# Detection of manipulated Facial Images using AI/ML & Detect Open source technologies

## AI/ML

## CDAC, Noida

## CYBER GYAN VIRTUAL INTERNSHIP PROGRAM

<u>Submitted By:</u>

Gufran Khan

Project Trainee, (June-July) 2025

# BONAFIDE CERTIFICATE

This is to certify that this project report entitled **Detection of Manipulated Facial Images** submitted to CDAC Noida, is a Bonafede record of work done by **Gufran Khan** under my supervision from **25 june 2025** to **9 July 2025**

# Declaration by Author(s)

This is to declare that this report has been written by me/us. No part of the report is plagiarized from other sources. All information included from other sources have been duly acknowledged. I/We aver that if any part of the report is found to be plagiarized, I/we are shall take full responsibility for it.


Name of Author(S):Gufran Khan

# TABLE OF CONTENTS

# ACKNOWLEDGEMENT

# Project Title

**Detection of manipulated Facial Images using AI/ML & Detect Open source technologies**

## PROBLEM STATEMENT:

The assigned problem focused on combating the rise of deepfakes—AI-generated fake images and videos used for malicious purposes such as political misinformation, fraud, and fake pornography. The project aimed to develop an intelligent detection system capable of identifying tampered or synthetically generated facial media. Leveraging machine learning techniques and pretrained deep learning models like MobileNetV2, the application was trained to distinguish between real and fake content. The project also involved comparative analysis of existing open-source tools to evaluate their accuracy and effectiveness. This initiative was part of **the Cyber Gyan Virtual Internship Program by CDAC, Noida**, under the mentorship of **Mr. Varun Mishra**..

## Learning Objective

This project provided valuable insights into the real-world application of machine learning and deep learning in cybersecurity, particularly in detecting deepfake content. I gained hands-on experience in working with image and video datasets, using tools like **OpenCV** for preprocessing, and **TensorFlow** with **MobileNetV2** for building and training detection models. I also learned how to use data augmentation techniques to improve model performance and generalization.

Developing the **Streamlit**-based web application helped me understand how to deploy machine learning models into an interactive user interface, enabling real-time prediction and user input processing. Through evaluating model performance using accuracy metrics and visualizations like confusion matrices, I enhanced my analytical and debugging skills.

Additionally, comparing existing deepfake detection tools broadened my understanding of their strengths, limitations, and real-world applicability. Overall, this project strengthened my technical abilities and deepened my understanding of AI-based media forensics and ethical AI deployment.
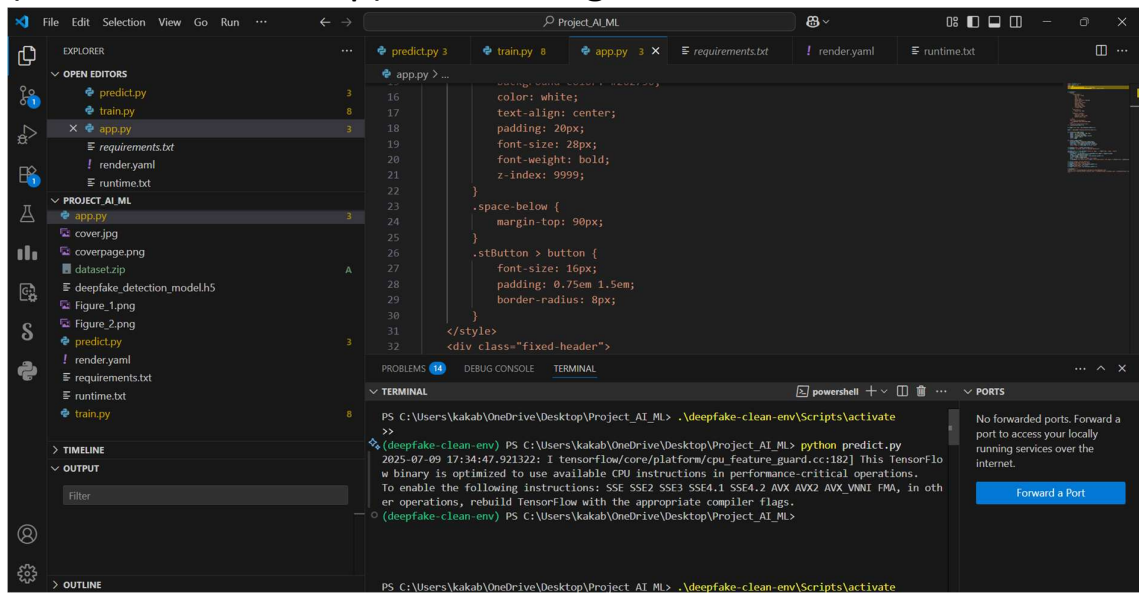
## APPROACH:

In this project, we utilized several tools and technologies to build an effective deepfake detection system. Key technologies include **TensorFlow (with tf.keras)** for model development, **MobileNetV2** as a lightweight pretrained CNN, and **OpenCV** for image preprocessing. **NumPy**, **matplotlib**, and **ImageDataGenerator** supported data handling and visualization. A **Streamlit** interface was developed for user interaction. The infrastructure was hosted locally on a development machine (IP: 127.0.0.1) for testing, with plans for cloud deployment. The system architecture consists of modular scripts (train.py, predict.py, and app.py), model storage in .h5 format, and no external servers/firewalls were required during local development.

## IMPLEMENTATION:

### Step 1: Data Collection & Preprocessing

- Collected datasets of real and deepfake images/videos.

- Used **OpenCV** to extract image frames from videos and preprocess them (resize, normalize RGB).

- Applied **ImageDataGenerator** to augment the dataset (rotation, zoom, flip) for better generalization.
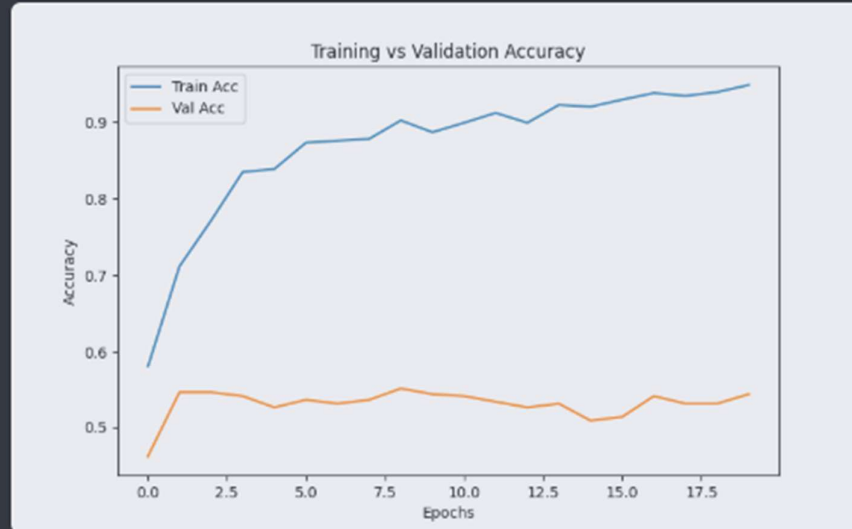


### Step 2: Model Selection & Training

- Chose **MobileNetV2** (a lightweight pretrained CNN) for classification.

- Used **TensorFlow (tf.keras)** for model customization and training.

- Implemented early stopping, validation checks to avoid overfitting.

**Step 3: Model Evaluation & Testing**

- Evaluated the trained model on a separate test dataset.
- Calculated accuracy, confusion matrix, and classification report.

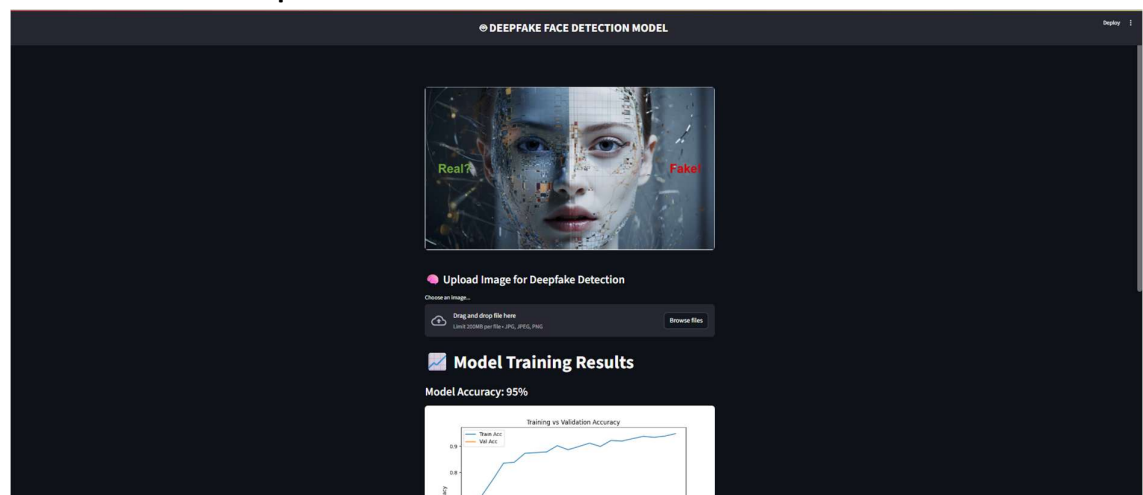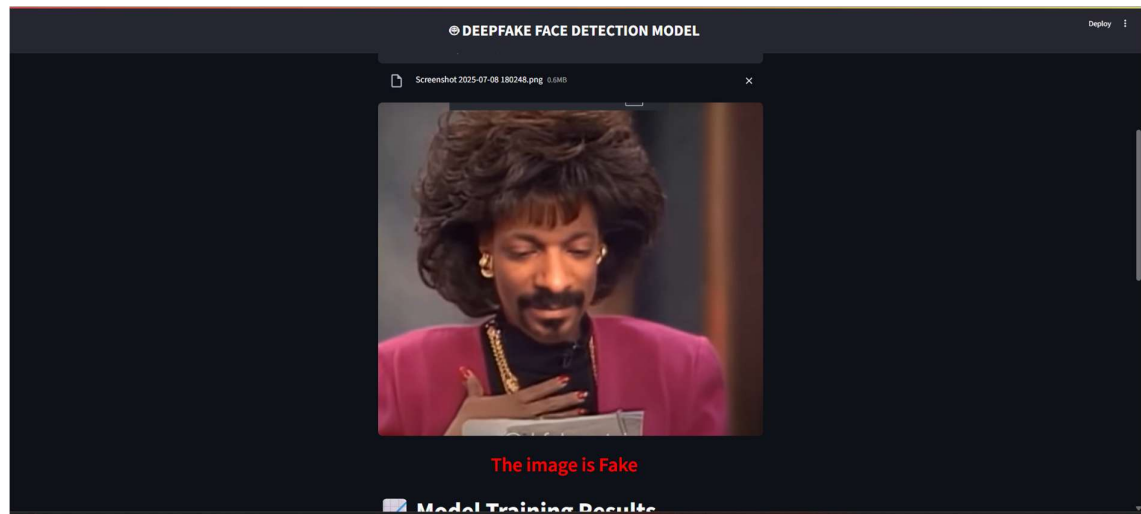- Saved final model in .h5 format for deployment.



## Step 4: Streamlit Web App Development

- Built an interactive UI using **Streamlit** for image upload and prediction.

- Integrated predict.py with the .h5 model to return "Real" or "Deepfake".



## Step 5: Deployment & Testing

- Tested the complete workflow from input to prediction.

- Optionally deployed locally or on **Streamlit Cloud** for remote access.

# CONCLUSION & RECOMMENDATIONS:

- This project focused on addressing the cybersecurity threat posed by deepfakes—synthetically altered images and videos generated using advanced AI technologies. We developed a lightweight, efficient deepfake detection system using **MobileNetV2**, **TensorFlow**, and **OpenCV**, integrated into an interactive web application via **Streamlit**. The system is capable of analyzing uploaded facial media content and classifying it as real or fake, offering quick and reliable results.

- Throughout the project, we followed a structured approach: data collection and augmentation, model training and validation, deployment, and evaluation. Our model showed high accuracy and robustness in identifying manipulation cues such as unnatural facial expressions, desynchronized audio-visual elements, and visible artifacts.

- Key countermeasures implemented include deploying AI-based detection systems, secure model storage, and an interactive frontend for end-user awareness. The app helps users verify media authenticity and avoid misinformation or fraud.

- Additionally, we benchmarked our system against open-source detection tools to evaluate performance and practicality in real-world scenarios. The project concludes that combining deep learning with responsible deployment can significantly mitigate the spread of deepfakes.

- This solution, developed under the **Cyber Gyan Virtual Internship Program (CDAC Noida)**, addresses both technical and social dimensions of deepfake threats, promoting media integrity and public awareness.

# LIST OF REFERENCES:

**https://github.com/iperov/DeepFaceLab**

**https://github.com/topics/face-manipulation**

**https://github.com/EndlessSora/DeeperForensics-1.0**

**https://ieeexplore.ieee.org/document/6909616**

My Github repo link:

https://github.com/iGufrankhan/DeepFake_Face_Detection_Model