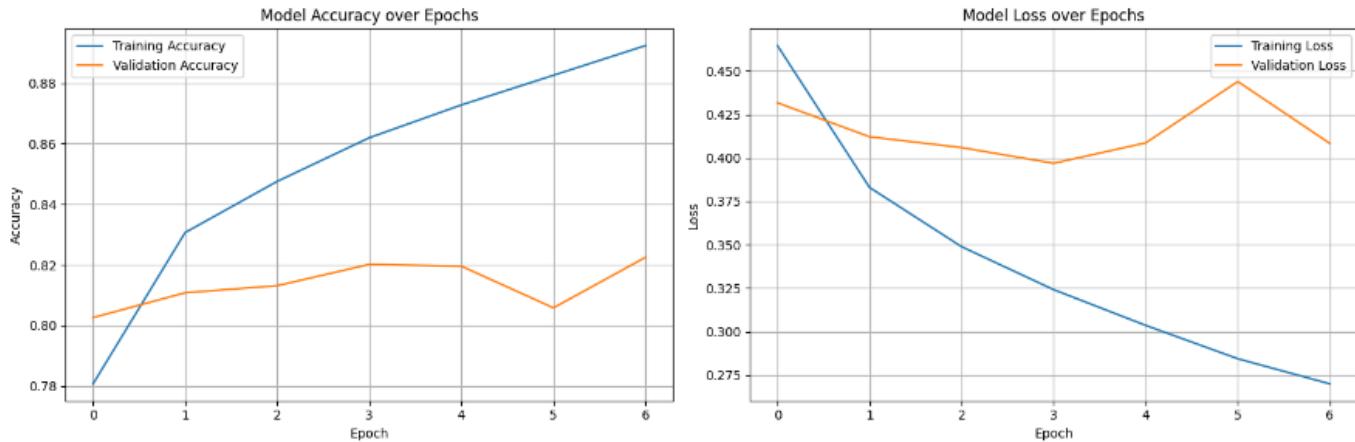


Figure 10A: NasNet Model Accuracy over Epochs

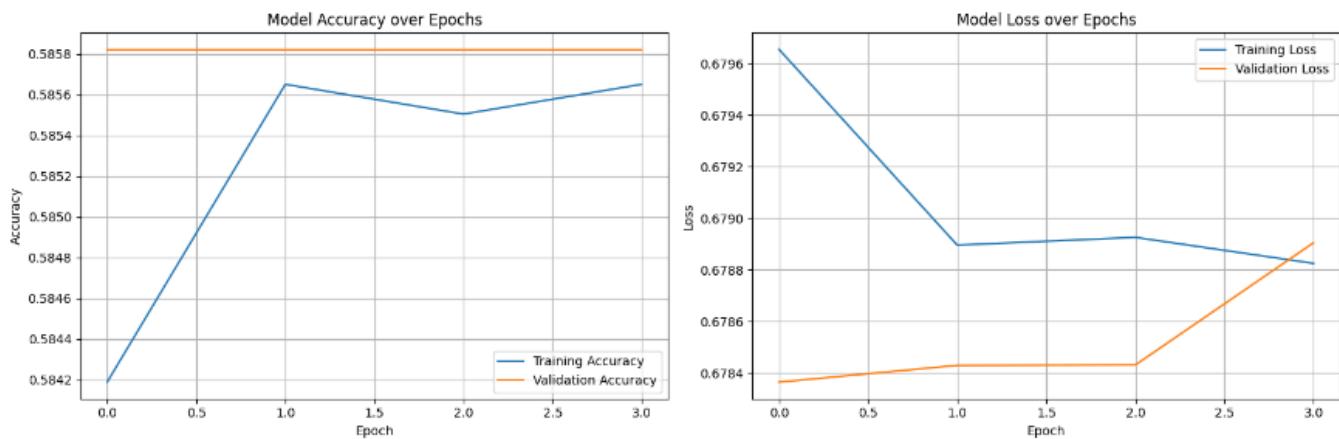


Figure 10B: EfficientNetBO Model Accuracy over Epochs

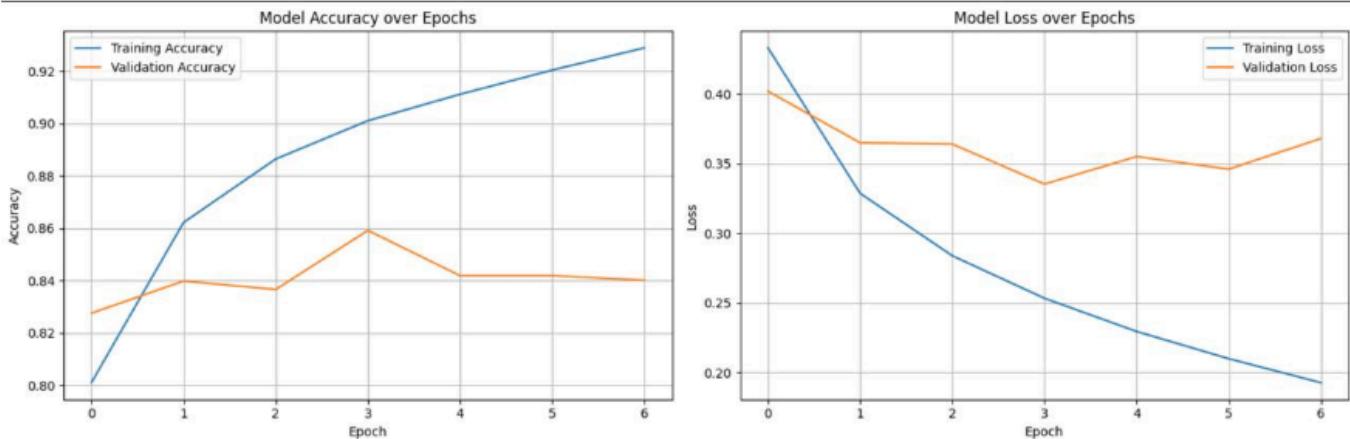
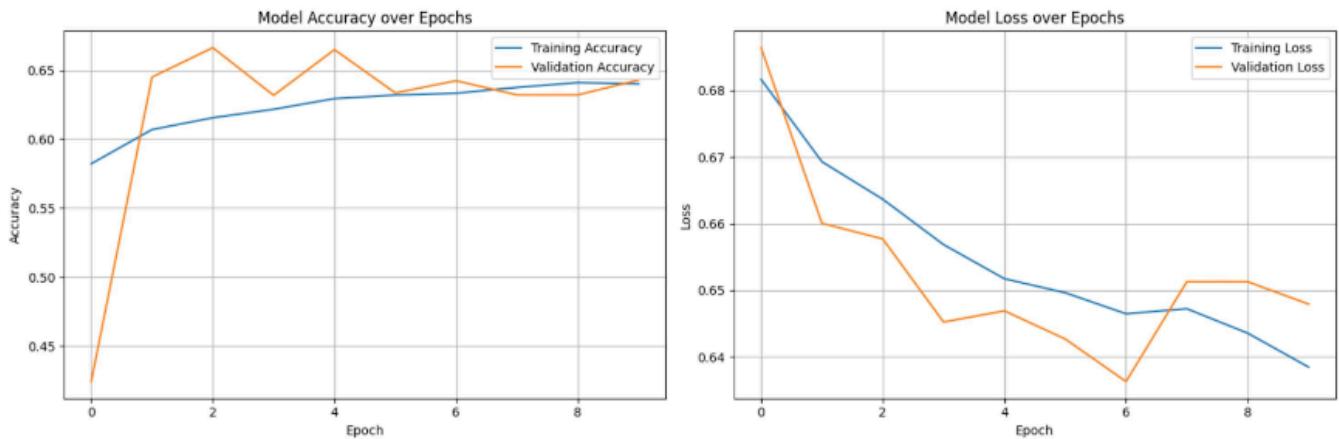


Figure 10C: Xception Model Accuracy over Epochs



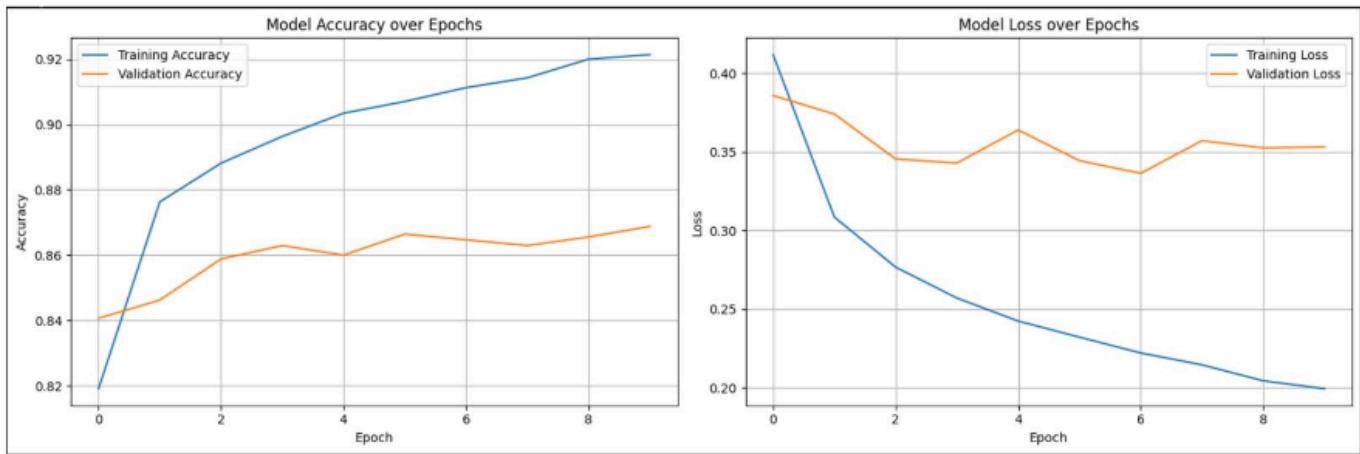
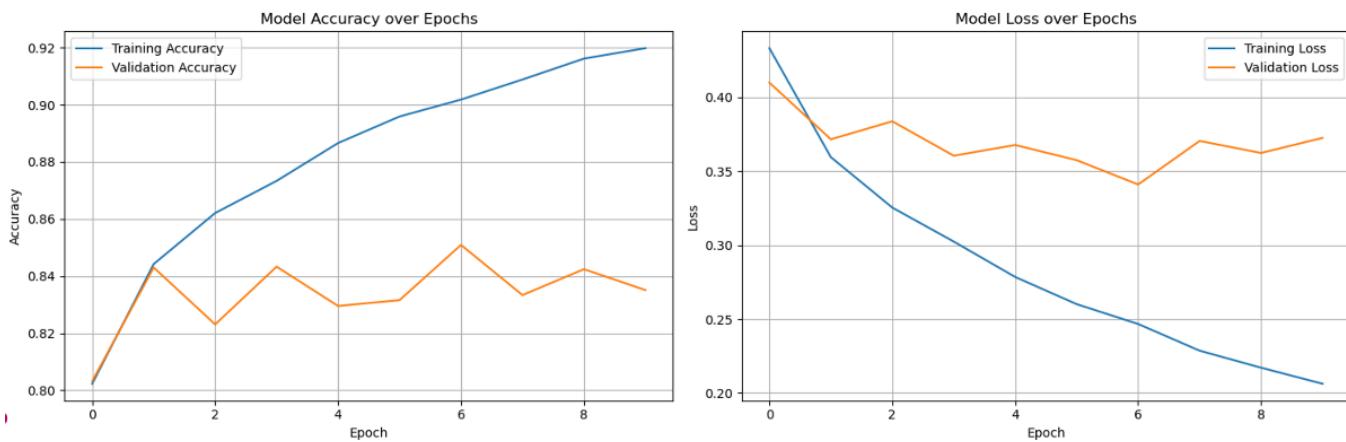
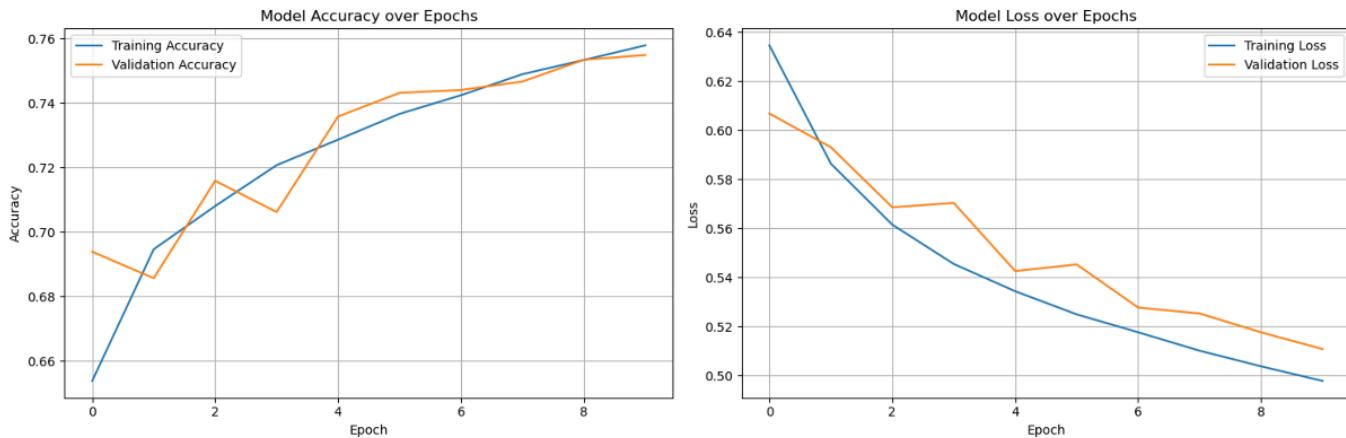
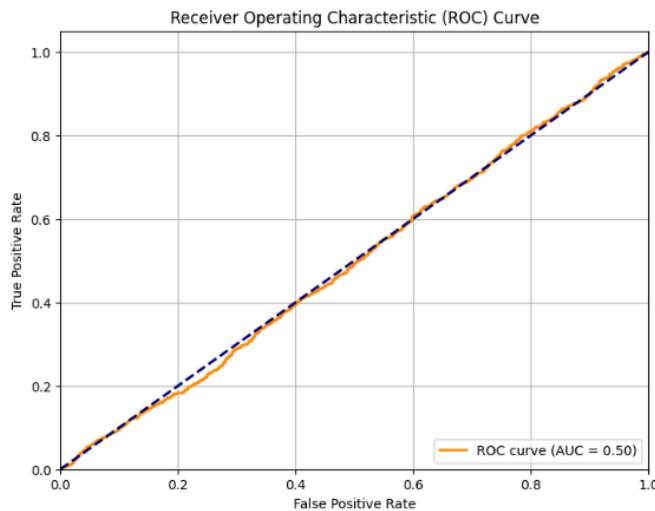
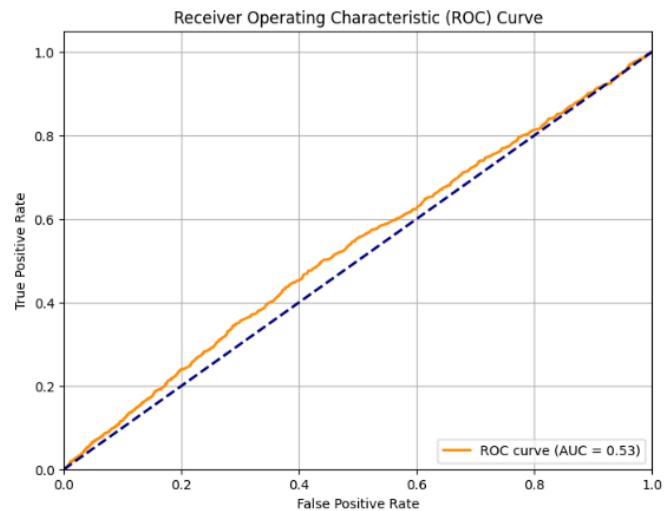
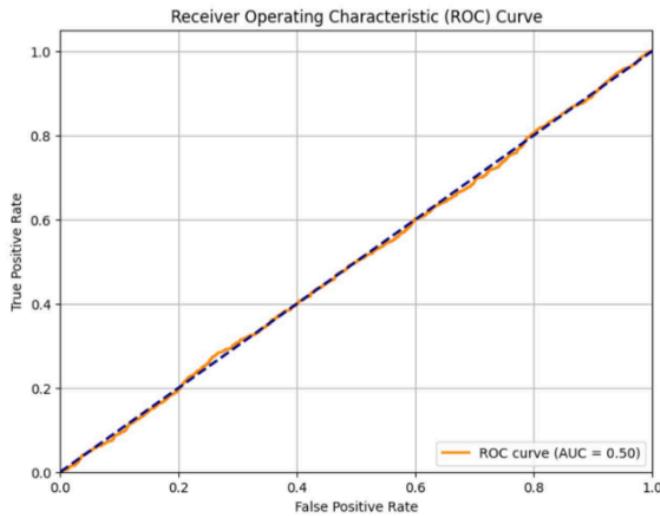
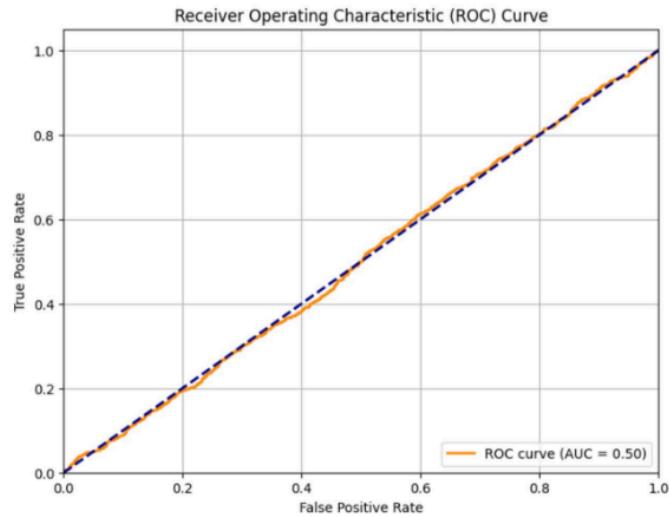
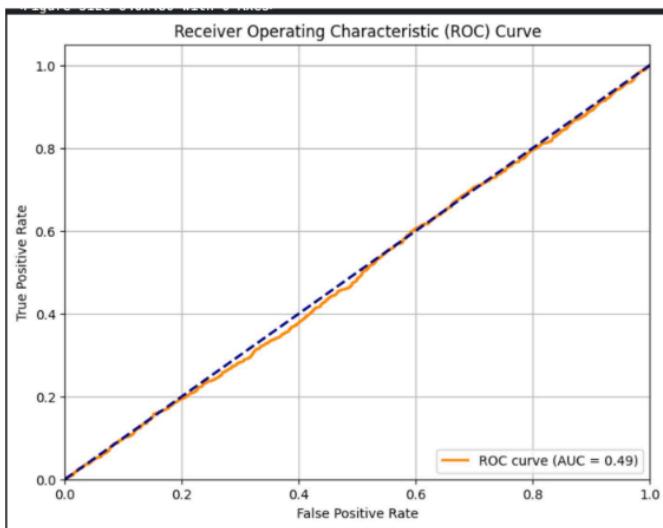
*Figure 10D: ResNet50 Model Accuracy over Epochs**Figure 10E: DenseNet Model Accuracy over Epochs**Figure 10F: InceptionResnetV2 Model**Figure 10G: VGG19 Model*

Figure 10 shows the progress of validation loss and accuracy with the number of epochs. For some models, we notice an increase in validation loss and decrease in validation accuracy while the same metrics improve in the opposite direction for the training dataset. This has been identified as a sign of potential overfitting. However, InceptionResNetV2 and VGG19 do not follow this trend but rather show either plateauing or slight improvements in validation loss and accuracy for all epochs during training.

Thus, these architectures may be best at generalizing, given their construction be better suited to adequately capturing complex feature representations without overfitting them to the training data.

*Figure 11A: NasNet ROC Curve**Figure 11B: EfficientNetBO ROC Curve**Figure 11C: Xception ROC Curve**Figure 11D: ResNet50 ROC Curve**Figure 11E: DenseNet ROC Curve*

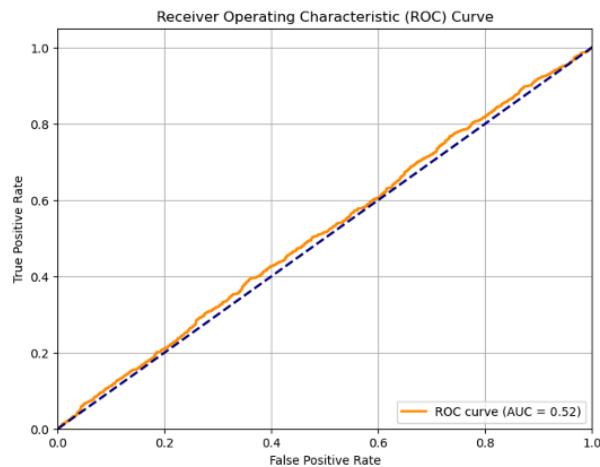


Figure 11F: InceptionResNetV2 ROC Curve

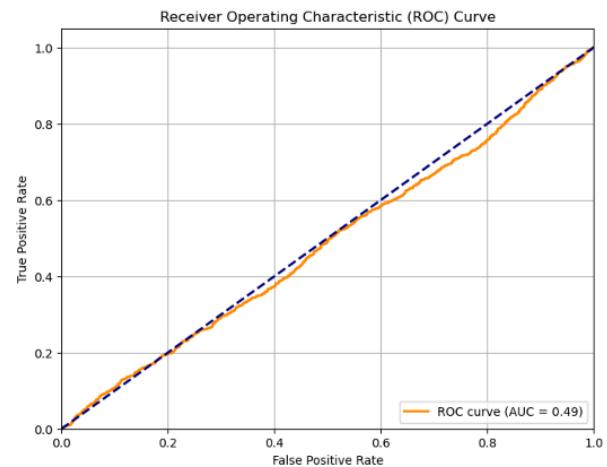


Figure 11G: VGG19 ROC Curve

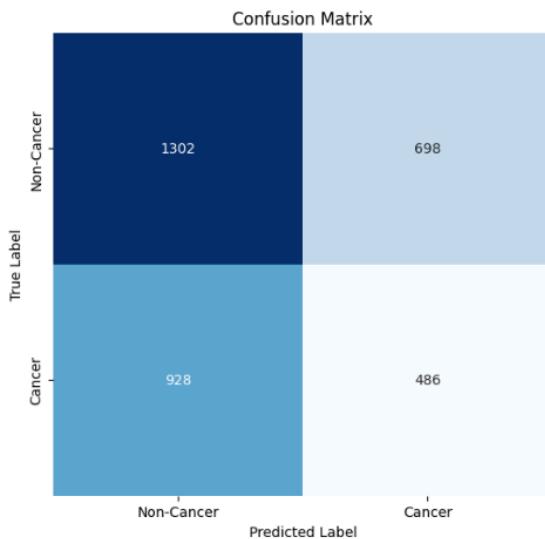


Figure 12A: NasNet Confusion Matrix

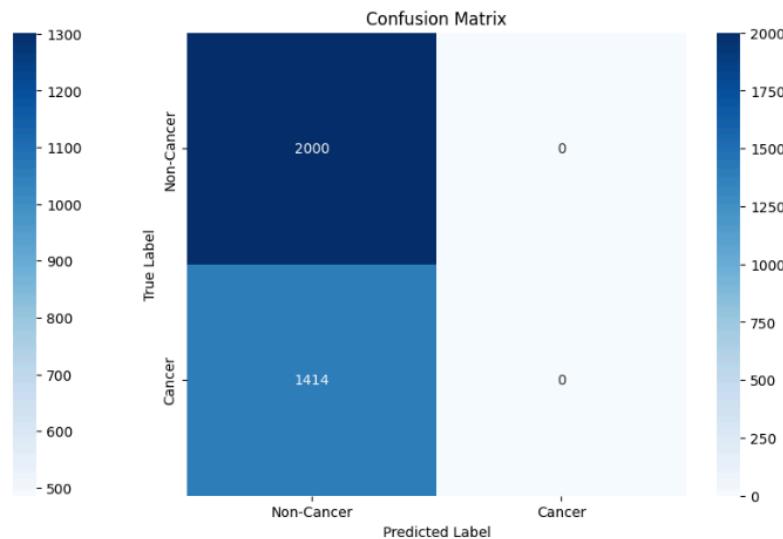


Figure 12B: EfficientNetBO Confusion Matrix

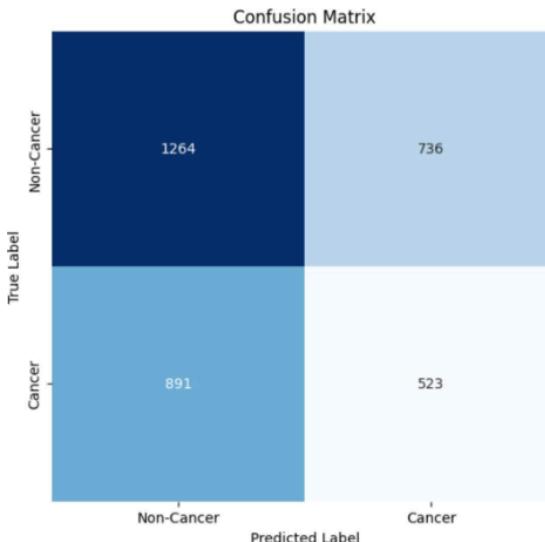


Figure 12C: Xception Confusion Matrix

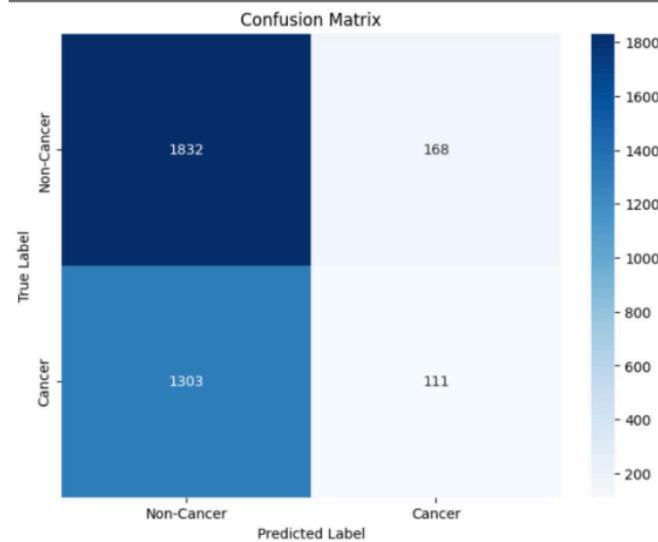


Figure 12D: ResNet50 Confusion Matrix

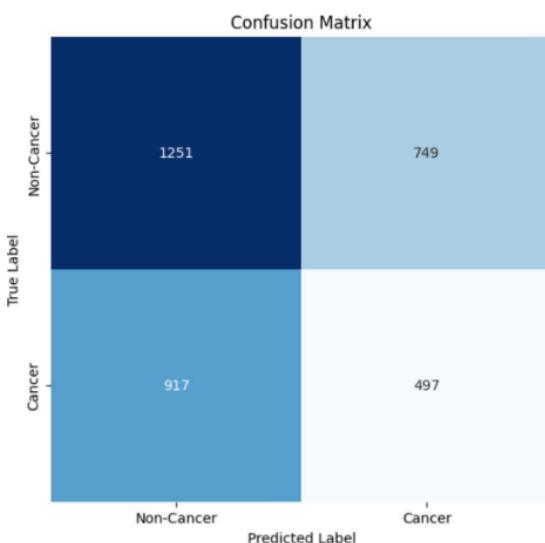


Figure 12E: DenseNet Confusion Matrix

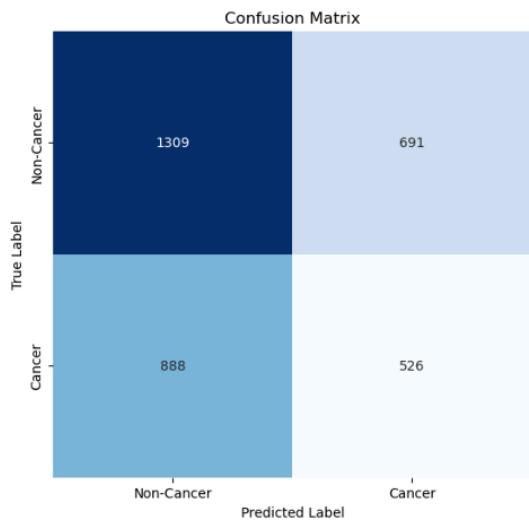


Figure 12F: InceptionResNetV2 Confusion Matrix

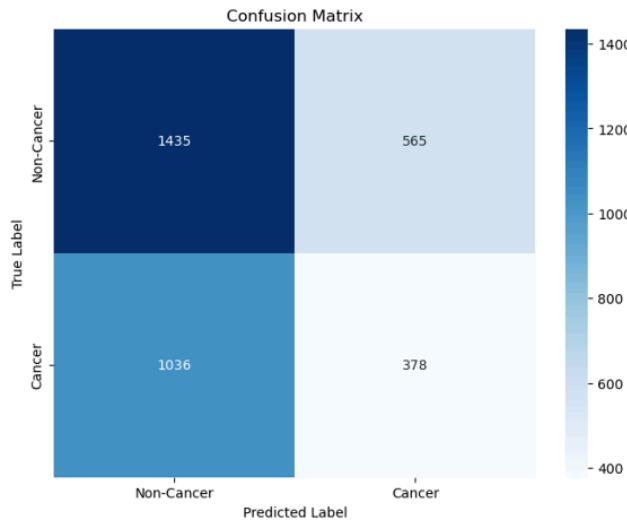


Figure 12G: VGG19 Confusion Matrix

Accuracy and loss are not sufficient in characterizing model performance in imbalance datasets. Further, in diagnostic settings, it is important to ensure good precision and recall due to the wide-ranging implication of mis-diagnosis. Being diagnosed with cancer causes severe mental and emotional distress - hence a low precision or high false positive rates can severely impact the well-being of the misdiagnosed patients. As mentioned earlier, timely diagnosis is the need of the hour, hence it is very important to reduce the number of false negatives and significantly improve our recall. As demonstrated by the AUC (Figure 11) for each model, even the models with high accuracy, have concerningly low AUC - this may be due to the heavy class imbalance in the dataset. The confusion matrix (Figure 12) also clearly indicates highly imprecise predictions. EfficientNetB0, for example, classifies everything as Benign - a clear sign of class imbalance. Even our model of choice - DenseNet categorizes a majority of Malignant Tumors as Benign (64.85%) and categorizes 37.45% of benign tumors as malignant. InceptionResNetV2 reflects more balanced results with increased precision and recall while compared with the other models; hence it is better at class imbalance handling as evident in its confusion matrix. In this context, InceptionResNetV2 has been selected as the primary model for our supervised system, although it is more demanding in terms of computation. We would prioritize precision and recall in the diagnosis of cancer over memory limitation. Even this model, however, is not without its limitations, as the imbalance in the dataset continues to affect its generalization efficiency.

These findings highlight the importance of going beyond accuracy and loss to evaluate models particularly for critical applications like cancer diagnosis. Robust metrics like AUC, precision, recall, and confusion matrices are essential tools for understanding model performance and identifying areas of improvement.

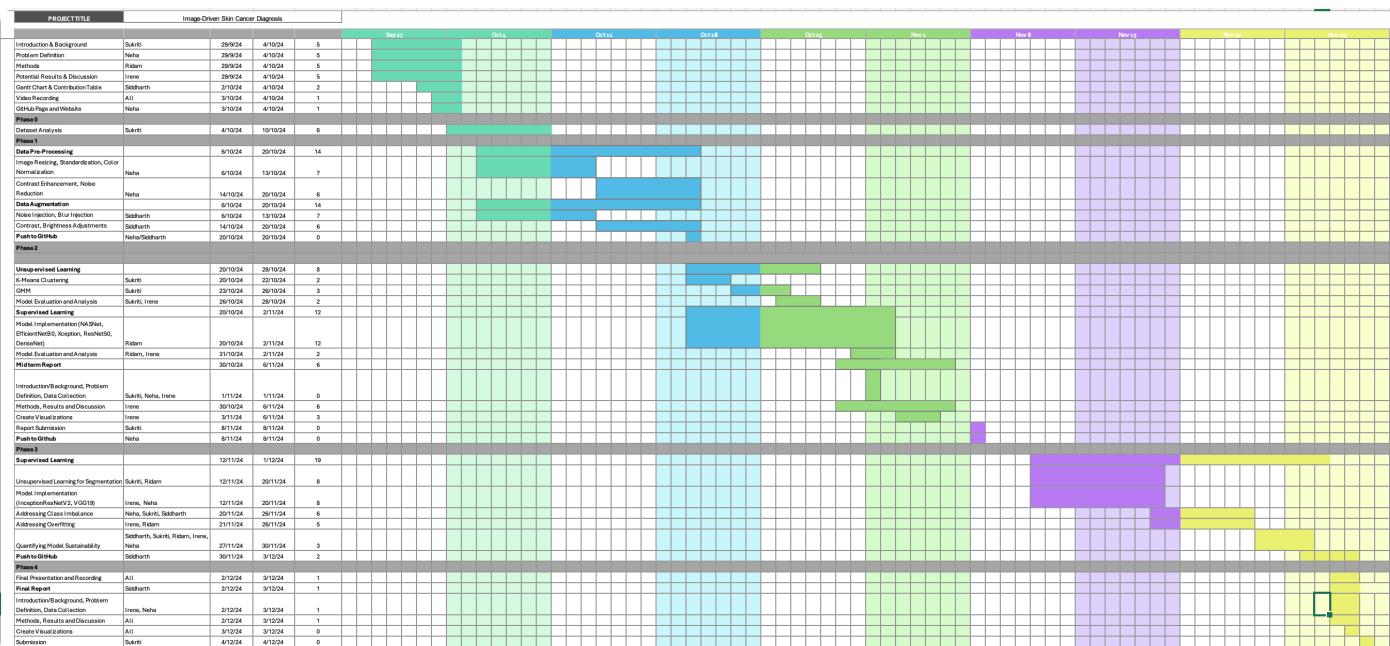
## Next Steps

To summarize, next steps include:

Addressing overfitting through strategies like early stopping based on validation loss, regularization or further data augmentation.

Focus on cost sensitive learning to improve precision and recall.

## Gantt Chart



## Contribution Table

## Contribution Table

PROJECT TITLE	Image-Driven Skin Cancer Diagnosis
MEMBER	FINAL CONTRIBUTIONS
Irene	Supervised Models, Analysis, Report, Visualizations, Website, Video
Neha	Data Preprocessing, Supervised Models, Analysis, Results, Website, Github
Ridam	Supervised Models, Unsupervised Models, Analysis, Video
Siddharth	Data Preprocessing, Report, Visualizations, Website
Sukriti	Unsupervised Models, Analysis, Website, Github, Video

All members of the team contributed equally to this project.

## References

---

- [1] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun, "Dermatologist-level classification of skin cancer with deep neural networks," *Nature*, vol. 542, no. 7639, pp. 115–118, 2017.
- [2] M. Naqvi, S. Q. Gilani, T. Syed, O. Marques, and H.-C. Kim, "Skin cancer detection using deep learning—a review," *Diagnostics*, vol. 13, no. 11, p. 1911, 2023.
- [3] S. S. Han, M. S. Kim, W. Lim, G. H. Park, I. Park, and S. E. Chang, "Classification of the clinical images for benign and malignant cutaneous tumors using a deep learning algorithm," *Journal of Investigative Dermatology*, vol. 138, no. 7, pp. 1529–1538, 2018.
- [4] Y. Liu et al., "A deep learning system for differential diagnosis of skin diseases," *Nature Medicine*, vol. 26, no. 6, pp. 900–908, 2020.
- [5] J. Savulescu, A. Giubilini, R. Vandersluis, and A. Mishra, "Ethics of artificial intelligence in medicine," *Singapore Medical Journal*, vol. 65, no. 3, p. 150, 2024.

\* Click on images to expand.