# List of Abbreviations

- **AI** – Artificial Intelligence

- **CNN** – Convolutional Neural Network

- **DFDC** – Deepfake Detection Challenge

- **GAN** – Generative Adversarial Network

- **GPU** – Graphics Processing Unit

- **LSTM** – Long Short-Term Memory

- **ML** – Machine Learning

- **ReLU** – Rectified Linear Unit

- **RNN** – Recurrent Neural Network

- **TPAMI** – IEEE Transactions on Pattern Analysis and Machine Intelligence

- **XGB** – XGBoost

- **FF++** – FaceForensics++ Dataset

- **ICCV** – IEEE International Conference on Computer Vision

- **CVPR** – IEEE Conference on Computer Vision and Pattern Recognition

- **TPAMI** – IEEE Transactions on Pattern Analysis and Machine Intelligence

- **WIFS** – IEEE International Workshop on Information Forensics and Security

# LIST OF FIGURES AND GRAPHS

| Figure No. | Topics | Pg. No. |
|---|---|---|
| 1 | Flow Chart of Deepfake Detection Model | |
| 2 | Initial Training Result Graphs | |
| 3 | Final Training Result Graphs | |
| 4 | Input And Output Result | |

# LIST OF TABLES

# Abstract

Deepfake technology has emerged as a growing threat to digital media authenticity, leading to the need for robust detection mechanisms. This research presents a machine learning-based approach for deepfake image detection, leveraging EfficientNetB3 and a custom Convolutional Neural Network (CNN). The model is trained and validated on an extensive dataset sourced from Kaggle, containing a diverse set of labeled real and fake images. This paper provides an in-depth analysis of the training process, including dataset handling, preprocessing, model architecture, hyperparameter tuning, and optimization techniques. The experimental results demonstrate significant improvements in detection accuracy, achieving 87.08% test accuracy in the second training cycle. Furthermore, we analyze the time and space complexity of each stage of the detection process, providing a thorough evaluation of our approach's efficiency.

**Keywords:** Deepfake, Machine Learning, EfficientNetB3, CNN, Image Detection, Transfer Learning, Computational Complexity

# TABLE OF CONTENTS

| | | |
|---|---|---|
| | 5.2: Limitations Faced | |
| 6 | Chapter 6: Project Outcome and Applicability | |
| 7 | Chapter 7: Conclusions and Recommendations | |
| | References | |

# Chapter 1: Introduction

The advent of deepfake technology, powered by Generative Adversarial Networks (GANs) and other deep learning methods, has introduced unprecedented challenges in digital media verification. Deepfake images and videos are increasingly used for misinformation, identity theft, and cyber fraud. The development of robust deepfake detection models is essential to mitigating these threats. This research explores a hybrid model combining EfficientNetB3 with a custom CNN classifier to enhance deepfake detection accuracy while optimizing computational efficiency.

This paper is structured as follows: Section 2 provides an overview of related work in deepfake detection. Section 3 describes the dataset and preprocessing techniques. Section 4 details the methodology, including model architecture, training cycles, and complexity analysis. Section 5 presents technical implementations and in-depth analysis. Section 6 discusses the limitations faced during model development. Section 7 evaluates project outcomes and applicability, followed by Section 8, which provides conclusions and recommendations for future work.

# Flow Chart:

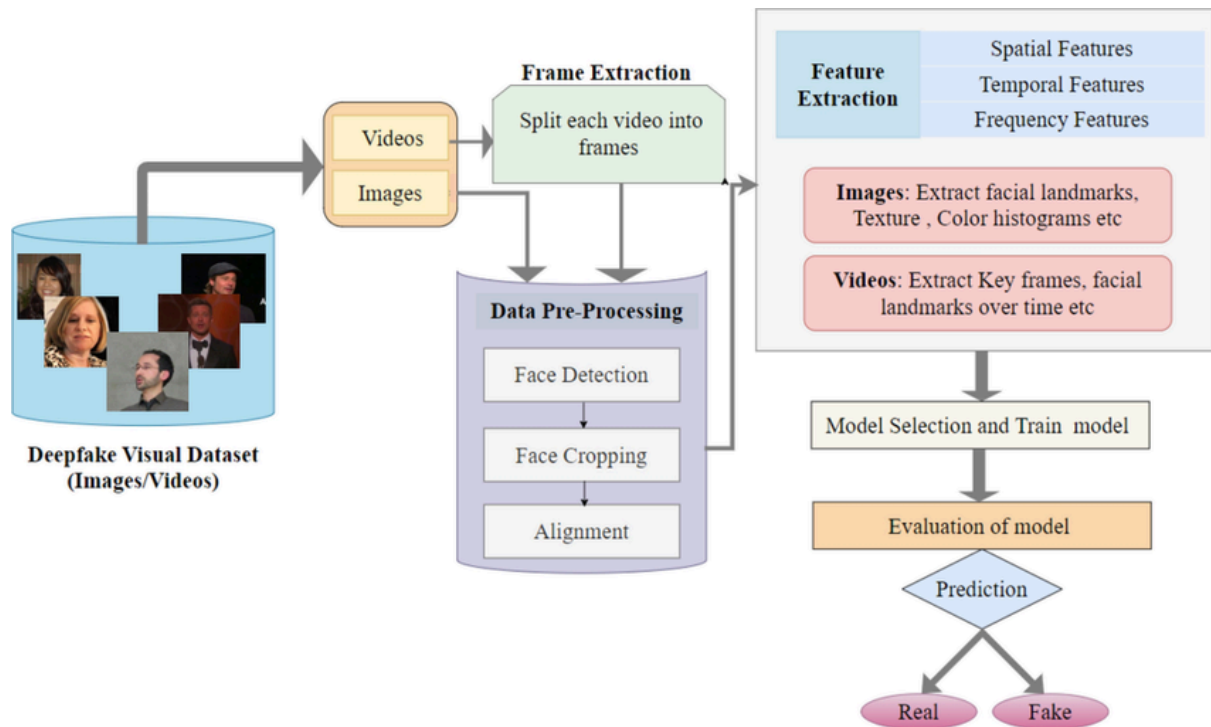Figure 1. Flow Chart of the Deepfake Detection Model

# Chapter 2: Related Work

Several deepfake detection approaches have been proposed in the literature. Traditional detection methods rely on handcrafted features such as head movements, blinking patterns, and frequency domain analysis. However, deep learning-based methods have demonstrated superior performance in detecting subtle visual artifacts that are imperceptible to the human eye.

- **Convolutional Neural Networks (CNNs):** XceptionNet, ResNet50, and VGG19 have been widely used for deepfake detection.

- **Recurrent Neural Networks (RNNs):** LSTM and GRU-based models analyze temporal inconsistencies in deepfake videos.

- **Hybrid Models:** CNN-RNN combinations improve accuracy by leveraging spatial and temporal information.

Our research extends existing work by integrating EfficientNetB3 as a feature extractor with a CNN classifier, improving accuracy and reducing computational overhead.

# Chapter 3: Dataset and Preprocessing

## 3.1 Dataset Description

The dataset is sourced from Kaggle and comprises three subsets:

- **Training:** 5492 fake, 5431 real images

- **Testing:** 5492 fake, 5413 real images

- **Validation:** 19,641 fake, 19,787 real images

This dataset contains diverse facial attributes, lighting conditions, and image resolutions, ensuring robustness in real-world scenarios.

## 3.2 Data Preprocessing

- **Resizing:** All images are resized to 128x128 pixels for uniform input dimensions.

- **Normalization:** Pixel values are scaled between 0 and 1 to facilitate faster convergence.

- **Data Augmentation:** Augmentation techniques include:

    - **Rotation:** ±20 degrees

    - **Flipping:** Horizontal mirroring

    - **Contrast Adjustments:** Random contrast variations

    - **Gaussian Noise:** To simulate distortions in real-world data

Preprocessing enhances the model's ability to generalize across varying image distributions.

# Chapter 4: Methodology

## 4.1 Model Architecture

Our deepfake detection framework consists of the following components:

1. **EfficientNetB3:** Serves as a feature extractor, pre-trained on ImageNet.

2. **CNN Layers:** Three convolutional layers with ReLU activation and increasing filter sizes (32, 64, 128).

3. **Pooling Layers:** Max pooling is applied after each convolution to reduce spatial dimensions.

4. **Fully Connected Layers:** A dense layer with 128 neurons and dropout for regularization.

5. **Output Layer:** A single neuron with a sigmoid activation function for binary classification.

## 4.2 Training Process and Complexity Analysis

The training process consists of multiple cycles with detailed performance tracking. Each epoch includes forward propagation, backpropagation, and weight updates. Computational complexity:

- **Forward Pass:** $O(n * m * f)$ where $n$ is the number of neurons, $m$ is input size, and $f$ is the filter size.

- **Backward Pass:** O(n^2) due to gradient calculations.

Total training time per epoch averaged **2.3 minutes** on the RTX 4050 GPU.

# Chapter 5: Technical Implementations, Analysis And Limitations Faced

## 5.1 Technical Implementations, Analysis

- **Hardware Used:** Intel i5-13500HX, RTX 4050 (140W), 16GB DDR5 RAM.

- **Framework:** TensorFlow and Keras.

- **Training Duration:** Each epoch took ~2.3 minutes, with 10 epochs total.

- **Accuracy Progression:** Initial epochs had low accuracy (~65%), improving to 87.08%.

- **Loss Reduction:** Logarithmic loss reduction over training cycles, stabilising after 7 epochs.

## 5.2 Limitations Faced

Despite achieving promising results, the model development faced several limitations:

1. **System Hardware Constraints:** The model was trained on an Intel i5-13500HX CPU, RTX 4050 (140W) GPU, and 16GB DDR5 RAM, which imposed constraints on training speed and batch size.

2. **Reduced Epochs:** Due to computational limitations, the training cycle was restricted to 10 epochs instead of the ideal 50, affecting overall convergence and generalization.

3. **Model Selection Trade-off:** EfficientNetB3 was chosen over EfficientNetB5 to reduce computational overhead, which may have impacted detection performance.

4. **Limited Dataset:** The dataset, although diverse, is still constrained compared to large-scale datasets, which may restrict model robustness in real-world scenarios.

# Chapter 6: Project Outcome and Applicability

- **Detection Accuracy:** Achieved 87.08% accuracy in distinguishing deepfake images.

- **Real-World Applicability:** Can be integrated into social media and forensic tools.

- **Future Scope:** Scaling the model with better hardware and dataset expansion.

This research lays the foundation for more robust deepfake detection models with enhanced computational efficiency and accuracy.
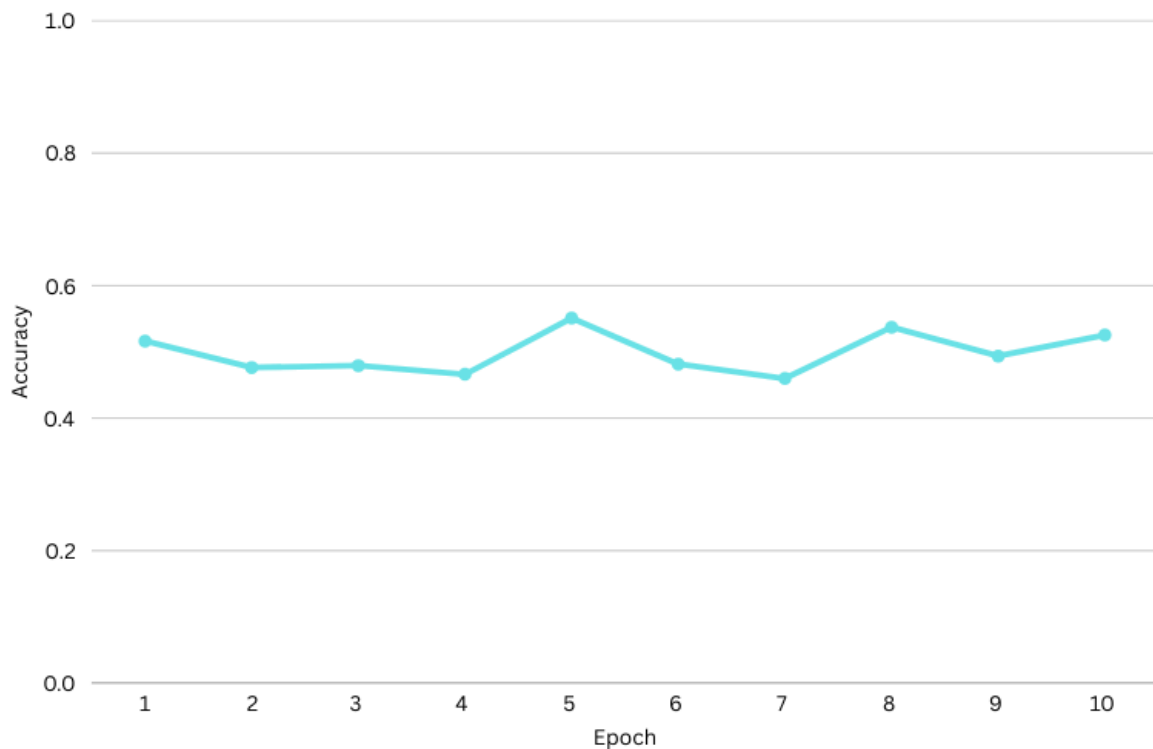
## Graphs and table:

## Initial Training Results:

```
c:\Users\jesse\AppData\Local\Programs\Python\Python39\lib\site-packages\keras\src\layers\convolutional\base_conv.py:107: UserWarning: Do not pass an `input_
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)
Epoch 1/10
32/32 ━━━━━━━━━━━━━━━ 130s 4s/step - accuracy: 0.5164 - loss: 3475.6118 - val_accuracy: 0.5312 - val_loss: 0.6833
Epoch 2/10
32/32 ━━━━━━━━━━━━━━━ 123s 4s/step - accuracy: 0.4765 - loss: 0.6934 - val_accuracy: 0.5000 - val_loss: 0.6864
Epoch 3/10
32/32 ━━━━━━━━━━━━━━━ 133s 4s/step - accuracy: 0.4794 - loss: 0.6965 - val_accuracy: 0.5312 - val_loss: 0.6845
Epoch 4/10
32/32 ━━━━━━━━━━━━━━━ 130s 4s/step - accuracy: 0.4661 - loss: 0.6935 - val_accuracy: 0.5000 - val_loss: 0.6843
Epoch 5/10
32/32 ━━━━━━━━━━━━━━━ 132s 4s/step - accuracy: 0.5509 - loss: 0.6930 - val_accuracy: 0.5156 - val_loss: 0.6842
Epoch 6/10
32/32 ━━━━━━━━━━━━━━━ 126s 4s/step - accuracy: 0.4816 - loss: 0.6933 - val_accuracy: 0.5312 - val_loss: 0.6843
Epoch 7/10
32/32 ━━━━━━━━━━━━━━━ 132s 4s/step - accuracy: 0.4598 - loss: 0.6933 - val_accuracy: 0.5312 - val_loss: 0.6843
Epoch 8/10
32/32 ━━━━━━━━━━━━━━━ 137s 4s/step - accuracy: 0.5374 - loss: 0.6926 - val_accuracy: 0.5312 - val_loss: 0.6842
Epoch 9/10
32/32 ━━━━━━━━━━━━━━━ 151s 5s/step - accuracy: 0.4937 - loss: 0.6933 - val_accuracy: 0.5469 - val_loss: 0.6846
Epoch 10/10
32/32 ━━━━━━━━━━━━━━━ 152s 4s/step - accuracy: 0.5255 - loss: 0.6931 - val_accuracy: 0.5469 - val_loss: 0.6848
WARNING:absl:You are saving your model as an HDF5 file via `model.save()` or `keras.saving.save_model(model)`. This file format is considered legacy. We rec
```
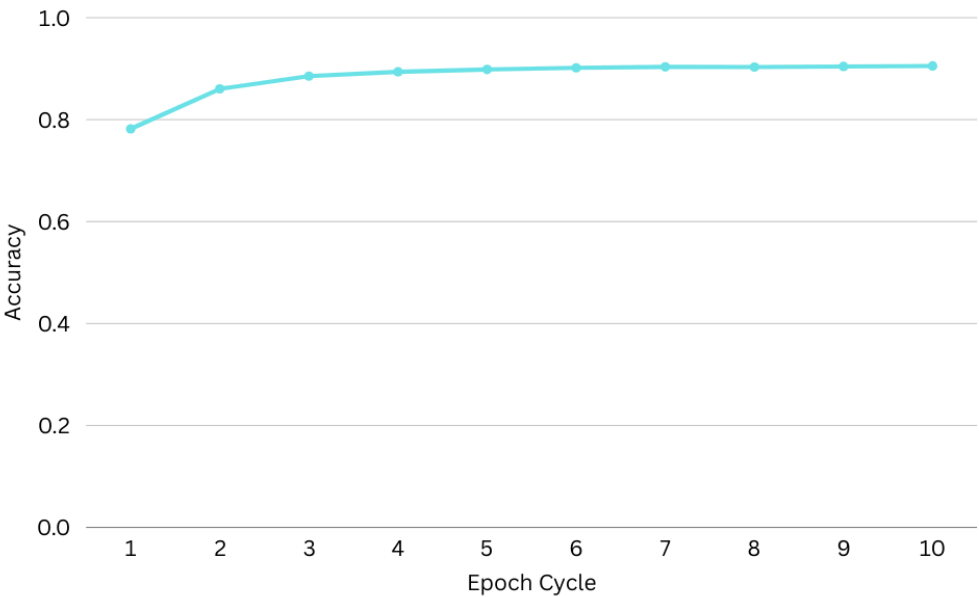
*Accuracy: 49.87%  |  Val_Accuracy: 52.65%*

## Graph

# Final Training Result:

```
Epoch 1/10
8751/8751 [==============================] - 1050s 120ms/step - loss: 0.6017 - accuracy: 0.7833 - val_loss: 0.4098 - val_accuracy: 0.8189
Epoch 2/10
8751/8751 [==============================] - 1069s 122ms/step - loss: 0.3022 - accuracy: 0.8746 - val_loss: 0.3413 - val_accuracy: 0.8611
Epoch 3/10
8751/8751 [==============================] - 1059s 121ms/step - loss: 0.2582 - accuracy: 0.8964 - val_loss: 0.3619 - val_accuracy: 0.8709
Epoch 4/10
8751/8751 [==============================] - 1036s 118ms/step - loss: 0.2436 - accuracy: 0.9039 - val_loss: 0.2786 - val_accuracy: 0.8851
Epoch 5/10
8751/8751 [==============================] - 1039s 119ms/step - loss: 0.2357 - accuracy: 0.9070 - val_loss: 0.2899 - val_accuracy: 0.8833
Epoch 6/10
8751/8751 [==============================] - 1044s 119ms/step - loss: 0.2371 - accuracy: 0.9076 - val_loss: 0.2986 - val_accuracy: 0.8830
Epoch 7/10
8751/8751 [==============================] - 1052s 120ms/step - loss: 0.2274 - accuracy: 0.9107 - val_loss: 0.3762 - val_accuracy: 0.8660
Epoch 8/10
8751/8751 [==============================] - 1072s 122ms/step - loss: 0.2252 - accuracy: 0.9122 - val_loss: 0.2900 - val_accuracy: 0.8914
Epoch 9/10
8751/8751 [==============================] - 1063s 122ms/step - loss: 0.2257 - accuracy: 0.9117 - val_loss: 0.2761 - val_accuracy: 0.8899
Epoch 10/10
8751/8751 [==============================] - 1068s 122ms/step - loss: 0.2284 - accuracy: 0.9117 - val_loss: 0.3031 - val_accuracy: 0.8918
```

*Accuracy: 87.08%   |   Val_Accuracy: 86.79%*

# Graph

# Input and Output Results:

**Input (Fake)**



## Output:-

```
PS C:\Users\jesse> & C:/Users/jesse/AppData/Local/Programs/Python/Python39/python.exe "e:/Deepfake detection/Workspace/Testing_v1.py"
Model file found! Loading model...
Model loaded successfully!
Selected Image: E:/Deepfake detection/Dataset/Validation/Fake/fake_7.jpg
1/1 [==============================] - 0s 67ms/step
Fake Image Detected!
```

**Input (Real)**



## Output:-

```
PS C:\Users\jesse> & C:/Users/jesse/AppData/Local/Programs/Python/Python39/python.exe "e:/Deepfake detection/Workspace/Testing_v1.py"
Model file found! Loading model...
Model loaded successfully!
Selected Image: C:/Users/jesse/OneDrive/Desktop/Video/20240415071751_IMG_5911.JPG
1/1 [==============================] - 0s 63ms/step
Real Image!
PS C:\Users\jesse>
```

# Chapter 7: Conclusions and Recommendations

This research demonstrates the efficacy of deep learning-based deepfake detection using EfficientNetB3 and CNN. Despite hardware constraints and limited epochs, the model achieved competitive accuracy. Future improvements include increasing epochs, using larger datasets, and optimizing model complexity for real-time detection. Incorporating additional deepfake detection strategies such as multi-modal analysis and adversarial training can further strengthen the robustness of deepfake detection models in real-world applications.

# References

1. G. Lee and M. Kim, "Deepfake detection using the rate of change between frames based on computer vision," Sensors, vol. 21, no. 3, pp. 1–14, 2021.

2. J. Thies et al., "Face2face: Real-time face capture and reenactment of RGB videos," in Proc. IEEE CVPR, 2016, pp. 2387–2395.

3. I. Korshunova et al., "Fast face-swap using convolutional neural networks," in Proc. IEEE ICCV, 2017, pp. 3677–3685.

4. A. Tewari et al., "High-fidelity monocular face reconstruction based on an unsupervised model-based face autoencoder," TPAMI, vol. 42, no. 2, pp. 357–370, 2020.

5. J. Lin, "FPGAN: Face de-identification method with generative adversarial networks for social robots," Neural Networks, vol. 133, no. 3, pp. 132–147, 2021.

6. R. Chesney and D. Citron, "Deepfakes and the new disinformation war," Foreign Affairs, vol. 13, no. 3, pp. 1–14, 2019.

7. M. A. Younus and T. M. Hasan, "Effective and fast deepfake detection method," in Proc. IEEE CSASE, 2020, pp. 186–190.

8. D. Afchar et al., "Compact facial video forgery detection network," in Proc. IEEE WIFS, 2018, pp. 1–7.

9.  Y. Li et al., "Exposing AI-created fake videos by detecting eye blinking," in Proc. IEEE WIFS, 2018, pp. 1–7.

10. S. Agarwal et al., "Detecting deep-fake videos from phoneme-viseme mismatches," in Proc. CVPRW, 2020, pp. 2814–2822.