



# DeepInsight:

## Navigating the Realms of Deepfake Detection

Barnana Ganguly (bg29596)

Abhijit Anil (aa94375)

Amey Ghate (ag84957)

Rathil Madihalli (rm63782)

Sankalp Kulkarni (sk57727)

Shrishty Mishra (sm82386)

# Table of Contents:

<b>1. Introduction to Deepfakes</b>	<b>3</b>
<b>2. How are they created?</b>	<b>4</b>
<b>3. Why are DeepFakes a matter of concern?</b>	<b>6</b>
<b>4. Methodology</b>	<b>8</b>
4.1 The DenseNet Model	9
4.1.1 DenseNet Model (Base Model)	11
4.1.2 DenseNet Model with Image augmentations	12
4.1.3 DenseNet Model with Grayscale Images	13
4.2 Custom CNN	14
4.3 Other Models	17
<b>5. Metrics</b>	<b>17</b>
<b>6. Results</b>	<b>19</b>
6.1 DenseNet Model (Base Model)	19
6.2 DenseNet Model with Image Augmentations	20
6.3 DenseNet Model with Grayscale Images	21
6.4 Custom CNN	22
<b>7. Comparison of Results</b>	<b>23</b>
<b>8. Performance Demo</b>	<b>26</b>
<b>9. Conclusion</b>	<b>27</b>
<b>10. Future Work/ Next Steps</b>	<b>28</b>
<b>11. References</b>	<b>29</b>



# 1. Introduction to Deepfakes

A DeepFake is a synthetic media in which a person in an existing image or video is replaced with someone else's likeness using artificial intelligence and machine learning technologies. The term "DeepFake" is a portmanteau of "deep learning" and "fake," indicating its reliance on deep learning technology. DeepFakes leverage advanced AI algorithms to manipulate or generate visual and audio content with a high potential to deceive.



*GIF 1.1: Let's Talk DeepFake!*

The technology is sophisticated enough to make it appear that individuals are saying or doing things they never did, making it both impressive and controversial.

## 2. How are they created?

Generative Adversarial Networks (GANs) are at the heart of DeepFake technology. Introduced by Ian Goodfellow and his colleagues in 2014, GANs are a class of AI algorithms used in unsupervised machine learning tasks.

The fundamental architecture of GANs is built upon the interaction between two distinct neural networks: the generator and the discriminator. This interaction is essentially a game, akin to a cat-and-mouse dynamic, where both networks compete against each other, leading to continuous improvements in the quality of the generated outputs.

They are composed of two neural networks, the generator and the discriminator, which are trained simultaneously through a competitive process.

1. **Generator:** The generator's role is to learn to create data that mimics some distribution of interest. It takes random noise as input and generates samples as close as possible to the real data. Initially, the generated samples are likely to be of poor quality, but the goal is for these samples to become indistinguishable from real data through the training process.
2. **Discriminator:** The discriminator acts as a judge, attempting to differentiate between genuine data (drawn from the training set) and fake data produced by the generator. It is essentially a binary classification model that outputs the probability of a given sample being real.

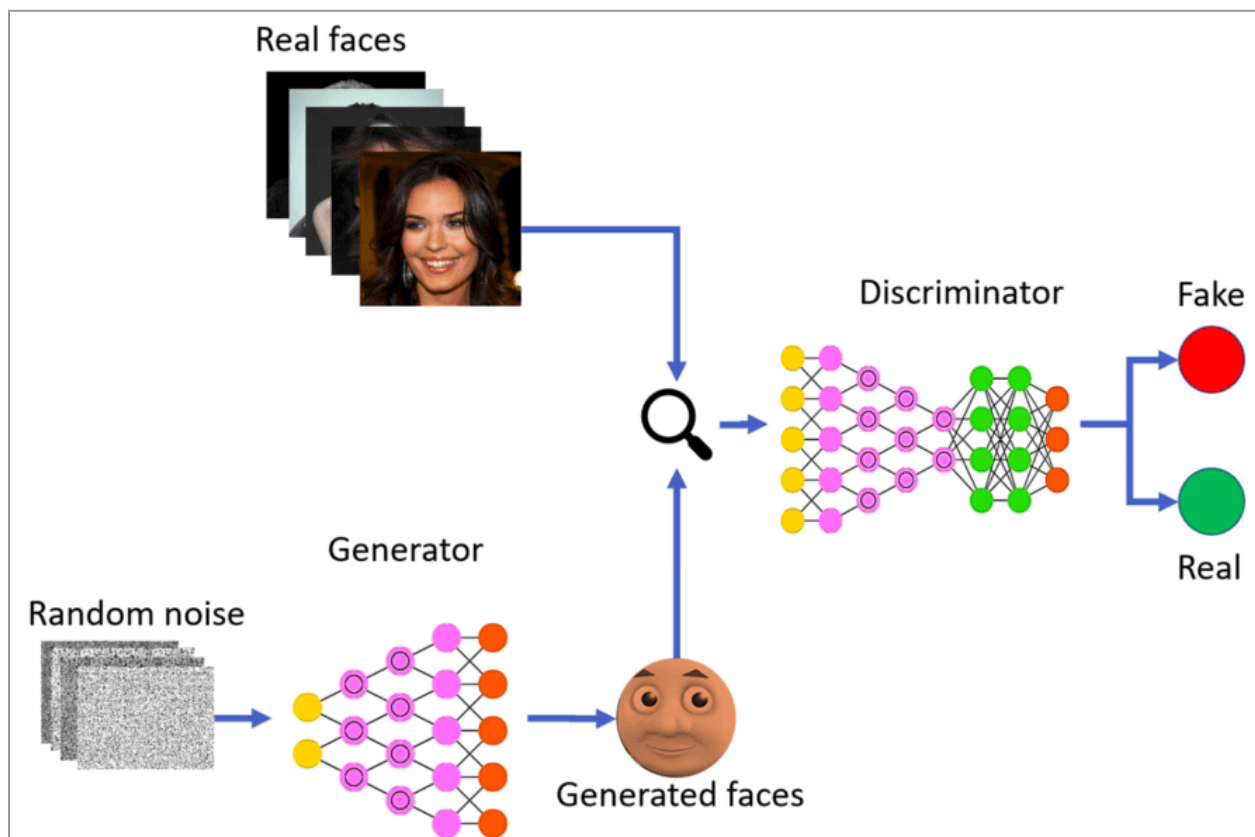


Figure 2.1: GAN

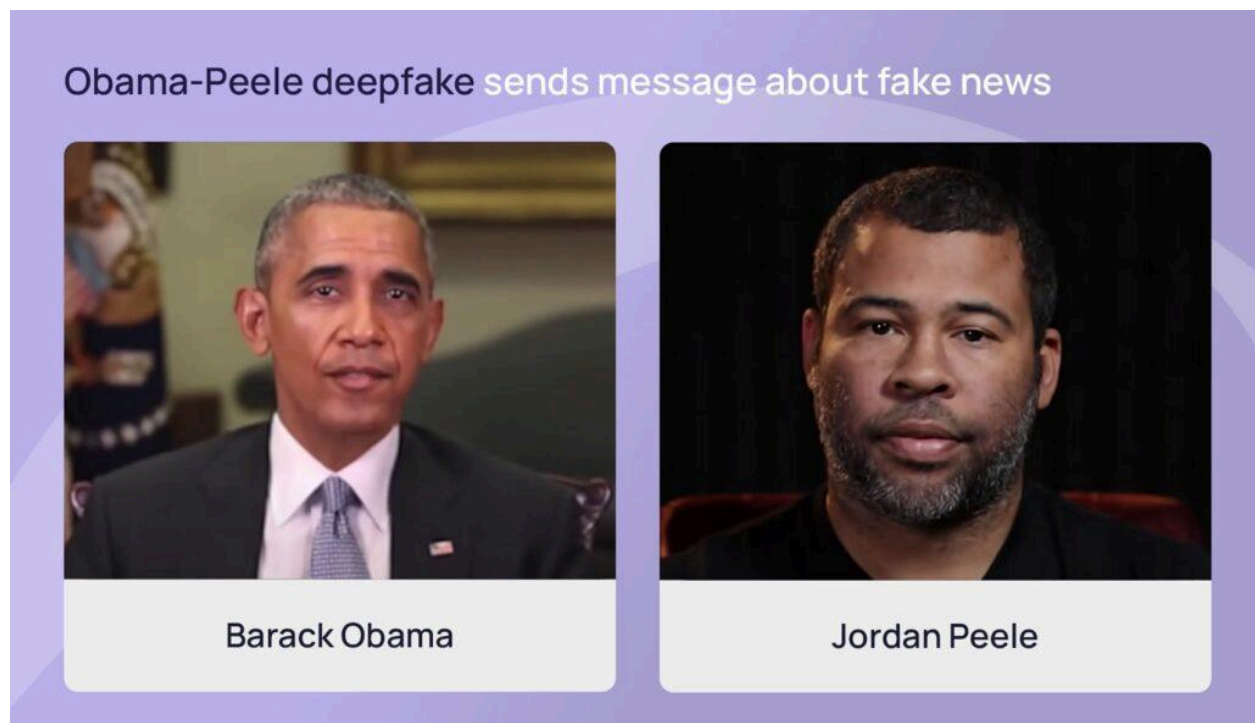
The theoretical foundation of GANs is rooted in a zero-sum game framework, where the generator's goal is to maximize the probability of the discriminator making a mistake. This is formally represented as a min-max game with a value function  $V(G, D)$  that both networks aim to optimize: the generator tries to minimize this function, while the discriminator tries to maximize it.

The interplay between these two networks drives them to improve continuously: the generator learns to produce more accurate fakes, while the discriminator becomes better at detecting them. This adversarial process enhances the quality of the generated outputs, making GANs incredibly effective for creating realistic DeepFakes.

GANs are not limited to creating DeepFakes; they are also used in various applications, including image generation, style transfer, image super-resolution, and more, showcasing their versatility and power in the field of AI and machine learning.

### 3. Why are DeepFakes a matter of concern?

DeepFakes represent a significant technological advancement but come with substantial risks and challenges that impact individuals, societies, and global stability. The core issue with DeepFakes lies in their ability to craft believable but entirely false realities, leading to deception, manipulation, and harm.



*Figure 3.1: DeepFakes Cause Concern*

Here are the principal concerns associated with DeepFakes:

**Misinformation and Fake News:** DeepFakes can propagate misinformation and fake news by producing highly realistic, yet false, videos or audio recordings. This enables malicious actors to fabricate events or statements, significantly manipulating public opinion, affecting elections, and eroding trust in media and public figures.

**Political Manipulation:** By creating counterfeit speeches or compromising visuals of politicians, DeepFakes pose a severe threat to the political sphere, capable of tarnishing reputations, influencing voter sentiments, and unsettling democratic processes.

**Personal and Social Harm:** Individuals can suffer greatly from DeepFakes that falsely portray them in undesirable situations, causing emotional distress, reputational damage, and social ostracization. This is particularly egregious in instances of non-consensual DeepFake pornography.

**Legal and Ethical Challenges:** The rise of DeepFakes prompts complex legal and ethical issues around consent, privacy, and free speech, creating a regulatory vacuum. As the technology advances, pinpointing liability and authenticity becomes increasingly challenging.

**Erosion of Trust:** DeepFakes contribute to a significant decline in trust towards digital content, complicating the ability to differentiate real from fake. This skepticism weakens the foundation of informed decision-making and societal unity, fostering cynicism and disconnection.

**Security Concerns:** DeepFakes also raises security alarms, particularly in fraud and identity verification contexts. The precision in replicating someone's appearance and voice can facilitate sophisticated phishing schemes and security breaches, undermining digital verification trust and posing risks to both individuals and organizations.

The introduction of DeepFake technology marks a new era with vast potential yet profound challenges. Despite its benefits, the misuse of DeepFakes demands immediate and comprehensive action from policymakers, technologists, and the community at large. Tackling the issues posed by DeepFakes will require an integrated approach, encompassing technological innovations, legal reforms, ethical standards, and increased public awareness to protect rights and preserve societal trust.

## 4. Methodology

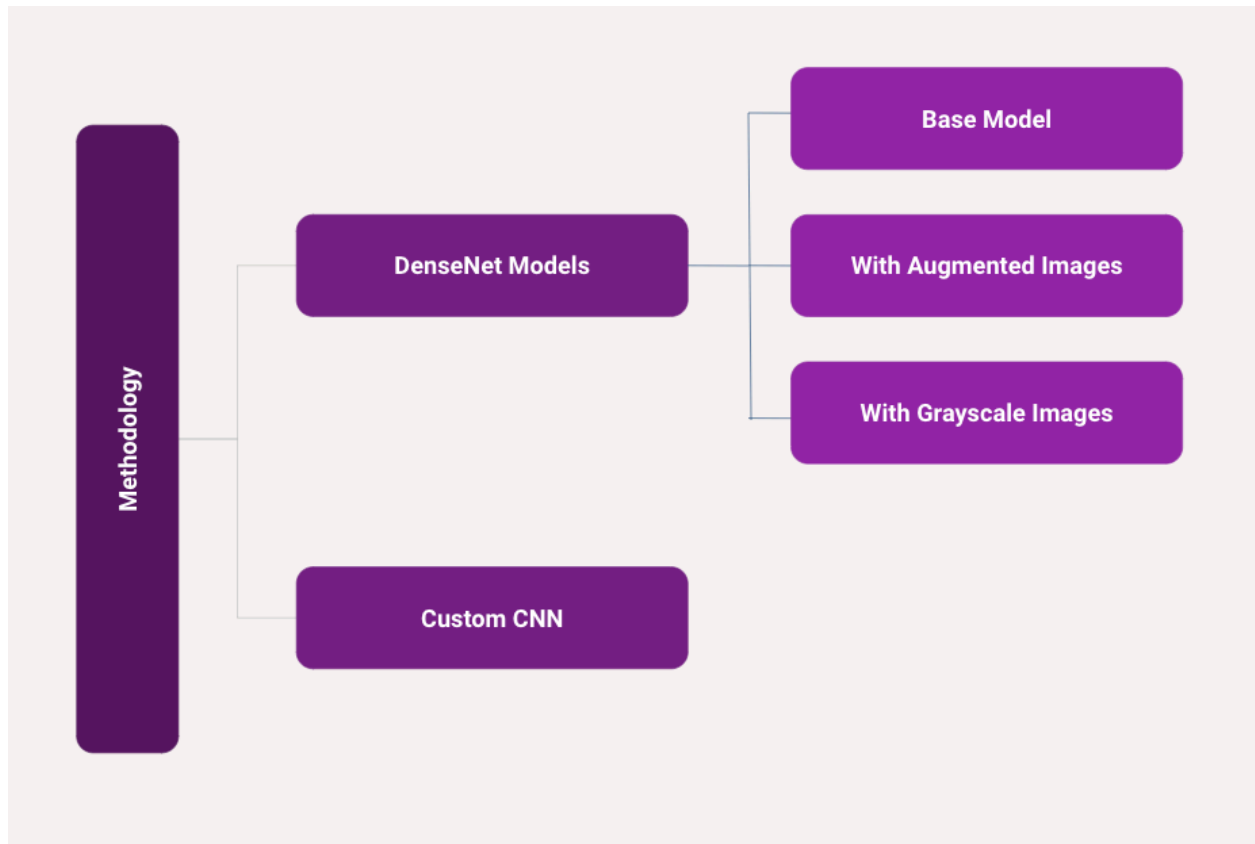


Figure 4.0.1: Methodology Flowchart

In this section, we delve into the comprehensive methodology adopted to develop and evaluate a series of models aimed at identifying DeepFake content. The overarching strategy encompasses the deployment of Dense Convolutional Network (DenseNet) architectures and a custom Convolutional Neural Network (CNN) model. Each model was meticulously crafted to address specific facets of the DeepFake detection challenge, leveraging unique characteristics and modifications to enhance performance and generalization capabilities across varied datasets.

The DenseNet models, known for their efficiency in feature extraction and depth, form the cornerstone of our approach. We explore three distinct variations of the DenseNet architecture to understand the impact of different enhancements and adjustments on the model's ability to distinguish between real and manipulated images. These variations include a base DenseNet model emphasizing raw feature extraction, an augmented DenseNet model incorporating data augmentation techniques to improve robustness, and a grayscale-focused DenseNet model exploring the influence of color information on detection accuracy.

Complementing the DenseNet architectures, we introduce a custom CNN model tailored to the specifics of DeepFake detection. This model amalgamates convolutional layers, pooling, batch

normalization, dropout, and global average pooling to construct a nuanced feature extraction and classification pipeline. The design philosophy behind this custom CNN emphasizes adaptability, efficiency, and the mitigation of overfitting, with a keen focus on achieving high accuracy in identifying DeepFake content.

## 4.1 The DenseNet Model

DenseNet-121, or Dense Convolutional Network with 121 layers, is a deep learning model that excels in image recognition tasks due to its unique architecture within the DenseNet family. It features dense connectivity, meaning each layer receives inputs from all preceding layers, which promotes maximum information flow and mitigates the risk of information loss as data passes through the network.

The DenseNet architecture is divided into dense blocks and transition layers.

Dense blocks consist of layers that are densely connected, and the transition layers reduce the feature map dimensions.

Composite functions including Batch Normalization, ReLU activation, and Convolution are followed by a pooling layer to decrease feature map sizes, enhancing computational efficiency.

We chose to use DenseNet because of its many advantages over traditional CNNs:

1. **Parameter Efficiency:** DenseNet requires fewer parameters than traditional CNNs, which makes it more resource-efficient and reduces the risk of overfitting.
2. **Improved Feature Propagation:** The dense connectivity allows features from earlier layers to be reused in subsequent layers, enhancing the model's capability to learn detailed features.
3. **Feature Reuse:** The architecture encourages feature reuse across layers, increasing efficiency and decreasing the necessary number of parameters.
4. **Better Gradient Flow:** DenseNet's structure facilitates gradient propagation during the training phase, which eases the training process.



You might be wondering why, given ResNet's popularity and proven effectiveness in various deep learning tasks, we opted for DenseNet over ResNet for our specific application.

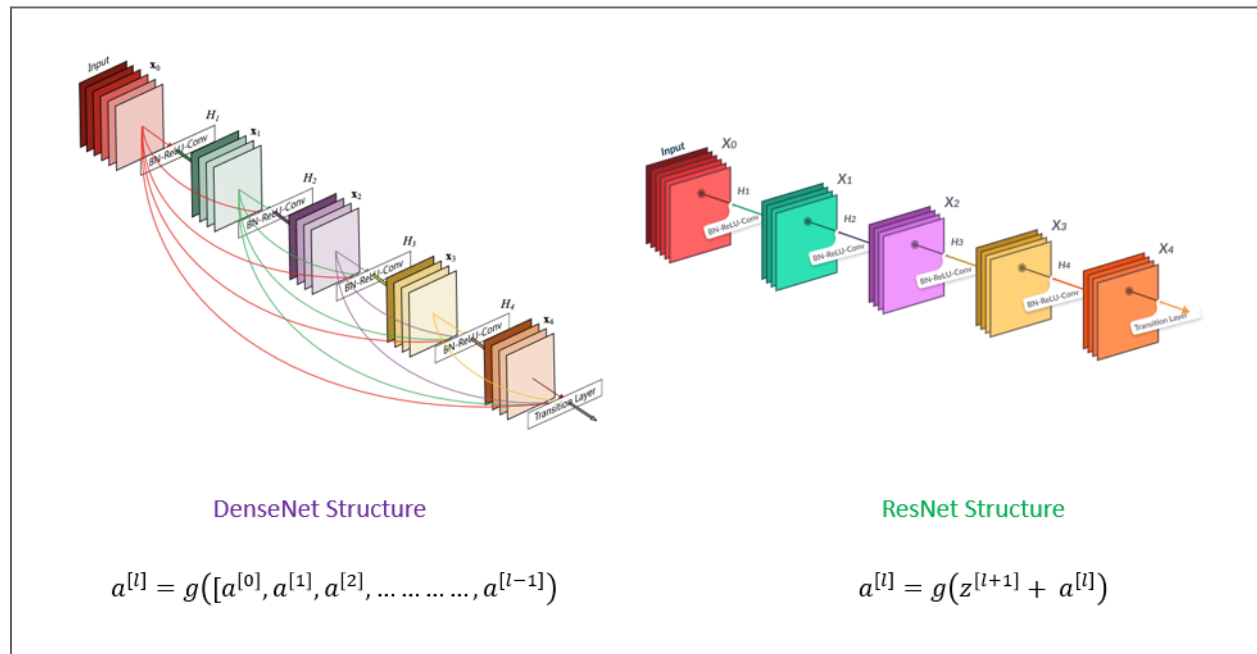


Figure 4.1.1: DeepNet vs ResNet

The choice of DenseNet over ResNet for our project was driven by a need for efficiency, effectiveness in handling diverse data types, and alignment with our specific goals.

DenseNet distinguishes itself from the popular ResNet by its architecture, which features direct connections between all layers, significantly reducing the model's parameter count while enhancing feature reuse. This not only facilitates better information flow and gradient propagation during training but also grants DenseNet superior adaptability to varied data complexities. Consequently, DenseNet emerges as the preferred option for our application, where the efficient use of computational resources is crucial. Its ability to achieve high performance with fewer parameters and its inherent flexibility offer clear advantages over the more parameter-heavy ResNet.

### 4.1.1 DenseNet Model (Base Model)

Model: "sequential"		
Layer (type)	Output Shape	Param #
=====		
densenet121 (Functional)	(None, 7, 7, 1024)	7037504
global_average_pooling2d (GlobalAveragePooling2D)	(None, 1024)	0
dense (Dense)	(None, 1)	1025
=====		
Total params: 7,038,529		
Trainable params: 6,954,881		
Non-trainable params: 83,648		

Figure 4.1.2: DeepNet Base Model Architecture

We developed a neural network-based model, capitalizing on the strengths of the DenseNet121 architecture, to address the challenge of DeepFake detection. This section outlines the core architectural components and strategies that define our model's capability to process and learn effectively.

At the heart of the model lies DenseNet121, selected for its impressive efficiency and depth as a feature extractor. Our approach deviates from typical usage by not employing pre-loaded weights, a decision that guarantees the extraction of features is directly pertinent to the specific task of discerning DeepFakes.

Subsequent to the DenseNet121 layer, we incorporated a Global Average Pooling 2D (GAP) layer. This critical component condenses each feature map into a single summary value. The employment of the GAP layer is crucial for distilling the feature information into a more digestible and focused form, highlighting the most salient features derived from the input images.

The culmination of our architecture is a dense output layer furnished with a single neuron. This layer is equipped with a sigmoid activation function designed to generate a probability score. This score is pivotal in the classification process, effectively discerning whether the input image is authentic or a DeepFake.

To ensure a balance between performance and computational demand, we trained our model with a batch size of 32 for 20 epochs, utilizing the Adam Optimizer to enhance the convergence rate. This careful calibration of training parameters and batch processing underscores our commitment to efficiency and performance in the model development process.

Our architectural strategy reflects a tailored approach, optimizing the DenseNet121's inherent strengths in feature extraction in tandem with the GAP layer for feature summarization, concluding with a precise output layer for classification. This methodical design is purpose-built for the specialized application of DeepFake detection.

### 4.1.2 DenseNet Model with Image augmentations

Model: "sequential"		
Layer (type)	Output Shape	Param #
=====		
densenet121 (Functional)	(None, 7, 7, 1024)	7037504
global_average_pooling2d (GlobalAveragePooling2D)	(None, 1024)	0
dense (Dense)	(None, 1)	1025
=====		
Total params: 7,038,529		
Trainable params: 6,954,881		
Non-trainable params: 83,648		
=====		

Figure 4.1.3: DeepNet Model Architecture (Image Augmentations)

In our development process, a crucial step was the implementation of data augmentation techniques to increase the robustness and generalization capabilities of our DenseNet-based model. This augmentation plays a significant role in ensuring that our model can effectively learn from the training dataset and accurately generalize to unseen images. The strategies for data augmentation we employed are as follows:

**Rescaling:** This technique normalizes pixel values across images to a 0-1 range. By ensuring a consistent input scale for the model, we minimize preprocessing disparities and set a uniform stage for model training.

**Rotation:** Our approach includes randomly rotating images within a predefined range, such as 20 degrees. This simulates real-world variations in camera angles and positions that the model may encounter in deployment.

**Horizontal Flip:** We applied horizontal mirroring to images, which is particularly effective for simulating various perspectives. This ensures that the model is not biased towards the orientation of the features in the training images.

**Shear Range:** Shear transformations are utilized to skew images, thereby enriching the dataset with a wider array of angles and perspectives. This enhances the model's ability to recognize features in a variety of visual contexts.

**Zoom Range:** Finally, random zooms are applied to pictures, allowing the model to focus on different parts of the image. This not only aids in recognizing features at various scales but also improves the model's resilience to scale variations in input data.

These data augmentation methods were meticulously chosen and integrated into our training pipeline, with the intent of exposing the model to a diverse spectrum of visual features. This diversity is key to developing a model that remains accurate and reliable when confronted with the myriad forms of media that it may need to evaluate in practical scenarios. Through these techniques, we aim to bolster the model's performance, making it a potent tool in the detection of DeepFakes.

### 4.1.3 DenseNet Model with Grayscale Images

Model: "sequential"		
Layer (type)	Output Shape	Param #
=====		
densenet121 (Functional)	(None, 7, 7, 1024)	7031232
global_average_pooling2d (GlobalAveragePooling2D)	(None, 1024)	0
dense (Dense)	(None, 1)	1025
=====		
Total params: 7,032,257		
Trainable params: 6,948,609		
Non-trainable params: 83,648		

Figure 4.1.4: DeepNet Model Architecture (Grayscale Images)

In our investigation of DeepFake detection, we explored the significance of color in distinguishing genuine imagery from manipulated content. To do so, we adapted our DenseNet model to process grayscale images, necessitating the conversion of all input data to grayscale using the 'color mode' setting.

While this adjustment represents a significant departure, we retained the original DenseNet architecture to ensure consistency in performance evaluation. Grayscale analysis serves not only as a procedural variation but also as an additional augmentation method, deepening our understanding of the model's robustness in the absence of color information.

Transitioning to grayscale acts as an experimental variable, compelling the model to rely solely on patterns, textures, and shapes. This enriches our comprehension of the model's adaptability and underscores the crucial role of color in image recognition and authentication.

This strategic modification, coupled with existing augmentation techniques, forms a comprehensive suite of experiments to assess the model's resilience against DeepFakes. Our initiative sets the stage for future investigations into the interplay between model complexity and performance.

## 4.2 Custom CNN

The proposed CNN is structured as a sequential model, a common framework for building neural networks in Keras, allowing for the layer-by-layer assembly of the network. The architecture is designed to progressively downsample the input image through a series of convolutional and pooling layers, extracting features while reducing spatial dimensions. This model turned out to be particularly effective for distinguishing between real versus fake images.



Layer (type)	Output Shape	Param #
batch_normalization (Batch Normalization)	(None, 224, 224, 3)	12
conv2d (Conv2D)	(None, 224, 224, 16)	448
max_pooling2d (MaxPooling2D)	(None, 112, 112, 16)	0
batch_normalization_1 (Batch Normalization)	(None, 112, 112, 16)	64
conv2d_1 (Conv2D)	(None, 112, 112, 32)	4640
max_pooling2d_1 (MaxPooling2D)	(None, 56, 56, 32)	0
batch_normalization_2 (Batch Normalization)	(None, 56, 56, 32)	128
dropout (Dropout)	(None, 56, 56, 32)	0
conv2d_2 (Conv2D)	(None, 56, 56, 64)	18496
max_pooling2d_2 (MaxPooling2D)	(None, 28, 28, 64)	0
batch_normalization_3 (Batch Normalization)	(None, 28, 28, 64)	256
dropout_1 (Dropout)	(None, 28, 28, 64)	0
conv2d_3 (Conv2D)	(None, 28, 28, 128)	73856
max_pooling2d_3 (MaxPooling2D)	(None, 14, 14, 128)	0
batch_normalization_4 (Batch Normalization)	(None, 14, 14, 128)	512
dropout_2 (Dropout)	(None, 14, 14, 128)	0
conv2d_4 (Conv2D)	(None, 14, 14, 256)	295168
max_pooling2d_4 (MaxPooling2D)	(None, 7, 7, 256)	0
batch_normalization_5 (Batch Normalization)	(None, 7, 7, 256)	1024
dropout_3 (Dropout)	(None, 7, 7, 256)	0
conv2d_5 (Conv2D)	(None, 7, 7, 512)	1180160
max_pooling2d_5 (MaxPooling2D)	(None, 3, 3, 512)	0

batch_normalization_6 (Batch Normalization)	(None, 3, 3, 512)	2048
dropout_4 (Dropout)	(None, 3, 3, 512)	0
global_average_pooling2d (GlobalAveragePooling2D)	(None, 512)	0
dense (Dense)	(None, 1)	513
=====		
Total params: 1,577,325		
Trainable params: 1,575,303		
Non-trainable params: 2,022		

Figure 4.2: Custom CNN Model Architecture

**Input Layer and Data Preprocessing:** The input to the network is an image of size 224x224 pixels with three color channels (RGB). The first layer of the model is a batch normalization layer, which normalizes the input layer by adjusting and scaling the activations. This step helps in speeding up the training process and reducing the number of epochs required to train the model.

**Convolutional and Pooling Layers:** Following the input layer, the model consists of several blocks, each comprising a convolutional layer followed by a max pooling layer. The convolutional layers use filters of increasing size (16, 32, 64, 128, 256, and 512) with a 3x3 kernel, 'relu' activation function, and 'same' padding. This configuration allows the model to learn increasingly complex and abstract features at each level. After each convolutional layer, a max pooling layer with a pool size of 2x2 is applied to reduce the spatial dimensions of the feature maps, thus reducing the computational complexity and helping to prevent overfitting.

Each convolutional block is also accompanied by batch normalization and dropout layers. The batch normalization layers, with a small epsilon value of 0.001, ensure that the activations remain within a range that promotes stable and fast training. The dropout layers, with a rate of 0.1, are crucial for regularizing the model, preventing overfitting by randomly setting a fraction of input units to 0 at each update during training.

**Global Average Pooling and Output Layer:** The final convolutional layer is followed by a global average pooling layer, which replaces the traditional fully connected layers found in many CNNs. This layer reduces the model's parameters, decreasing the risk of overfitting while maintaining the spatial hierarchy of the feature maps. The output of the global average pooling layer feeds into a dense layer with a single unit and a 'sigmoid' activation function, which classifies the image as real or fake.

This CNN architecture is designed with the specific task of DeepFake detection in mind, leveraging a series of convolutional, pooling, batch normalization, and dropout layers to efficiently and effectively learn discriminative features from images. The use of global average pooling instead of fully connected layers helps in reducing the model's complexity and the risk of overfitting, making it well-suited for the nuanced task of distinguishing between real and fabricated images. Through training on a dataset comprising real and fake faces, this model aims to achieve high accuracy in identifying DeepFake content, contributing valuable tools in the fight against digital misinformation.

## 4.3 Other Models

In addition to DenseNet and our custom CNN architecture, we experimented with leveraging EfficientNet, a state-of-the-art convolutional neural network architecture known for its efficiency and scalability. Despite its promising attributes, our trials with EfficientNet yielded suboptimal results in the context of DeepFake detection. A preliminary analysis suggested that to significantly enhance the model's performance, it would be necessary to scale up the network by increasing its depth and width, thereby augmenting its parameter count. However, such adjustments would demand considerably more computational resources, exceeding the capabilities of our available GPU hardware. This limitation constrained our ability to effectively utilize EfficientNet within our project's scope, compelling us to focus on optimizing and fine-tuning our custom CNN model to achieve the desired balance between performance and computational feasibility.

If we can secure access to more powerful GPUs, we intend to delve into EfficientNet's extensive capabilities to significantly improve our DeepFake detection efforts.

## 5. Metrics

Before we dive into results, let's look at the main metrics we will be using to evaluate our models:

1. **Accuracy** is the proportion of true results among the total number of cases examined, meaning it measures how often the model is correct overall.  
For instance, the DenseNet Model with Grayscale Images has an accuracy of 0.98, which means it correctly identified 98% of the images as either real or fake.
2. **Precision** is the ratio of true positive observations to the total positive observations predicted by the model, indicating how reliable the model is when it predicts something is a deepfake.

For example, the DenseNet Base Model's precision of 0.97 suggests that when it predicts an image is a deepfake, it's correct 97% of the time.

3. **Recall** is the ratio of true positive observations to all actual positives, measuring how well the model catches all the actual deepfakes.  
For example, the DenseNet Base Model has a recall of 0.97, meaning it identifies 97% of all actual deepfakes in the dataset.
4. The **F1-Score** is the harmonic mean of Precision and Recall, providing a balance between the two. It looks at both the reliability and completeness of the model's deepfake detection.  
For instance, the DenseNet Base Model's F1-Score of 0.97 suggests it is both reliable and thorough in detecting deepfakes.

Strong scores in Precision, Recall, and the F1 score demonstrate high effectiveness—key indicators of its capability to discern real from fake content. A balanced F1 score reflects the model's accuracy and consistency, while individual attention to Precision and Recall is critical due to the serious implications of misidentifying deepfakes. High Precision helps prevent mislabeling genuine content as fake, protecting the content's credibility, while High Recall ensures fake content does not slip through, crucial for maintaining security and trust.

The focus on these metrics is particularly vital in deepfake detection, where both false positives and false negatives have significant real-world consequences—from misinformation spread to unwarranted censorship. Thus, while the F1 score is a useful summary statistic, optimizing Precision and Recall separately ensures our model's performance is not only technically sound but also ethically and pragmatically aligned with the stakes of accurately detecting deepfakes in digital communications.

# 6. Results

## 6.1 DenseNet Model (Base Model)

ROC AUC Score: 0.9975173699999998					
AP Score: 0.9973757115620031					
	precision	recall	f1-score	support	
0	0.94	0.99	0.97	10000	
1	0.99	0.94	0.96	10000	
accuracy			0.97	20000	
macro avg	0.97	0.97	0.97	20000	
weighted avg	0.97	0.97	0.97	20000	

Figure 6.1: DenseNet Base Model Performance

In the application of DenseNet for DeepFake detection, our model has achieved remarkable results, demonstrating exceptional efficacy in distinguishing between genuine and manipulated content. The ROC AUC score reached an impressive 0.9975, indicative of the model's high true positive rate and low false positive rate. Additionally, the AP score closely paralleled this performance with a value of 0.9974, suggesting the model's precision in ranking its predictions across various thresholds is highly reliable.

The model also exhibited excellent precision and recall, with scores of 0.94 and 0.99 for authentic images (class '0'), and 0.99 and 0.94 for DeepFake images (class '1') respectively. These metrics are complemented by f1-scores of 0.97 and 0.96, which further confirm the balanced detection capabilities of our DenseNet-based algorithm. With an overall accuracy of 0.97 and consistently high macro and weighted averages for precision, recall, and f1-score, the results underscore the robustness and reliability of our approach in the domain of DeepFake detection.



## 6.2 DenseNet Model with Image Augmentations

ROC AUC Score: 0.968640695					
AP Score: 0.966609199565477					
	precision	recall	f1-score	support	
0	0.96	0.81	0.88	10000	
1	0.83	0.97	0.90	10000	
accuracy			0.89	20000	
macro avg	0.90	0.89	0.89	20000	
weighted avg	0.90	0.89	0.89	20000	

Figure 6.2: DenseNet Model with Image Augmentations Performance

Implementing image augmentations in our DenseNet-based DeepFake detection model has yielded promising outcomes. The model achieved a ROC AUC score of 0.9686, demonstrating a strong ability to differentiate between genuine and fabricated images. The AP score, reflecting the precision of the model across various decision thresholds, was also notable at 0.967. These metrics suggest that the inclusion of image augmentations has enhanced the model's generalization capabilities, enabling it to maintain robust performance even under varied input conditions.

In terms of precision, recall, and f1-score, the model has shown a solid performance, particularly in identifying DeepFake images (class '1'), with precision at 0.83, recall at 0.97, and an f1-score of 0.90. This indicates a high sensitivity to DeepFake detection. For authentic images (class '0'), the model also performed well, with a precision of 0.96, recall of 0.81, and an f1-score of 0.88. The overall accuracy of the model stood at 0.89, with both macro and weighted averages for precision, recall, and f1-score at 0.90 and 0.89, respectively.

### 6.3 DenseNet Model with Grayscale Images

ROC-AUC Score: 0.9980781299999999				
AP Score: 0.9978361434792907				
	precision	recall	f1-score	support
0	0.99	0.97	0.98	10000
1	0.97	0.99	0.98	10000
accuracy			0.98	20000
macro avg	0.98	0.98	0.98	20000
weighted avg	0.98	0.98	0.98	20000

Figure 6.3: DenseNet Model with Grayscale Images Performance

Utilizing grayscale images for our DenseNet-based DeepFake detection model has demonstrated exceptional performance. The model reached an impressive ROC AUC score of approximately 0.9980, showcasing an outstanding ability to distinguish between authentic and manipulated images. Furthermore, the AP score stands at roughly 0.9978, signifying an excellent model precision across different decision thresholds. These results indicate that the model is highly effective and reliable in identifying DeepFake content when trained on grayscale images.

In terms of precision, recall, and f1-score, the model presents an almost perfect performance. For authentic images (class '0'), the model has achieved a precision of 0.99 and a recall of 0.97, which results in an f1-score of 0.98. This denotes that the model is nearly flawless in correctly labeling authentic images. For DeepFake images (class '1'), the results are similarly outstanding, with a precision of 0.97, a recall of 0.99, and an f1-score of 0.98. The high recall indicates the model's exceptional sensitivity in detecting almost all DeepFake instances.

The overall accuracy of the model is 0.98, which is remarkable and suggests that the model makes correct predictions most of the time. The macro average and weighted averages for precision, recall, and f1-score are all at 0.98, indicating a consistent and high-level performance across both classes.

These metrics, particularly the elevated f1-scores, imply that converting images to grayscale as a preprocessing step might have contributed to reducing the complexity of the task, allowing the model to focus more effectively on the structural features pertinent to DeepFake detection.

## 6.4 Custom CNN

ROC AUC Score: 0.9919095450000001					
AP Score: 0.9914724911124952					
		precision	recall	f1-score	support
	0	0.92	0.98	0.95	10000
	1	0.97	0.92	0.94	10000
accuracy				0.95	20000
macro avg		0.95	0.95	0.95	20000
weighted avg		0.95	0.95	0.95	20000

Figure 6.4: Custom CNN Model Performance

Employing a custom Convolutional Neural Network (CNN) has produced outstanding results in our DeepFake detection project. The model has achieved an ROC AUC score of 0.9919, showcasing its exceptional capability to differentiate between genuine and altered images. The Average Precision (AP) score closely mirrors this excellent performance with a score of 0.9914, indicating the model's high precision across various decision-making thresholds.

The evaluation of the model's precision, recall, and f1-score highlights its robustness in classifying both authentic and DeepFake images. For class '0', representing authentic images, the model obtained a precision of 0.92 and an impressive recall of 0.98, culminating in an f1-score of 0.95. The detection of DeepFakes, denoted as class '1', was also highly accurate, with the model achieving a precision of 0.97, recall of 0.92, and an f1-score of 0.94. With an overall model accuracy of 0.95 and uniform macro and weighted averages for precision, recall, and f1-score at 0.95, these results affirm the effectiveness of the custom CNN in accurately identifying and mitigating DeepFake content.

## 7. Comparison of Results

Model	Accuracy	Precision	Recall	F1-Score
DenseNet Model (Base Model)	0.97	0.97	0.97	0.97
DenseNet Model with Image augmentations	0.89	0.90	0.89	0.89
DenseNet Model with Grayscale Images	0.98	0.98	0.98	0.98
Custom CNN	0.95	0.95	0.95	0.95

Table 7.1: Model Comparison

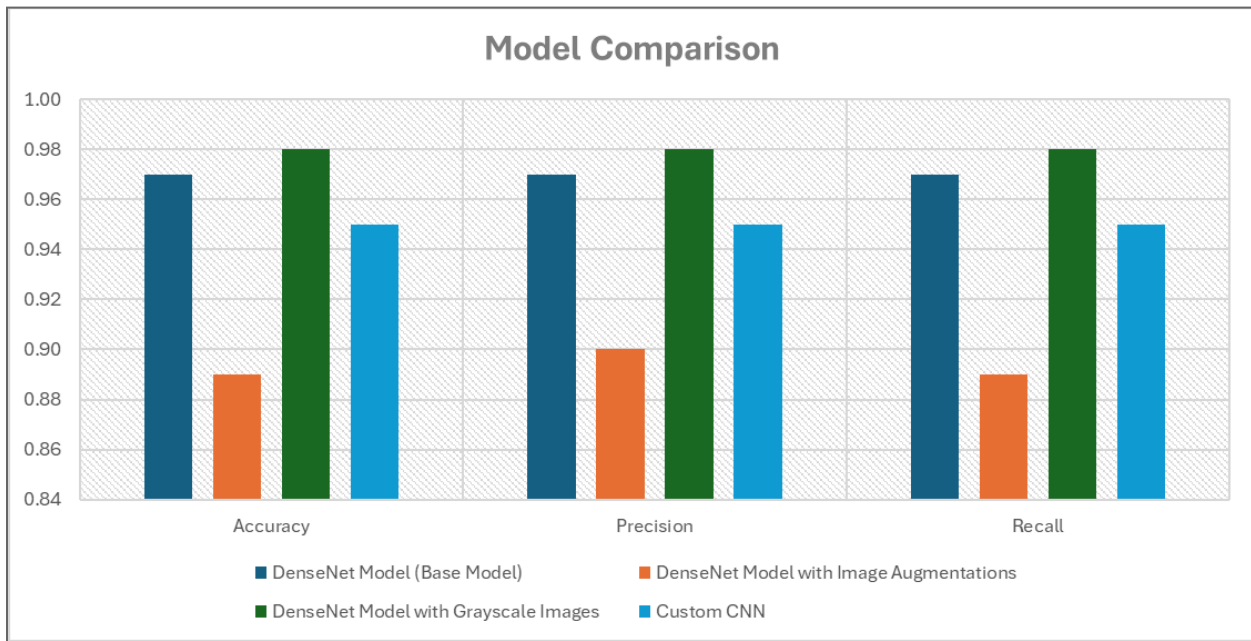


Figure 7.1: Model Comparison

In our study, we evaluated the DenseNet architecture, which demonstrated notable performance, affirming its efficiency and effectiveness as mentioned previously. The integration of data augmentation led to a decrease in accuracy, a result that aligns with expectations since augmented data inherently complicates the training process. We believe that extending the training duration could potentially enhance these results. Our investigations also indicated that converting images to grayscale actually improved the DenseNet model's ability to classify GAN-generated images, suggesting that color information is not critical for detection in this context.

The Custom Model developed for our project also yielded positive outcomes, finding a balance between efficiency and performance on augmented data. It surpassed the DenseNet model in terms of performance with augmented datasets.

The relative performance of the models, including a visual comparison through the ROC plot (Figure 7.2) , offers insight into their efficacy in detecting DeepFake images.

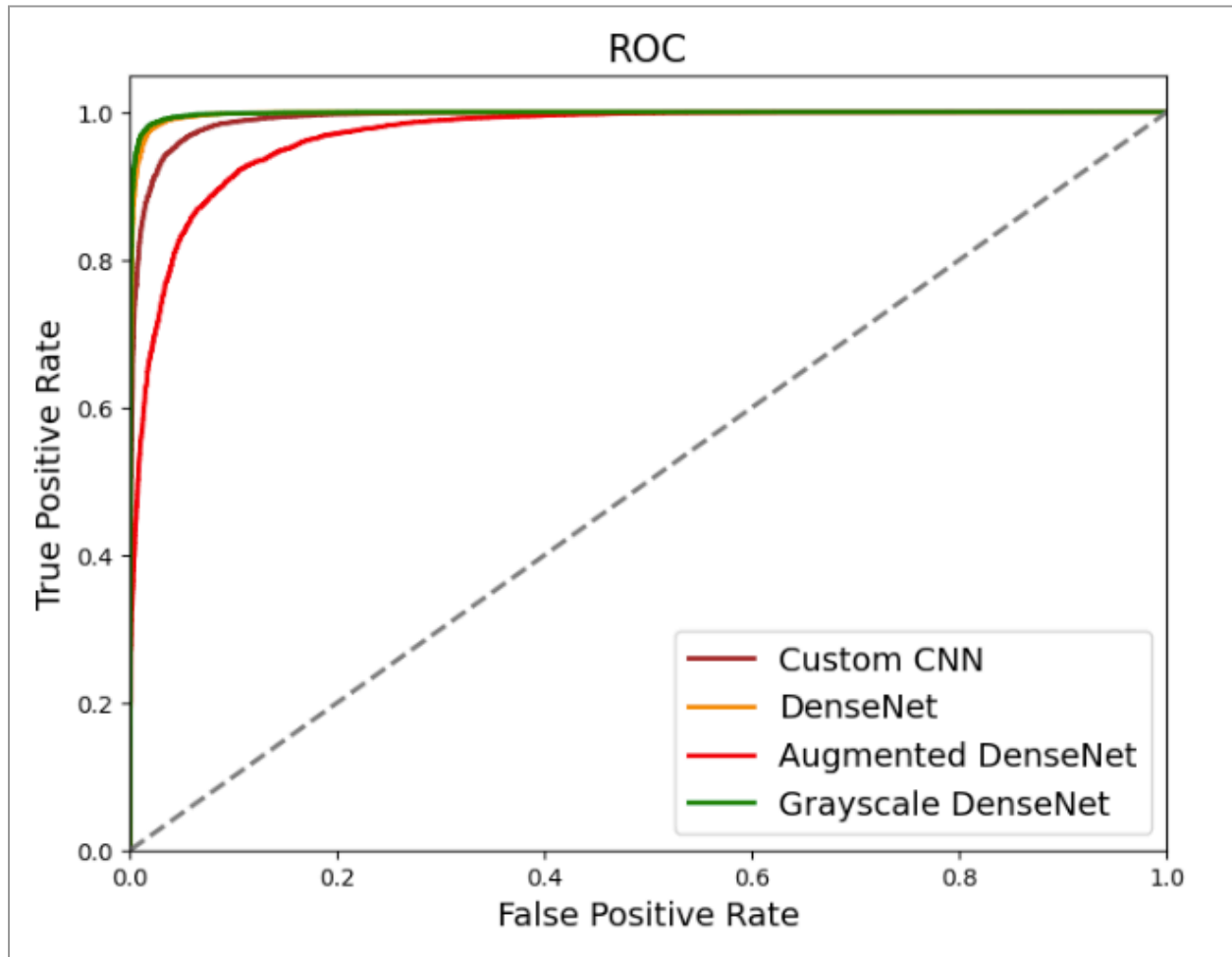


Figure 7.2: ROC Plot



The scatter plots of the principal components for each model, with the exception of DenseNet with Grayscale, also provide valuable insights into the capability of these models to distinguish between real and fake images. These plots demonstrate the effectiveness of the feature extraction process, as seen through the lens of principal component analysis. Specifically, the first two principal components, often referred to in analysis, showcase distinct clusters for real versus fake images in most models analyzed.



Figure 7.3: Scatter Plot of PCA components for each model

These clusters signify the models' ability to capture and accentuate the differences between real and fake images, making it possible to apply classification techniques like SVM (Support Vector Machine) effectively. A clear separation between clusters in these plots is a strong indicator of a model's potential in accurately classifying and thus identifying DeepFakes.

## 8. Performance Demo

### Identify Deepfakes:

For practical evaluation of our DeepFake detection models, we provide a Jupyter notebook named Performance\_Demo.ipynb. This enables users to upload images, load our DenseNet and custom CNN models to swiftly obtain the model's verdict on their authenticity. The screenshots below illustrate the notebook's user-friendly model performance evaluation.

#### Upload a Fake Image:



**Actual Label: Fake**

1/1 [=====] - 0s 82ms/step

**Custom\_CNN\_model: Fake**

1/1 [=====] - 0s 316ms/step

**DenseNet\_model: Fake**

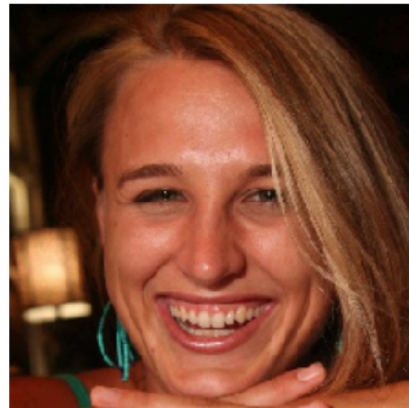
1/1 [=====] - 0s 249ms/step

**DenseNet\_Augmented\_model: Fake**

1/1 [=====] - 0s 194ms/step

**DenseNet\_Grayscale\_model: Fake**

#### Upload a Real Image:



**Actual Label: Real**

1/1 [=====] - 0s 71ms/step

**Custom\_CNN\_model: Real**

1/1 [=====] - 0s 282ms/step

**DenseNet\_model: Real**

1/1 [=====] - 0s 257ms/step

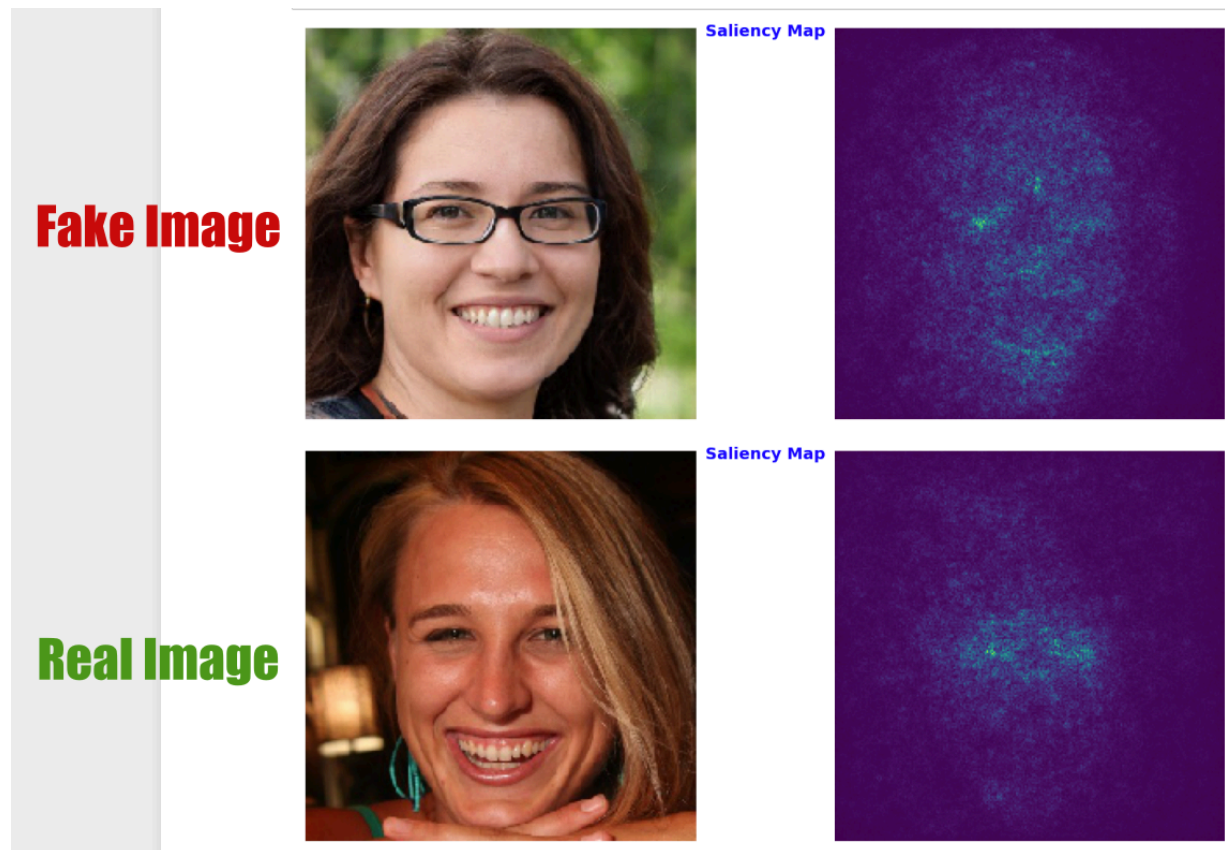
**DenseNet\_Augmented\_model: Real**

1/1 [=====] - 0s 245ms/step

**DenseNet\_Grayscale\_model: Real**

Figure 8.1: Model Performance/Evaluation

## Saliency Maps:



*Figure 8.2: Saliency Maps*

In the demonstration, we also explored the application of saliency maps on images. Saliency maps serve as a valuable tool for understanding which regions of an image capture the attention of the algorithm when distinguishing between real and fake images. Our observation revealed that, in the case of a fake image, the algorithm predominantly focuses on the eyes, nose, and chin areas. This is contrasted with the real image, where the algorithm's focus is primarily on the eye area.

It's important to note that saliency maps offer deep insights into the decision-making processes of algorithms. By highlighting the specific areas that the algorithm considers significant, saliency maps can help developers and researchers identify potential biases, understand the strengths and weaknesses of the model, and guide improvements in algorithm design. This can be especially useful in applications where understanding the 'why' behind a decision is as critical as the decision itself, such as in security, authentication, and deepfake detection.

## 9. Conclusion

Our team has harnessed the capabilities of the DenseNet121 convolutional neural network to forge a model adept at identifying DeepFakes, pivotal in combating digital misinformation. Our tailored approach eschews pre-loaded weights, ensuring the feature extraction is intricately aligned with the task of discerning authentic from manipulated content. By integrating a Global Average Pooling 2D (GAP) layer post-DenseNet121, we've honed in on the most crucial features from input images, thereby streamlining our model's processing prowess.

Additionally, the incorporation of our Custom Convolutional Neural Network (CNN) has further enhanced our DeepFake detection capabilities. This model, designed specifically for this task, efficiently extracts discriminative features from images, contributing to the overall robustness and accuracy of our approach.

Both CNN and DenseNet models architecture peaks with a dense output layer featuring a sigmoid activation function, which adeptly computes a probability score to ascertain the authenticity of images. This critical functionality anchors the model's ability to classify images as real or synthetic with impressive precision.

Training over 20 epochs with a batch size of 32 and utilizing the Adam Optimizer, we struck an equilibrium between performance and computational demands, ensuring our model delivers on accuracy without excessive resource consumption. Notably, our grayscale-infused DenseNet model demonstrated superior performance over other variations, including our custom CNN and the base model with image augmentations. This underscores the potential minimal role of color in the accurate detection of DeepFakes, suggesting a significant focus on texture and structure.

Our research delineates the delicate interplay between model complexity and its efficacy in generalizing across diverse datasets, noting that augmentations can sometimes diminish performance by introducing counterproductive variance.

In essence, our concerted effort has yielded a robust and precise model that not only benchmarks DenseNet121's effectiveness in detecting DeepFakes but also showcases the valuable contribution of our Custom CNN. This work propels the conversation forward concerning the refinement of deep learning models for enhanced digital security. The insights gained from the effectiveness of grayscale images in this context beckon further exploration and optimization of models. Moreover, our work accentuates the necessity for a composite approach, encompassing legal, ethical, and policy frameworks, to effectively counter the multifaceted challenges posed by DeepFake technology.

## 10. Future Work/ Next Steps

Due to previously limited computational resources, our project's ability to enhance model robustness and generalization in DeepFake detection was constrained. These limitations impacted the scale and sophistication of the data processing and model training we could achieve. However, with the prospect of accessing increased computational power, we are poised to significantly advance our efforts in four key areas:

1. **Advanced Architectural Implementation:** With more computational resources, we can explore and implement more complex and computationally intensive models such as EfficientNet and other cutting-edge architectures. These models are known for their superior performance in image classification tasks due to their scalable architecture, allowing for better generalization and robustness in detecting DeepFakes across various datasets.
2. **Diverse Dataset Utilization:** Greater computational power will enable us to work with larger and more varied datasets, including high-resolution images and videos from multiple sources. This diversity is crucial for training our models to recognize a wide range of DeepFake techniques and ensuring they are robust against new and evolving methods.
3. **Enhanced Data Augmentation:** We plan to significantly expand our data augmentation techniques. By leveraging more advanced and diverse augmentation methods, such as simulating different lighting conditions, adding noise, and applying various transformations, our model's ability to generalize and accurately identify DeepFakes under different scenarios will be greatly improved.
4. **Intensive Training Regimes:** With the availability of more computational resources, we can afford to run longer, more intensive training regimes. This includes employing techniques such as adversarial training, where the model is continuously challenged to improve its detection capabilities by being exposed to an ever-evolving set of DeepFake examples generated by an adversarial network.

In conclusion, by addressing these key areas with the aid of enhanced computational resources, we want to set a strong foundation for significantly advancing the state of DeepFake detection technology, ensuring it remains effective and resilient against the rapidly evolving landscape of digital misinformation.



# 11. References

[1] DeepFake-Detection

[2] Awesome Face Forgery Generation and Detection

[3] Comparison of Deepfake Detection Techniques through Deep Learning by Maryam Taeb and Hongmei Chi

[4] Combining EfficientNet and Vision Transformers for Video Deepfake Detection

[5] Deepfake-image-detection

[6] DeepFake-Learning

**Note:** For References, click to visit the link