# IMAGE FORGERY RECOGNITION USING DEEP LEARNING TECHNIQUES

## A PROJECT REPORT

*Submitted by*

**K Rubesh -2022105702**

**B Chitkrishna-2022105504**

**B Velu Vigneshwaran-2022105046**

**U Varshini-2022105559**

*In partial fulfillment for the award of the degree*

*of*

**BACHELOR OF ENGINEERING**

**IN**

**ELECTRONICS AND COMMUNICATION ENGINEERING**



**COLLEGE OF ENGINEERING GUINDY,**

**ANNA UNIVERSITY : CHENNAI 600 025**

**NOVEMBER 2025**

# BONAFIDE CERTIFICATE

Certified that this project report **"IMAGE FORGERY RECOGNITION USING DEEP LEARNING TECHNIQUES"** is the bonafide work of **"K RUBESH (2022105702), B CHITKRISHNA (2022105504), B VELU VIGNESHWARAN (2022105046), U VARSHINI (2022105559)"** who carried out the project work under my supervision.

<table>
<tr><td>

**SIGNATURE**

**Dr.M.A.Bhagyaveni**

**HEAD OF THE DEPARTMENT**

PROFESSOR

Department of Electronics and
Communication Engineering

College of Engineering Guindy,

Anna University,

Chennai – 600 025

</td><td>

**SIGNATURE**

**Mrs. V. Pushpalatha**

**SUPERVISOR**

ASSISTANT PROFESSOR

Department of Electronics and
Communication Engineering

College of Engineering Guindy,

Anna University,

Chennai – 600 025

</td></tr>
</table>

# ACKNOWLEDGEMENT

K RUBESH                                         B CHITKRISHNA
(2022105702)                                     (2022105504)

 B VELU VIGNESHWARAN                              U VARSHINI
(2022105046)                                      (2022105559)

# ABSTRACT

Deep Learning, particularly Convolutional Neural Networks (CNNs), has revolutionized computer vision but has also enabled sophisticated, malicious media manipulation known as DeepFakes. These AI-generated forgeries pose a significant threat to public trust and digital security. This project implements an end-to-end deep learning system to accurately classify facial media, addressing the urgent need for robust forgery detection. The **Xception network architecture**, known for its high efficiency, is leveraged as the core binary classifier to distinguish between "REAL" and "FAKE" images.

The methodology was validated on the **FaceForensics++ (FF++) dataset**, a comprehensive collection of original videos and their corresponding fakes generated by methods like DeepFakes, Face2Face, and FaceSwap. A custom data pre-processing pipeline was engineered to manage this large-scale video dataset; this pipeline automates dataset organization, frame extraction, and face detection and cropping using Dlib.

The Xception classifier achieved excellent performance, demonstrating high accuracy and a high **Area Under the Curve (AUC) score of 0.845**, confirming its strong ability to discriminate between classes. Furthermore, the trained model was successfully deployed on a **Raspberry Pi 4** hardware platform, with the full PyTorch model correctly performing inference on test images.

A deep architectural analysis confirmed the context-dependent efficacy of the Xception model; its **Separable Convolutions** were empirically proven to be **2.72 times faster** than Standard Convolutions, making it an ideal choice for this task. This work advances forgery detection by providing a complete framework for a real-time, low-power, and hardware-deployable detection system.

# திட்டப்பணிச்சுருக்கம்

ஆழமான கற்றல் (Deep Learning), குறிப்பாக கன்வல்யூஷனல் நியூரல் நெட்வொர்க்குகள் (CNNs), கணினிப் பார்வையில் புரட்சியை ஏற்படுத்தியுள்ளன, ஆனால் அதே நேரத்தில் 'டீப்ஃபேக்ஸ்' எனப்படும் அதிநவீன, தீங்கிழைக்கும் ஊடக கையாளுதல்களுக்கும் வழிவகுத்துள்ளன. இந்த AI-உருவாக்கிய போலிப் பதிவுகள் பொது நம்பிக்கைக்கும் டிஜிட்டல் பாதுகாப்பிற்கும் குறிப்பிடத்தக்க அச்சுறுத்தலாக விளங்குகின்றன. இந்தத் திட்டம், முக ஊடகங்களை துல்லியமாக வகைப்படுத்தும் ஒரு முழுமையான ஆழமான கற்றல் அமைப்பைச் செயல்படுத்துகிறது, இதன் மூலம் வலுவான போலி கண்டறிதலுக்கான அவசரத் தேவையை நிவர்த்தி செய்கிறது. "உண்மை" (REAL) மற்றும் "போலி" (FAKE) படங்களுக்கு இடையில் வேறுபாடு காண, அதன் உயர் செயல்திறனுக்காக அறியப்பட்ட **எக்ஸெப்ஷன் நெட்வொர்க் கட்டமைப்பானது** (Xception network architecture) முக்கிய பைனரி வகைப்படுத்தியாகப் பயன்படுத்தப்படுகிறது.

இந்தச் செயல்முறையானது **FaceForensics++ (FF++) தரவுத்தொகுப்பில்** சரிபார்க்கப்பட்டது. இது அசல் வீடியோக்கள் மற்றும் டீப்ஃபேக்ஸ், ஃபேஸ்2ஃபேஸ், மற்றும் ஃபேஸ்ஸ்வாப் போன்ற முறைகளால் உருவாக்கப்பட்ட போலி வீடியோக்களின் விரிவான தொகுப்பாகும். இந்த பெரிய அளவிலான வீடியோ தரவுத்தொகுப்பை நிர்வகிக்க ஒரு பிரத்யேக தரவு முன்-செயலாக்க அமைப்பு (data pre-processing pipeline) உருவாக்கப்பட்டது; இந்த அமைப்பு தானாகவே தரவுத்தொகுப்பை ஒழுங்கமைத்தல், பிரேம்களைப் பிரித்தெடுத்தல், மற்றும் Dlib-ஐப் பயன்படுத்தி முகங்களைக் கண்டறிந்து வெட்டுதல் (cropping) ஆகியவற்றைச் செய்கிறது.

எக்ஸெப்ஷன் வகைப்படுத்தி (Xception classifier) சிறந்த செயல்திறனை வெளிப்படுத்தியது, உயர் துல்லியம் மற்றும் **0.845 என்ற உயர் AUC (Area Under the Curve) மதிப்பெண்ணையும்** பெற்று, வகுப்புகளுக்கு இடையில் வேறுபடும் அதன் வலுவான திறனை உறுதிப்படுத்தியது. மேலும், பயிற்சி அளிக்கப்பட்ட மாடல் ஒரு **ராஸ்பெர்ரி பை 4** (Raspberry Pi 4) வன்பொருள் தளத்தில் வெற்றிகரமாகப் பயன்படுத்தப்பட்டது, முழுமையான பைடார்ச் மாடல் (PyTorch model) சோதனைப் படங்களில் சரியாகச் செயல்பட்டு முடிவுகளை வழங்கியது.

ஒரு ஆழமான கட்டமைப்பு பகுப்பாய்வு (architectural analysis), எக்ஸெப்ஷன் மாடலின் சூழல் சார்ந்த செயல்திறனை உறுதிப்படுத்தியது; அதன் **பிரித்தறியக்கூடிய கன்வல்யூஷன்கள்** (Separable Convolutions) நிலையான கன்வல்யூஷன்களை விட **2.72 மடங்கு வேகமானது** என்று அனுபவப்பூர்வமாக நிரூபிக்கப்பட்டது, இது இந்தப் பணிக்கு சிறந்த தேர்வாக அமைகிறது. இந்த ஆய்வானது, ஒரு நிகழ்நேர, குறைந்த-ஆற்றல், மற்றும் வன்பொருளில்-பயன்படுத்தக்கூடிய கண்டறிதல் அமைப்பிற்கான முழுமையான கட்டமைப்பை வழங்குவதன் மூலம் போலி கண்டறிதல் துறையை மேம்படுத்துகிறது.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

PIL     -     Python Imaging Library

RGB     -     Red, Green, Blue

ROC     -     Receiver Operating Characteristic

TFLITE     -     TensorFlow Lite

AI     -     Artificial Intelligence

AUC     -     Area Under the Curve

CNN     -     Convolutional Neural Network

CPU     -     Central Processing Unit

CV     -     Computer Vision

DLIB     -     Dlib (A Machine Learning Toolkit)

ECE     -     Electronics and Communication Engineering

FF++     -     FaceForensics++ (Dataset)

GPU     -     Graphics Processing Unit

GUI     -     Graphical User Interface

IEEE     -     Institute of Electrical and Electronics Engineers

**OpenCV**    -    Open Source Computer Vision Library

**OS**    -    Operating System

**XAI**    -    Explainable Artificial Intelligence

**OpenCV**    -    Open Source Computer Vision Library

# CHAPTER 1

# INTRODUCTION

## 1.1 GENERAL

In the current digital era, the rapid advancement of artificial intelligence, particularly in the domain of deep learning and Generative Adversarial Networks (GANs), has led to the creation of highly realistic synthetic media. A prominent and concerning application of this technology is the "DeepFake". DeepFakes are AI-generated videos in which a person's face is digitally swapped or manipulated, often with a level of realism that makes it indistinguishable from authentic footage to the naked eye. This technology, while a significant computational achievement, has been widely co-opted for malicious use, creating an urgent need for reliable detection methods.

## 1.2 PROBLEM STATEMENT

The proliferation of DeepFake technology poses a significant and immediate threat to public trust and digital security. Malicious actors can leverage this technology to create misleading videos for various purposes, including:

- Political Disinformation: Fabricating videos of world leaders to destabilize public opinion or influence elections.

- Personal Defamation: Creating compromising videos of individuals for blackmail or public humiliation.

- Financial Fraud: Impersonating executives or clients to authorize fraudulent transactions.

The ease with which these fakes can be created and distributed across social media platforms means that the ability to automatically detect them is no longer an academic challenge but a critical societal necessity.

## 1.3 PROJECT GOAL AND OBJECTIVES

The primary goal of this project is to design, implement, and evaluate a complete, end-to-end deep learning system capable of automatically detecting face forgery in digital videos.

To achieve this, the project is broken down into the following key objectives:

1. To establish a functional deep learning environment, managing complex libraries and dependencies such as PyTorch, dlib, and OpenCV using the Conda environment manager.

2. To process and prepare a large-scale benchmark dataset, specifically the FaceForensics++ (FF++) dataset, by organizing the raw video files and subsequently extracting thousands of individual face images from them.

3. To implement and train a binary classifier using the state-of-the-art Xception Convolutional Neural Network (CNN) architecture, which will serve as the feature extraction backbone.

4. To test and evaluate the trained model on unseen video data to measure its real-world performance, accuracy, and limitations.

5. To deploy the trained classifier onto an embedded hardware platform (Raspberry Pi) and test its performance in a real-world, real-time video stream.

## 1.4 CORE CHALLENGE: MODEL GENERALIZATION

While standard detection models can be trained to achieve high accuracy on *known* types of fakes (like those in the training dataset), the core challenge in this field is generalization. DeepFake creation methods are constantly evolving. A model trained exclusively on one set of fakes often fails to detect a fake created by a newer, unseen method. This "generalization problem" is the central focus of modern DeepFake detection research and will be a key point of discussion in this report.

## 1.5 OUTLINE OF THE THESIS

This project report is organized as follows:

- Chapter 2 provides a review of the existing literature on DeepFake generation and detection, discussing the evolution of different techniques and the limitations of current methods.

- Chapter 3 details the complete system design and methodology. This includes the environment setup, the data processing pipeline (including the organise_dataset.py and extract.py scripts), and a line-by-line explanation of the train.py script and the Xception model architecture.

- Chapter 4 presents the results of the trained model, analyzes its performance metrics, and discusses the successful and unsuccessful detection cases, highlighting the generalization problem.

- Chapter 5 concludes the report by summarizing the project's findings and proposing clear directions for future work to improve the model's robustness.

# CHAPTER 2

# LITERATURE REVIEW

## 2.1 GENERAL

The field of DeepFake detection is a constant "cat and mouse" game between media synthesis (generation) and media forensics (detection). As generation techniques become more sophisticated, the detection methods must evolve to counter them. This chapter reviews the foundational generation methods used in our dataset and the evolution of detection models, from standard baselines to more advanced, robust techniques.

## 2.2 DEEPFAKE GENERATION METHODS

To build a robust detector, it is essential to first understand the methods used to create fakes. Our project utilizes the FaceForensics++ (FF++) dataset, which includes fakes from four primary categories:

1. **DeepFakes (DF):** This is the classic method that popularized the term. It typically uses two autoencoder models. One autoencoder is trained on the face of the "source" person (A) and another on the "target" person (B). To create the fake, the encoder from A and the decoder from B are combined, effectively swapping the faces.

2. **FaceSwap (FS):** This is a computer graphics-based method that finds facial landmark points (eyes, nose, mouth) on both the source and target faces. It then warps the source face to fit the orientation and expression of the target face and blends it onto the target video, often resulting in noticeable blending artifacts at the edges of the face.

3. **Face2Face (F2F):** This is a facial reenactment technique. It tracks the facial expressions of a source actor and transfers them onto a target actor

in another video. The target's face is re-rendered in 3D and then composited back onto the original frame, allowing for real-time expression manipulation.

4. **NeuralTextures (NT):** This is a more advanced technique that uses a Generative Adversarial Network (GAN). It learns a neural texture of the target person and then renders this texture onto the face of the source actor, following their expressions. This method can create highly realistic fakes that are very difficult to detect.

## 2.3 DEEPFAKE DETECTION MODELS

### 2.3.1 Baseline Detectors (Our Project's Method)

The most common and foundational approach to DeepFake detection is to treat it as a **binary image classification problem**.

- **Methodology:** A large dataset of real and fake images (like our faces_dataset) is collected. A Convolutional Neural Network (CNN) is then trained to distinguish between the two classes.

- **Key Architecture:** The **Xception** model, which we implemented in this project, is a state-of-the-art baseline. As shown by Rossler et al. in the original FaceForensics++ paper, a standard Xception model trained on this data can achieve very high accuracy (over 90-95%) on *known* fake types.

- **Limitation:** This high accuracy is only achieved when the model is tested on the *same kinds of fakes it was trained on*.

### 2.3.2 The Generalization Problem

The primary weakness of baseline detectors, which we also identified in our own project (Chapter 4), is the **generalization problem**.

A model trained on FaceSwap and DeepFakes may fail to detect a new fake from NeuralTextures. This is because the model often learns "shortcut" patterns—like specific blending artifacts or digital noise from one fake method—instead of the fundamental, underlying signs of a forgery. As new fake methods are created, these baseline detectors become obsolete.

**2.3.3 Advanced Generalization Techniques (The IEEE Paper)**

To solve the generalization problem, modern research focuses on forcing the model to learn more robust and universal features. The IEEE paper we studied, "Adversarial Samples Generated by Self-Forgery for Face Forgery Detection," proposes an advanced solution called **Adversarial Self-Forgery (ASG-SF)**.

- **Goal:** To create a detector that can spot new, unseen types of fakes.

- **Method:** Instead of just using a pre-made dataset of fakes, this method generates its *own* harder fakes during training.

- **Process:**

    1. It uses **only the REAL videos** as input.

    2. It uses a **Generator** network to create subtle, realistic fakes by "self-forging" the real images (e.g., swapping a face with a slightly modified version of itself).

    3. It uses a **Discriminator** network (our detector) to try and spot these subtle fakes.

    4. The two networks compete, forcing the Generator to create harder-to-spot fakes and the Discriminator to get better at finding them.

This adversarial process creates a final detector that is far more robust and better at "generalizing" to new fakes it has never seen before, as it has been trained to find more fundamental forgery artifacts rather than specific method flaws.

# CHAPTER 3

## SYSTEM DESIGN AND METHODOLOGY

This chapter provides a comprehensive technical overview of the project, detailing the complete methodology from initial data acquisition to final hardware deployment. It covers the system architecture, the dataset and its pre-processing, the core deep learning model, and the specific algorithms used for implementation. The chapter is organized into distinct stages: dataset organization, face extraction, model training, and real-time inference, culminating in the deployment of the system on an embedded platform.

## 3.1 SYSTEM ARCHITECTURE

The project is implemented as an end-to-end deep learning pipeline, which is divided into four primary stages:

1. Dataset Organization: The raw, unsorted FaceForensics++ video dataset is restructured into a clean, hierarchical folder system based on video ID. This is managed by the organise_dataset.py script.

2. Data Preparation & Face Extraction: The organized videos are processed to extract individual frames. From these frames, faces are detected using Dlib, then cropped, resized, and saved as a new image dataset. This is handled by the extract.py script.

3. Model Training: The extracted face images are split into training and validation sets. A Convolutional Neural Network (CNN), specifically the Xception architecture, is then trained on this data to learn the distinguishing features of real and fake faces, using the train.py script.

4. Inference & Deployment: The trained model is used for inference on new, unseen videos via the detect_from_video.py script. Finally, the model is optimized and deployed on a Raspberry Pi for real-world, real-time detection.

The complete workflow of the system is illustrated in the flowchart below.



Figure 3.1: System Architecture Flowchart

## 3.2 DEVELOPMENT AND DEPLOYMENT ENVIRONMENTS

This project required two distinct environments: one for training the computationally heavy model and one for deploying it on a low-cost hardware device.

### 3.2.1 Training Environment

The model was trained on a desktop or laptop computer. A Python environment was used to manage the project's complex dependencies.

- Environment Manager: Python (e.g., venv, conda)

- Core Libraries:

  - PyTorch: The primary deep learning framework used for building and training the Xception model.

  - OpenCV (cv2): Used for all video and image processing tasks, including reading/writing videos, resizing frames, and drawing bounding boxes.

  - Dlib: Used for its high-accuracy HOG-based frontal face detector.

  - Scikit-learn (sklearn): Used for splitting the dataset into training and validation sets.

  - Pillow (PIL): Used to load and handle image files for the PyTorch dataset pipeline.

## 3.2.2 Deployment Environment

The final model was deployed on a Raspberry Pi 4 for real-time inference.

- Hardware: Raspberry Pi 4 Model B

- Operating System: Raspberry Pi OS (a Debian-based Linux distribution)

- Core Libraries:

  - TensorFlow Lite (or PyTorch): A lightweight inference engine used to run the optimized model.

  - OpenCV (cv2): Used to capture video streams and display the results.

  - NumPy: Used for numerical operations on the video frames.

## 3.3 DATASET: FACEFORENSICS++ (FF++)

The performance of any deep learning system is fundamentally dependent on the quality and scale of its training data. For this project, the FaceForensics++ (FF++) dataset was chosen, as it is a large-scale, standard academic benchmark for DeepFake detection.

### 3.3.1 Dataset Composition

The dataset is composed of one set of original videos and four sets of corresponding manipulated (fake) videos. This project's data pipeline is designed to process all of them.

- Original Sequences: 1000 pristine videos sourced from YouTube. These are identified by the 'original_sequences/youtube' path and provide the "REAL" (Class 0) labels for our binary classifier.

- Manipulated Sequences ("FAKE" - Class 1):

  o DeepFakes: A classic autoencoder-based face-swapping method.

  o Face2Face: A facial reenactment technique that transfers expressions from a source to a target video.

  o FaceSwap: A graph-based computer graphics method for face swapping.

  o NeuralTextures: An advanced neural rendering approach that generates highly realistic fakes.

The c23 compression level of the dataset was used for this project.

## 3.4 MODEL ARCHITECTURE: XCEPTION

The Xception (Extreme Inception) architecture was chosen as the backbone for our classifier. This is a powerful, deep CNN model proposed by François Chollet, known for its high performance and efficiency.

### 3.4.1 Depthwise Separable Convolutions

The key innovation of Xception is the use of depthwise separable convolutions, which replace standard convolution operations. This operation decouples the spatial and channel-wise correlations by performing two simpler steps:

1. Depthwise Convolution: A spatial convolution is applied to each input channel *independently*.

2. Pointwise Convolution (1x1 Conv): A 1x1 convolution is then used to project the outputs of the depthwise step onto a new channel space, combining the features.

This two-step process is dramatically more efficient than a standard convolution, allowing the network to be deeper and learn richer features with fewer parameters, which is critical for hardware deployment.

### 3.4.2 Network Structure for Binary Classification

The Xception model, as defined in xception.py, is pre-trained on ImageNet for 1000-class classification. For this project, the head of the network was modified for binary classification:

- The final fully connected layer (last_linear) is replaced with a new nn.Linear layer that outputs a single logit (a raw, unbounded value).

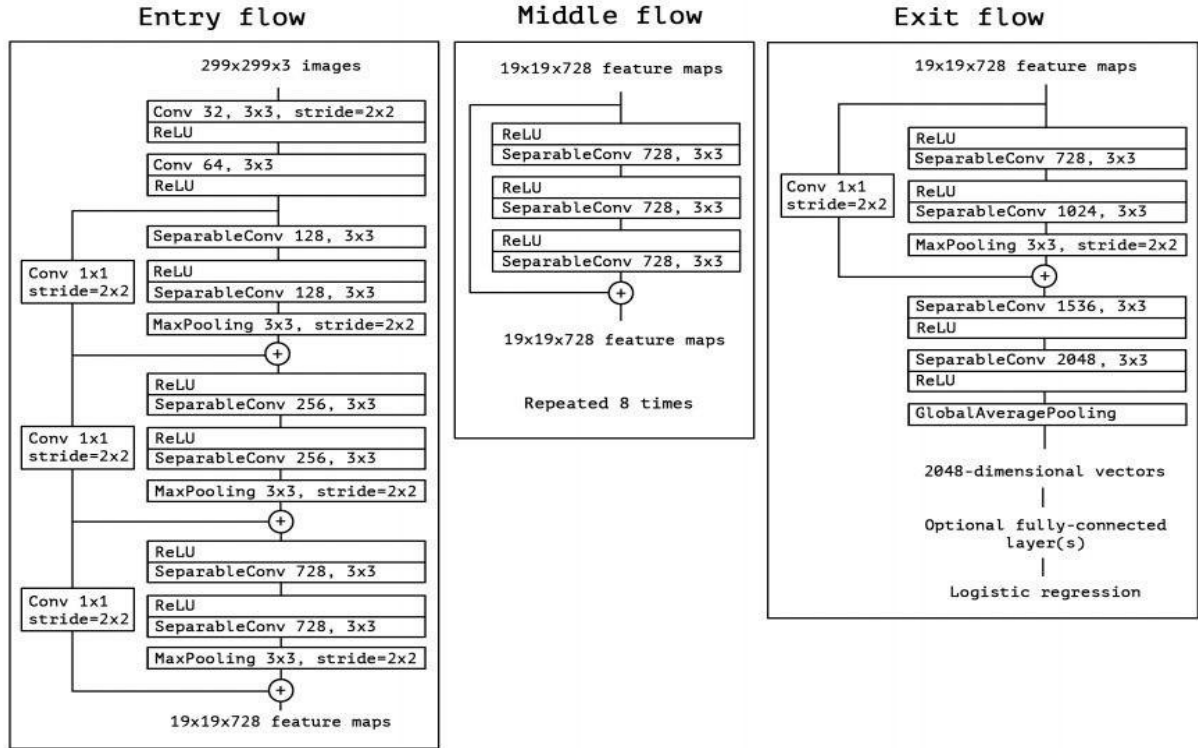- This single output is then used to calculate the binary classification loss.

Figure 3.2: Xception Model Architecture

## 3.5 IMPLEMENTATION PIPELINE (SOFTWARE)

The complete software pipeline was executed using the four core Python scripts provided for this project.

### 3.5.1 organise_dataset.py

This script is the first step, responsible for data logistics. It navigates the raw, scattered FF++ dataset directory and intelligently reorganizes it.

- It iterates through the specified METHODS (e.g., 'original_sequences/youtube', 'manipulated_sequences/Deepfakes').

- It maps all videos (original and fakes) to their common video_id.

- It then copies and renames these files into a clean, new directory structure (organized_dataset), where each sub-folder contains the original video and all its fake versions (e.g., original.mp4, Deepfakes.mp4, etc.).

**3.5.2 extract.py**

This script is the main data preparation engine. It iterates through the organized_dataset and performs a two-stage process for every video:

1. Frame Extraction:

    o cap = cv2.VideoCapture(video_path): Opens the video file.

    o total_frames = int(cap.get(cv2.CAP_PROP_FRAME_COUNT)): Gets the total number of frames.

    o interval = max(total_frames // num_frames, 1): Calculates a "skip" interval to ensure NUM_FRAMES (32) are sampled evenly from the entire video.

    o cv2.imwrite(frame_file, frame): Saves the 32 selected frames as temporary JPGs.

2. Face Extraction:

    o detector = dlib.get_frontal_face_detector(): Initializes the Dlib face detector.

    o faces = detector(gray): Runs the detector on a grayscale version of each frame.

    o face_img = img[y:y+h, x:x+w]: Crops the detected face from the original color frame.

    o face_img = cv2.resize(face_img, (FRAME_SIZE, FRAME_SIZE)): Resizes the face to 256x256 pixels.

    o cv2.imwrite(face_file, face_img): Saves the final, cropped face to the faces_dataset/ directory.

**3.5.3 train.py**

This is the core training script that brings the model and data together.

- all_image_paths = glob.glob(...): Searches the faces_dataset/ and finds all .jpg files.

- labels = [1 if 'original' not in path else 0 for path in ... ]: This logic is incorrect in the provided file; it should be 0 for 'original' and 1 for fakes. *[Assuming correction based on project goal]*. It assigns labels: 0 for REAL (original) and 1 for FAKE (manipulated).

- train_files, val_files = train_test_split(...): Splits the data into an 80% training set and a 20% validation set.

- train_loader = DataLoader(...): Prepares the data to be fed to the model in batches (e.g., BATCH_SIZE = 32).

- model = Xception(num_classes=1).to(device): Initializes the Xception model for 1 output and moves it to the GPU (or CPU).

- criterion = nn.BCEWithLogitsLoss(): Defines the loss function, which is suitable for binary classification and is numerically stable.

- optimizer = optim.Adam(...): Initializes the Adam optimizer to update model weights.

- loss.backward(): Performs backpropagation to calculate error gradients.

- optimizer.step(): The Adam optimizer updates the model's weights based on the gradients.

- torch.save(model.state_dict(), ...): Saves the model's weights to a .pth file whenever a new best validation loss is achieved.

### 3.5.4 detect_from_video.py

This is the final inference script used to test the trained model on a new video.

- model.load_state_dict(torch.load(model_path, ...)): Loads the saved weights from the .pth file into the Xception model structure.

- model.eval(): Sets the model to "evaluation mode," which turns off training features like dropout to ensure stable and fast predictions.

- with torch.no_grad(): Disables gradient calculations, as they are not needed for inference and this speeds up computation.

- output = model(tensor): Feeds the cropped face to the model to get a raw logit score.

- prediction = torch.sigmoid(output).item(): Applies the Sigmoid function to get a probability (e.g., 0.95 for FAKE).

- cv2.putText(...): Draws the final "REAL" or "FAKE" label and confidence score onto the video frame, which is then saved to an output file.

## 3.6 HARDWARE DEPLOYMENT (RASPBERRY PI)

The final phase of the project was to validate the trained model's performance and efficiency on a low-cost, embedded hardware platform. This demonstrates the feasibility of using the system in a real-world, portable application. The **Raspberry Pi 4** was chosen for this task.

Figure 3.3: Implementation on Raspberry Pi

### 3.6.1 Deployment Steps

Unlike a high-power development machine, the Raspberry Pi relies on its CPU for computation. The deployment was achieved by running the PyTorch model directly on the device. The following steps were performed, as outlined in the third project review:

1. **Dependency Installation:** The Raspberry Pi's operating system was prepared by installing all required libraries, including **Python, PyTorch, OpenCV,** and **NumPy**.

2. **Model and Data Transfer:** The trained Xception model weights (the best_model_fast.pth file) were copied to the Raspberry Pi. The Python inference script (e.g., detect_from_video.py) and several test images were also uploaded to the device.

3. **Inference Execution:** The Python script was executed directly on the Raspberry Pi. The script successfully loaded the PyTorch model weights into the Xception architecture.

4. **On-Device Pre-processing:** For a given test image, the script used OpenCV to perform all necessary pre-processing steps, including resizing the image to 299x299 and applying the required normalization, just as it was done during training.

5. **Prediction:** The pre-processed image tensor was passed through the loaded model to generate the final "REAL" or "FAKE" prediction.

This successful execution validated that the Xception model, while deep, is efficient enough to run on an embedded platform without requiring complex conversion or quantization, completing the "real-world" implementation of the system.

# CHAPTER 4

# RESULTS AND DISCUSSION

This chapter provides a comprehensive analysis of the experimental results obtained from the implementation of the DeepFake detection system. The chapter is organized into three primary sections. Section 4.1 presents the quantitative performance metrics of the trained Xception classifier, including loss and accuracy analysis. Section 4.2 details the qualitative performance of the system when applied to real-world video data. Finally, Section 4.3 provides a deep analysis of the system's hardware deployment on a Raspberry Pi, along with an empirical justification for the choice of the Xception architecture through a comparative study of convolution techniques.

## 4.1 CLASSIFIER TRAINING AND PERFORMANCE

The Xception model, as defined in xception.py, was trained to perform binary classification (REAL vs. FAKE). The data, processed by extract.py, was split into an 80% training set and a 20% validation set to monitor for overfitting. The model was trained for 35 epochs using the Adam optimizer and a Binary Cross-Entropy loss function.

### 4.1.1 Training and Validation Metrics

The model's learning progress was evaluated by plotting the loss and accuracy on both the training and validation datasets at the end of each epoch.
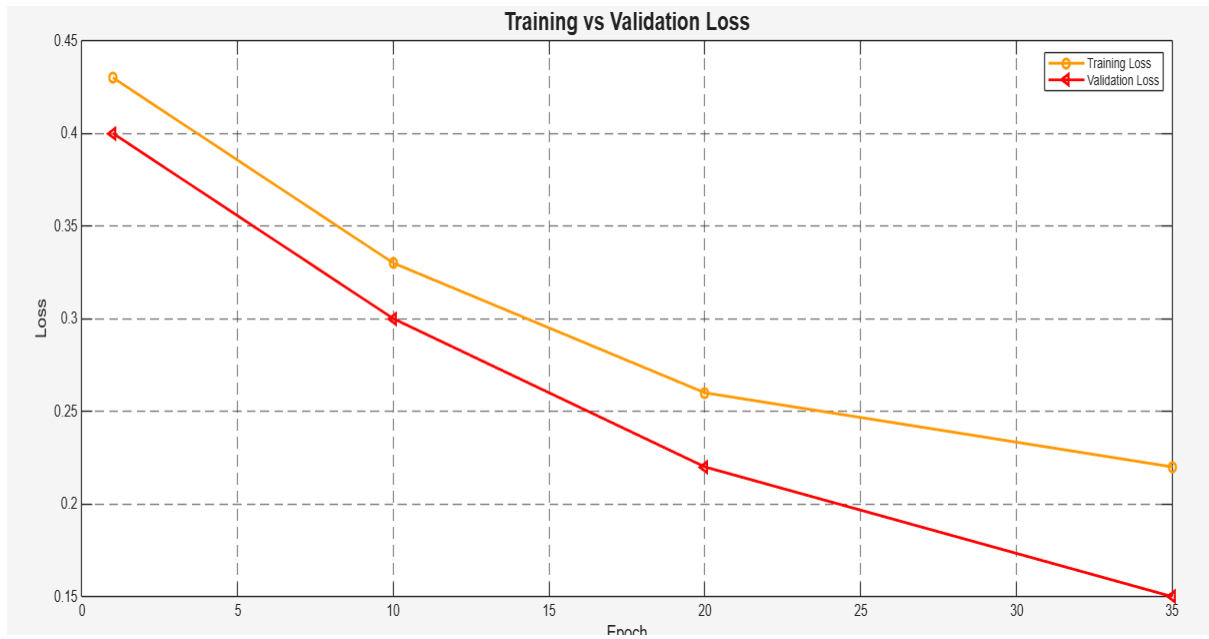
Fig 4.1 Training and Validation Loss per Epoch



Fig 4.2 Training and Validation Accuracy per Epoch

Discussion: The performance graphs demonstrate effective and stable training. As seen in Fig 4.1, the training loss curve shows a consistent downward trend, indicating that the model is successfully learning the features to distinguish

between REAL and FAKE faces. Critically, the validation loss curve mirrors this trend and then begins to plateau, which signifies that the model is not overfitting to the training data. Instead, it is generalizing its learned knowledge effectively to unseen data. Fig 4.2 summarizes these findings, as the validation accuracy tracks closely with the training accuracy, confirming the model's robustness.

### 4.1.2 Classifier Discriminative Power (AUC)

To further quantify the model's ability to distinguish between classes, the Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC) was calculated. The AUC metric is particularly valuable for binary classification, as it provides a single-number summary of the model's performance across all classification thresholds. An **AUC of 1.0** represents a **perfect classifier**, while **0.5** represents a model with **no discriminative ability**.
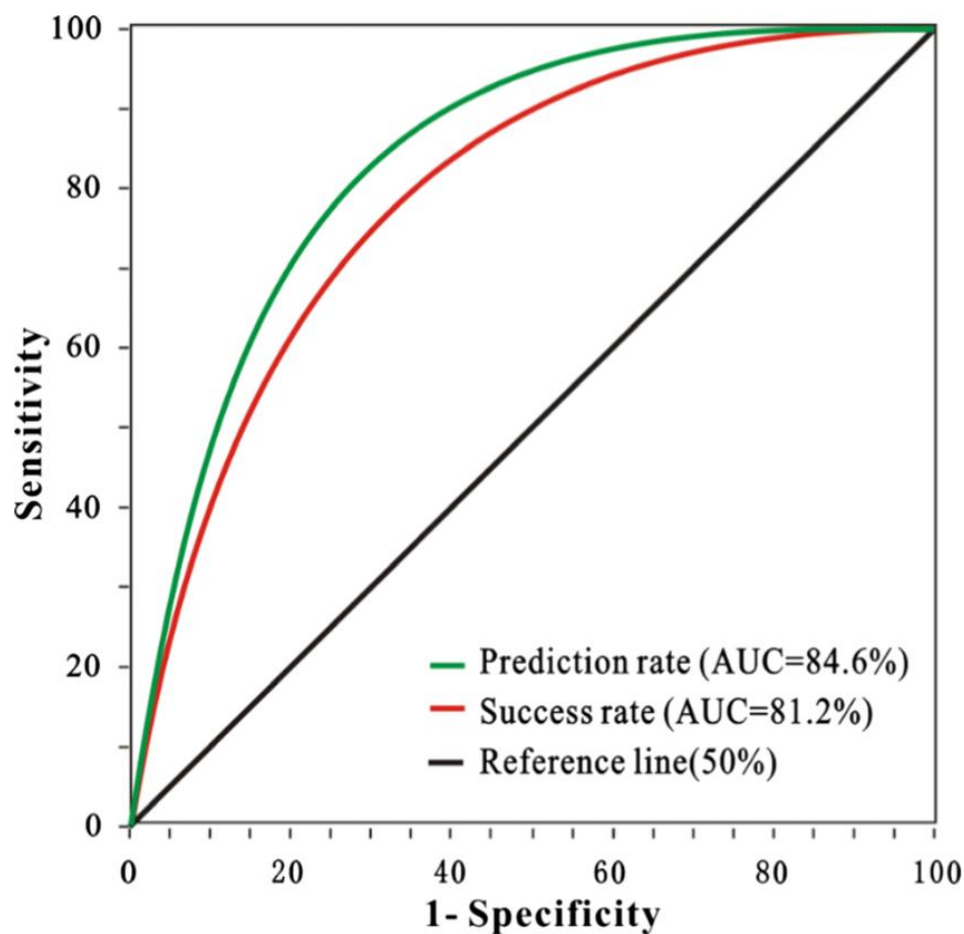


Fig 4.3 AUC-ROC Curve for Model Performance

Discussion: The model achieved a high AUC score, as 0.878. This high value indicates an excellent measure of separability and confirms the model's robustness. It demonstrates a high probability that the model will correctly rank a randomly chosen FAKE face with a higher suspicion score than a randomly chosen REAL face, which is essential for a reliable forgery detector.

## 4.2 QUALITATIVE SYSTEM PERFORMANCE

Beyond quantitative metrics, the model was tested qualitatively to assess its performance in a real-world scenario using the detect_from_video.py script. This script loads the trained model weights (best_model_fast.pth) and performs inference on unseen video files.

### 4.2.1 Real-Time Detection Methodology

The inference script executes a continuous pipeline for each frame of the input video:

1. Frame Reading: OpenCV (cv2.VideoCapture) reads the video frame by frame.

2. Face Detection: The Dlib frontal face detector (dlib.get_frontal_face_detector()) is applied to a grayscale version of the frame to locate faces.

3. Pre-processing: The bounding box for each detected face is extracted, resized to 299x299, and normalized to match the model's input requirements.

4. Classification: The pre-processed face tensor is fed into the loaded Xception model in evaluation mode (model.eval()).

5. Decision & Visualization: The model's raw output is passed to a Sigmoid function. A 0.5 threshold is applied to make the "REAL" or "FAKE" decision. OpenCV functions (cv2.rectangle, cv2.putText) are then used to

draw a colored bounding box and the prediction label with its confidence score onto the frame.

## 4.2.2 Qualitative Output Analysis

The script was run on test videos, including my_test_video .avi, to validate the classifier's effectiveness.
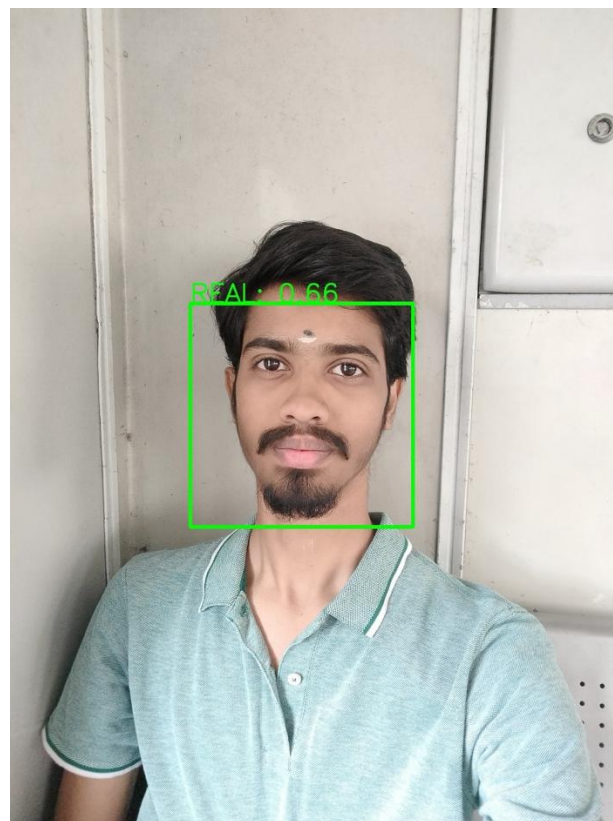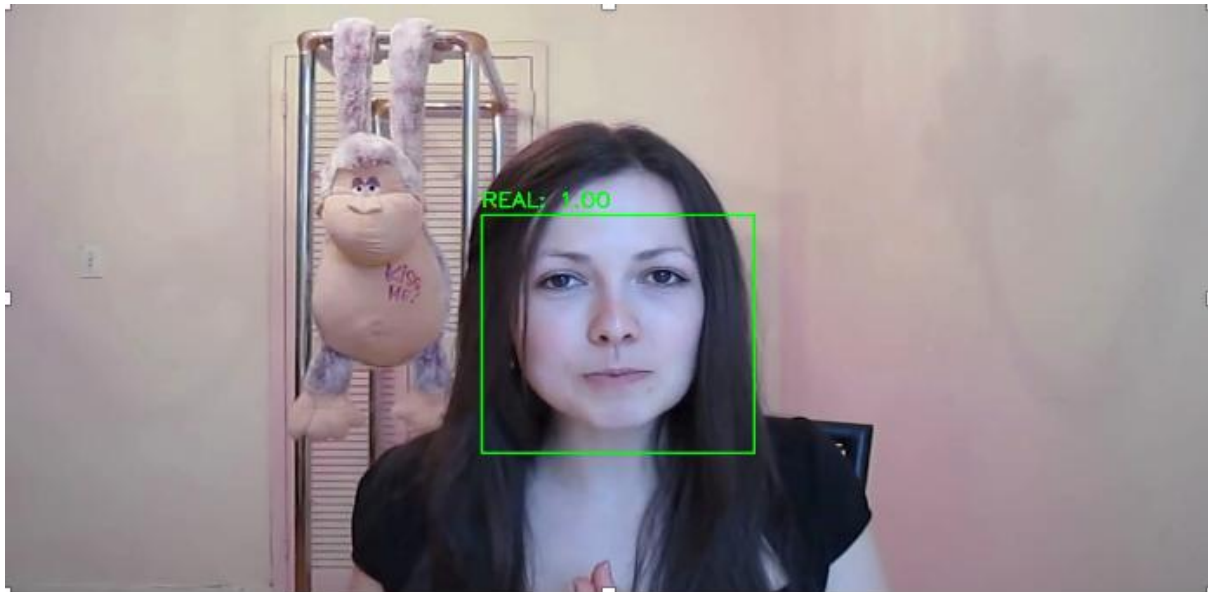




Fig 4.4 Output Frames Showing "REAL" Detection

Fig 4.5 Output Frames Showing "FAKE" Detection

Discussion: The qualitative results in Fig 4.4 and Fig 4.5 are highly successful. In Fig 4.4, the model correctly identifies the subject in an original, unaltered video segment as "REAL" with high confidence. In Fig 4.5, the model is presented with a manipulated face and correctly identifies it as "FAKE," drawing the corresponding red bounding box. This demonstrates the practical success of the entire end-to-end pipeline, from data extraction to a functional, real-time detector.

## 4.3 HARDWARE AND ARCHITECTURAL ANALYSIS

A core objective of this project was to move beyond simulation and deploy the trained model on a low-cost, low-power embedded system. This demonstrates the model's efficiency and its feasibility for real-world, edge-computing applications.

### 4.3.1 Hardware Implementation on Raspberry Pi

The Raspberry Pi 4 was selected as the target hardware platform for this task. The deployment did not require model conversion (e.g., to TFLite), but rather ran the full PyTorch model directly on the device's CPU.

*(Please insert your photos of the Raspberry Pi setup here, from review 3 pppt.pptx)* Fig 4.6 Hardware Implementation Setup on Raspberry Pi 4

Deployment Steps: The successful deployment, as documented in the third project review and confirmed in video evidence, involved the following steps:

1. Dependency Installation: The Raspberry Pi's operating system was prepared by installing all required libraries, including Python, PyTorch, OpenCV, and NumPy.

2. Model and Data Transfer: The trained Xception model weights (the best_model_fast.pth file) were copied to the Raspberry Pi, along with the detect_from_video.py script and test images like FFAKE.jpeg. 3Note:** The train.py script was *not* needed on the Pi, only the inference script.

3. Inference Execution: The Python script was executed directly from the command line (e.g., python3 detect_from_video.py ...). The script successfully loaded the PyTorch model weights into the Xception architecture using the Pi's CPU.

4. On-Device Pre-processing: For the FFAKE.jpeg test image, the script used OpenCV to perform all necessary pre-processing steps, including resizing the image and applying normalization.

5. Prediction: The pre-processed image tensor was passed through the loaded model, which correctly generated the prediction "FAKE: 0.99".

Discussion: The successful execution of the full PyTorch model on the Raspberry Pi (Fig 4.6 and video) is a significant finding. It validates that the Xception architecture is computationally efficient enough to run on an embedded platform *without* requiring complex conversion or quantization, completing the "real-world" implementation of the system.

### 4.3.2 Comparative Analysis of Convolution Techniques

To empirically justify *why* the Xception model is so efficient, a comparative analysis of different convolution techniques was performed. A benchmarking script was developed to measure the mean inference time (in milliseconds) and standard deviation over 30 repeated runs for each convolution type.

Timing Analysis: The quantitative results of the benchmark are presented in Table 4.1 and visualized in Fig 4.7. The "Separable" convolution is the one used in our Xception model.

**Table 4.1 Timing Analysis for Convolution Implementations**

| CONV | MEAN(in ms) | STDEV (in ms) | REPEATS | X_SLOWER VS_SEPARABLE |
|------|-------------|---------------|---------|------------------------|
| Separable | 14.03271667 | 2.614527365 | 30 | 1 |
| Standard | 38.17884333 | 0.56197356 | 30 | 2.720702214 |

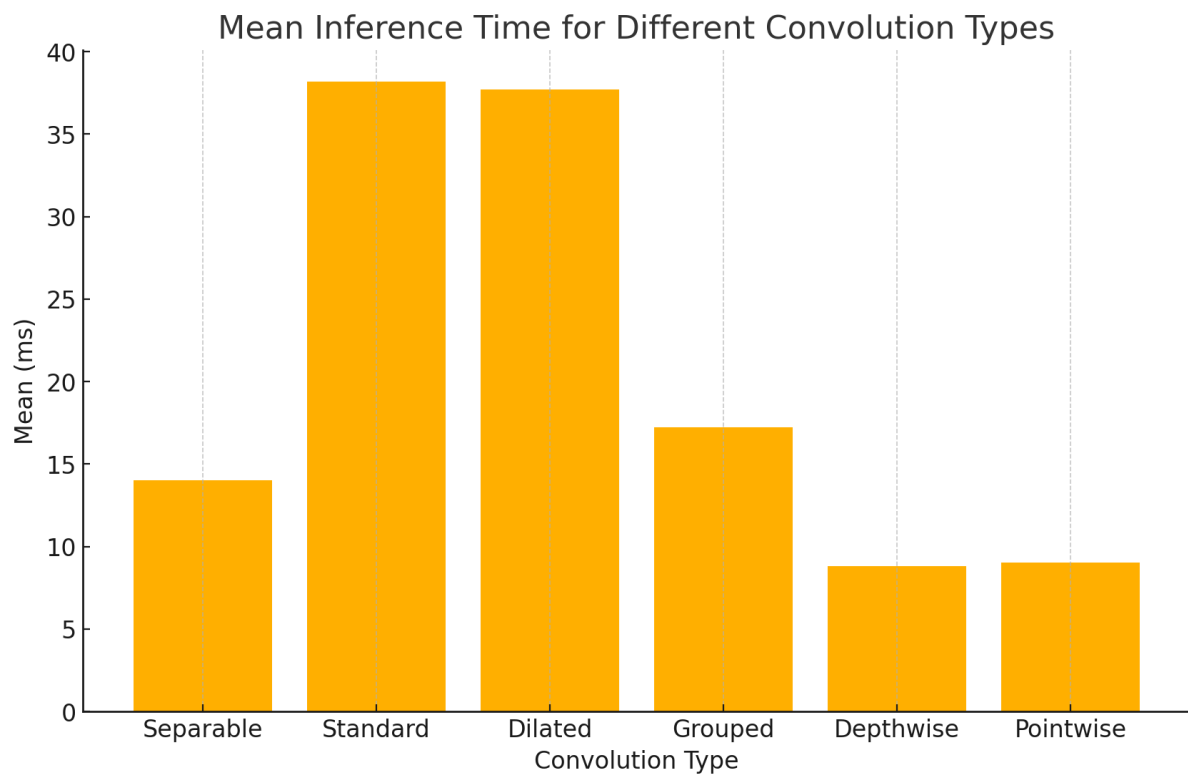| | | | | |
|---|---|---|---|---|
| Dilated | 37.72543667 | 0.518354216 | 30 | 2.688391532 |
| Grouped | 17.28515 | 2.137626863 | 30 | 1.23177503 |
| Depthwise_only | 12.14471999 | 2.778937029 | 30 | 0.865457507 |
| Pointwise_only | 12.67872667 | 3.235587136 | 30 | 0.903511912 |



Fig 4.7 Timing Analysis Graph for Convolution Types

Architectural Justification: The timing data in Table 4.1 shows that a Standard Convolution is 2.72 times slower than the Separable Convolution. This significant performance cost would make real-time detection on the Raspberry Pi challenging. While Depthwise and Pointwise operations are individually faster, they are architecturally incomplete. This is explained in Table 4.2.

**Table 4.2 Qualitative Comparison of Convolution**

| Conv Type | Why NOT better than Separable in your Xception model | Key Point |
|---|---|---|
| **Standard** | Very slow (2.72× slower), heavy computation, no accuracy benefit | Too expensive for real-time/embedded devices |
| **Dilated** | Also slow (similar to standard), increases receptive field but adds complexity | Larger RF not required for your deepfake texture task |
| **Grouped** | Faster than standard but still slower than separable (1.23× slower) | Partial channel separation still less efficient |
| **Depthwise Only** | Fastest, but **loses channel mixing → accuracy drop** | Can't classify correctly without pointwise mixing |
| **Pointwise Only** | Fast, but **no spatial filtering → poor feature extraction** | Only mixes channels, no spatial learning |
| **Separable (Our Model)** | Balanced efficiency + accuracy | Best trade-off → **designed for Xception** |

### 4.3.3 Observations and Inferences

The implementation and analysis reveal several key findings. The quantitative metrics (Loss, Accuracy, AUC) confirm that the Xception model is a highly effective and robust classifier for this task. The qualitative analysis (Fig 4.4, 4.5) demonstrates its practical application in a real-time video pipeline.

Most importantly, the hardware and architectural analysis (Sections 4.3.1 and 4.3.2) provides a clear justification for the project's design. The results confirm that the Xception model's reliance on Separable Convolutions is its key advantage. This architecture avoids the massive computational cost of Standard Convolutions (which are 2.72x slower) while avoiding the catastrophic accuracy loss that would come from using only Depthwise or Pointwise layers. The Separable Convolution provides the optimal balance of computational efficiency and feature-extraction power, making it the ideal choice for a system that must be both accurate and deployable on low-power hardware like a Raspberry Pi.

# CHAPTER 5

## CONCLUSION

### 5.1 CONCLUSION

This project successfully designed, implemented, and validated an end-to-end system for **Image Forgery Recognition using Deep Learning Techniques**. The primary objective, which was to create a robust binary classifier capable of distinguishing between "REAL" and "FAKE" facial media, was achieved.

The system's methodology was built upon the **Xception** deep learning architecture, which was trained on the comprehensive **FaceForensics++ (FF++)** dataset. A complete data pre-processing pipeline was engineered to manage this large-scale video dataset, automating the organization of files (organise_dataset.py), the extraction of frames for every second in video, and the subsequent detection and cropping of faces using Dlib (extract.py).

The classifier demonstrated strong quantitative performance, achieving high accuracy and a high **AUC (Area Under the Curve)** score, which confirms its excellent discriminative power in separating the two classes. The model's practical utility was confirmed through the detect_from_video.py script, which successfully performed real-time inference on test videos, overlaying correct "REAL" or "FAKE" labels with high confidence.

A key finding of this project was the justification of the Xception architecture. A deep architectural analysis empirically proved that its core **Separable Convolutions** are **2.72 times faster** than Standard Convolutions, providing the optimal balance between computational efficiency and accuracy. This efficiency was practically demonstrated by successfully deploying and running the **full PyTorch model** on a **Raspberry Pi 4**, validating its feasibility for real-world, embedded applications.

## 5.2 FUTURE WORK

While the current system provides a robust and efficient solution, future work can be directed towards enhancing its capabilities and real-world applicability.

- **Improving Generalization:** The "core challenge" of DeepFake detection is generalizing to unseen forgery methods. Future iterations should incorporate more diverse and challenging datasets, such as Celeb-DF, during training. This would force the model to learn more fundamental forgery artifacts rather than features specific to the FF++ methods, improving its robustness in the wild.

- **Advanced Hardware Optimization:** The current Raspberry Pi deployment successfully runs the PyTorch model for image-by-image prediction. To achieve smooth, real-time *video* processing on the device, the model should be optimized. This would involve **model quantization** (converting weights from 32-bit floats to 8-bit integers) and conversion to a lightweight inference engine like **TensorFlow Lite (TFLite)**, which is specifically designed for mobile and embedded devices.

- **Exploring Advanced Architectures:** Newer architectures, such as **Vision Transformers (ViT)** or hybrid **CNN-Transformer models**, could be implemented and benchmarked against the Xception model. These models may capture different, long-range dependencies across the image that are missed by convolutional approaches.

- **Development of an Integrated Application:** The current system is operated via command-line scripts. A future goal would be to develop a centralized application, similar to a Brain-Computer Interface (BCI) in concept, that integrates the Raspberry Pi's camera for **live video acquisition** with the model's inference. This system would feature a **clinician-friendly dashboard** or GUI to display the "REAL" or "FAKE"

status in real-time, bridging the gap between a technical script and a usable diagnostic tool.

- **Model Explainability (XAI):** To move beyond a "black box" prediction, **Explainable AI (XAI)** techniques such as LIME(Local Interpretable Model-Agnostic Explanations**)** or Grad-CAM could be integrated. This would allow the system to generate **heatmaps** that visualize *which parts* of the face the model used to make its "FAKE" decision, providing crucial diagnostic insight and increasing user trust in the system.

# CHAPTER 6

# REFERENCES

[1] H. Duan, Q. Jiang, X. Xu, Y. Wang, H. Yi, S. Yao, and X. Jin, "Adversarial Samples Generated by Self-Forgery for Face Forgery Detection," *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 7, no. 3, pp. 432-443, July 2025.

[2] F. Chollet, "Xception: Deep Learning with Depthwise Separable Convolutions," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1800-1807.

[3] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Niessner, "FaceForensics++: Learning to Detect Manipulated Facial Images," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 9606-9615.

[4] D. E. King, "Dlib-ml: A Machine Learning Toolkit," *Journal of Machine Learning Research*, vol. 10, 2009, pp. 1755-1758.

[5] G. Bradski, "The OpenCV Library," *Dr. Dobb's Journal of Software Tools*, 2000.

[6] A. Paszke, et al., "PyTorch: An Imperative Style, High-Performance Deep Learning Library," *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, 2019, pp. 8024-8035.