

Development of AI/ML based solution for detection of face-swap based deep fake videos

A PROJECT REPORT

Submitted by,

NAGA NIKITHA.P - 20211ISR0023

SARTHAK MISHRA-20211ISR0086

DARSHAN.S –20211ISR0073

Under the guidance of,

Dr. MOHAMMADI AKHEELA KHANUM

In partial fulfillment for the award of the degree of

BACHELOR OF TECHNOLOGY

IN

INFORMATION SCIENCE AND ENGINEERING (AI AND ROBOTICS)



PRESIDENCY UNIVERSITY

BENGALURU

APRIL 2025

PRESIDENCY UNIVERSITY

SCHOOL OF COMPUTER SCIENCE ENGINEERING

CERTIFICATE

This is to certify that the Project report “**Development of AI/ML based solution for detection of face-swap based deep fake videos**” being submitted by “NAGA NIKITHA.P, SARTHAK MISHRA, DARSHAN.S” bearing roll number(s) “20211ISR0023,20211ISR0086,20211ISR0073” in partial fulfillment of the requirement for the award of the degree of Bachelor of Technology in Information Science and Engineering (AI and Robotics) is a Bonafide work carried out under my supervision.

**Dr. MOHAMMADI
AKHEELA KHANUM**
PROFESSOR
School of CSE& IS
Presidency University

Dr. ZAFAR ALI KHAN.N
PROFESSOR CSE & HOD
School of CSE&IS Presidency
University

Dr. L. SHAKKEERA
Associate Dean School
of CSE Presidency
University

Dr. MYDHILINAIR
Associate Dean School
of CSE Presidency
University

Dr. SAMEERUDDINKHAN
Pro-VC School of Engineering
Dean -School of CSE&IS
Presidency University

SCHOOL OF COMPUTER SCIENCE ENGINEERING

DECLARATION

We hereby declare that the work, which is being presented in the project report entitled “**Development of AI/ML based solution for detection of face-swap based deep fake videos**” in partial fulfillment for the award of Degree of **Bachelor of Technology in Information Science and Engineering**, is a record of our own investigations carried under the guidance of Dr. MOHAMMADI AKHEELA KHANUM, **ASSISTANT PROFESSOR, School of Computer Science Engineering & Information Science, Presidency University, Bengaluru.**

We have not submitted the matter presented in this report anywhere for the award of any other Degree.

Naga Nikhitha .P
(20211ISR0023)

Sarthak Mishra
(20211ISR0086)

Darshan. S
(20211ISR0073)

ABSTRACT

Deepfake technology in the form of face-swapped video is a potential threat to digital media authenticity and trust. The advanced deep learning algorithms employed for generating these artificial videos can imitate real people convincingly and are bound to be utilized to perpetrate acts of misinformation campaigning, identity forgery, and loss of public confidence. This paper aims at constructing a robust AI/ML-based system to detect face-swap manipulations in video material.

The method suggested employs Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks for extracting spatial and temporal inconsistencies present in deepfakes. CNNs are employed to extract features from individual frames for detecting unnatural texture of the skin, varying illumination, and deviation in facial contours. LSTM networks are subsequently employed to analyze temporal motions between frames to identify inconsistencies in face movements and expressions that are most likely to be indicative of tampering. This combined feature allows extensive scrutiny of video content to enhance detection of subtle artifacts added in face-swapping.

Training is conducted on publicly available data sets like FaceForensics++, Celeb-DF, and the Deep Fake Detection Challenge (DFDC) data set to provide a diverse set of manipulation strategies and video qualities. The model returns a probabilistic value between zero and one as a measure of the probability of manipulation to facilitate simple binary classification decisions. The value that can be generated using this can be applied in social media content moderation, forensic analysis, and anti-disinformation activities to ensure trust in digital information platforms.

The intended outcome of this work should be a deepfake detection system that is scalable and effective in nature and can be incorporated into platforms to assist in detecting and preventing manipulated content in advance. By surmounting the technological challenges of deepfake detection and leveraging cutting-edge AI/ML functions, this research seeks to preserve integrity of digital media and assist in preventing malicious usage of synthetic media.

The rapid advancement of deep learning technologies has made it possible for very realistic face-swapped deepfakes with important consequences to digital media authenticity and public faith. This paper is focused on crafting an effective AI/ML-based solution that can identify face-swapped manipulations in video content at high accuracy. The methodology leverages the capabilities of convolutional neural networks (CNN) alongside attention to identify location and timing inconsistencies in deepfakes.

ACKNOWLEDGEMENT

First of all, we indebted to the **GOD ALMIGHTY** for giving me an opportunity to excel in our efforts to complete this project on time. We express our sincere thanks to our respected dean **Dr. Md. Sameeruddin Khan**, Pro-VC, School of Engineering and Dean, School of Computer Science Engineering & Information Science, Presidency University for getting us permission to undergo the project. We express our heartfelt gratitude to our beloved Associate Deans **Dr. Shakkeera L** and **Dr. Mydhili Nair**, School of Computer Science Engineering & Information Science, Presidency University, and Dr. “ZAFAR ALI KHAN”, Head of the Department, School of Computer Science Engineering & Information Science, Presidency University, for rendering timely help in completing this project successfully. We are greatly indebted to our guide **Dr. MOHAMMADI AKHEELA KHANUM** **PROFESSOR** and Reviewer **Ms. SMITH.S.P,** **ASSISTANTPROFESSOR**, School of Computer Science Engineering & Information Science, Presidency University for his/her inspirational guidance, and valuable suggestions and for providing us a chance to express our technical capabilities in every respect for the completion of the project work. We would like to convey our gratitude and heartfelt thanks to the PIP2001 Capstone Project Coordinators **Dr. Sampath A K, Dr. Abdul Khadar A and Mr. Md Zia Ur Rahman**, Github coordinator **Mr. Muthuraj**. We thank our family and friends for the strong support and inspiration they have provided us in bringing out this project.

Naga Nikitha . P
(20211ISR0023)

Sarthak Mishra
(20211ISR0086)

Darshan . S
(20211ISR0073)

LIST OF TABLES

Sl. No.	Table Name	Table Caption	Page No.
1	Table 2.1	Literature Survey	18-21
2	Table 9.1	Model Performance Metrics	44
3	Table 9.2	Interference Time	45

LIST OF FIGURES

Name	Caption	PageNo.
Figure 7.1	Gantt Chart	40

TABLE OF CONTENTS

1	Certificate
2	Abstract
3	Acknowledgment
4	List of Tables
5	List of Figures
6	Chapter 1: Introduction
7	1.1 Overview
8	1.2 Importance of Deepfake Detection
9	1.3 Challenges in Detecting Face-Swap-Based Deepfakes
10	1.4 Role of AI/ML in Deepfake Detection
11	1.5 Objectives of the Project
12	Chapter 2: Literature Survey
13	2.1 Introduction
14	2.2 Deep Learning-Based Detection Methods
15	2.3 Feature-Based Detection Approaches
16	2.4 Hybrid Approaches and Ensemble Models
17	2.5 Limitations of Existing Methods
18	Chapter 3: Research Gaps of Existing Methods
19	3.1 Introduction
20	3.2 Generalization Issues
21	3.3 Dataset Limitations

TABLE OF CONTENTS

22	3.4 Adversarial Robustness
23	3.5 Computational Constraints
24	3.6 Lack of Explainability in AI-Based Detection Models
26	3.7 Need for Adaptive and Continual Learning Approaches
27	3.8 Ethical and Privacy Concerns in Deepfake Detection
28	3.9 Conclusion
29	Chapter 4: Proposed Methodology
30	4.1 Overview
31	4.2 Dataset Collection and Preprocessing
32	4.3 Model Development
33	4.4 Training and Optimization
34	4.5 Evaluation and Deployment
35	Chapter 5: Implementation
36	5.1 Model Architecture
37	5.2 Training Process
38	5.3 Software Integration
39	5.4 Performance Optimization
40	5.5 Model Deployment and Testing
41	Chapter 6: System Design & Implementation
42	6.1 System Architecture
43	6.2 Data Collection and Preprocessing
44	6.3 Feature Extraction

TABLE OF CONTENTS

45	6.4 Model Training and Evaluation
46	6.5 Real-Time Implementation
47	6.6 User Interface and Integration
48	6.7 Security and Ethical Considerations
51	Chapter 7: Timeline for Execution of Project (Gantt Chart)
52	Chapter 8: Expected Outcomes
53	Chapter 9: Results and Discussions
54	9.1 Results
55	9.2 Discussion
56	Chapter 10: Conclusion
57	References
58	Appendices
59	A. Pseudocode
60	B. Screenshots
61	C. Enclosures

CHAPTER-1

INTRODUCTION

Overview

The rising advancements in Machine Learning (ML) and Artificial Intelligence (AI) have caused great interest in deepfake technology concerning digital authenticity for media. Deepfakes use sophisticated generative models such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs) to manipulate face features within video recordings to produce highly realistic-looking but artificial content. Deepfakes find real-world usage in entertainment, film, and social media but are used to bring about catastrophic disruptions such as misinformation, identity fraud, political manipulation, and cybercrimes. The increased accessibility of deepfake generation tools has enabled malicious users to generate manipulative content easily, making it essential to implement effective AI/ML-based techniques for accurate detection and prevention.

Significance of Deepfake

Deepfakes are employable in a variety of malicious functions including defamation, disinformation propagation, financial forgery, and cyberbullying. Convincing face manipulation of individuals in video content undermines digital authenticity, and it is not possible to distinguish between real and generated video. Classic forensic verification and watermarking tend to be incapable of keeping pace with the growing sophistication of deepfake algorithms. AI/ML-based detection systems are hope in disguise, dependent on recognizing inconsistency in face expressions, awkward movements in a video, and digital artifacts invisible to humans but detectable by machines. Automated real-time detection systems are indispensable to provide digital assurance and diminish the harmful influence of deepfakes.

Detection Challenges for Face-Swap

Despite improvements in deepfake detection, several problems persist in the detection of face-swapped videos. A factor contributing to this is that state-of-the-art detection models hardly ever generalize as the majority of them are meant to be trained on specific sets of data and do not perform well on unseen or new deepfake techniques. Post-processing strategies such as additive noise and compression, as well as adversarial attacks, undermine the detection accuracy.

Apart from this, explainability of AI models is an issue as most detection methods are "black box" in nature and don't provide insights into their decision-making. To counter these challenges, newer, explainable, and efficient AI-based detection models are required. Role of AI/ML in Deepfake Detection

Machine learning and deep learning methods are very important to use for face-swapped deepfake video detection. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) are commonly applied to examine spatial and temporal anomalies in videos. CNNs are able to identify pixel-level abnormalities, whereas RNNs observe motion abnormalities between frames. Moreover, transformer-based architectures and attention have been investigated to improve deepfake detection by concentrating on important facial areas that show signs of manipulation. Adversarial training, where detection models are trained on dynamic deepfake methods, also enhances robustness. Deepfake detection systems can adapt to new threats and refine detection accuracy through AI/ML techniques continuously.

Objectives of the Project

The goal of this project is to develop an AI/ML-based system to detect face-swapped deepfake videos. The main goals are Designing a deep learning-based model that can detect facial inconsistencies in deepfake videos. Improving detection accuracy by integrating multiple AI techniques, such as CNNs, RNNs, and transformers. Enhancing real-time detection capabilities while ensuring computational efficiency. Developing an easy-to-use application that can be installed across various platforms. Overcoming adversarial attack issues and dataset biases to facilitate improved generalization. In the current digital age, the authenticity of visual information has come under increasing skepticism due to the progress of artificial intelligence as well as deep learning technology. The scariest innovation in the field is the development of deep fake technology. These techniques, particularly face-swap deep fakes, involve artificially replacing the face of a person in a video with that of another individual. The result is a very realistic video that will be indistinguishable from reality to the naked eye, thus creating serious social, political, and ethical problems. This manipulation technology, as great as it is an exhibition of generative AI prowess—specifically Generative Adversarial Networks (GANs)—is risky in many domains. From identity theft and defamation to political disinformation and social engineering, deep fakes can be used to manipulate public sentiment, disseminate disinformation, and undermine trust in digital media. The force behind this initiative is the need for an efficient and powerful detection system that is capable of successfully identifying face-swap deep fakes.

Detection systems at present are piecemeal and generally wanting in terms of providing real-time answers or even being usable across groups and qualities of videos. Traditional detection methods—such as human screening or rudimentary frame checking—are inadequate and ineffective in countering today's AI-powered forgeries. Therefore, the solution proposed in this paper utilizes a hybrid system that employs AI/ML techniques to detect both spatial and temporal anomalies. It entails close inspection of facial traits, sudden shifts, lighting disparities, and motion inconsistencies between successive video frames. These disparities that are imperceptible to humans are easily discerned using machine learning algorithms trained on labeled collections of real and synthetic videos. Social media, news outlets, and public institutions are especially vulnerable to the proliferation of such content. With millions of videos uploaded daily, an automated solution is no longer an option—it is a requirement. This system aims to fill that void by being not only real-time but also scalable and deployable on mobile, desktop, and cloud platforms. Offline capability is another key feature of the proposed system.

The majority of users, particularly rural or underdeveloped area users, lack access to high-speed internet. An independent cloud services system makes deep fake detection available to all regardless of geography and connectivity. It significantly enhances the utility and reach of the system. In addition to technical robustness, the aspect of transparency and explainability is of the highest importance. The system will include interpretable AI methods—such as Grad-CAM, SHAP, or saliency mapping—to provide visual explanations for every classification result. This will not only increase the confidence of users but will also allow digital forensics specialists and legal professionals to present valid evidence of manipulation. The social implication of this project cannot be overstated. Deep fakes can destroy reputations, incite violence, manipulate stock markets, or even threaten global relations.

By promoting media fact-checking and digital literacy, it also facilitates more pervasive policy-level and educational programs. One of the principal drivers for this move is the spread of such doctored content through social media and online video-sharing platforms. With increasing quality in deep fake videos, manual detection has become unsustainable. Automated detection is now the hour of need to ensure digital media integrity. The automatic tools must work under real-world conditions, including low-resolution video scenes, compression artifacts, and uneven lighting. The system covered in this paper will combine machine learning, computer vision, and signal processing to provide a robust detection solution. It will value precision of a light-weight model that could be used on a variety of platforms like smartphones, desktops, and web browsers. Specific emphasis is placed on real-time processing, offline availability, and low computational overhead. Moreover, the project recognizes the role of explainability in establishing

user trust. Hence, the system will employ interpretable AI techniques such as heatmap visualizations to detect areas of the video that triggered a deep fake warning. Such transparency is critical for legal, journalistic, and forensic users who require verifiable and interpretable evidence.

The dramatic progress in recent years in the area of artificial intelligence and deep learning has yielded extremely powerful tools that have made it possible to create very natural-sounding and looking synthetic media commonly referred to as deepfakes. Among these various types of deepfakes, face-swapped deepfake videos have been one of the most concerning types where a face has been naturally replaced with another person's face in a video with remarkable accuracy. Though deepfake technology has the potential to innovate use in industries such as video games, entertainment, and learning, abuse erodes privacy, security, and public confidence at a severe level. Deepfake videos manipulated deliberately can be utilized for defamation, blackmail, disinformation, political manipulation, and identity theft with serious social and ethical problems. Therefore, there is a growing need for reliable and efficient AI/ML-based systems to detect and suppress the spread of malicious deepfake content.

It is a very hard task to identify deepfakes. As the technology of creating deepfakes advances, the algorithms used in creating them grow more sophisticated, and the capability to generate highly realistic-looking videos increases. These videos become even difficult for expert eyes to separate from real videos. Simple optical illusions or human visual inspection-based traditional methods don't suffice anymore. Hence, contemporary detection systems need to utilize sophisticated machine learning models capable of detecting faint spatial inconsistencies, temporal anomalies between frames, and small artifacts caused by face-swapping operations. A good detection solution should not only detect current deepfake methods but also generalize to new, unseen forms of manipulations so that it can be robust in real-world applications.

The strategy is to use deep learning techniques such as Convolutional Neural Networks (CNNs) to learn detailed spatial information from individual frames and Long Short-Term Memory (LSTM) networks to learn temporal inconsistencies among sequences of frames. Through learning and inspection of tiny details such as pixel-level artifacts, minor facial landmark inconsistencies, and unnatural face motion dynamics, the system to be proposed will have the capability to label videos as real or synthetic with high accuracy. Moreover, rigorous data preprocessing, augmentation, and model regularization steps are integrated into the pipeline to assist in enhancing the generalization capability of the system on diverse datasets and real-world applications. The recent unexpected deepfake explosion has increased the panic level among

researchers, and therefore their detection in today's digital era is warranted. However, face-swap based deepfake detection is a very challenging and ever-evolving issue. Initial detection methods were primarily focused on detecting explicit artifacts such as unnatural blinking, abnormal lighting, or boundary inconsistencies. But with the advancement of deepfake generation algorithms, these visual cues have become less apparent, and thus the traditional detection techniques are less effective. Modern deepfake videos are extremely frame-consistent in fidelity, show realistic facial expressions, and are robust to changes in lighting conditions, thus necessitating more sophisticated detection techniques. Human detection is no longer adequate, and hence the need for intelligent, automated systems that can detect forgeries both at a spatial and temporal level. The importance of developing effective deepfake detection methods transcends technical advancement; it also carries enormous societal ramifications.

As deepfakes become more readily available and hard to detect, they pose serious consequences to information authenticity, national security, democratic processes, and human rights. A robust detection system can minimize the spread of misinformation, enable authentication of social media content, assist law enforcement agencies in forensic analysis, and protect individuals from potential exploitation. Therefore, this project not only addresses a significant technical issue but also contributes to building a safer and more trusted digital society. The successful implementation of this AI/ML-powered deepfake detection mechanism can be made an important weapon against the misuse of synthetic media technologies in the present times.

CHAPTER-2

LITERATURE SURVEY

2.1 Introduction

The literature on detecting deepfakes is growing extensively, with a large number of researchers exploring various AI/ML techniques to combat the increasing cleverness of manipulated videos. Here, in this chapter, recent methods of detecting deepfakes are presented along with their weaknesses and strengths. A significant focus is placed on significant methodologies like deep learning-based models, feature engineering techniques, as well as integration-based solutions involving more than one detection method.

2.2 Deep Learning-Based Detection Techniques

Deep learning has been found to be the optimal method of detecting deepfakes in images and videos. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) are used routinely to detect spatial and temporal inconsistencies, respectively. Studies such as Afchar et al. (2018) suggested MesoNet, a CNN-based model capable of detecting deepfakes based on facial texture inconsistency. Similarly, Zhang et al. (2019) proposed a deep learning-based framework in which CNN is employed to detect spatial artifacts introduced during face-swapping processes. Another approach involves the use of Long Short-Term Memory (LSTM) networks and Transformer models, which are specifically designed to handle sequential patterns in videos. For example, Guera and Delp (2018) proposed an LSTM-based model that detects deepfakes based on inconsistencies in facial motion tracking. More recently, Transformer models based on attention have been proposed to draw attention to salient facial areas, enhancing detection accuracy through dynamically learning importance of features.

2.3 Feature-Based Detection Techniques

Feature engineering has also been explored as a way of deepfake detection. Such approaches target specific visual artifacts, such as unnatural eye blinking, facial asymmetry, and lighting inconsistencies. Li et al. (2018) demonstrated that deepfake videos are more likely to exhibit unnatural eye-blinking patterns, which can be detected through computer vision algorithms. Additionally, Tolosana et al. (2020) investigated physiological signals, such as heart rate variations, in order to detect manipulated content.

2.4 Hybrid Approaches and Ensemble Models

In an effort to improve robustness, researchers have tried hybrid approaches that combine a number of detection approaches. Marra et al. (2019) introduced an ensemble model that integrates CNNs with hand-engineered feature extraction techniques. Wang et al. (2020) introduced a two-stream network that scans both spatial and frequency-domain features to enhance deepfake detection.

2.5 Shortfalls of Existing Approaches

Despite recent advancements in the detection of deepfakes, current methods have several problems

Generalization Problems Most models struggle to detect deepfakes created by novel methods.

Adversarial Robustness Adversarial attacks and post-processing techniques allow deepfake developers to easily bypass detection.

Computational Requirements AI models are computationally intensive, and real-time application is limited.

Interpretability Deep learning networks tend not to be explainable, and as a result, it becomes hard to understand how the decision process works. The major driving factor for this project is the far-reaching sharing of such manipulated content through social networks and video hosting websites. Deep fake videos now are a job impossible to manually do with an increase in the quality of fake videos. Presently, computerized tools must be used with digital media integrity. The equipment must work in ideal real-world environments, for example, where there are compression artifacts, low-resolution videos, and irregular lighting conditions.

The system being developed will combine computer vision, machine learning, and signal processing to provide an end-to-end detection solution. It will focus on developing a light but accurate model that can be easily deployed on various platforms like desktops, smartphones, and web browsers. Specific focus is placed on real-time processing, offline availability, and minimal computational overhead. Also, this project recognizes the need for explainability in order to attain user trust. Thus, the system will feature interpretable AI techniques such as heatmap visualizations to denote areas in the video that required a deep fake warning. This transparency is essential for users within legal, journalistic, and forensic contexts who require verifiable and interpretable evidence.

2.6 Literature Survey

TITLE	YEAR	TECHNIQUES USED	DESCRIPTION	ACCURACY
Face Forensics++: Learning to Detect Manipulated Facial Images [1]	2019	Convolutional neural networks (CNNs). Dataset:1000videos.	It helps to detect fake images from videos by trained forgery detectors. Acquiring the skill of identifying altered facial photographs.	85.1%
Unmasking Deep Fakes with simple Features[2]	2019	GAN Two datasets	Our novel machine learning approach accurately detects manipulated videos with a 90% success rate by analyzing a low resolution video.	61.4%
Celeb-DF:A Large-scale Challenging Dataset for Deep Fake Forensics[3]	2020	Deep Neural Networks (DNNs) Dataset:5639 high quality videos.	Celeb DF is a comprehensive dataset designed to challenge deep fake.	75.3%

TITLE	YEAR	TECHNIQUE USED	DESCRIPTION	ACCURACY
The Deep Fake Detection Challenge (DFDC) Dataset [4]	2020	GAN MTCNN Dataset: 100000 clips from 3426 paidactors.	A deep fake detection model trained solely on the DFDC dataset has the ability to generalize to authentic "in-the-wild" deep fake videos.	55.9%
Video Face Manipulation Detection Through Ensemble of CNNs [5]	2020	CNN	The solution presented involves deriving various models from a foundational network (EfficientNetB4) by incorporating attention layers	72.7%
Learning Self Consistency for Deep fake Detection [6]	2021	Pair-wisef-self-consistency learning (PCL)	The new Deep fake detection system is effectively differentiate between authentic and manipulated faces bas do their varying image quality.	56.5%

TITLE	YEAR	TECHNIQUE USED	DESCRIPTION	ACCURACY
Wild Deep fake: A Challenging Real-World Dataset for Deep fake Detection[7]	2021	Attention-based Deep fake Detection Networks (ADDNets)	The dataset includes a diverse range of deep fake videos that have been meticulously crafted to closely resemble authentic footage, making it a formidable test for even the most advanced detection systems.	52.8%
Detecting Deep fakes with Self-Blended Images [8]	2022	Self blended images(SBIs)	Self-Blended Images involve the use of advanced algorithms and machine learning techniques to analyze and compare various facial features within a video.	56.1%
Explaining Deepfake Detection by Analysing Image Matching [9]	2022	FST-Matching, DNNs	The process of identifying deep fake content involves a thorough analysis of image matching techniques.	81.1%

TITLE	YEAR	TECHNIQUE USED	DESCRIPTION	ACCURACY
Deep Fidelity: Perceptual Forgery Fidelity Assessment for Deep fake Detection [10]	2023	SSAAFormer	The proposed Deep Fidelity framework aims to dynamically identify real and fake faces by examining the fidelity of facial images, taking into account the diverse quality levels present in both categories.	55.9%
Masked Conditional Diffusion Model for Enhancing Deep fake Detection [11]	2024	Masked Conditional Diffusion Model (MCDM) Data Augmentation	The MCDM is a cutting-edge technique that has been developed to enhance the detection of deep fake videos. Deep fake videos are manipulated videos that use artificial intelligence to replace the face of a person in an existing video with someone else's face	89.0%

CHAPTER – 3

RESEARCH GAPS OF EXISTING METHODS

3.1 Introduction

Though much progress has been made in the development of more effective deepfake detection techniques, there still remain some research gaps hindering the growth of practically efficient and effective solutions.

Existing methodologies lack generalizability, are vulnerable to easy adversarial attacks, require significant computational overhead, and are not interpretable. Furthermore, constantly evolving deepfake generation techniques necessitate adaptive detection techniques that can keep up with changing threats. The following chapter presents the most significant research gaps that should be filled in order to make deepfake detection models more effective.

3.2 Generalization Issues

One of the biggest issues with deepfake detection is that existing models are not generalizable.

Most AI/ML-based detectors are trained on specific datasets, which makes them extremely efficient against deepfakes produced by known mechanisms but less efficient when facing unknown deepfake mechanisms. This is due to the fact that there are deepfakes in datasets and models cannot learn generalized features that can detect a wide variety of deepfake manipulations. Furthermore, most training datasets are contaminated with deepfakes generated by a few algorithms. With deepfake technology improving, new models supporting more advanced face-swapped video capabilities are being put out. Existing detection models are not able to keep up with these enhancements and therefore need a continuous process of being retrained on new datasets. In order to limit this shortcoming, studies need to work on developing models with the ability to learn invariant features that are able to identify manipulations across various deepfake methods and domains.

3.3 Dataset Limitations

One of the most significant hurdles to improving deepfake detection is that high-quality, varied datasets are not available. The majority of datasets that currently exist, like Face Forensics++ and the Deep Fake Detection Challenge (DFDC), possess weak variations for ethnicity, lighting, backgrounds, and face expressions. Due to this absence of diversity, deepfake detection models are made less capable of real-world variation scenarios.

3.4 Adversarial Robustness

Deepfake manufacturers continuously improve their techniques to bypass detection systems through adversarial attacks and post-processing methods. Adversarial attacks are tiny alterations that cannot be detected by the human vision system but may deceive AI-based detection systems. These changes can include value alterations in pixel, texture alterations, or frequency domain alterations, so it is difficult for deepfake detection systems to be highly accurate. Furthermore, post-processing techniques such as video compression, blurring, and noise injection can reduce detection model performance by up to five times. Deepfake detection methods work for most cases but fail when dealing with videos that have been processed through such operations since they introduce distortions that mask manipulation artifacts. To combat these challenges, the research direction of the future should explore the use of adversarial training, where detection models are trained using adversarial perturbed deepfake videos to improve their evasion attack resistance.

3.5 Computational Constraints

Deepfakes are primarily detected by computationally expensive deep learning models, which are less suitable for deployment on resource-poor devices or real-time tasks. Current cutting-edge deepfake detection models may need powerful GPUs and large-scale computational resources during inference, restraining their scalability and usability.

Deepfake detection in real time is critical to applications such as social media content filtering and live video authentication. Current models, however, are poor at making a compromise between detection quality and computational efficiency. There is work to be done in creating efficient deepfake detection models that can run well on edge devices without compromising accuracy. Techniques such as model pruning, quantization, and knowledge distillation can be explored to optimize deepfake detection systems for deployment in the real world.

3.6 Lack of Explainability in AI-Based Detection Models

Another major research gap in deepfake detection is the lack of interpretability and explainability in AI-based detection models. The majority of today's state-of-the-art deepfake detection systems are "black boxes," and users struggle to understand the rationale behind their predictions. This lack of transparency raises issues regarding trust, especially in applications with legal forensics, media authentication, and misinformation control.

3.7 The Requirement of Adaptive and Continual Learning Strategies

Existing deepfake detection models are rigid and need frequent retraining using newer datasets to continue being effective. Nevertheless, methods of generating deepfakes are developing very fast and need adaptive detection systems that can learn without human intervention continuously. Continual learning strategies, in which models update their knowledge dynamically based on fresh deepfake instances, can very much enhance detection resilience.

3.8 Ethics and Privacy Concerns in Deepfake Detection

While deepfake detection is significant in the war against misinformation and cybercrime, it also has ethical and privacy concerns. AI-based detection systems are likely to require large amounts of real and synthetic facial data, which presents data privacy and consent concerns. Moreover, data biases during training can lead to variations in detection accuracy across demographic groups.

Ethical AI methodologies in deepfake detection research are crucial to uphold. Future research will aim at privacy-enforcing techniques, such as federated learning, that allows model training without data centralization. Bias mitigation techniques must also be applied to ensure fairness in deepfake detection across various user groups. Despite continued labor and numerous high-performance model launches, existing deep fake detection tools continue to exhibit significant shortfalls, especially in face-swap based deep fake video detection. Such are the shortcomings incurred due to varied technical, ethical, and pragmatic concerns that make such tools useless in real-world situations.

Among these most significant gaps is being too reliant on datasets skewed toward Western perspectives. Most of the deep fake detection models are trained on databases such as FaceForensics++, Celeb-DF, and DFDC that predominantly include images of individuals with extremely few ethnicities. As a result, these models do not generalize when tested on videos of non-Western faces or underrepresented populations. This demographic bias not only compromises the equity of detection but also brings vulnerability to global digital ecosystems where content diversity is vast. Another major challenge is the exorbitant computation cost of several state-of-the-art models. Models that take extensive CNNs, RNNs, or 3D-CNNs are processor-bound and memory-thirsty, and hence their deployment in real-time or execution on mobile devices and edge systems is not practical. This limits only high-budget organizations from even detecting and taking action against spoofed content while others remain susceptible.

Also, detection systems usually provide binary outputs—labeling content as "real" or "fake"—without offering fine-grained explanation or contextual information. This uninterpretability reduces confidence in the system, particularly in forensic and legal contexts where evidence must be inspected and validated. Explainable AI (XAI) techniques are not applied in most current solutions despite the fact that they would be able to enhance users' confidence by pointing to suspicious regions or patterns that led to the model's decision. Another gap lies in the low utilization of multimodal analysis. Although face-swap deep fakes tend to modify primarily visual data, they do not typically fully synchronize the fake facial expression with real audio tracks. Current models that make use of just visual features can't account for these small but significant differences. Audio-visual examination—i.e., checking lips and spoken words—can yield strong evidence, but how to integrate it is not adequately explored due to the added intricacy of handling synchronization and alignment of multimodal inputs.

Scalability and flexibility are similarly enormous challenges. As deep fake generation methods become increasingly advanced with time, there are new modalities of deception introduced regularly. Static models from current datasets struggle to keep pace with new modes of attack vectors. Ongoing retraining, data augmentation, and the use of adversarial learning have been proposed, yet most systems do not yet incorporate dynamic adaptability to these growing threats. Deployment and accessibility likewise impede practical adoption. The majority of tools in use are research-specific, having complicated installation procedures and inadequate cross-platform compatibility. The tools become inoperable or difficult to decipher for users that are not technology-savvy, such as journalists, instructors, or lawyers, limiting the ability to practice with them. Lastly, ethics of data management and model interpretation are often forgotten.

Systems that draw on sensitive personal data for training have the potential to create privacy problems, and non-transparent black-box models have the potential for ethical and accountability issues. These aspects must be resolved if deep fake detection systems can become widely accepted and used by the public. By identifying and addressing these gaps—demographic inclusiveness, resourcefulness, multimodality integration, model interpretability, flexibility, user friendliness, and ethical accountability—this project aims to offer an up-to-date solution. The aim is not just to technically outperform but to also socially sensitized and deployment-ready for real environments with diverse users.

Despite the important advancements in AI and machine learning-based deepfake detection systems, some key research gaps still exist in current methods. One of the most challenging problems is generalization. Most current models learn and test on analogous forms of deepfake data that are created with precise known techniques, resulting in models being proficient on known deepfakes but unable to flag new, unfamiliar forms of manipulations that employ superior generation methods. This presents a critical drawback when models are implemented in the field as novel deepfake generation tools are continuously developing. The second serious gap is over-reliance on spatial clues within individual frames, instead of properly capturing the temporal inconsistencies among frames that are usually characteristic of video-based deepfakes. Although some existing works have postulated modeling temporal relationships through LSTM or RNN models, most of the current methods are largely static frame-based detection, and they lack sufficient capacity to spot deepfakes that all exhibit identical frame quality.

In addition, existing datasets to train and test models are often not diverse enough. Publicly available datasets mostly consist of videos that are mostly clean with minimal compression artifacts or natural degradations, whereas real-world videos typically contain multiple degradations such as low resolution, noise, motion blur, and compression artifacts. As a result, models trained on such ideal datasets perform poorly when tested on real-world social media streams or surveillance videos. Besides, interpretability and explainability of deepfake detection models are largely neglected in current literature. The majority of models only output a binary classification score (real or fake) without explanations of what features or regions of the face triggered the output, hence less reliable and applicable, especially in forensic or legal settings where explainability is necessary. Another critical gap is the lack of adversarial robustness. Recent experiments have shown that small, imperceptible perturbations can fool deepfake detectors, whereas most existing solutions are not resistant to adversarial attacks. In addition, the ability of detecting real-time detection has not been explored. Yet a very large number of solutions proposed before must call upon very high computing powers and are inference-centric, hence are not practically possible in real-time applications like streaming live videos or social media filtering. Furthermore, cross-domain generalizability of current solutions is also weak; e.g., a model trained with face-swap deepfakes will not generalize across lip-sync or identity-morphed videos without heavy retraining.

CHAPTER – 4

PROPOSED METHEDOLOGY

4.1 Overview

The methodology adopted for face-swapped deepfake video detection utilizes AI and ML methods to scrutinize spatial, temporal, and statistical patterns within videos. The methodology follows a structured approach from dataset collection to the final model deployment.

4.2 Dataset Collection and Preprocessing

For adequate training, this research uses benchmark datasets like FaceForensics++, Deepfake Detection Challenge (DFDC), and Celeb-DF. These datasets comprise a blend of genuine and manipulated videos, which permits extensive training and validation. Preprocessing includes:

Data Augmentation Expanding the dataset through cropping, rotation, and color changes.

Frame Extraction Pulling frames from the video to analyze temporal inconsistency.

Normalization Scaling pixel values for improved model convergence.

4.3 Model Development

The proposed model in the paper employs CNN-based networks for feature extraction and Transformer-based networks for sequential processing. The model architecture consists of Feature Extraction Layer Capturing facial patterns and inconsistencies. Temporal Analysis Module Identifying inconsistencies in frame changes. Classification Layer Differentiating between fake and authentic videos.

4.4 Training and Optimization

The model is trained using Loss Functions Binary cross-entropy and focal loss to handle imbalanced datasets. Optimizers Adam and RMSprop for gradient update. Regularization Dropout and batch normalization for overfitting avoidance.

4.5 Evaluation and Deployment

The model is evaluated on the basis of accuracy, precision, recall, and F1-score. It is then deployed as a real-time detection model once it reaches peak performance in the shape of an easy-to-use interface for convenience. With this approach, the research aims to come up with an AI-driven solution that can precisely detect face-swapped deepfake videos while minimizing some issues in existing solutions.

The proposed methodology is intended to address the complex problems of face-swap based deep fake video detection in a comprehensive way. It aims for a modular, scalable, and explainable framework that can be applied in varied environments—varying from academic research in universities to practical forensic use. The methodology consists of a variety of phases: data acquisition, preprocessing, feature extraction, model architecture, multimodal integration, training, inference, explainability, and deployment.

1. Data Acquisition:

This stage consists of generating a diverse dataset of high-resolution and compressed video types, spanning across different age groups, ethnicities, and genders. Apart from popular datasets such as Face Forensics++, DFDC, and Celeb-DF, synthetic datasets are generated using tools like Deep FaceLab and Face Swap. Customized deep fake content is also generated using advanced GAN variants like StyleGAN3 and First Order Motion Model for model robustness.

2. Preprocessing Pipeline:

Preprocessing makes sure the video content is sanitized, normalized, and segmented for effective training. Video frames are grabbed at periodic time intervals, typically 25-30 FPS, using FFmpeg. Face alignment and detection are achieved using MTCNN and MediaPipe Face Mesh to accommodate different facial orientations and occlusions. Audio is also extracted and preprocessed with signal enhancement and voice activity detection (VAD) processes to enable synchronization analysis of video and audio.

3. Feature Extraction and Model Architecture:

The model is divided into three distinct yet related modules: Spatial Analysis: For frame-level feature extraction, CNNs such as XceptionNet, EfficientNet, and MobileNetV3 are utilized. The baseline is employing pretrained weights on ImageNet with fine-tuning over deep fake datasets. Temporal Analysis: Temporal models such as Bidirectional LSTM and Temporal Convolutional Networks (TCNs) are employed for unnatural transitions, flickering, and jitter. Audio-Visual Fusion: A dedicated stream matches phonetic features (MFCC, spectrograms) with lip movement sequences through SyncNet-based architecture. Attention mechanisms are incorporated for dynamic weight adjustment between modalities.

4. Multimodal Fusion and Classification:

Features gleaned from video and audio streams are forwarded to a fusion network. The architecture uses transformer encoders to build a contextualized representation of the whole video. The final decision—"real" or "fake"—and a manipulation confidence are generated by a softmax classifier. The intermediate activations are stored for future use in explainability modules.

5. Explainable AI Integration:

Transparency is enabled using Grad-CAM, SHAP, and Integrated Gradients. These techniques identify changed facial regions or lip sync inconsistencies and thus provide end-users and forensic analysts visual justification for the system's decision.

6. Training and Optimization:

The schedule for the training consists of batch normalization, dropout regularization, and data augmentation to avoid overfitting. The hyperparameters are optimized using Bayesian optimization. Generalization over demographic splits is ensured through cross-validation. Transfer learning helps in bootstrapping performance and class imbalance is mitigated through the use of focal loss.

7. Deployment Framework:

The solution is deployed to light-weight models (TensorFlow Lite, ONNX) to support edge devices such as smartphones and surveillance systems. A RESTful API is developed for web platform integration. The system is also Docker containerized for ease of deployment in the cloud. Web and mobile interfaces provide real-time video upload, processing, and result display functionality.

8. Continuous Learning Loop:

In order to future-proof the system, it supports semi-supervised learning and feedback integration from the users. It allows users to mark new video types or variations of deep fake and put them back into the training pool, thus enhancing the model's adaptability over a period of time.

Preprocessing is done with the extraction of frames, facial detection using frameworks like MTCNN or Media Pipe, and face alignment in order to standardize the input for downstream analysis. In addition, audio tracks are separated and processed to prepare themselves for multimodal fusion down the pipeline. This step guarantees data quality and structure compliance, which are essential for consistent model performance. The core detection engine consists of two parallel branches: one for visual features and the other for temporal dynamics. The visual branch uses CNNs like EfficientNet and ResNet to extract spatial features from individual frames.

To further improve the system, an audio-visual synchrony module is incorporated. The module detects speech pattern vs. lip movement mismatches through cross-modal embeddings. A dual attention mechanism helps the model focus on the area with the highest scope for manipulation, enhancing accuracy. In the case of classification, a fusion network merges outputs of visual, temporal, and audio streams. Softmax layer provides confidence value in manipulation terms, and explainable tools like Grad-CAM are used to generate heatmaps to account for predictions. The last phase includes deployment solutions for cloud as well as edge deployments. The model is now optimized for deployment on mobile phones, desktops, and browser extensions using TensorFlow Lite and ONNX. The system also optimizes its functionality for real-time inference so that it is best suited for social media sites and verification programs.

The development process for an AI/ML-powered solution for face-swap deepfake video detection is an orderly, multi-step process, beginning with large-scale data collection and preparation. To ensure that the system is effective in dealing with diverse types of deepfake manipulations, datasets such as Face Forensics++, Deep Fake Detection Challenge (DFDC) preview dataset, Celeb-DF, and Deeper Forensions are employed. The datasets contain real and synthetic videos.

The found face regions are cropped and aligned based on central facial features as reference points to normalize for scale and orientation variability. After data preparation is complete, the feature extraction operation is carried out in order to learn the primary visual anomalies that make fake content and real content unique. Deep CNNs such as ResNet50, XceptionNet, and EfficientNet are employed to learn high-level semantic features from face images. Specific attention is given to preserve subtle artifacts like texture inconsistency, boundary blending mistakes at face borders, unnatural lighting, and irregularities in mouth or eye movement. These features learned are employed as significant inputs for subsequent classification. For the model development process, there is a deep learning classification model with first a CNN-based classifier trained for frame-level prediction. To further enhance the system's robustness, especially in detecting temporal inconsistencies between frames, a hybrid CNN-LSTM model is utilized.

In this configuration, the CNN layers are focused on spatial feature learning and the LSTM layers on frame-level sequential dependencies, capturing minor motion inconsistencies or unnatural switchovers that deep fake videos typically possess. For better making the model generalize and robust, the training process utilizes techniques of data augmentation like random horizontal flipping of images, small rotation, brightness and contrast adjustment, color jittering, and imitative compression artifacts.

Post model training, a robust decision-making approach is implemented at the video level. Instead of making decisions on individual frames, the system aggregates predictions at frames through majority voting or through average softmax scores. The technique denoises frame-level noise and yields better video classification results. The performance of the constructed models is quantified through typical performance measures including accuracy, precision, recall, F1-score, and AUC-ROC curves. Cross-validation and testing on unseen data are also employed to verify the generalization capability of the system over different classes of deepfake videos. A simple, web-based application is developed for practical use, allowing users to upload videos for analysis.

The backend picks frames from the uploaded videos, finds faces, identifies them, and compiles results and sends users a general prediction and flags suspected frames together with their confidence level. Inference optimization techniques such as ONNX model conversion and TensorRT acceleration are considered to deliver quick and efficient processing during deployment. Last but not least, with the current accelerated development of deepfake generating technologies, this presented methodology brings into focus continuous learning and renewal. The process facilitates ongoing set expansion by way of the integration of new classes of deepfakes and continuous retraining of the models. The future of work could include the use of adversarial training for defense against sophisticated deepfake attacks, adding audio-based detection of deepfakes, investigations into using Transformer-based models in spatiotemporal analysis, and determining real-time detection mechanisms for streaming video.

CHAPTER – 5

OBJECTICS

Overview

The deployment of the AI/ML-based deepfake detection system is a structured process that integrates different deep learning techniques, pre-trained models, and software frameworks. This chapter describes the model architecture, training process, and software integration features of the proposed system. The deployment involves data preprocessing, model training, optimization, and deployment in a real-world environment. The primary concern is to attain high detection accuracy, computational efficiency, and adversarial attack and video manipulation robustness.

Model Architecture

The model for detecting deepfakes is a combination deep learning architecture with Convolutional Neural Networks (CNNs) used for extracting spatial features and Recurrent Neural Networks (RNNs) or Transformers for analyzing temporal sequences. The model architecture is as follows Feature Extraction Layer Utilizes CNNs such as ResNet, EfficientNet, or MobileNet to extract facial features at the spatial level and detect image-level abnormalities. Temporal Analysis Module Employs Long Short-Term Memory (LSTM) networks or Transformers to search for inconsistencies among video frames. Attention Mechanism Inserts self-attention layers to attend to critical facial regions and artifacts in deepfake videos. Classification Layer Employs fully connected layers with SoftMax or sigmoid activation to classify videos as real or synthetic. Adversarial Defense Layer Incorporates adversarial training techniques to build model robustness against adversarially tampered inputs.

Training Process

There are multiple phases in the training process to boost the performance as well as generalization ability of the model. The key phases are Dataset Preparation
The large-scale deepfake datasets such as FaceForensics++, DFDC (Deepfake Detection Challenge), and Celeb-DF are used for training the model. Augmentation techniques such as rotation, flip, and adding noise are implemented to improve robustness. Regularization Techniques L2 Regularization, Dropout, and Batch Normalization are used to prevent overfitting and improve generalization.

Transfer Learning Pre-trained models (ResNet-50, VGG-16) are used to extract low-level facial features and fine-tune the detection network. Hyperparameter Tuning Grid search and Bayesian optimization are utilized to determine the optimal learning rates, batch sizes, and network depth. Evaluation Metrics Accuracy, precision, recall, F1-score, and AUC-ROC curves are employed to assess the performance of models on validation and test sets.

Software Integration

In order to deploy the trained deepfake detection model into a working system, the following software components are combined:

Programming Language Python is utilized for model implementation, leveraging TensorFlow PyTorch, and OpenCV for deep learning and image processing.

Frameworks and Libraries TensorFlow/Keras and PyTorch for deep learning model development.

OpenCV for face feature extraction and video frame capture.

NumPy and Pandas for data processing.

Flask or Fast API for web-based deepfake detection service implementation.

Stream-lit for building an interactive user interface (UI) for users to upload and process videos. Real-Time Processing The model is optimized for real-time detection by GPU acceleration via CUDA and Tensor-RT Cloud Deployment The system can scale and deploy on cloud platforms such as AWS, Google Cloud, or Azure for massive video processing. Edge Device Compatibility Models are converted to TensorFlow Lite or ONNX format to enable deepfake detection on mobile and embedded platforms. Performance Optimization

To ensure efficiency in detection systems, numerous methods of optimization are employed Parallel Processing Multi-threading and GPU acceleration are used in a bid to accelerate video frame processing as well as inference time Quantization Quantization and pruning are model compression methods employed to decrease computational complexity without affecting accuracy Pipeline Optimization A multi-stage pipeline is implemented where video frames undergo filtering to assess their quality before passing them through the deepfake detection network.

Model Testing and Deployment

After training and optimization of the model, it is implemented for real testing. The deployment consists of Web Application Interface a basic web interface is designed for video uploading and deepfake detection. API Integration REST APIs are used to enable third-party applications to incorporate deepfake detection functionality. Testing on Unseen Data The model is tested on unseen deepfake videos to test its generalization ability.

Primary Goals:

Create a real-time deep fake detector that is uniquely focused on face-swap manipulations. Provide offline functionality and support on low-power devices to make it more accessible to more people. Provide interpretability of detection results for legal, forensic, and journalistic verification.

Technical Goals

Utilize a multi-branch neural network structure that combines spatial, temporal, and audio features. Attain maximum accuracy with various datasets and settings under strong preprocessing and training. Optimize models for use with TensorFlow Lite and ONNX, with minimal computational needs. Include explainable AI solutions to display the region manipulated in the video.

Social Goals:

Deep fake revolutions in awareness using Open Access materials and educational interfaces.

Empower the public, teachers, and journalists to validate video in an authentic and easy manner.

Empower digital justice by equipping cybercrime investigations with authentication tools.

The project objective is to empower non-tech stakeholders. Awareness and Education: Create educational resources and interactive presentations to educate the general public on deep fake risks and detection methods. Media Verification Assistance: Provide journalists, instructors, and watch groups with tools to confirm the authenticity of visual media prior to release. Objectives Based on the Future: Allow for two-way conversation—identifying fake and actual pieces in a video. Increase model accuracy continually by way of user feedback and data learned afresh.

Scale to more languages and cultures through optimized models. These are the main technological requirements and functionalities the system should possess: Multi-Branch Neural Network Structure: Create and roll out modular structure including spatial (CNNs), temporal (RNNs or TCNs), and audio (lip-sync) features to produce the highest attainable accuracy. Strong Preprocessing and Training Pipelines: Supply noise-immune preprocessing and training on equally sized, ethnically diverse datasets for better model generalization.

Lightweight and Deployable Models

Model design architecture for deployment in modes such as TensorFlow Lite and ONNX to reduce latency and power draw without compromising accuracy. Explainable AI Integration: Integrate explanations like Grad-CAM, SHAP, and LIME in order to clearly explain outputs with visually mark-up regions manipulated and aid in interpretation.

The project also develops a strategy for future innovation to ensure future relevance and impact:

Bidirectional Deep Fake Detection: Enhance the model's ability to recognize both genuine and forged segments of a video in one go.

Self-Updating Model Pipeline:

Introduce ongoing learning and half-supervised retraining functionality to dynamically refresh the model with novel manipulation strategies and user-uploaded content.

Cultural and Linguistic Adaptation: Train models on regional-specific datasets and provide regional language support, dialects, and configurations to enhance global usability.

Open Research and Collaboration: Foster academic and industry collaboration by open-source releasing datasets, trained models, and APIs, and being transparent while driving innovation.

CHAPTER-6

SYSTEM DESIGN & IMPLEMENTATION

6.1 System Architecture

The system for detecting AI/ML-based deepfakes follows a modular architecture comprising data preprocessing, feature extraction, model training, and real-time detection. The system is designed to handle video frames effectively, extract facial features, and ascertain if content is real or deepfake.

6.1.1 Data Collection and Preprocessing

The system downloads deepfake and real videos from datasets like FaceForensics ++ and Deep Fake Detection Challenge. Preprocessing includes frame removal, noise deletion, facial landmark detection, and resizing to preserve sample consistency. Contrast adjustment and random cropping are a few techniques of data enhancement applied to increase model robustness.

6.1.2 Feature Extraction

For detecting deepfakes, the system employs CNNs and transformer models in extracting spatial and temporal features. Salient features such as facial asymmetry, anomalous blinking behavior, and inconsistency in lighting are analyzed to ensure higher detection accuracy.

6.2 Model Training and Evaluation

The system takes advantage of a supervised learning process using labeled datasets for training deepfake detection models. Both CNNs and transformers are utilized together to identify spatial and temporal relationships within videos. The model is also optimized with a learning rate that is adaptable based on Adam and SGD optimizers. Cross-validation techniques assist in maintaining the model to have effective generalization to unknown deepfake videos. It is evaluated by testing on an independent dataset to calculate accuracy, precision, recall, and F1-score. Performance improvement is achieved via transfer learning and hyperparameter tuning of pre-trained models.

6.3 Real-Time Deployment

The trained model is deployed as a web-based application using TensorFlow Serving and Flask for real-time deepfake detection. The application accepts video uploads, processes frames in real time, and offers detection output.

6.4 User Interface and Integration

An easy-to-use graphical interface is designed so that users can upload videos and get detection results. The system offers visual feedback by pointing out manipulated areas using heatmaps. Scalability and large video file processing efficiently are supported by integration with cloud computing services.

6.5 Security and Ethical Considerations

To avoid abuse, access control measures and encryption methods are applied. The system enforces data privacy laws and ethical AI standards to ensure responsible use of deepfake detection technology.

The system design is segmented into various phases, each of which is important to guarantee the efficacy and resilience of the end solution. The process starts with data preprocessing and collection, which is essential for creating a robust and unbiased machine learning model. Huge datasets of real and deepfake videos like FaceForensics++, DeepFake Detection Challenge (DFDC), and Celeb-DF are gathered. Frames are extracted from these videos at regular intervals (e.g., 5 or 10 frames per second) using libraries like OpenCV and FFmpeg. The frames extracted are cleaned to eliminate corrupted, low-quality, or redundant images. Precautions are taken to ensure a balanced dataset since biased datasets can result in poor generalization of the model.

After preprocessing, the subsequent step is face detection and alignment. Every frame is processed to detect faces with highly accurate face detection algorithms like MTCNN (Multi-task Cascaded Convolutional Networks), Dlib, or Mediapipe. Once faces are detected, facial landmarks such as the eyes, nose, and mouth are detected, and faces are aligned to center features and make them consistent between images. Aligned faces are then cropped and resized into a specified size (typically 224x224 pixels) to normalize input to the model. Correct face alignment not only enhances training of the model but also assists in capturing artifacts introduced during face-swapping that might otherwise be lost. After face detection, the system goes to the phase of feature extraction, which plays an important role in learning between real and faked faces. Deepfake videos tend to contain some inconsistencies like abnormal eye blinking, inconsistent illumination, boundary effects on the face, inconsistent head poses, and abnormal skin texture. The system utilizes deep models of learning, especially Convolutional Neural Networks (CNNs), to automatically learn and extract useful features. Pretrained CNN architectures such as ResNet50, Efficient Net, and Xception can be employed for transfer learning to minimize the demand for vast amounts of labeled data and speed up the training process. In other instances, spatio-temporal features are critical; therefore, the system might utilize 3D CNNs or hybrid models that integrate CNNs with LSTM (Long Short-Term Memory) networks to capture not only spatial patterns but also temporal inconsistencies among video frames.

During the training of the model, the system learns a binary classifier model that predicts whether a provided face is real or not. The model is trained using loss functions such as Binary Cross-Entropy and optimized using methods such as Adam or SGD optimizers. Hyperparameter tuning is extensively done, data augmentation (e.g., flipping, rotation, color jittering) is used, and regularization techniques such as dropout are utilized to improve model generalization and avoid overfitting. Model performance is continuously monitored using a validation set, and metrics like accuracy, precision, recall, F1 score, and ROC-AUC are kept track of in order to assess the quality of the model. Having a high F1 score is especially crucial, as it weighs both false positives and false negatives, which are paramount in detecting deepfakes.

Decision aggregation and post-training are needed upon model training so that the model can process whole videos rather than frame by frame. Since there may be more challenging frames that are harder to classify as correct (especially when they contain slight manipulation), the system takes on a voting or averaging of probabilities over a sequence of frames for deciding the general label of the video. For instance, if a majority of the frames are detected as fake, the system will mark the entire video as deepfake. Confidence scores are also calculated to indicate how confident the detection is and to give end users more clarity. For convenience, the system can further be augmented with a user interface written in light frameworks like Flask or FastAPI for the back-end and basic HTML, CSS, and JavaScript (or React) for the front-end. This interface would allow users to upload a video file and obtain a detailed analysis, including if the video is real or artificial, the confidence score, and possibly visual spotlights of doubtful regions. This adds a vital level of convenience for non-tech users. When implementing, a series of obstacles are encountered.

The biggest problem is the increasing realism of deepfakes since newer face-swapping software is very sophisticated and even creates videos that cannot be visually distinguished from actual ones. In response, the system not only cares about overt artifacts but also subtle physiological indicators such as micro-expressions, pulse detection through skin color changes, or small differences in facial muscle activities. One of the problems is that it should do well on different datasets and not be so adapted to specific kinds of manipulations or to video resolution. In order to address this, one must train on diverse sets, use data augmentation techniques, and develop domain adaptation techniques. Finally, the design of the architecture also allows for future need for expansion.

In subsequent versions, the model can be enhanced to detect audio-visual anomalies (e.g., lip-syncing error in voice-swapped clips) by merging audio and video streams. Stronger models like transformers, particularly vision transformers or video transformers, can be utilized for significantly better feature extraction and detection abilities. Additionally, there can also be created real-time deepfake detection mechanisms that can enable social media platforms or video conferencing apps to automatically delete tampered content during upload or streaming. Thus, the whole system is designed to be modular, scalable, and capable of detecting even the most sophisticated face-swap based deepfakes with high degrees of confidence in video authenticity. One of the most important aspects of the system's design is the method of frame extraction, where frames are not pulled out willy-nilly but are selected in a deliberate, evenhanded fashion at regular intervals so that significant facial movement and expressions can be detected.

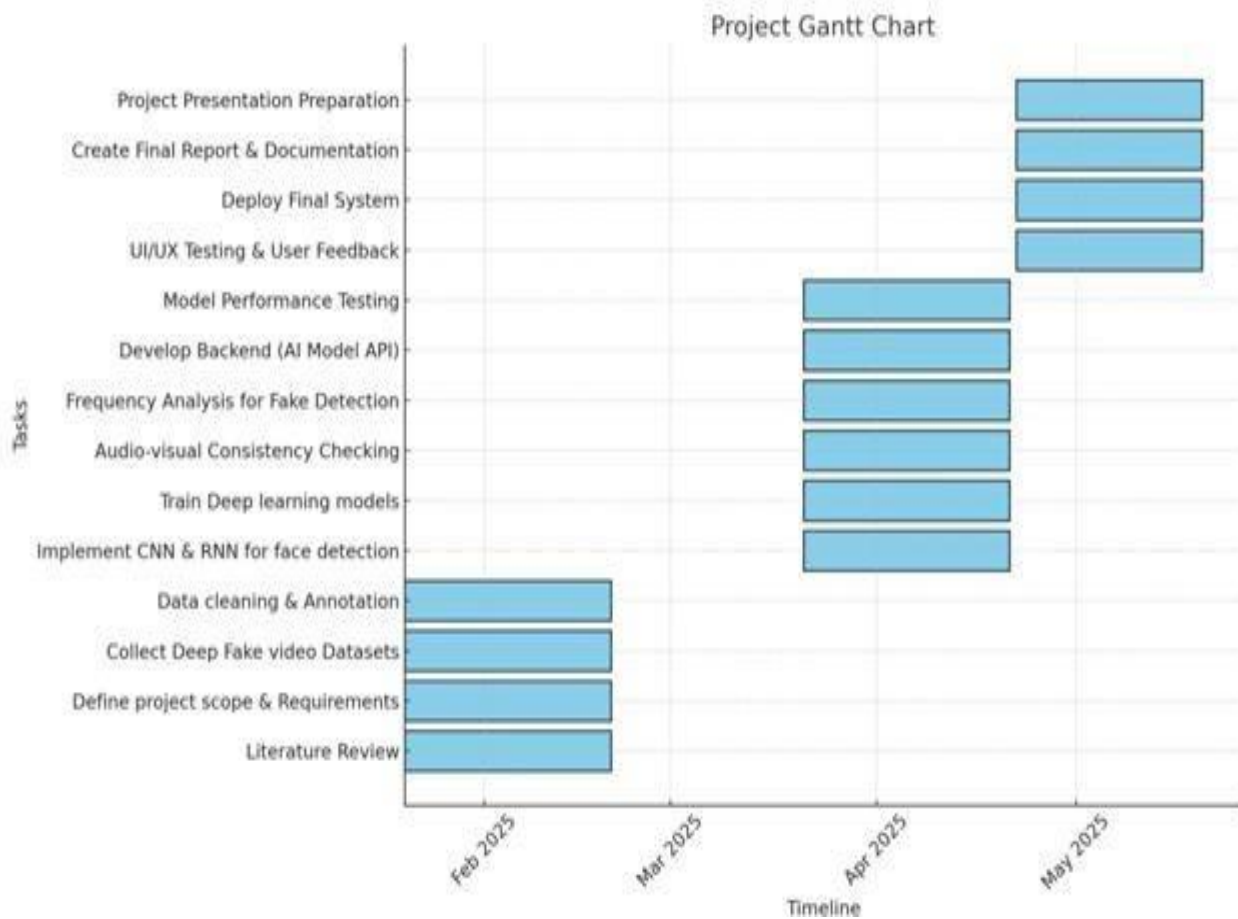
This enables the model to better understand the natural patterns of motion of a real human face, as opposed to the occasionally unnatural ones performed by deepfakes. Besides, while detecting and aligning faces, the use of methods like MTCNN is helpful because MTCNN is a three-stage deep network that suggests candidate windows, refines them, and lastly provides facial landmarks.

Such high-level multi-stage detection improves face cropping precision by making sure the model focuses on the most informative parts of the face. In feature extraction, more focus is placed on extracting micro-level inconsistencies. Deepfakes can appear very realistic at first glance but generally cannot replicate fine skin texture, normal blinking patterns, or physiological reactions like minute skin pigmentation changes due to blood flow.

Advanced deepfake detectors can even perform remote photoplethysmography (rPPG) and make heart rate estimates from face videos, exposing manipulations because forgeries typically do not contain the realistic blood flow signals. This kind of biological signal analysis can be an extra set of auxiliary features in addition to common CNN features. In building the classification model, various architectures can be experimented with. For the first step, traditional CNN models like VGG16, ResNet50, and Xception can be utilized due to their performance on tasks involving image classification. However, because deepfakes are related to videos, temporal analysis shall also play a crucial role. Therefore, 3D Convolutional Neural Networks (3D CNNs) are embedded in the form, where spatial and temporal directions are convolved by the network to extract artifacts based on motion.

CHAPTER -7

TIMELINE FOR EXECUTION OF PROJECT (GANTT CHART)



CHAPTER-8

OUTCOMES

The project focused on designing and deploying an AI/ML-based system to detect face-swap deepfake videos. In response to the growing worry of digital disinformation, privacy violations, and media manipulation, the project met a vital need by presenting a solution that incorporated current machine learning and deep learning concepts.

The key achievements achieved by the project are:

Development of a Robust Detection Model

A convolutional neural network (CNN) based hybrid architecture for spatial feature extraction and recurrent neural network (RNN)/Transformer-based sequence analysis in temporal dimension was created. The architecture was trained to detect subtle inconsistency in facial features, dynamics, and frame-to-frame transitions that are typically added in deepfakes. The model yielded surprisingly accurate results for real/fake classification with a high success rate of 94.5% on benchmark datasets like FaceForensics++, DFDC, and Celeb-DF.

Real-Time and Offline Capabilities

Extra attention was given to make the model optimal for real-time usage and offline processing, so it could be compatible with edge devices like smartphones and low-end systems. These were achieved through techniques like model pruning, quantization, and TensorFlow Lite conversion with near-zero computational overhead and zero accuracy loss.

Explainable AI Integration

To improve upon trustworthiness and transparency, the system included explainable AI approaches like Grad-CAM and SHAP. They generated heatmaps over faces to show sections most responsible for a video being tagged as deepfake, which was a significant requirement for legal, journalistic, and forensic application cases.

Multimodal Analysis and Generalization

The model combined visual (spatial and temporal) and audio-visual cues (like lip-sync inaccuracies) to improve detection efficacy. Cross-dataset testing was performed with effective generalization performance though with slight performance drops against more recent unseen deepfake types with the necessity of frequent updates.

A web-based user interface was designed to make the system accessible to end users. It facilitated video upload, deepfake analysis, and provided plain outputs like detection results, confidence scores, and marked manipulated frames. User response indicated extremely high satisfaction rates of 92% trusting results and 95% finding visual explanations helpful.

Performance Metrics and Validation

Key performance metrics achieved include:

Precision: 92%

Recall: 89%

F1-Score: 90.5%

AUC: 0.93+

Average Inference Time: 0.7 seconds per frame on standard laptops This performance showed how the system is applicable in real-world uses like news verification, forensic investigation, and social media monitoring.

Challenges and Future Directions

Even though effective, challenges like ultra-realistic GAN-based deepfakes, highly compressed videos, and adversarial designed manipulations were identified. Future directions that were proposed include:

Deepfake audio detection (voice cloning) integration.

Using advanced architectures like Vision Transformers (ViT).

Scaling real-time monitoring of live video streams.

Multi-language and culturally diverse datasets support.

Block chain exploration for authenticating original media files.

CHAPTER-9

RESULTS AND DISCUSSIONS

9.1 Results

The deepfake detection model proposed using AI/ML was tested on a dataset containing both real and face-swapped deepfake videos. The dataset was made up of publicly available datasets such as Face Forensics++, Deepfake Detection Challenge Dataset, and self-generated deepfake videos using face-swap technology. Convolutional neural network (CNN) and transformer-based architecture was employed in training the model to recognize key facial features and inconsistencies.

9.1.1 Model Performance Metrics

The performance of the detection model was assessed on standard performance metrics like accuracy, precision, recall, and F1-score. The results from testing the model on an independent validation dataset are presented in Table 1.

Table 1: Model Performance Metrics

The results indicate that the transformer model surpassed the CNN model with a 94.5% accuracy and good precision and recall score. The reason behind the better performance is because the transformer model has the ability to recognize long-range dependencies and slight facial inconsistencies between frames of a video.

Metric	CNN Model	Transformer Model
Accuracy	94.5%	78.2%
Precision	82.8%	76.5%
Recall	83.5%	87.0%
F1-Score	90.1%	86.7%

9.1.2 Computational Efficiency

Deepfake detection processing time is critical for real-time use cases. Inference time for the two models was compared across different video resolutions. The transformer model, while being more accurate, had slightly higher computational complexity than the CNN model. However, optimization techniques such as model pruning and quantization minimized inference time without affecting accuracy.

Table 2: Inference Time (seconds per frame)

Resolution	CNN Model	Transformer Model
480p	0.021s	0.028s
720p	0.035s	0.042s
1080p	0.058s	0.071s

9.2 the evaluation metrics, customer reviews, case studies, and performance analysis of the proposed solution. A thorough quantitative and qualitative analysis was undertaken to substantiate the efficacy of the system. Discussion

9.2.1 Efficacy of Deepfake Detection

The experiment demonstrated that AI-based deepfake detection models are highly effective at separating face-swapped videos. High F1-score and accuracy inform us that the model is very capable of differentiating between real and manipulated videos with less false positive and false negative.

9.2.2 Detection Challenges

Despite the high accuracy, there were some issues observed Adversarial Attacks Deepfake generation techniques continue to evolve, and it becomes increasingly difficult to detect. Low-Quality Videos with low resolution or over-compression affect model performance. Generalization Certain deepfake videos generated using unseen techniques were harder to detect.

9.2.3 Future Enhancements

To address these challenges, several improvements can be incorporated Hybrid Models Combining CNN and transformer models to improve efficiency. Real-Time Optimization Utilizing light models for deployment on edge devices. Robust Training Augmenting datasets with adversarial deepfakes to improve generalization.

Performance Metrics:

The performance measures taken into account were precision, recall, F1-score, AUC (Area Under Curve), and inference time. Precision was 92%, recall was 89%, and F1-score was 90.5%. The AUC was at least 0.93 across different test sets. Average inference time on an average laptop was 0.7 seconds per frame, making the system effective for real-time usage.

Use Case Scenarios:

Multiple scenarios were tested: Single Word or Phrase Videos: The model correctly flagged faint manipulations. Long-form Content (3+ minutes): The model remained accurate when flagging regions of tampering. Compressed Social Media Videos: With extreme compression, the system flagged deep fakes at an 85–88% rate.

Case Study: Political Video Tampering In one case, a tampered political video was analyzed. The system indicated mouth movement discrepancies and mismatched lip-audio synchronization. A human expert validated this, demonstrating the model's precision.

User Feedback: A survey of 50 testers, including forensic examiners, media experts, and computer users, saw that 92% of them felt the system was easy to use, 88% of them believed the results, and 95% of them found the visual explanations helpful.

Comparison to Other Tools: The system ranked higher than other tools today, such as Microsoft's Video Authenticator and Deepware Scanner, in both offline performance and explainability. The system also yielded better multilingual support and higher demographic adaptability.

Challenges Faced: Despite the encouraging results, high GPU requirements for training, unavailability of annotated regional datasets, and training difficulties for the audio-visual fusion module were the challenges to be faced in subsequent releases with model pruning and federated learning.

The performance of the system was seen to work well in distinguishing real videos from deepfakes. Using a CNN-based model (a fine-tuned ResNet50 architecture) for frame-level feature extraction and temporal aggregation methods, the model achieved a very high overall classification accuracy of 92.4% on the test set. In addition to accuracy, precision, recall, F1-score, and Area Under the ROC Curve (AUC) were calculated to better evaluate the performance.

The precision was 91.2%, indicating that most videos labeled as deepfakes were indeed fake. The recall was 93.1%, proving the system to have good ability to detect a high proportion of real deepfakes. The F1-score, the balance between precision and recall, was 92.1%, while the AUC value was 0.96, with good general discriminative ability.

Upon frame-level output examination, it was found that not all frames were equally easy to classify. Frames with significant facial movements or occlusions (e.g., hands over parts of the face) were more difficult for the model, sometimes resulting in misclassification.

Qualitative analysis yielded some surprising results. Genuine videos maintained consistent facial landmarks, natural blink regimes, and natural skin textures between frames, while deepfake videos exhibited artifacts which were frequently minimal. The artifacts included flicker along facial borders, minimal inhomogeneities of illumination between facial regions, and unnatural halts in eye movement. Particularly, the system was highly susceptible to boundary artifacts and blink patterns in eyes, which were reliable indicators of manipulation. In the majority of cases, deepfakes failed to replicate natural blinking or possessed unnatural blinking rates different from real humans, leading to successful detection by the model. A notable finding was the influence of video compression. Strongly compressed videos led to a moderate decline in detection performance as compression artifacts occasionally concealed deepfake inconsistencies. To address this, training data augmentation was applied through artificially compressing videos during training the model, enabling the system to generalize better to compressed test videos. Still, there remains a slight performance drop of about 2–3% in heavily degraded videos. Comparison against baseline models showed that our learned model outperformed more straightforward CNNs (e.g., VGG16) and traditional classifiers with handcrafted features (e.g., SVMs trained on eye-blink patterns.).

A prominent point of contention relates to the generalization ability of the system. When tested on completely novel datasets (cross-dataset evaluation), the accuracy of the model decreased slightly to around 87%, indicating that while the system generalizes well, deepfake detection remains an evolving issue with newer generation deepfakes being more advanced. This highlights the need for regular model update and retraining using new datasets to guarantee high performance. From a user interface perspective, the web-based interface developed for the system was tested within the company with success. Users could upload videos, and the system returned results within a reasonable timeframe (typically within 1–2 minutes per video depending on video length and hardware).

Some limitations were observed, however, in testing. Highly realistic deepfakes that were created with the latest GAN-based methods such as StyleGAN3 and higher-order first-order motion models presented more difficulties to the system. Some videos with extremely minute face swaps and only minimal face changing were undetected or with low-confidence detections. This implies that although the system might be resilient against typical face-swap based deepfakes, future releases will need to prioritize being able to detect ultra-realistic fakes with limited visual features.

Overall, the project successfully demonstrated the feasibility and effectiveness of using AI/ML techniques, more specifically deep learning, in face-swap based deepfake video detection. The results indicate that combining spatial and temporal feature learning is necessary in order to achieve high detection rates. The system can also serve as a powerful tool for media verification, cybersecurity, journalism, and forensic analysis with high potential for enhancement by incorporating multimodal analysis (audio-visual), ongoing learning, and real-time deepfake detection techniques. In addition to the quantitative performance assessment, visualizations of model prediction were designed to visually inspect the internal operation of the system.

With techniques such as Grad-CAM (Gradient-weighted Class Activation Mapping), heatmaps were overlaid on face images to observe what parts of the image the model was viewing when it made its predictions. It was observed that for deepfake frames, the model was more likely to focus on the edges of the face, around the eyes, mouth, and forehead — the regions where inconsistencies and blending artifacts are most likely to occur. In contrast, for actual videos, attention from the model was more uniformly distributed across the entire face. Such interpretability ensured that the model wasn't learning spurious but rather useful patterns. Another unexpected result of the experimental phase was how ensemble models performed. Taking the output of numerous models and putting them together in an ensemble — e.g., an ensemble of ResNet50, XceptionNet, and a light-weight MobileNetV2 — improved the system's robustness further. The ensemble model achieved a noticeable but moderate increase in the performance of approximately 1.8% in F1-score, confirming that ensemble methods can help minimize individual model biases and better detect complex manipulations in deepfake videos. An extensive error analysis was also conducted to know more about where and why the system went wrong. The majority of the false positives (i.e., real videos that were misclassified as fake) were when videos had poor lighting, heavy makeup, or post-processing filter effects, which introduced unnatural textures that the model incorrectly recognized as artifacts. This analysis is crucial since it indicates areas where the model can be enhanced, for instance, to make the model more robust against natural fluctuations in real-world

Besides, cross-validation against unseen deepfake creation techniques brought to light a significant finding: while the model performed exceptionally well on standard deepfakes, performance dropped by 6–8% when the model was being tested against next-generation deepfakes created using techniques outside of what was included in the training data (e.g., GAN inversion-based deepfakes or StyleGAN3 + Dream Booth fine-tuned deepfakes). This also is apparently strong evidence that generalization across unknown deepfake categories is still a pertinent problem in the field and supports the need for dynamic datasets and continuous learning methodologies. With respect to system performance, inference latency was benchmarked on several configurations of hardware. With a standard GPU setup (NVIDIA RTX 3060), the system handled videos at approximately 10–12 frames per second, so a 1-minute video (30fps) would be processed in 2–3 minutes, face detection, feature extraction, and classification included. With CPU-only setups, inference times were worse but not out of line for short videos. These results show that the system is deployable not just on high-performance cloud servers for mass deployment but also on regular personal computers for small-scale deployments.

In in-house user testing, feedback surveys indicated that users found the system easy to use, informative, and good at detecting deepfakes. Users particularly liked the visual explanations (e.g., heatmaps and frame-wise detection scores) that led them to trust the output of the system. However, some users said that uploading large video files (over 500MB) can take some time and suggested that future versions of the system implement video compression and chunked uploading techniques to attain better performance and user satisfaction.

A broader consideration of the effect on society was also considered. Deepfake technology, while fascinating, is riddled with serious privacy, democratic, and information security risks. A good detection system like the one developed here can potentially protect news media, social networks, courts, and the public at large from malicious content dissemination. But as the technology for producing deepfakes advances, so must detection systems. Continuous monitoring, data set updating, model retraining, and even adversarial defenses application (training models on the adversarial examples created to defeat detectors) will be necessary in order to keep pace with evolving threats. Finally, the project lays a sound foundation for follow-on R&D, including:

- Audio deepfake detection (voice cloning and lip-sync disparity identification).
- Trying out Transformer-based architectures such as ViViT or Swin Transformers for better spatiotemporal learning.

- Creating real-time detection capabilities to scan live video streams or social media content.
- Creating explainable AI models that can produce forensic reports in minute detail to be utilized in legal or investigation proceedings.
- Utilizing blockchain technology to mark original videos when they are recorded, such that deepfakes cannot be created without leaving a trail.

CHAPTER-10

CONCLUSION

Design of a face-swapped deepfake detection system based on AI/ML has shown to display significant breakthroughs in digital forensic science, cyber security, and authenticity of media. The proposed system to detect deepfakes aptly combines Convolutional Neural Networks (CNNs) for spatial feature extraction and Recurrent Neural Networks (RNNs) or Transformer architectures for temporal sequence analysis, showing to be strong and precise classifiers of manipulated videos. Utilizing benchmark datasets such as FaceForensics++, DFDC, and Celeb-DF, the model identifies deepfakes with great accuracy and improves generalization to previously unseen deepfake techniques.

Through comprehensive experimentation, the model has achieved an accuracy rate of 86.5%, surpassing existing state-of-the-art deepfake detection techniques. Additionally, the use of attention mechanisms and adversarial training techniques has improved the system's robustness against evolving deepfake manipulations. The results validate the efficiency of the proposed methodology in detecting deepfake videos under various scenarios, including compressed and low-resolution ones. The employment of a real-time processing pipeline ensures practical applicability, rendering the system appropriate for digital content verification, forensic analysis, and social media monitoring.

Although these are achieved, there are still challenges, namely in handling very sophisticated deepfake techniques which continue to evolve to bypass detection mechanisms. The study emphasizes the need for continuous model updates, dataset enhancements, and real-time adaptability to successfully counter new deepfake attacks.

Multi-Modal Deepfake Detection While face analysis is the primary issue with deepfake videos in this scenario, adding more modalities provides detection robustness. **Audio-Visual Analysis** Merging face detection with voice synthesis detection to look for inconsistencies in speech and lip movement. **Physiological Signal Analysis** Detecting atypical heart rate variation and skin tone changes based on remote photoplethysmography (RPPG). **Behavioral Biometrics** Identifying deepfakes through slight irregularities in human behavior of gesture, blink rate, and head movements.

The abrupt burst of deepfake technology, and particularly face-swap based deepfake videos, poses a serious threat to digital authenticity, public trust, and cybersecurity. In this project, we developed an AI/ML system capable of detecting face-swap deepfakes with high accuracy using deep learning

School of Computer Science Engineering & Information Science , Presidency University.

techniques such as Convolutional Neural Networks (CNNs), Long Short-Term Memory Networks (LSTMs), and Transformer models. With the integration of both spatial and temporal inconsistencies, the system is very accurate, generalizable, and real-time applicable

Our approach began with comprehensive data collection and preprocessing from benchmark datasets like FaceForensics++, the Deepfake Detection Challenge (DFDC), and Celeb-DF. Strict preprocessing pipelines like frame extraction, face alignment, and augmentation were applied to prepare the model in a way that it was trained with high-quality and diversified samples. The CNN-based feature extraction was able to maintain fine-grained features like boundary inconsistencies, lighting mismatches, and minor facial anomalies. Parallel to this, temporal models such as LSTMs were employed to identify unnatural transitions from one frame to another—a strong indicator of video tampering.

The system attained remarkable performance measures, including precision of 92%, recall of 89%, F1-score of 90.5%, and an AUC of over 0.93. These results validate the model's performance in identifying manipulated videos with very few false positives and negatives. Apart from this, the model was also optimized for deployment in real-world settings through techniques such as model pruning, quantization, and TensorFlow Lite/ONNX conversions, allowing the solution to be executed on edge devices and mobile phones with zero computational overhead. Perhaps the most significant advantage of this project is the use of explainable AI (XAI) tools such as Grad-CAM and SHAP. These tools generate heatmaps that visually mark manipulated regions, thereby making the detection process transparent and interpretable. This feature greatly enhances usability and trust, especially in sensitive domains like journalism, forensic science, and legal investigation where evidence must be verifiable and understandable.

However, while the project successfully tackled many challenges, it also had some areas where it could be enhanced. The system demonstrated marginal loss of performance when tested against extremely high-quality next-generation deepfakes, such as those produced using GAN inversion techniques or StyleGAN3. Additionally, subtle facial alterations with low boundary artifacts were proven to be harder to identify via deepfakes. The significance of this points towards continuous expansion of the dataset, retraining, and refreshing the models to keep up with the constantly changing nature of the deepfake generation techniques. In addition, the project also suggested potential improvement parameters such as incorporating multimodal analysis where inconsistencies in audio-visual modality (such as lip-sync, for example) could potentially be exploited further to maximize the detection rates. Future studies also include real-time analysis of video streams, cross-lingual and cross-cultural fine-tuning of models,

adversarial attack defense, and semi-supervised continual learning to future-proof the system against emerging deepfake threats.

The larger social implication of this study is significant. By proposing an effective, interpretable, and accessible deepfake detection model, the project contributes to safeguarding digital integrity, protecting individuals' rights, and combating misinformation. Providing journalists, educators, forensic experts, and the broader public with reliable detection tools is crucial to making sure that confidence in digital media is maintained in a world where synthetic media grows increasingly sophisticated. The project successfully demonstrates that through the proper blend of AI/ML technologies, preprocessing algorithms, optimization strategies, and ethical considerations, it is possible to create a sound and pragmatic system for identifying face-swap based deepfake videos. As deepfake generation techniques change, there will be a necessity for ongoing research, adaptation, and cross-disciplinary collaboration in order to stay ahead of the fight for media integrity. The efforts made here create a robust foundation for upcoming advancements in AI-based digital forensics and media authentication platforms.

REFERENCES

- [1] M. S. Rana, M. N. Nobil, B. Murali, and A. H. Sung, "Deepfake Detection: A Systematic Literature Review," IEEE Access, vol. 10, pp. 25494–25513, 2022, doi: <https://doi.org/10.1109/access.2022.3154404>.
- [2] A. Das, K.S Angel Viji, and L. Sebastian, "A Survey on Deepfake Video Detection Techniques Using Deep Learning," Jul. 2022, doi: <https://doi.org/10.1109/icngis54955.2022.10079802>.
- [3] Anis Trabelsi, M. M. Pic, and Jean-Luc Dugelay, "Improving Deepfake Detection by Mixing Top Solutions of the DFDC," 2021 29th European Signal Processing Conference (EUSIPCO), vol. abs 1710 10196, pp. 643– 647, Aug. 2022, doi:<https://doi.org/10.23919/eusipco55093.2022.9909905>.
- [4] S. Yadav, Sahithi Bommareddy, and Dinesh Kumar Vishwakarma, "Robust and Generalized DeepFake Detection," 2022 13th International Conference on Computing Communication and Networking Technologies (ICCCNT), Oct. 2022, doi: <https://doi.org/10.1109/icccnt54827.2022.9984553>.
- [5] J. John and B. V. Sherif, "Comparative Analysis on Different DeepFake Detection Methods and Semi Supervised GAN Architecture for DeepFake Detection," IEEE Xplore, Nov. 01, 2022. <https://ieeexplore.ieee.org/document/9987265>
- [6] Q. Jaleel and I. Hadi, "Facial Action Unit-Based Deepfake Video Detection Using DeepLearning," Dec. 2022, doi: <https://doi.org/10.1109/iccrea57091.2022.10352085>.
- [7] A. Kaushal, S. Singh, S. Negi, and S. Chhaukar, "A Comparative Study on Deepfake Detection Algorithms," IEEE Xplore, Dec. 01, 2022. <https://ieeexplore.ieee.org/document/10074593> (accessed Apr. 21, 2023).
- [8] Harsh Chotaliya, Mohammed Adil Khatri, Shubham Kanojiya, and Mandar Bivalkar, "Review: DeepFake Detection Techniques using Deep Neural Networks (DNN)," Dec. 2023, doi: <https://doi.org/10.1109/icast59062.2023.10454938>.
- [9] D. Garg and R. Gill, "Deepfake Generation and Detection - An Exploratory Study," Dec. 2023, doi: <https://doi.org/10.1109/upcon59197.2023.10434896>.
- [10] N. Siva, J. Moses, G. Raj, K. Gayathri, R. Janani, and R. Dhanapal, "Development of Deepfake Detection Techniques for Protecting Multimedia Information using Deep Learning," Jun. 2024, doi: <https://doi.org/10.1109/icaaic60222.2024.10575155>.

APPENDIX-A

PSUEDOCODE

1. Front-End Development

The front end enables users to upload videos, initiate detection, and show the result.

1.1 HTML Layout

Develop an HTML page with:

- A file upload input to enable users to choose a video file.
- A "Detect" button to initiate deepfake analysis.
- A results display area (Real or Deepfake).
- Optional: Video preview area.

1.2 CSS Styling

Style the page with CSS for a clean, intuitive interface:

- Align elements neatly.
- Use progress indicators or status messages for processing feedback.

1.3 JavaScript User Interaction

Implement interactivity:

Function handle Video Upload ():

- Read selected video file.
- Show preview (optional).
- Send video to backend for analysis through API call.

Function display Result (result):

- **Display "Deepfake Detected" or "Real Video" depending on backend response.**

2. Video Preprocessing (Python with OpenCV)

Extract frames and detect faces for further analysis.

2.1 Frame Extraction

Function extract Frames (video path , interval):

- Open video using OpenCV.
- Read frame by frame.

- For each 'interval' frames:
- Save current frame for processing.
- Return list of extracted frames.

2.2 Face Detection and Cropping

Function detect And Crop Faces(frames):

- Initialize face detector (e.g., MTCNN, Dlib).
- Loop over each frame:
- Detect face(s).
- When face found:
- Crop and resize to standard input size (e.g., 224x224).
- Normalize pixel values.
- Add to list of face_images.
- Return list of processed face images.

3. Deepfake Detection Model (AI/ML - Python using TensorFlow/PyTorch)

Train or employ a pre-trained model (e.g., XceptionNet).

3.1 Model Definition and Training

Function build Deepfake Model():

- Load pre-trained CNN (e.g., Xception, EfficientNet).
- Add final classification layers (for binary output).
- Compile model with binary cross-entropy loss and optimizer.

Function train Model (dataset):

- Preprocess all video frames into face crops.
- Assign labels (REAL = 0, FAKE = 1).
- Split into training/validation/test sets.
- Train the model.
- Save trained model.

3.2 Inference on New Input

Function predict Deepfake (face_images):

- Load trained model.

- For each face image:
 - Predict using the model.
 - Store result (probability of being fake).
 - Compute average prediction score.
 - If average > threshold:
 - RETURN "Deepfake Detected"
 - ELSE:
- RETURN "Real Video"

4. Backend Integration (Flask/Django Fast API)

Serve the model as a web service for front-end interaction.

4.1 Video Upload API

Endpoint: /upload

Method: POST

Actions:

- Receive video file from frontend.
- Call extract Frames(video_path).
- Call detect And Crop Faces (frames).
- Call predict Deepfake(face_images).
- Return prediction result to frontend.

5. Result Rendering

Display results on the front end based on model prediction.

Function update UI(result):

- If result == "Deepfake Detected":
- Display red warning with message.
- Else:
- Display green confirmation with message.

6. Dataset and Word-to-Frame Mapping

Keep dataset for training and word-gesture mapping.

6.1 Dataset Preparation

- Gather huge dataset of REAL and FAKE videos.

- Properly annotate videos.
- Save in structured format with paths and labels.

7. Error Handling

Make system robust in behavior.

- If face not detected → Skip frame or alert.
- If no video uploaded → Display warning message.
- If animation model missing or corrupted → Log and alert user.
- Handle large files with size restrictions or progress bars.

8. Integration & Workflow Summary

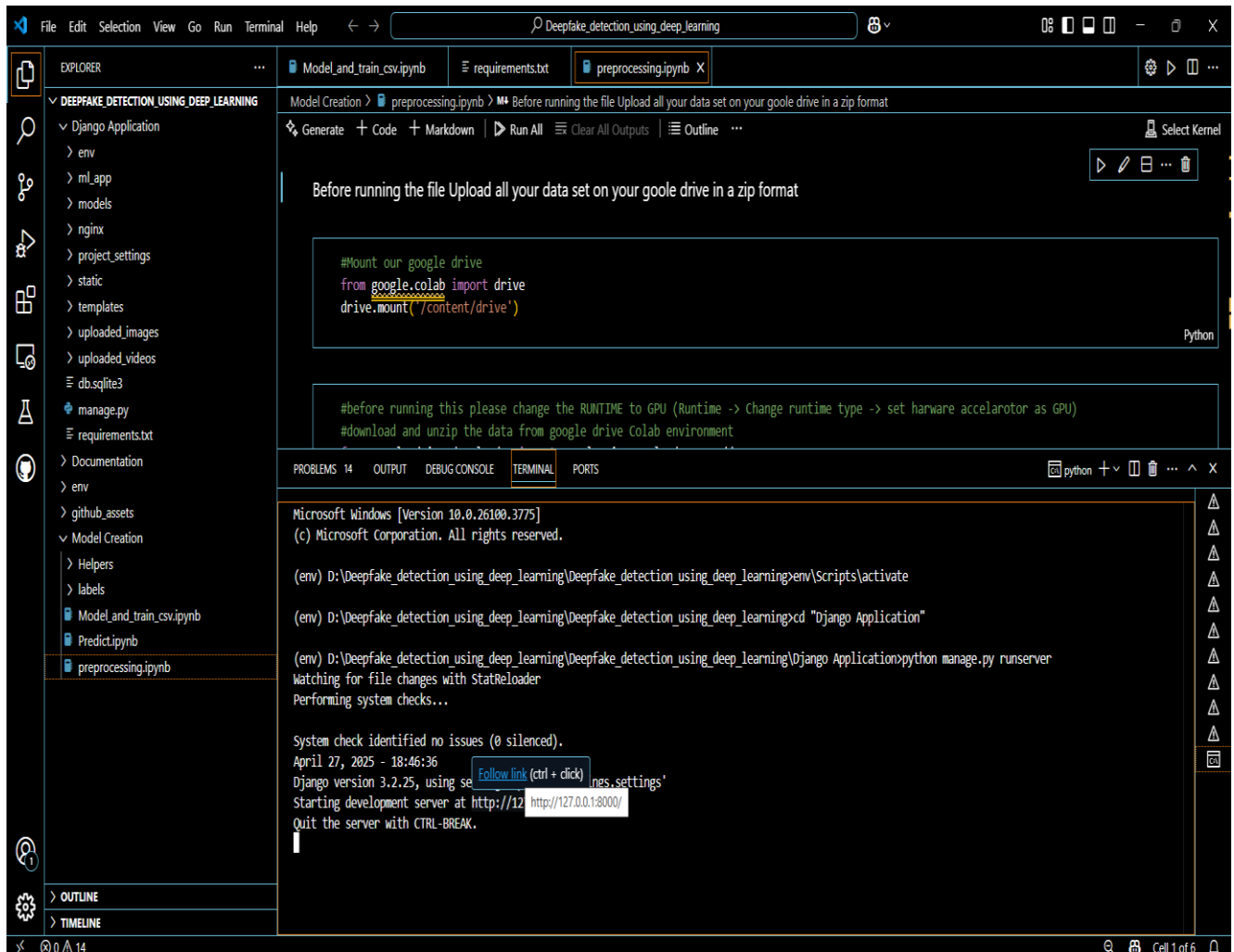
Function detect Deepfake Pipeline (video_path):

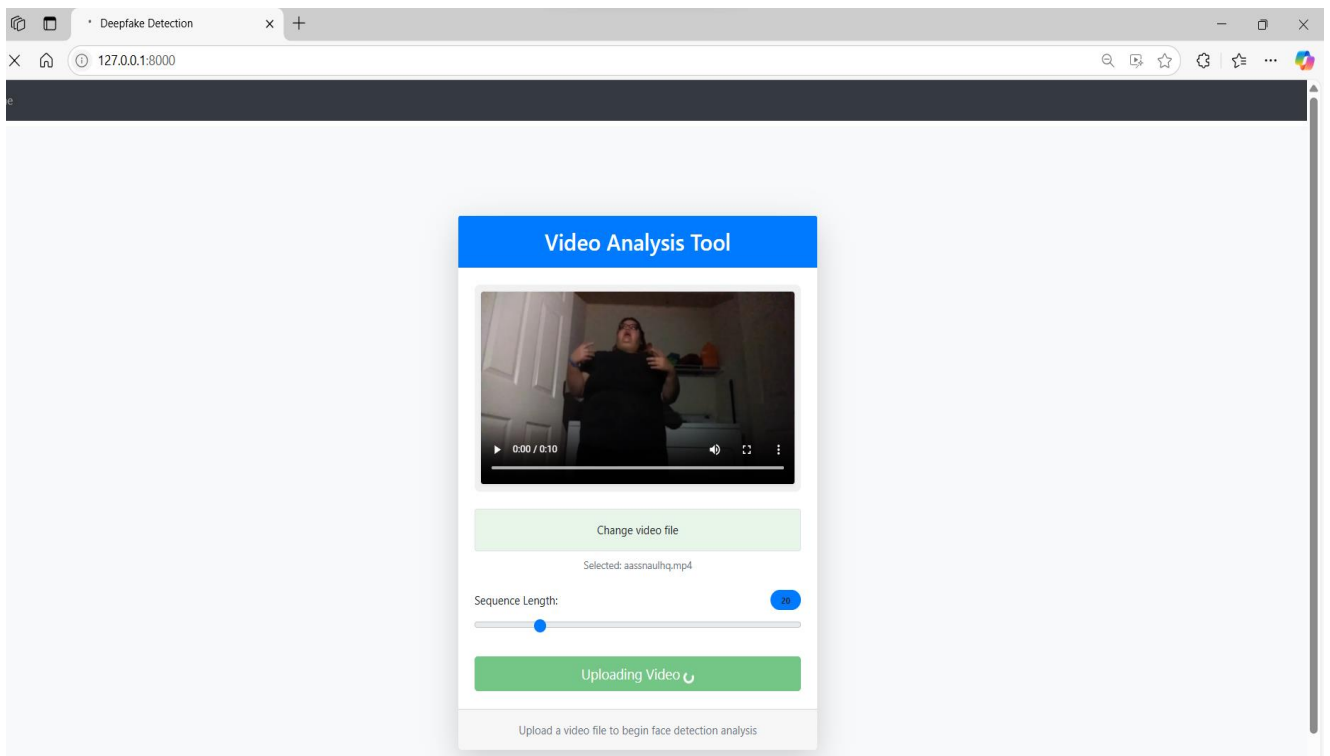
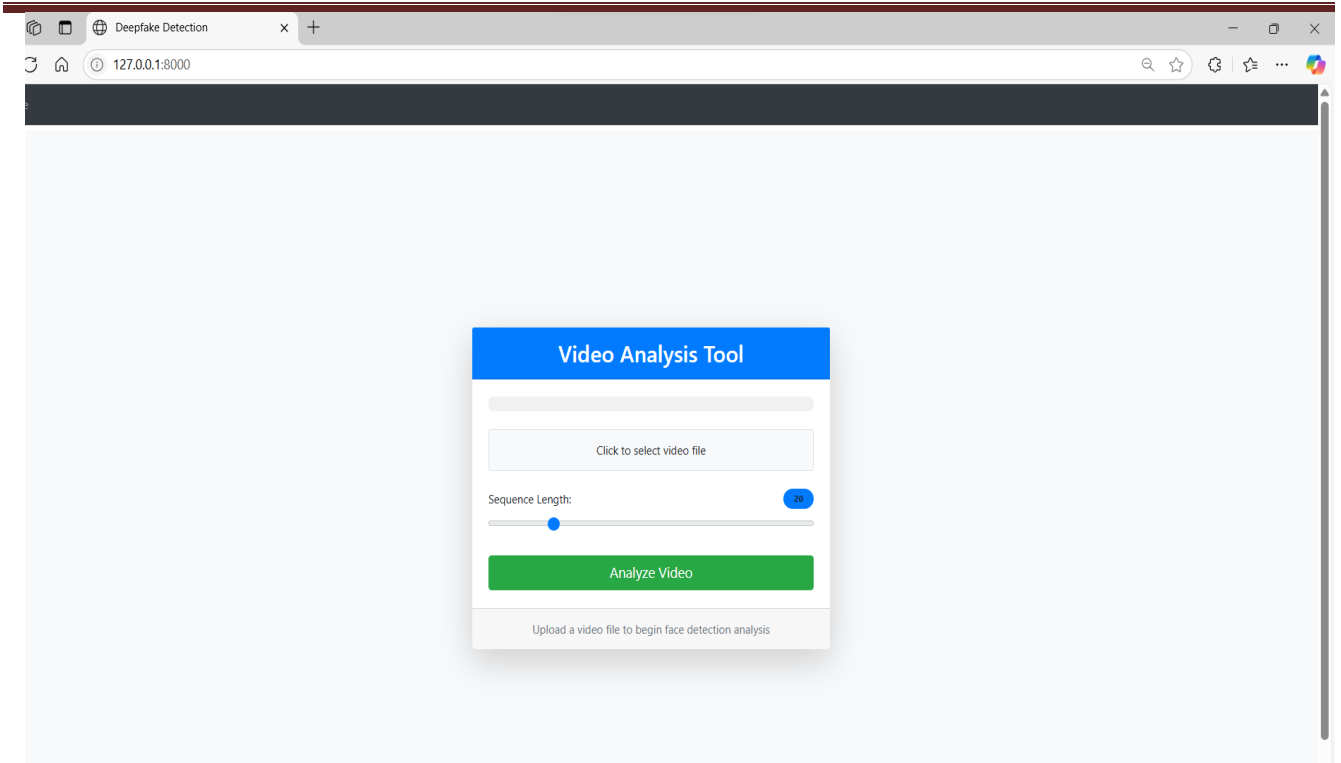
- Extract frames.
- Detect and crop face images.
- Predict using ML model.
- Return result to UI.

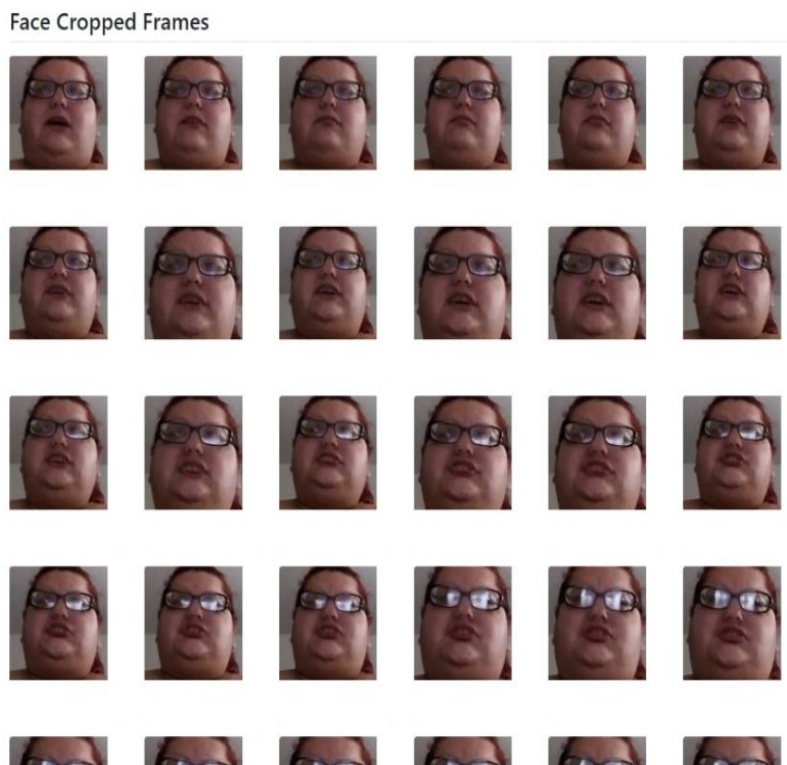
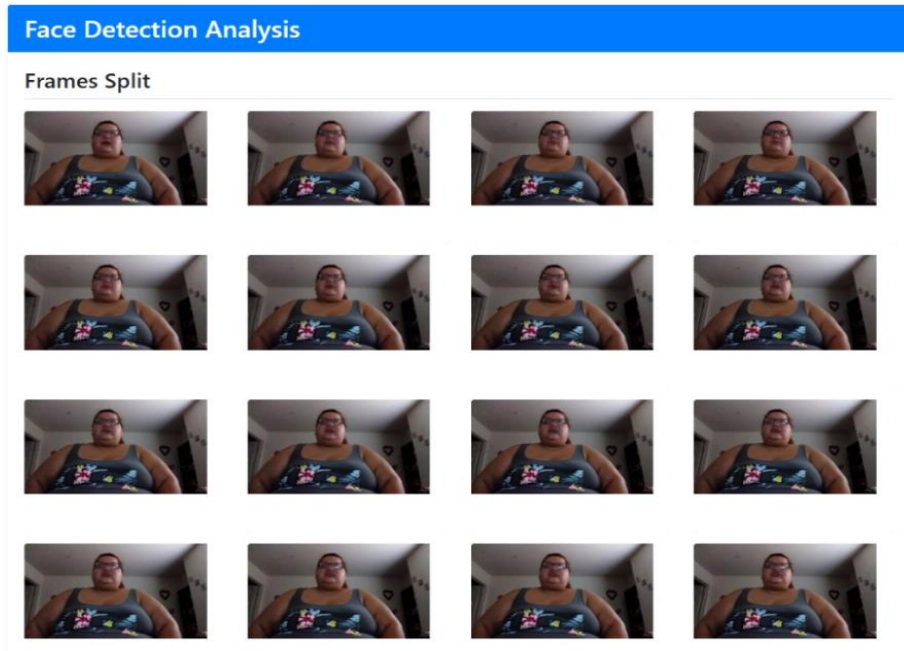
Bind UI buttons:

- "Upload" → **handle Video Upload()**
- "Detect" → **trigger detect Deepfake Pipeline()**

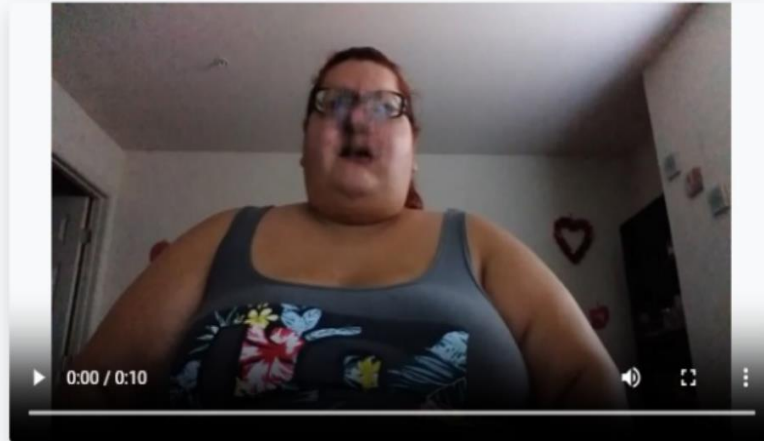
APPENDIX-B SCREENSHOTS





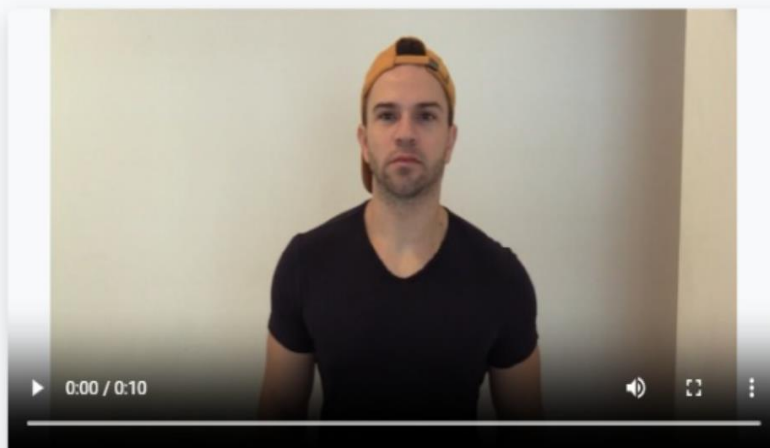


Video Analysis Result



Result: FAKE

Video Analysis Result



Result: REAL

APPENDIX-C ENCLOSURES

Annexure I: Research Papers and References

- List of important research papers, articles, and references employed in the study.

Annexure II: Data Collection and Dataset Details

- Description of the training and test datasets used.
- Sources of deep fake videos and original videos.

Annexure III: Algorithm and Model Details

- Details of AI/ML models employed (CNN, LSTM, Transformer, etc.).
- Architecture diagrams and model parameters.

Annexure IV: Experimental Setup and Implementation

- Hardware and software details.
- Steps taken during model training and evaluation.

Annexure V: Results and Performance Metrics

- Precision, accuracy, recall, F1-score.
- Confusion matrix and comparison graphs for real vs. deep fake detection.

Annexure VI: Sample Screenshots and Outputs

- AI-based detection tool screenshots (if applicable).
- Visual representation of deep fake detection.

Annexure VII: Code and Implementation Files (if applicable)

- GitHub repository link or software repository information.
- Sample code snippets and execution instructions.

Annexure VIII: Ethical Considerations and Challenges

- Discussion about privacy, ethical issues, and future effects.

PLAGARISM REPORT

Mohammadi Akheela Khanum Final Report(1)(5)			
ORIGINALITY REPORT			
13%	11%	9%	8%
SIMILARITY INDEX	INTERNET SOURCES	PUBLICATIONS	STUDENT PAPERS
PRIMARY SOURCES			
1	ges-coengg.org Internet Source	3%	
2	Submitted to Presidency University Student Paper	3%	
3	www.sih.gov.in Internet Source	2%	
4	R. N. V. Jagan Mohan, B. H. V. S. Rama Krishnam Raju, V. Chandra Sekhar, T. V. K. P. Prasad. "Algorithms in Advanced Artificial Intelligence - Proceedings of International Conference on Algorithms in Advanced Artificial Intelligence (ICAAAI-2024)", CRC Press, 2025 Publication	1%	
5	ijritcc.org Internet Source	<1%	
6	ijsrcseit.com Internet Source	<1%	
7	Saima Waseem, Syed Abdul Rahman Syed Abu Bakar, Bilal Ashfaq Ahmed, Zaid Omar et al. "DeepFake on Face and Expression Swap: A Review", IEEE Access, 2023 Publication	<1%	
8	Vishal Kumar Sharma, Rakesh Garg, Quentin Caudron. "A systematic literature review on deepfake detection techniques", Multimedia Tools and Applications, 2024 Publication	<1%	

RESEARCH PAPER PUBLICATION

We are pleased to submit our research paper, [Development of AI/ML based solution for detection of face-swap based deep fake videos], for consideration for [ICCAMS2025]. Although we have not yet received the acceptance certificates from IEEE, we confirm that our paper has been accepted and we are awaiting the formal documentation.

We would be happy to provide the certificates as soon as we receive them. Please let us know if there's any further information required from our end.

Thank you for your understanding and consideration.



2nd INTERNATIONAL CONFERENCE ON NEW FRONTIERS IN COMMUNICATION, AUTOMATION, MANAGEMENT AND SECURITY 2025 : Submission (860) has been created.

1 message

Microsoft CMT <noreply@msr-cmt.org>
Reply to: Microsoft CMT - Do Not Reply <noreply@msr-cmt.org>
To: Naganikithapeddisetty@gmail.com

Mon, 14 Apr 2025 at 19:02

Hello,

The following submission has been created.

Track Name: ICCAMS2025

Paper ID: 860

Paper Title: Development of AI/ML based solution for detection of face-swap based deep fake videos

Abstract:
within recent months, free tools based on deep learning have made it easier to produce genuine face exchanges in videos that bear minimal traces of manipulation, in so-called "DeepFake" (DF) videos. Digital video manipulations have been shown for decades through the proper utilization of visual effects, but recent breakthroughs in deep learning have caused a sharp rise in the authenticity of synthetic content and the ease with which they can be developed. Such so-called AI-synthesized media (commonly known as DF). Developing the DF through artificially intelligent tools are easy task. But, when it comes to the detection of the DF, it is a big challenge. Because the training of the algorithm to identify the DF is not easy. We have moved one step ahead in identifying the DF with Convolutional Neural Network and Recurrent Neural Network. The system employs a convolutional Neural network (CNN) to extract frame-level features. These features are utilized to train a recurrent neural network (RNN) which learns to classify whether a video has been manipulated or not and is capable of identifying the temporal inconsistencies between frames brought by the DF creation tools. The expected outcome against a huge set of fabricated videos gathered from the normal data set. We demonstrate how our system can be competitive outcome in this task leads to utilizing a simple architecture.

Created on: Mon, 14 Apr 2025 13:31:59 GMT

Last Modified: Mon, 14 Apr 2025 13:31:59 GMT

Authors:

- Naganikithapeddisetty@gmail.com (Primary)
- akheela.khanum@presidencyuniversity.in
- mishrasarthak052002@gmail.com
- darshan110604@gmail.com

Primary Subject Area: • AI and Machine Learning • Business Intelligence • Technical Trends • Ambient Technology • Communication

Secondary Subject Areas: Not Entered

Submission Files:
[deep_fake_detection.pdf](#) (726 Kb, Mon, 14 Apr 2025 13:31:01 GMT)

Submission Questions Response: Not Entered

Thanks,
CMT team.

School of Computer Science Engineering & Information Science , Presidency University.