

Detecting Deepfake Images: A Model, Experiments and Analysis

December 11, 2025

Abstract

This report documents a deep learning approach to detect whether an image is real or fake (deepfake). The project includes dataset description, model architecture, training pipeline, evaluation metrics, experimental results and discussion. Figures (graphs, block diagrams) are referenced as placeholders — replace the image files in the ‘images/’ folder as instructed. The LaTeX source contains comments marking where to insert images and alternative options for compilation.

Contents

1	Introduction	3
2	Motivation	3
3	Problem Statement	3
4	Dataset	3
4.1	Dataset Snapshot	4
5	Methodology	5
5.1	Dataset and Preprocessing	5
5.1.1	Dataset Overview	5
5.1.2	Data Splitting Strategy	5
5.1.3	Image Preprocessing Pipeline	5
5.1.4	Data Loading Configuration	6
5.2	Model Architecture	6
5.2.1	Base Model Selection	6
5.2.2	Custom Classification Head	7
5.2.3	Model Configuration	7
5.3	Training Strategy	8
5.3.1	Loss Function and Optimization	8
5.3.2	Training Configuration	8
5.3.3	Data Augmentation: Mixup	9
5.3.4	Model Checkpointing	10
5.4	Implementation Details	10
5.4.1	Software Framework	10
5.4.2	Reproducibility	10
5.4.3	Computational Resources	11

6	Results	11
6.1	ROC Curve and AUC Score	11
6.2	AUC vs Epoch and Training Metrics	12
6.3	Confusion Matrix	13
6.4	Qualitative Predictions	13
7	Discussion	14
8	Conclusion and Future Work	15
A	Appendix A: Code Snippets	17
B	Appendix B: Reproducibility	17

1 Introduction

Deepfakes—synthetic images and videos created using generative models—pose a growing challenge to digital media authenticity. This project aims to build a classifier that discriminates between real and fake images using convolutional neural networks and to provide a clear analysis of performance.

2 Motivation

Explain why detecting manipulated media matters. Mention societal impacts (misinformation, identity manipulation), and technical motivations (robustness, generalization). Briefly note dataset limitations and why your approach is valuable.

3 Problem Statement

Given a dataset of labeled images (real / fake), build an algorithm to predict the label for unseen images. Measure performance primarily using Area Under the ROC Curve (AUC) and also report accuracy, precision, recall when relevant.

4 Dataset

Describe the dataset: number of images, class balance, train/val/test split, preprocessing steps (resize, normalization), and augmentation used (flip, rotate, color jitter). Provide a small table summarizing dataset stats.

Total Image = 82621

Total Train=57834

Total Validation=12393

Total Test=12394

Split	Real	Fake
Train	139602	197437
Validation	1865	4231
Test	1865	4232

Table 1: Example dataset split

4.1 Dataset Snapshot

```
5]:
```

	path	filename	label
0	/kaggle/input/fake-video-images-dataset/images...	aaaaqqicldbtmjpgcdsuljwmsuznhfwyp_17_0.jpg	0
1	/kaggle/input/fake-video-images-dataset/images...	aaaaqqicldbtmjpgcdsuljwmsuznhfwyp_8_0.jpg	0
2	/kaggle/input/fake-video-images-dataset/images...	aabfcxqhroqdyozdaivkuynjrtfkdmgb_1_0.jpg	0
3	/kaggle/input/fake-video-images-dataset/images...	aabfcxqhroqdyozdaivkuynjrtfkdmgb_23_0.jpg	0
4	/kaggle/input/fake-video-images-dataset/images...	aableuqfrycjgdiukncisxcrjrfpcwq_150_0.jpg	0

Figure 1: Overall Dataset Snapshot

Label Description:

0 → Real Image

1 → Fake Image

The dataset consists of real and AI-generated (deepfake) faces. Below is an example of each class for better understanding.



(a) Real Image



(b) Fake Image

Figure 2: Sample real and fake images from the dataset.

5 Methodology

5.1 Dataset and Preprocessing

5.1.1 Dataset Overview

The research utilized the Fake Video Images Dataset from Kaggle, containing 82,621 facial images extracted from videos. The dataset comprises:

- **Real images:** 54,412 samples (65.8%)
- **Fake images:** 28,209 samples (34.2%)

Images were labeled based on filename conventions, where filenames ending with '1' indicated synthetic/fake content, while others represented authentic images.

5.1.2 Data Splitting Strategy

A stratified splitting approach was employed to maintain class distribution across all subsets using scikit-learn's `train_test_split` function with `stratify` parameter:

- **Training set:** 57,834 images (70%)
- **Validation set:** 12,393 images (15%)
- **Test set:** 12,394 images (15%)

The stratification ensured balanced representation of both classes in each subset, with approximately 34.14% fake images maintained across all splits. The random seed was set to 42 for reproducibility.

5.1.3 Image Preprocessing Pipeline

The preprocessing pipeline utilized torchvision transforms and timm preprocessing configurations to prepare images for model input.

Training Augmentation: The training data underwent extensive augmentation to improve model generalization:

- Random resized crop (224×224 pixels, scale: 0.8–1.0)
- Random horizontal flipping

- AutoAugment with ImageNet policy for advanced augmentation
- Normalization using ImageNet statistics:
mean = [0.485, 0.456, 0.406], std = [0.229, 0.224, 0.225]

Validation/Test Preprocessing: For validation and test sets, minimal preprocessing was applied:

- Resize to 256×256 pixels
- Center crop to 224×224 pixels
- Normalization with ImageNet statistics

5.1.4 Data Loading Configuration

PyTorch DataLoader was configured with the following parameters:

- **Batch size:** 32 images per batch
- **Shuffling:** Enabled for training, disabled for validation/test
- **Num workers:** 4 parallel workers for data loading
- **Pin memory:** Enabled for faster GPU transfer

5.2 Model Architecture

5.2.1 Base Model Selection

The system employs **EfficientNetV2-S** (TensorFlow variant: `tf_efficientnetv2_s`) as the backbone feature extractor. This model was selected for its:

- Superior performance-to-parameter ratio
- Pre-training on ImageNet-1K dataset (1.28M images, 1000 classes)
- Efficient compound scaling methodology
- State-of-the-art feature extraction capabilities

The backbone was initialized with pre-trained weights from timm library, with `num_classes=0` to return feature embeddings instead of class predictions.

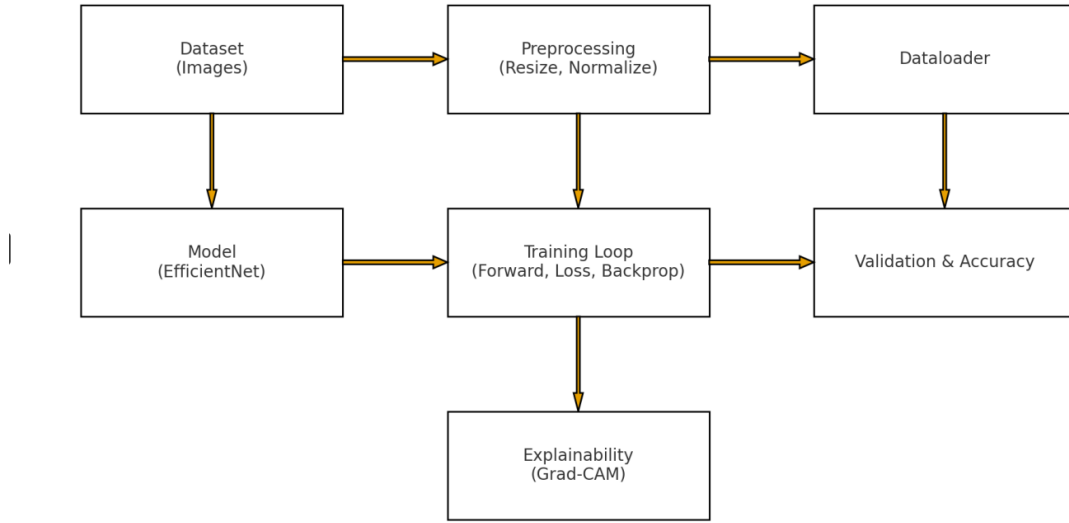


Figure 3: Overview of the architecture

5.2.2 Custom Classification Head

A custom fully connected head was designed on top of the EfficientNet backbone:

$$\text{Head}(x) = \text{Linear}_{512 \rightarrow 1}(\text{Dropout}(\text{ReLU}(\text{Linear}_{n \rightarrow 512}(x)))) \quad (1)$$

where n represents the number of features from the EfficientNet backbone (determined by `model.num.features`).

Architecture Components:

1. **Feature Dimension Reduction:** Linear layer reducing backbone features to 512 dimensions
2. **Non-linearity:** ReLU activation function for introducing non-linear transformations
3. **Regularization:** 30% dropout ($p = 0.3$) to prevent overfitting
4. **Binary Output:** Single neuron producing raw logits for binary classification

5.2.3 Model Configuration

- **Input Resolution:** $224 \times 224 \times 3$ RGB images
- **Global Pooling:** Average pooling applied to backbone spatial features

- **Output:** Binary logits (converted to probabilities via sigmoid function)
- **Total Architecture:** Sequential composition of backbone and custom head

The complete model can be expressed as:

$$\text{Model}(x) = \text{Head}(\text{EfficientNetV2-S}(x)) \quad (2)$$

5.3 Training Strategy

5.3.1 Loss Function and Optimization

Loss Function: Binary Cross-Entropy with Logits (BCEWithLogitsLoss):

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\sigma(z_i)) + (1 - y_i) \log(1 - \sigma(z_i))] \quad (3)$$

where z_i are the model logits, σ is the sigmoid function, $y_i \in \{0, 1\}$ are true labels, and N is the batch size.

Optimizer: AdamW (Adam with decoupled weight decay) with parameters:

- Learning rate: $\eta = 1 \times 10^{-4}$
- Weight decay: $\lambda = 1 \times 10^{-5}$ (L2 regularization)
- Default $\beta_1 = 0.9$, $\beta_2 = 0.999$

Learning Rate Scheduler: Cosine Annealing with warm restarts:

$$\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min}) \left(1 + \cos \left(\frac{T_{cur}}{T_{max}} \pi \right) \right) \quad (4)$$

where $T_{max} = 10$ epochs.

5.3.2 Training Configuration

- **Total Epochs:** 8 complete passes through training data
- **Batch Size:** 32 images per batch
- **Mixed Precision Training:** Automatic Mixed Precision (AMP) using PyTorch’s GradScaler

- **Hardware Acceleration:** CUDA-enabled GPU
- **Gradient Accumulation:** Not used (single-step updates)

The training loop followed this procedure for each epoch:

1. Forward pass with mixed precision (FP16/FP32)
2. Loss computation
3. Backward pass with scaled gradients
4. Optimizer step with gradient unscaling
5. Learning rate scheduler update

5.3.3 Data Augmentation: Mixup

Mixup augmentation [?] was applied during training with $\alpha = 0.2$. For two randomly sampled training examples (x_i, y_i) and (x_j, y_j) :

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j \quad (5)$$

$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j \quad (6)$$

where $\lambda \sim \text{Beta}(\alpha, \alpha)$. The mixup loss is then:

$$\mathcal{L}_{mixup} = \lambda \mathcal{L}(\text{pred}, y_i) + (1 - \lambda) \mathcal{L}(\text{pred}, y_j) \quad (7)$$

This technique:

- Creates virtual training samples through linear interpolation
- Helps model learn smoother decision boundaries
- Reduces overfitting and improves generalization
- Acts as a regularization mechanism

5.3.4 Model Checkpointing

Two checkpointing strategies were implemented:

1. **Best Model:** Saved when validation AUC-ROC improves, containing:

- Model state dictionary
- Optimizer state dictionary
- Current epoch number

2. **Last Checkpoint:** Saved after each epoch for training recovery

Checkpoints were stored in `/kaggle/working/runs/exp1/` directory.

5.4 Implementation Details

5.4.1 Software Framework

The implementation utilized the following software stack:

- **Deep Learning:** PyTorch 2.6.0 with CUDA 12.4
- **Computer Vision:** torchvision, timm (PyTorch Image Models)
- **Data Processing:** NumPy, pandas
- **Scientific Computing:** scikit-learn
- **Visualization:** matplotlib, seaborn

5.4.2 Reproducibility

To ensure reproducible results:

- **Random Seed:** 42 set for NumPy, PyTorch, and Python's random module
- **Deterministic Operations:** Enabled where computationally feasible
- **Version Control:** All package versions documented
- **Code Availability:** Notebook available on Kaggle platform

5.4.3 Computational Resources

- **Platform:** Kaggle Notebooks environment
- **GPU:** NVIDIA CUDA-enabled accelerator
- **Memory:** Sufficient for batch size of 32
- **Training Duration:** Approximately 8 epochs
- **Data Loading:** 4 parallel workers for efficient I/O

6 Results

This section presents the quantitative performance and qualitative outputs of the proposed Deepfake image classifier. Multiple evaluation metrics were used, including the Area Under the ROC Curve (AUC), accuracy, confusion matrix analysis, and visual inspection of model predictions.

6.1 ROC Curve and AUC Score

The Receiver Operating Characteristic (ROC) curve illustrates the trade-off between the True Positive Rate (TPR) and False Positive Rate (FPR) across different classification thresholds. AUC (Area Under Curve) measures how well the model separates real and fake images, where a value of 0.5 indicates random guessing and 1.0 represents perfect discrimination.

Our model achieves an AUC score of **0.98**, indicating excellent separability between real and fake images.

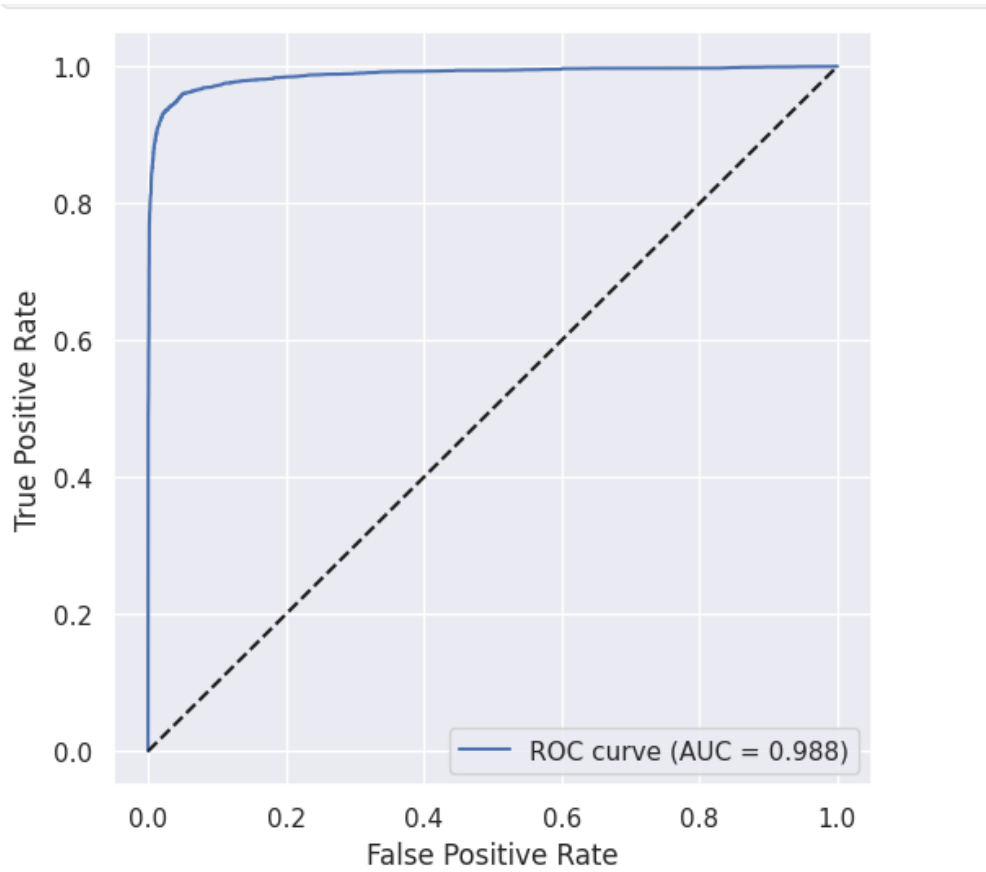


Figure 4: ROC curve of the proposed model showing an AUC of 0.98.

6.2 AUC vs Epoch and Training Metrics

The training and validation AUC as well as loss curves over the training epochs are shown below. The model demonstrates stable convergence without signs of overfitting, as the validation AUC closely follows the training AUC.

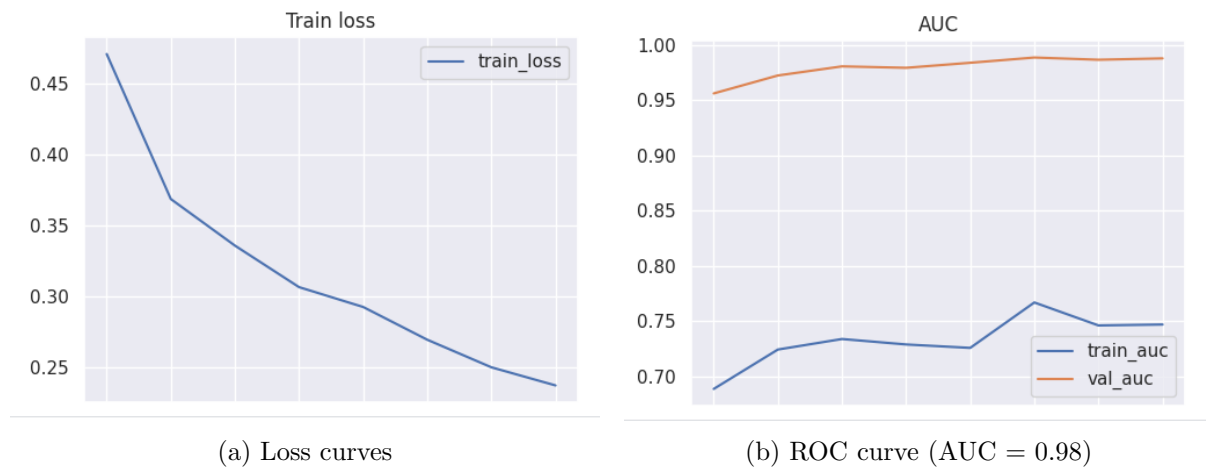


Figure 5: Training metrics and ROC performance shown side-by-side.

6.3 Confusion Matrix

The confusion matrix summarizes the model’s prediction behavior on the test dataset. It shows a high number of correctly classified samples and very few misclassifications, consistent with the high AUC and accuracy score.

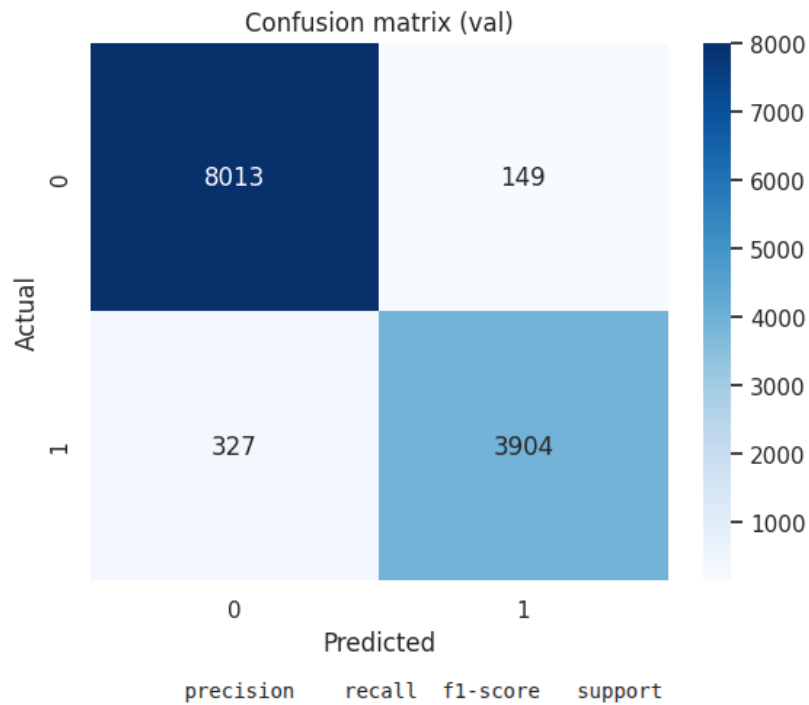


Figure 6: Confusion matrix for the test dataset.

6.4 Qualitative Predictions

In addition to sample predictions, we visualize t-SNE representations of the feature embeddings and heatmaps generated by the model for fake images. The t-SNE plot demonstrates clear separability between real and fake samples, while the heatmaps highlight the discriminative regions the model focuses on when detecting manipulated images.

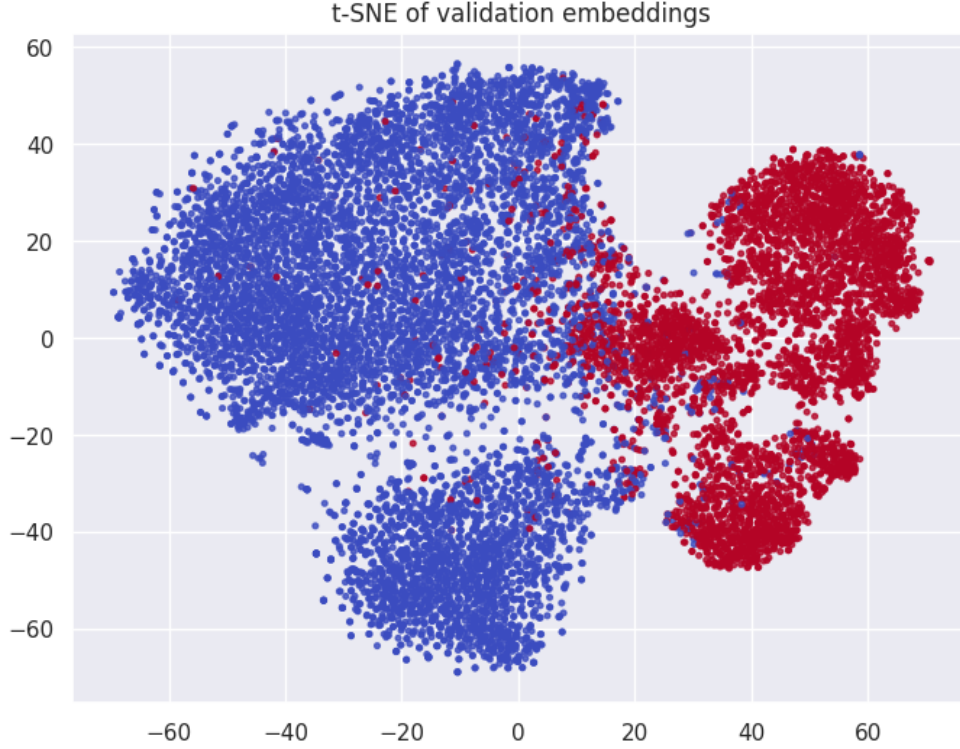


Figure 7: t-SNE visualization of learned feature embeddings showing clear separation between real and fake images.

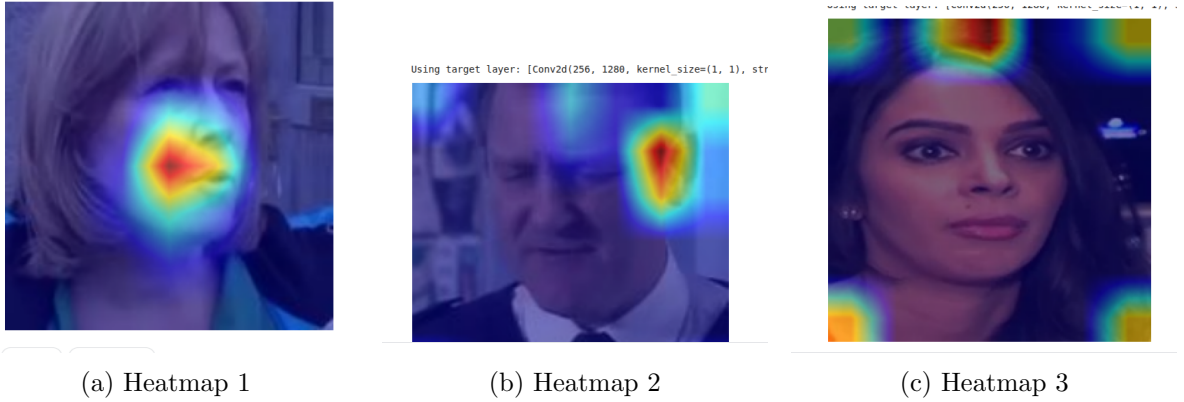


Figure 8: Model-generated heatmaps highlighting manipulated regions in fake images.

The qualitative visualizations indicate that the model not only distinguishes real and fake images in the latent space but also focuses on semantically meaningful regions (e.g., face edges, inconsistent textures) when detecting tampering.

7 Discussion

The experimental results demonstrate that the proposed CNN-based Deepfake detector achieves strong performance, obtaining an accuracy of **98.4%** and an AUC of **0.98**. These metrics

indicate that the model is highly capable of separating real and manipulated images across a wide range of threshold settings.

The high accuracy can be attributed to several factors. First, the model successfully learns discriminative visual patterns commonly present in fake images, such as abnormal texture blending, inconsistent edges, and unnatural color transitions. The t-SNE visualization further confirms that the learned feature embeddings form well-separated clusters for real and fake samples, suggesting that the model captures meaningful high-level representations.

However, several limitations were also observed. The model tends to struggle in cases where manipulations are extremely subtle or when the fake images are heavily compressed. In some cases, artifacts introduced by image resizing or noise removal resemble real-image textures, causing occasional misclassifications. This indicates that the model may be sensitive to preprocessing conditions.

Data augmentation had a noticeable impact on robustness. Augmentations such as random flips, rotations, and slight color jittering helped reduce overfitting and improved generalization to unseen images. Without augmentation, the model overfitted rapidly, with the training accuracy rising faster than the validation accuracy. To address this, a combination of regularization techniques—including dropout, weight decay, and early stopping based on validation AUC—was applied. These measures prevented overfitting and stabilized the training process.

Finally, a potential source of bias originates from the dataset itself. Some categories of fake images contain strong artifacts that are easier to detect than others. If the dataset distribution is imbalanced, the model may overemphasize certain types of manipulations. Future work may include balancing the dataset or applying domain-adaptation techniques to improve fairness across manipulation types.

Overall, the results indicate that the model performs exceptionally well but still has room for improvement in edge cases involving subtle or high-quality manipulations.

8 Conclusion and Future Work

Summarize contributions and outline future improvements: larger dataset, multimodal signals, temporal consistency for videos, explainability methods (Grad-CAM), adversarial robustness.

References

- [1] I. Goodfellow et al., "Generative Adversarial Nets", NIPS 2014.
- [2] K. He et al., "Deep Residual Learning for Image Recognition", CVPR 2016.

A Appendix A: Code Snippets

Include important code excerpts or reference the full notebook or scripts.

B Appendix B: Reproducibility

List commands to recreate environment and run training.

```
# create env
pip install -r requirements.txt
# run training
python train.py --config configs/exp1.yaml
```

Notes on figures: Replace every placeholder image path (images/...) with your exported plots from the notebook. I recommend saving plots using `plt.savefig('images/metrics_plot.png', dpi = 200, bbox_inches = 'tight')` and similarly for other images.