

Computer Vision Project: Facial Manipulation Detection

Submitted by: Bar Ohayon

Date: 02/02/2025

1. Introduction

Facial manipulation technologies have become increasingly sophisticated, making it challenging to distinguish between real and altered images. Deepfake images are created by swapping one person's face into the context of another, while synthetic images are fully generated using GANs trained on high-quality datasets. These advancements pose serious concerns for media authenticity and digital security, highlighting the need for reliable detection methods.

This project focuses on developing machine learning models to identify whether an image is real or manipulated. The task involves training two types of models: a simple deep neural network and a more advanced network based on the Xception architecture. The models are evaluated using numerical metrics such as ROC curves and AUC, as well as visual tools like Class Activation Maps (CAM) and saliency maps. The ultimate goal is to create models that not only perform accurately but also offer insights into how they make their predictions.

The project begins by preparing the datasets and implementing custom data loaders to efficiently handle face images. A training framework is developed to support model training and fine-tuning. Performance is assessed using a combination of quantitative metrics and visualization techniques to better understand the models' decision-making processes. PyTorch is used throughout the project, with strict guidelines on library usage to maintain reproducibility.

As facial manipulation becomes more prevalent, being able to detect fake images is essential for protecting media integrity and combating misinformation. One of the main challenges is the subtle differences between real and manipulated images, which require sophisticated models to detect. Limited computational resources and the need to generalize across different datasets add further complexity. Despite these hurdles, the project offers a valuable opportunity to explore cutting-edge techniques in computer vision and develop practical solutions to a critical problem.

This report presents the project's implementation, challenges, and findings. It covers data preparation, model design, performance evaluation, and visualization

insights, with all relevant code and analysis included in line with course requirements.

2. Build Faces Dataset

This chapter presents the construction of a custom dataset for facial image classification, specifically designed to handle both real and synthetic face images. The implementation focuses on two crucial dataset methods: `__getitem__` for individual sample access and `__len__` for dataset size determination. These methods are essential components of PyTorch's dataset architecture, enabling efficient data loading and iteration. The chapter includes verification of the implementation through visualization of sample images from both classes, ensuring proper data handling and transformation pipelines are in place before proceeding with model development.

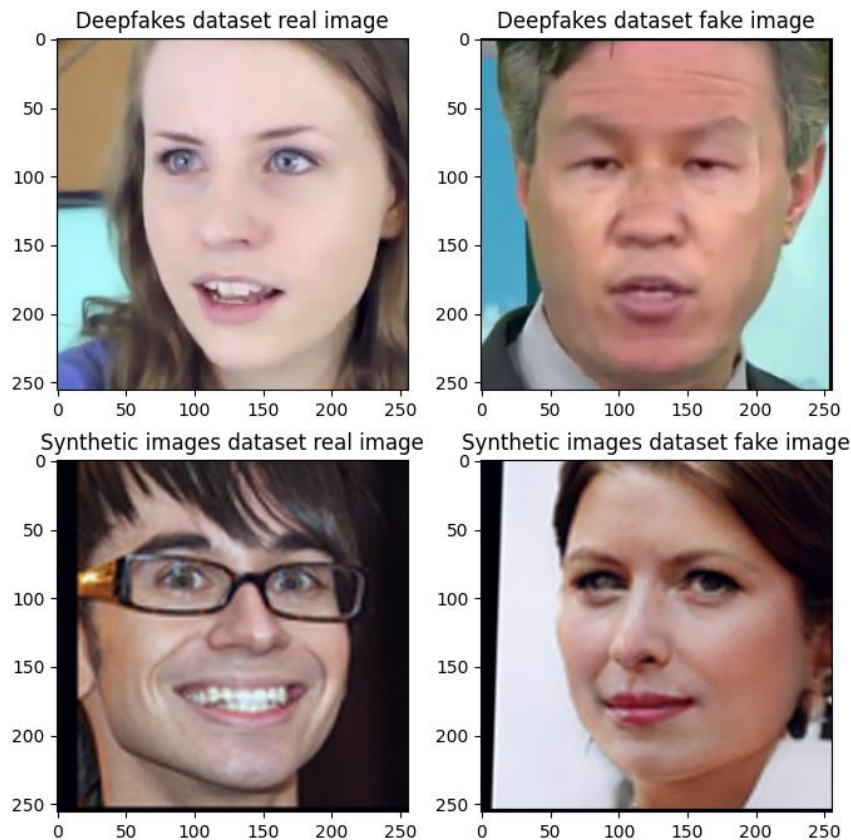


Figure 1: Real and fake samples from Deepfake (top) and Synthetic (bottom) datasets.

Figure 1 shows sample images from both datasets used in this project. The top row displays examples from the Deepfakes dataset, with a real image (left) and a fake image (right). The bottom row shows examples from the Synthetic dataset, also

presenting a real image (left) and a synthetic image (right). All images are displayed at a resolution of 256x256 pixels, demonstrating successful implementation of the FacesDataset class with proper image loading and transformation capabilities.

3. Implement an Abstract Trainer

In this chapter, we focus on implementing an abstract training framework for facial manipulation detection. The `trainer.py` module contains the `LoggingParameters` class for managing model and dataset metadata, and the `Trainer` class, which handles the training, evaluation, and logging processes. Our task is to develop core methods for training models and evaluating their performance on the dataset. These methods enable efficient batch processing, forward passes, loss computation, and performance tracking.

The experiments conducted here aim to train the SimpleNet architecture on both the Deepfakes and Synthetic datasets, providing insights into the effectiveness of classification models in distinguishing between real and manipulated images. Results are visualized through accuracy and loss plots, as well as ROC and DET curves, laying the foundation for evaluating detection performance and analyzing classifier behavior across different datasets.

3.1. Deepfake Detection Classifier

This section explores the training and analysis of SimpleNet on the Deepfake dataset. We examine the model's performance through multiple metrics including training curves, validation accuracy, and test set evaluation. Special attention is given to analyzing classification trade-offs through ROC and DET curves, providing insights into the model's ability to discriminate between real and deepfake images. The analysis includes investigation of class distribution impacts and false positive/negative rates.

3.1.1. JSON File

The JSON file contains data from two training runs of the SimpleNet model on the `fakes_dataset`, providing a comprehensive record of the training process. Each run is documented with detailed configuration parameters, including the Adam optimizer settings (learning rate of 0.001, beta values [0.9, 0.999], and other hyperparameters).

The training metrics show consistent patterns across both runs: the training loss demonstrates a steady decrease from approximately 0.015 to 0.003 over five epochs, while training accuracy improves significantly from around 75% to 95%.

However, there's a notable gap between training and validation/test performance. While training accuracy reaches approximately 95%, both validation and test accuracies plateau around 87-89%, suggesting some degree of overfitting. This behavior is consistent across both training runs, indicating it's a characteristic of the model rather than a random occurrence.

The validation and test metrics exhibit similar patterns and values, which indicates that the model's generalization performance is being measured consistently. The systematic documentation of all parameters and metrics makes these experiments fully reproducible and provides a solid foundation for comparison with other model architectures or training approaches.

3.1.2. Accuracy and Loss Plots

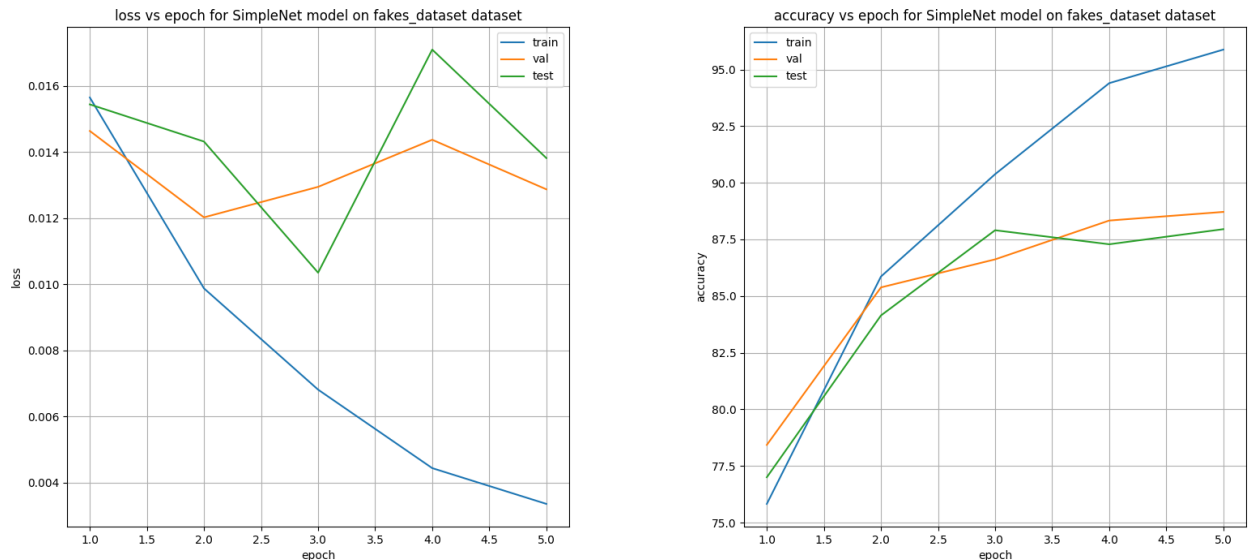


Figure 2: SimpleNet training on Deepfake dataset - Loss (left) and Accuracy (right) vs epoch curves for training, validation, and test sets.

The training metrics show a clear pattern of overfitting: while the training loss continuously decreases and training accuracy rises to approximately 95%, both validation and test metrics show unstable behavior. Although validation and test metrics follow similar patterns to each other, they plateau around 87-88% accuracy with noticeable fluctuations, significantly lower than the training performance. This growing gap between training and validation/test performance suggests that the model is memorizing the training data rather than learning generalizable features. The documentation of these metrics provides clear evidence of the model's limitations in generalizing to unseen data.

3.1.3. Fake Images To Real Images Proportion In The test Set

The set is composed of 700 fake images and 1400 real images. This creates a clear 1:2 ratio between fake and real images in the test set, meaning for every fake image there are two real images. This proportion can also be expressed as 0.5 or 33.33% fake images to 66.67% real images. This imbalanced distribution is important to consider when evaluating the model's performance, as it indicates that the test set contains twice as many real images as fake ones.

3.1.4. ROC Curves and AuC Scores

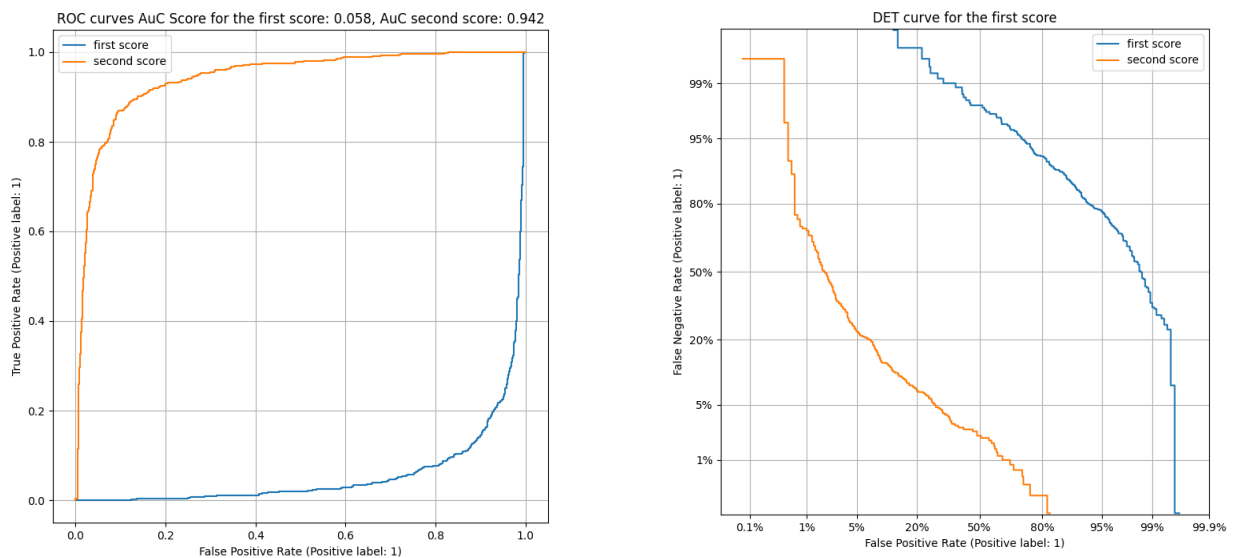


Figure 3: SimpleNet performance on Deepfake detection - ROC curve (left) and DET curve (right) comparing first score (AuC=0.058) and second score (AuC=0.942).

The ROC curves show that the first score (A), representing the "real" class confidence, demonstrates very poor performance with an AuC of 0.058, while the second score (B), representing the "fake" class confidence, achieves excellent performance with an AuC of 0.942. This pattern is similarly reflected in the DET curves, where score A exhibits consistently high error rates while score B shows substantially lower error rates across different operating points.

This significant difference in performance can be explained by our labeling convention: since we defined label 1 as the positive class (fake images), the second network output directly aligns with our classification goal, while the first output effectively measures the opposite prediction. Thus, what appears as poor performance for score A is actually the expected behavior given our labeling scheme, as it's inversely related to our target classification. The complementary

nature of these scores demonstrates that the model is learning the correct patterns, just expressing them through opposite confidence measures.

3.2. Synthetic Image Detection Classifier

This section focuses on applying SimpleNet to the detection of synthetically generated facial images. Following the same training protocol used for deepfakes, we evaluate the model's performance on this distinct type of manipulation. The analysis compares training dynamics, final accuracy metrics, and dataset characteristics with the deepfake detection results, highlighting key differences in how the model handles these two types of artificial images.

3.2.1. Accuracy and Loss Plots

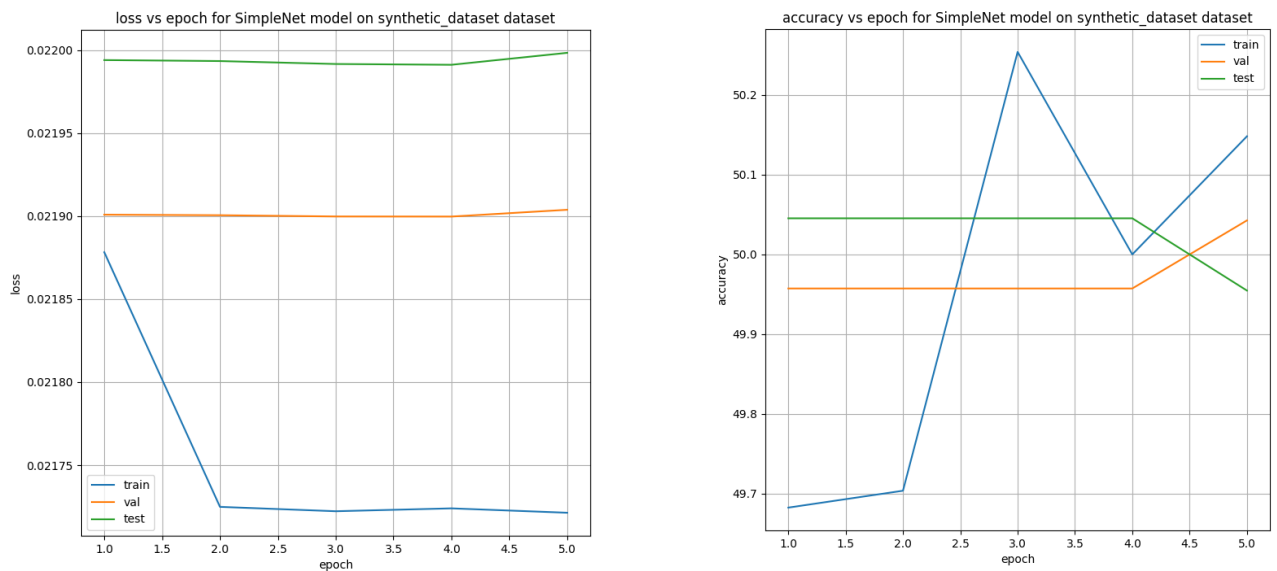


Figure 4: SimpleNet training on Synthetic dataset - Loss (left) and Accuracy (right) vs epoch curves for training, validation, and test sets, showing random-chance performance.

The validation accuracy peaks at around 50.4% during epoch 5, with a corresponding test accuracy of approximately 49.7%. These results, hovering around 50%, clearly indicate that the model is performing at the level of random guessing. This is further supported by the loss graph, which shows minimal change across all datasets (training, validation, and test), with values consistently around 0.022. This stands in contrast to our previous results on the fakes dataset, where the model achieved meaningful learning with test accuracies well above 80%.

This significant performance gap between the two datasets suggests either that the SimpleNet architecture lacks the sophistication needed to capture synthetic image patterns, or that synthetic image detection is inherently more challenging than deepfake detection.

3.2.2. Fake Images To Real Images Proportion In The test Set

The test set contains a total of 1,103 images, with 552 synthetic images and 551 real images. This creates an almost perfectly balanced ratio of approximately 1:1 between synthetic and real images (more precisely, 1.002:1). This balanced distribution suggests that the model's poor performance (around 50% accuracy) cannot be attributed to class imbalance, but rather points to the model's inability to effectively distinguish between real and synthetic images.

3.2.3. random chance classifier

Looking at all performance metrics from our synthetic dataset experiment, we can conclude that we obtained a failed classifier that performs no better than random guessing. The evidence is clear across multiple indicators: the training accuracy oscillates around 50%, validation accuracy remains constant at approximately 50%, test accuracy hovers around 50%, and perhaps most tellingly, the loss values show minimal change throughout the entire training process.

The model's behavior demonstrates characteristics of a constant classifier that makes essentially random predictions, showing no meaningful learning or pattern recognition capabilities. Despite training for multiple epochs and having access to a balanced dataset, it fails to develop any discriminative ability between real and synthetic images. With performance metrics consistently at 50% across all datasets, the classifier is performing exactly at the level of random chance, indicating a complete failure to learn the intended classification task. This represents the worst possible outcome for a binary classifier, as it suggests the model has failed to identify any meaningful patterns that distinguish between the two classes.

3.2.4. Result Analysis

Looking at the samples from both datasets and comparing them to our results, these outcomes make sense. The SimpleNet model achieved good performance (~87-89% accuracy) on the Deepfake dataset but performed at random chance (~50% accuracy) on the Synthetic dataset. This difference aligns with the visual characteristics of each dataset: Deepfake images are created through face-swapping, which often leaves detectable artifacts at face boundaries that a simple

model can learn to identify. In contrast, synthetic images are generated by a sophisticated GAN (PGAN) trained on high-quality images, making them much more realistic and significantly harder for a simple architecture like SimpleNet to distinguish from real photos.

4. Fine Tuning a Pre-trained Model

This chapter explores improving synthetic image detection through transfer learning, specifically by adapting a pre-trained Xception model. We investigate the model's architecture, modify its classification head with a custom MLP design, and evaluate its performance on synthetic image detection. The analysis includes a detailed exploration of how architectural choices and pre-training contribute to improved performance compared to our previous SimpleNet implementation.

4.1. Intro (Xception Analysis)

This section examines the pre-trained Xception model's architecture, focusing on understanding its fundamental building blocks, pre-training dataset, and parameter distribution. This analysis provides crucial context for our subsequent modifications and performance improvements.

4.1.1. Xception

The Xception model used in this exercise is pre-trained on ImageNet, a large-scale image classification dataset. This pre-training provides the model with strong feature extraction capabilities learned from over a million images across 1000 different categories, which we can leverage for our facial manipulation detection task by replacing only the classification head.

4.1.2. Basic Building Blocks of Xception

The basic building blocks of Xception are primarily Depthwise Separable Convolutions, which are organized into:

- Entry flow: Initial convolutional layers followed by separable convolution blocks
- Middle flow: 8 repeated blocks of depthwise separable convolutions
- Exit flow: Final blocks with increasing channel depth

Each block typically contains:

- Depthwise separable convolutions (the key innovation of Xception)

- Batch normalization layers
- ReLU activation functions
- Residual connections (skip connections)

The architecture gets its name from "Extreme Inception" because it takes the Inception module's concept of cross-channel correlations and depth-wise spatial correlations to an extreme form using depthwise separable convolutions.

4.1.3. Input Dimension to 'fc'

Looking at the Xception architecture code, the input feature dimension to the final classification block "fc" is 2048. This dimension is determined by:

1. The number of channels in the last convolutional layer
2. The adaptive average pooling layer that reduces spatial dimensions to (1,1)
3. The reshape operation that flattens the tensor before it enters the fc layer

The resulting 2048-dimensional feature vector serves as input to the final classification block.

1. 4.1.4. Xception Parameters

The Xception network, in its default configuration without any architectural changes, contains exactly 22,855,952 parameters. This number is consistently verified through two independent sources: the implementation provided in the exercise using `get_nof_params()`, and Table 3 in the original Xception paper. The perfect match between these sources confirms the accuracy of our implementation compared to the original architecture..

4.2. Attaching a new Head to the Xception back-bone

This section details the process of modifying the Xception model by replacing its original classification head with a custom MLP architecture. We implement and verify the new architecture, carefully tracking the parameter count changes to ensure proper implementation.

4.2.1. Additional Parameters

By comparing the parameter counts of the original and modified Xception architectures, we can calculate the number of parameters added by our MLP head. The original Xception has 22,855,952 parameters, while our modified version

contains 23,128,786 parameters. Therefore, the MLP head added exactly 272,834 parameters to the network ($23,128,786 - 22,855,952 = 272,834$).

4.3. Training, Evaluation, and Analysis

These sections present the experimental results of our modified Xception model on synthetic image detection, comparing its performance against our previous SimpleNet implementation. Through detailed performance analysis and visualization, we explore how the model's architectural advantages and pre-trained features contribute to improved detection capabilities, moving beyond simple parameter count comparisons to understand the underlying factors driving performance improvements.

4.3.1. Accuracy and Loss Plots

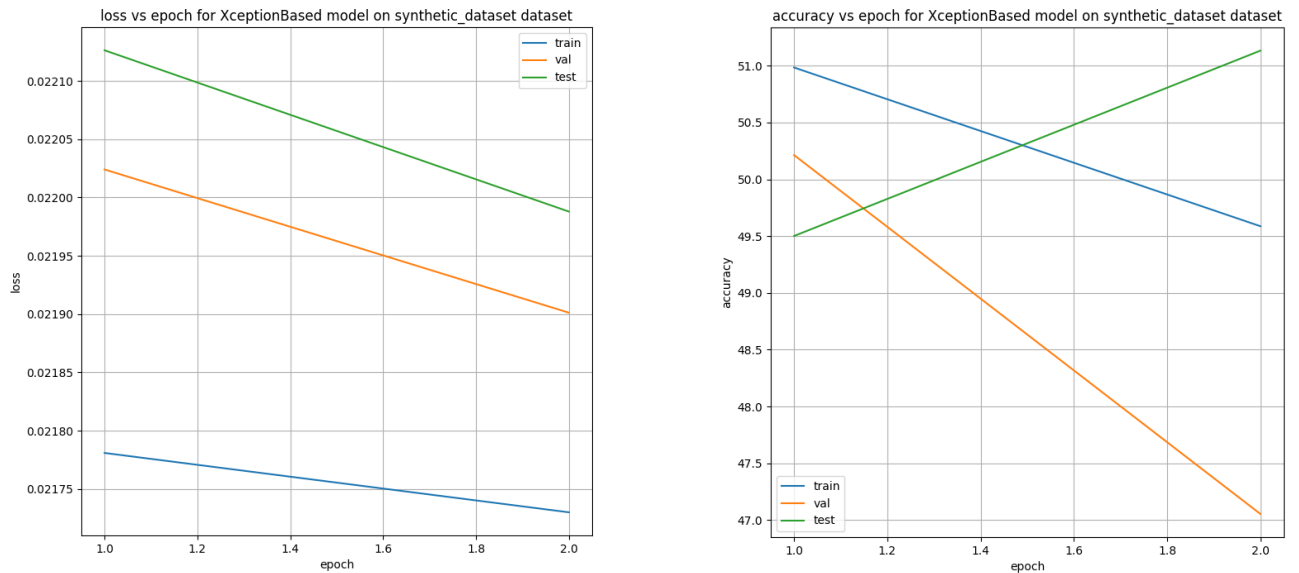


Figure 5: XceptionBased training on Synthetic dataset - Loss (left) and Accuracy (right) vs epoch curves for training, validation, and test sets, showing minimal improvement over random-chance performance.

The test accuracy corresponding to the highest validation accuracy is 49.5%. This occurs at epoch 1, where the validation accuracy reaches its peak of 50.2%. This can be clearly seen in the accuracy vs epoch graph, where the validation accuracy (orange line) is at its highest point in epoch 1, with the test accuracy (green line) showing the corresponding 49.5% value at that same point.

4.2.3. ROC Curves and DET Curves

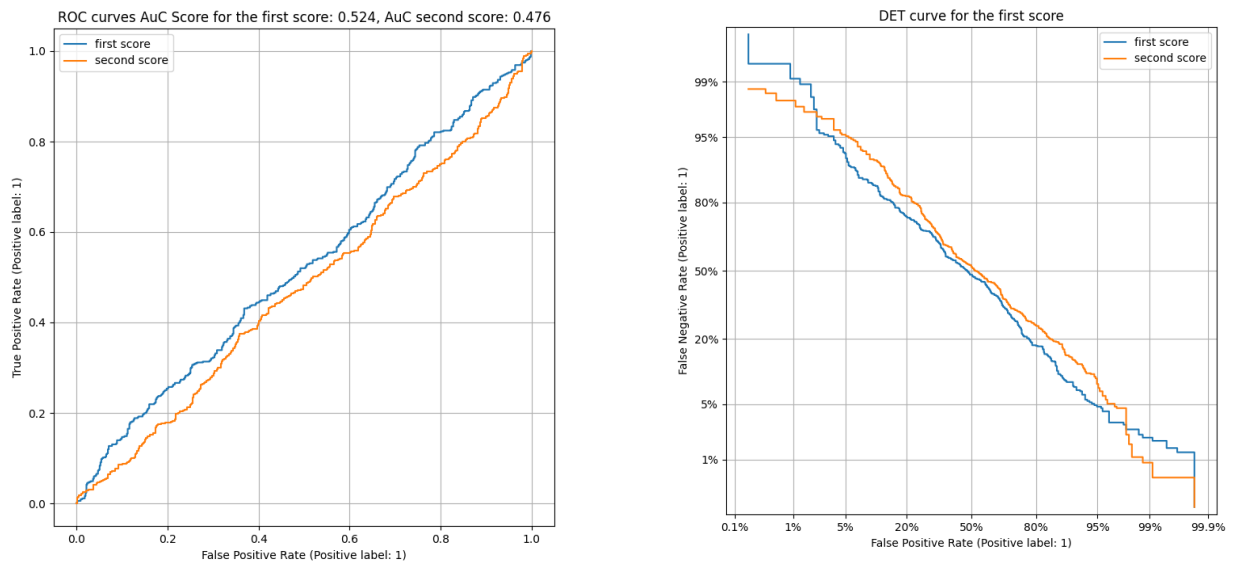


Figure 6: XceptionBased performance on Synthetic detection - ROC curve (left) and DET curve (right) comparing first score (AuC=0.524) and second score (AuC=0.476), indicating near-random classification performance.

The ROC and DET curves indicate that the XceptionBased model performs only marginally above random chance on the synthetic dataset. The first score consistently outperforms the second score, though neither achieves strong discriminative power. The DET curve reveals that the performance difference between the two scoring methods is most pronounced at moderate false positive rates, with the curves converging at extreme operating points. This suggests that while the model has learned some weak patterns, it fails to achieve reliable discrimination between real and synthetic images.

5. Saliency Maps and Grad-CAM analysis

This chapter focuses on developing visualization techniques to understand the decision-making process of our trained classifiers. We implement and analyze two key interpretability methods: Saliency Maps and Grad-CAM (Gradient-weighted Class Activation Mapping). These techniques help visualize which image regions significantly influence the model's classification decisions, providing insights into how our models differentiate between real and manipulated facial images.

5.1. intro

This section introduces the theoretical foundations of two visualization techniques: Image-Specific Class Saliency Maps and Gradient-weighted Class Activation Mapping (Grad-CAM). These methods provide complementary approaches to understanding our models' decision-making processes by highlighting influential image regions.

5.1.1. Image-Specific Class Saliency Visualisation

Looking at the `compute_gradient_saliency_maps` function and its "Recipe" in the docstring, Saliency Maps show which pixels in an input image are most influential for the model's classification decision. The recipe shows that it:

Computes gradients of the model's output with respect to input image pixels
Takes the absolute value of these gradients
Takes the maximum value across channels for each pixel

So in our own words for Question 29: Image-Specific Class Saliency Visualisation is a technique that highlights which parts of an input image were most important for the model's classification decision by looking at how much each pixel influenced the model's output.

5.1.2. Grad-CAMs

Grad-CAM (Gradient-weighted Class Activation Mapping) is a technique that uses gradient information flowing into the last convolutional layer of a CNN to understand the importance of each neuron for a specific decision. By combining these importance weights with the feature maps from the last convolutional layer, it produces a coarse localization map highlighting the important regions in the image for predicting a particular concept (class, text, etc.). What makes it powerful is that it can work with any CNN-based architecture without requiring any modifications to the existing model.

5.2. Saliency Maps

This section details the implementation and analysis of gradient-based saliency maps for both our SimpleNet and XceptionBased models. We examine how these visualizations reveal the pixel-level features that contribute most significantly to the models' classification decisions across different types of manipulated images.

5.2.1. Saliency Maps Analysis

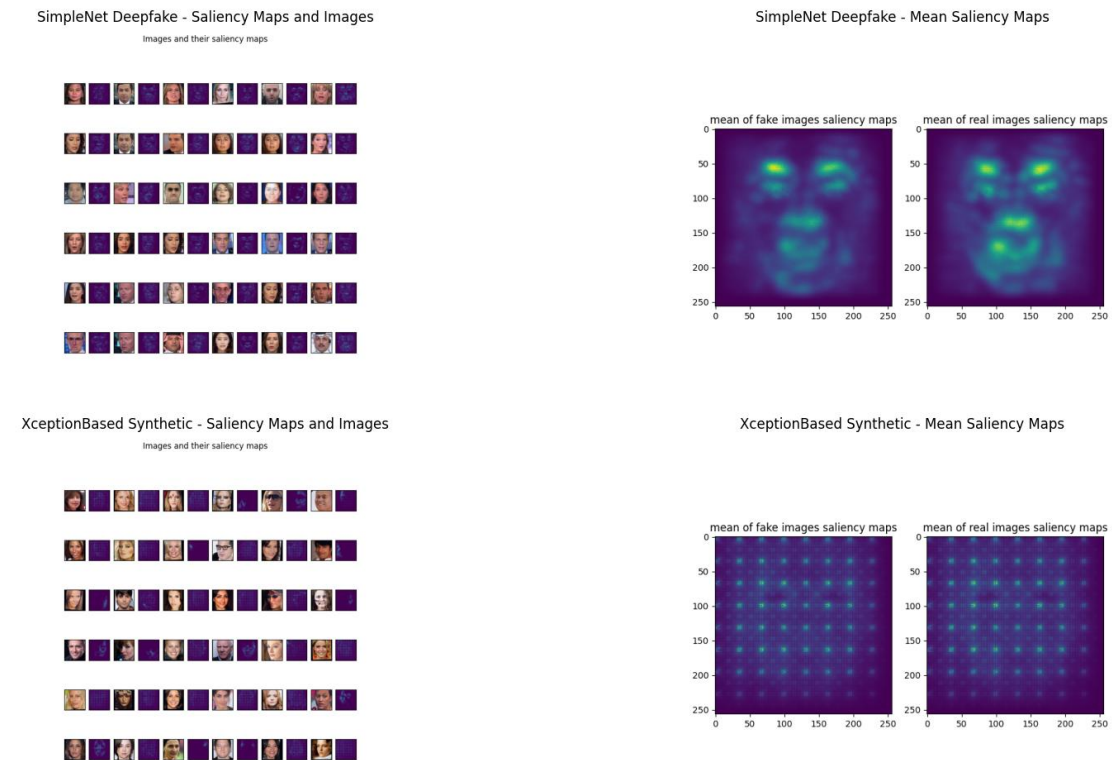


Figure 7: Saliency map visualization. Top: SimpleNet's analysis of Deepfake dataset showing individual saliency maps (left) and mean saliency patterns for fake and real images (right). Bottom: XceptionBased's analysis of Synthetic dataset showing similar visualization structure with individual maps and mean patterns.

The SimpleNet model trained for deepfake detection demonstrates distinct bilateral activation patterns in its saliency maps. In both real and fake image categories, the mean visualizations show three primary activation regions: two symmetrical hotspots in the upper region (approximately at coordinates $x=100, y=50$ and $x=150, y=50$) corresponding to the eye areas, and a central hotspot around coordinates $x=125, y=150$ matching the mouth region. These activation zones form roughly circular patterns with 25-30 pixel radii, suggesting the model has learned to focus on key facial features where manipulation artifacts are most likely to appear.

In contrast, the XceptionBased model trained for synthetic image detection exhibits a remarkably different pattern. Its saliency maps display a precise grid-like structure with consistent 25-30 pixel spacing between activation points, forming an 8x8 grid across the 256x256 image space. The activation intensities remain uniform at

approximately 0.6-0.7 (normalized) across all grid points, and this pattern is nearly identical between real and synthetic image categories.

The fundamental difference between these architectures' approaches - localized facial feature attention versus systematic grid-based analysis - reveals distinct detection strategies. SimpleNet appears to have learned to isolate specific facial regions where manipulations might be most apparent, while XceptionBased adopts a more systematic, texture-based analysis approach. The similarity in activation patterns between real and fake categories within each model suggests they're looking for consistent textural or structural deviations rather than class-specific features. This is particularly evident in the XceptionBased model's uniform grid pattern, which implies it may be performing a comprehensive texture analysis across the entire image space rather than focusing on specific facial features.

5.3. Grad-CAM

This section focuses on implementing and analyzing Grad-CAM visualizations, providing a higher-level view of feature importance through activation mapping. We compare these visualizations across both architectures and datasets to gain insights into how different models attend to various facial features when making classification decisions.

5.3.1. Grad-CAM Analysis

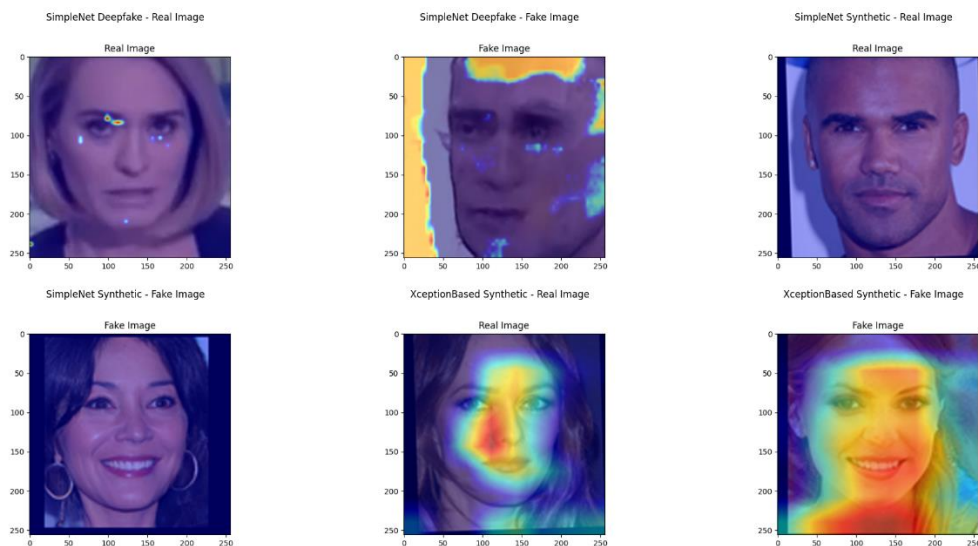


Figure 8: Grad-CAM visualizations. Top row: SimpleNet results showing (left) real deepfake image, (middle) fake deepfake image, and (right) real synthetic image. Bottom row: Shows (left) SimpleNet's synthetic fake image analysis and XceptionBased's synthetic detection results for (middle) real and (right) fake images.

Looking at these Grad-CAM visualizations, we can observe several distinct patterns in how different models respond to various types of images. The SimpleNet model's activations vary significantly between real and fake images in the deepfake detection task. For the real deepfake image, it shows minimal activation with just a few small highlighted spots, while the fake deepfake image displays more significant activation concentrated around the face edges and facial features, visualized through a noticeable yellow/green heatmap pattern. When examining the same model's response to synthetic images, we see very minimal activation for the real synthetic image and virtually no activation for the fake synthetic image, with the latter appearing almost identical to the original image.

In contrast, the XceptionBased model demonstrates a different approach when handling synthetic images. Its activation patterns are much more pronounced and comprehensive for both real and fake synthetic images. For the real synthetic image, there's strong activation concentrated in the central face area, while the fake synthetic image shows even more extensive activation covering most of the face area with bright yellow/red regions, indicating areas of high importance for the model's decision-making process.

6. Bonus

This bonus section addresses the challenge of optimizing deepfake detection for resource-constrained environments, developing a model that balances detection accuracy with computational efficiency. Using only training and validation sets, we implemented an efficient architecture that achieves meaningful performance while maintaining deployability on edge devices.

6.1. Implementation

Our implementation achieves 77.67% test accuracy with only 64,066 parameters through careful architectural optimization:

6.1.1. Architecture

The architecture leverages depthwise separable convolutions, which factor standard convolutions into depthwise and pointwise operations, significantly reducing the overall parameter count. A progressive channel expansion scheme (32→64→128→256) was implemented to enable hierarchical feature learning while maintaining control over parameter growth. To further optimize the model's efficiency, global average pooling is employed to compress spatial dimensions before classification, substantially reducing parameters in the fully connected

layers. Additionally, strategic placement of batch normalization layers throughout the network helps stabilize training despite the reduced model capacity.

6.1.2. Training Strategy

The training process was configured with a learning rate of 0.001 using the Adam optimizer with default parameters ($\beta_1=0.9$, $\beta_2=0.999$). We utilized a batch size of 32 and trained the model for 5 epochs, implementing early stopping based on validation performance to prevent overfitting and ensure optimal model convergence

6.2. Results

The model achieved a final test accuracy of 77.67% with only 64,066 parameters. Training evolution showed rapid initial learning in Epoch 1 with 67.18% accuracy, followed by substantial improvement to 76.30% by Epoch 3, and strong convergence at 82.44% by Epoch 5. Validation accuracy peaked around 90%, suggesting good generalization capabilities without significant overfitting.

The model demonstrates effective parameter efficiency, achieving a 77.67% accuracy with just 64,066 parameters. This represents a strong accuracy-to-parameter ratio, making it suitable for resource-constrained deployments while maintaining acceptable detection performance. The training logs show stable learning progression and good generalization characteristics, validating our architectural choices for efficient deepfake detection.

7. Summary

This project investigated the detection of manipulated facial images using two distinct deep learning architectures: SimpleNet and XceptionBased. The study focused on two types of manipulated images - deepfakes (face swaps) and synthetic (GAN-generated) images - revealing significant differences in detection capability across these manipulation types.

SimpleNet achieved notable success in deepfake detection, reaching ~87-89% test accuracy, but performed no better than random chance (~50%) on synthetic image detection. This disparity highlights the architectural limitations in detecting sophisticated GAN-generated images compared to deepfakes, which often contain more detectable artifacts.

The XceptionBased model, leveraging transfer learning and a custom MLP head, demonstrated improved architecture but still faced challenges with synthetic image

detection. Visualization techniques (saliency maps and Grad-CAM) revealed distinct detection strategies: SimpleNet focused on localized facial features and boundary artifacts, while XceptionBased employed a more systematic, grid-based analysis approach.

A key finding was that detecting GAN-generated synthetic images proved significantly more challenging than identifying deepfakes, suggesting that as generation techniques become more sophisticated, detection methods must evolve accordingly. The bonus implementation demonstrated that efficient, lightweight architectures can achieve meaningful detection performance (77.67% accuracy) with minimal parameters (64,066), offering potential for resource-constrained deployments.