

TEBLE OF CONTENTS

Declaration	i
Acknowledgement	ii
Vision & Mission of Dept	iii
Index	iv
List of Figures	vi
List of Snapshots	vii
Abstract	viii

INDEX

CONTENT	PAGE NUMBER
1. INTRODUCTION	1-5
1.1. Introduction to Deep Learning	1-3
1.2. Introduction to Topic	4-5
2. LITERATURE SURVEY	6-8
3. SYSTEM ANALYSIS	9-12
3.1 Existing System:	9
3.2 Disadvantages of the existing system	9-10
3.3 Proposed System	10-12
3.4 Advantages of proposed system	12
4. SYSTEM REQUIREMENTS	13
4.1. Hardware Requirements	13
4.2. Software Requirements	13

5. SYSTEM DESIGN	14-22
5.1. System Architecture	14-17
5.2. Data Flow Diagram	17-19
5.3. Use Case Diagram	20
5.4. Sequence Diagram	21
5.5. Activity Diagram	22
6. IMPLEMENTATION	23-42
7. RESULT	43
8. CONCLUSION	44
REFERENCES	45-46

LISTS OF FIGURES

SL.NO	FIGURES	PAGE NUMBER
1	DEEPFAKE VIDEO OF BARACK OBAMA.	5
2	DEEPFAKE VIDEO OF NICHOLAS CAGE	5
3	PRE-PROCESSING FLOW CHART	11
4	OVERALL PROCESS FLOWCHART	12
5	SYSTEM ARCHITECTURE	14
6	LEVEL 0 DIAGRAM	17
7	LEVEL-1 DIAGRAM	18
8	LEVEL-2 DIAGRAM	19
9	USE CASE DIAGRAM	20
10	SEQUENCE DIAGRAM	21
11	ACTIVITY DIAGRAM	22
12	DEEPFAKE MODEL RESULTS	43

LIST OF SNAPSHOTS

SL.NO	FIGURES	PAGE NUMBER
1	HOME PAGE	38
2	UPLOADING VIDEO	38
3	UPLOADING REAL VIDEO	39
4	FACE CROPPING IN REAL VIDEO	39
5	REAL VIDEO OUTPUT	40
6	UPLOADING FAKE VIDEO	40
7	FACE CROPPING IN FAKE VIDEO	41
8	FAKE VIDEO OUTPUT	41
9	PRESSING UPLOAD BUTTON WITHOUT SELECTING VIDEO	42

ABSTRACT

In recent months a machine learning based free software tool has made it easy to create believable face swaps in videos that leaves few traces of manipulation, in what are known as “deepfake” videos. Scenarios where these realistic fake videos are used to create political distress, blackmail someone or fake terrorism events are easily envisioned. Recent advancements in deep learning generative models have raised concerns as they can create highly convincing counterfeit images and videos. This poses a threat to people’s integrity and can lead to social instability. To address this issue, there is a pressing need to develop new computational models that can efficiently detect forged content and alert users to potential image and video manipulations. Our system uses a convolutional neural network (CNN) to extract frame-level features. These features are then used to train a recurrent neural network (RNN) that learns to classify if a video has been subject to manipulation or not. We evaluate our method against a large set of deepfake videos collected from multiple video websites. We show how our system can achieve competitive results in this task while using a simple architecture

Keywords – Machine learning, Deepfake, Internet, Cloud

CHAPTER 1

INTRODUCTION

1.1 Introduction to Deep Learning

Deep learning, a revolutionary subset of machine learning, has emerged as a transformative force in the field of artificial intelligence (AI), exhibiting unparalleled capabilities in solving complex problems and mimicking human cognitive functions. At its core, deep learning leverages artificial neural networks, inspired by the intricate structure of the human brain, to process vast amounts of data and uncover intricate patterns. This introduction aims to delve into the fundamental concepts, historical evolution, and diverse applications of deep learning, providing readers with a comprehensive understanding of this dynamic and rapidly evolving discipline.

The roots of deep learning can be traced back to the mid-20th century when researchers first began exploring the idea of creating computational models inspired by the human brain. However, it was not until the 21st century that deep learning witnessed a renaissance, fueled by the convergence of three crucial factors: 1. the accumulation of massive datasets, 2. the availability of powerful computational resources, and 3. groundbreaking advancements in neural network architectures. This confluence of factors paved the way for the resurgence of deep learning, marking the beginning of its ascent to prominence.

A cornerstone of deep learning is the neural network, a computational model composed of interconnected nodes, or artificial neurons, organized into layers. The architecture of these networks is designed to resemble the intricate connections within the human brain, with each layer extracting progressively more abstract features from the input data. The depth of these networks, characterized by multiple hidden layers, distinguishes deep learning from traditional machine learning approaches, allowing it to automatically learn hierarchical representation of data.

The training process of deep neural networks involves exposing the model to vast amounts of labeled data, allowing it to adjust its internal parameters through a process known as backpropagation. This iterative learning process enables the network to fine-tune its parameters and optimize its ability to make accurate predictions or classifications.

The efficiency of deep learning in learning complex representations directly from raw data has empowered it to excel in various domains, ranging from computer vision and natural language processing to speech recognition and healthcare.

One of the remarkable achievements of deep learning is its triumph in the realm of computer vision. Convolutional Neural Networks (CNNs), a specialized class of deep neural networks, have demonstrated unparalleled prowess in image recognition tasks, surpassing human performance in certain benchmarks. This breakthrough has paved the way for advancements in facial recognition, object detection, and autonomous vehicles, redefining the possibilities in fields that heavily rely on visual data.

Natural Language Processing (NLP) is another domain where deep learning has made significant strides, enabling machines to comprehend, generate, and interact with human language. Recurrent Neural Networks (RNNs) and Transformer architectures have played a pivotal role in enhancing language understanding, enabling applications such as machine translation, sentiment analysis, and chatbots to reach new levels of sophistication. The ability of deep learning models to grasp contextual nuances and semantic intricacies has fueled a paradigm shift in the way we interact with technology.

The impact of deep learning is not confined to specific domains; its influence extends to healthcare, finance, robotics, and beyond. In healthcare, deep learning models have demonstrated exceptional performance in medical image analysis, disease diagnosis, and drug discovery. Despite its remarkable achievements, deep learning is not without challenges. The insatiable appetite for data, the computational demands of training large models, and the interpretability of complex neural networks pose ongoing research questions.

Financial institutions leverage deep learning for fraud detection, risk assessment, and algorithmic trading, harnessing its capacity to discern intricate patterns within vast datasets. In robotics, deep learning facilitates advancements in autonomous navigation, object manipulation, and human-robot interaction, ushering in a new era of intelligent machines. The insatiable appetite for data, the computational demands of training large models, and the interpretability of complex neural networks pose ongoing research questions.

Key Features of Deep Learning:

These key features collectively contribute to the power and versatility of deep learning, allowing it to excel in a wide range of applications and continue its transformative impact on various industries.

- **Recurrent Neural Networks:** Tailored for sequential data, maintaining memory of past inputs. Suited for tasks with temporal dependencies, such as natural language processing.
- **Neural Network Architecture:** Deep learning relies on neural networks with multiple layers (deep architectures). The depth allows these networks to automatically learn hierarchical representations of data, capturing intricate patterns and features.
- **Feature Hierarchy and Abstraction:** Deep learning models can automatically extract hierarchical features from raw data, starting from low-level features and progressing to more abstract and complex representations.
- **End-to-End Learning:** Deep learning systems are capable of learning directly from raw data, eliminating the need for manual feature engineering. This end-to-end learning approach simplifies the model development process.
- **Scalability:** Deep learning architectures can scale to handle large amounts of data and compute resources. This scalability is crucial for training complex models on massive datasets, contributing to their impressive performance.
- **Adaptability to Diverse Data Types:** Deep learning is versatile and applicable to various types of data, including images, text, speech, and structured data.
- **Neural Network Architecture:** Deep learning relies on neural networks with multiple layers (deep architectures). The depth allows these networks to automatically learn hierarchical representations of data, capturing intricate patterns and features.
- **Backpropagation:** Training method involving the iterative adjustment of weights based on prediction errors.
- **Activation Functions:** Non-linear transformations applied to node outputs. Enhances the model's capacity to capture and model complex patterns.
- **Adaptability to Diverse Data Types:** Deep learning is versatile and applicable to various types of data, including images, text, speech, and structured data.

1.2 Introduction to Topic

Deepfake detection using deep learning involves leveraging advanced neural network architectures to identify manipulated or synthetic content in multimedia, particularly in images and videos. Deepfakes are created using deep neural networks, making it challenging to distinguish them from authentic content. The detection process typically involves training deep learning models on datasets containing both real and manipulated media.

The key steps in deepfake detection using deep learning begin with dataset preparation. This entails collecting a diverse dataset that includes both authentic and deepfake media, ensuring coverage across various scenarios, lighting conditions, and facial expressions to enhance model generalization. The next crucial step is choosing a deep learning architecture suitable for the task. For image-based deepfakes, Convolutional Neural Networks (CNNs) are commonly employed, while Recurrent Neural Networks (RNNs) may be chosen for video-based detection. Advanced architectures such as Variational Autoencoders (VAEs) or Generative Adversarial Networks (GANs) designed for detecting anomalies or inconsistencies can also be considered.

The training phase involves optimizing the model for accuracy in distinguishing between real and deepfake content. Techniques like transfer learning or fine-tuning on pre-trained models are utilized to enhance performance, especially when dealing with limited data. Feature extraction from the media is another critical aspect, focusing on identifying patterns indicative of deepfake generation. Both spatial and temporal features are explored for video-based detection.

To improve the model's ability to generalize, data augmentation techniques are applied to artificially increase the size of the training dataset. After training, validation on a separate dataset is conducted to ensure the model's performance on unseen data. Rigorous testing follows to assess the model's robustness against different types of deepfakes. Post-processing methods are then implemented to refine the model's predictions and reduce false positives or negatives. Continuous improvement is emphasized through regular updates to the model with new data, enabling it to adapt to evolving deepfake techniques. Researchers and practitioners also explore ensemble methods or combine multiple models to achieve enhanced detection accuracy. By employing deep learning techniques in deepfake detection, they contribute to ongoing efforts to stay ahead of the rapidly

evolving landscape of synthetic media, addressing and mitigating the potential misuse of deepfake technology.



Fig.1.2.1: Deepfake Video of Barack Obama.

A deepfake video of Barack Obama is a fake video created using AI to make it look like he's saying or doing something he didn't actually do. It's made by manipulating existing footage or generating new content to mimic his appearance and voice, raising concerns about spreading misinformation and manipulating public opinion.



Fig.1.2.2: Deepfake Video of Nicholas Cage

A deepfake video of Nicholas Cage is a manipulated video created using AI to make it seem like he's doing or saying something he didn't actually do. It's made by altering existing footage or generating new content to mimic his appearance and voice, raising concerns about spreading misinformation and manipulating public perception.

CHAPTER 2

LITERATURE SURVEY

- 1) **Title:** Deepfake Video Detection Using Recurrent Neural Network.

Author: David Guera, Edward J. Delp.

Journal: 2021 IEEE 8th International Conference on Industrial Engineering.

Year: 2021

Training and using RNNs for deepfake detection can be computationally expensive, especially when dealing with high-resolution videos. This can limit real-time or near-real-time detection capabilities. RNNs might struggle to generalize well to new and unseen deepfake techniques or variations. They may be more suitable for specific types of deepfakes but not versatile across different scenarios. RNNs may generate false positives, flagging legitimate videos as deepfakes, which can lead to user frustration and mistrust of the detection system. The methodology section meticulously outlines the data collection process, preprocessing steps, and the architecture of the RNN model tailored for this task.

- 2) **Title:** Deepfake Crection And Detection Using Deep Learning.

Author: Handy A.Khalil, Shady A.Maged.

Journal: 2021 International Mobile, Intelligent and Ubiquitous Conference .

Year: 2021

The Deepfake creation and detection represent two facets of the evolving landscape of artificial intelligence, specifically within the realm of deep learning. In the context of deepfake creation, a methodological paper delves into the intricacies of leveraging advanced neural network architectures to seamlessly generate highly realistic synthetic media, often replacing a person's likeness in videos with another's. These deepfake algorithms, driven by generative models like Generative Adversarial Networks (GANs) or Variational Autoencoders (VAEs), learn to mimic facial expressions, gestures, and voice patterns with remarkable accuracy, posing a significant challenge to the authenticity of digital content. Conversely, deepfake detection relies on the application of sophisticated deep learning techniques to discern genuine from manipulated media.

3) Title: DeepVision:Deepfakes detection using human eye blinking pattern**Author:** TackHyun Jung, SangWon Kim, KeeCheon Kim.**Journal:** 2021 IEEE Access Volume 8.**Year:** 2020.

The paper introduces DeepVision, a new method for detecting GAN-generated Deepfakes by analyzing changes in blinking patterns. Unlike traditional pixel-based approaches, DeepVision considers factors like physical conditions and cognitive activities. Utilizing insights from medicine, biology, brain engineering, and machine learning, DeepVision's heuristic method achieves an 87.5% accuracy rate in identifying Deepfakes across various video types. Human eye blinking occurs unconsciously and spontaneously, the algorithm is unpredictable, bolstering its security. An attacker will find it difficult to use GAN's model to disable DeepVision. The methodology section elucidates the intricate process of extracting and analyzing blinking patterns from video data, likely employing a combination of computer vision techniques and machine learning algorithms.

4) Title: Deepfake Detection through Deep Learning.**Author:** Deng Pan, Lixian Sun, Rui Wang, Xingjian Zhang, Richard O. Sinnott.**Journal:** 2020 IEEE/ACM International Conference on Big Data Computing,**Year:** 2020.

The paper explores the use of neural networks, specifically Xception and MobileNet, to automatically detect deepfake videos. These technologies aim to address the societal challenges posed by the controversial use of deepfake technology, including issues like election biasing. The results indicate a high accuracy range of 91-98% across different datasets and deepfake technologies. Additionally, the paper introduces a voting mechanism that enhances detection by aggregating results from multiple methods, surpassing reliance on a single approach. The methodology section meticulously delineates the data acquisition process, preprocessing techniques, and the architecture of deep learning models employed for deepfake detection. neural network selection, training methodologies, and performance evaluation metrics are discussed to provide clarity on the experimental framework.

5) Title: DeepFake Videos Detection Based on Texture Features**Author:** Bozhi Xu, Jiarui Liu, Jifan Liang, Wei Lu, Yue Zhang.**Journal:** 2019 International Conference on Computer and IS (ICCIS).**Year:** 2019

Texture-based approaches can sometimes produce false positives, flagging authentic videos as DeepFakes due to variations in lighting, camera quality, or compression artifacts, which may be similar to those found in manipulated videos. DeepFake creators can exploit the weaknesses of these detectors by using adversarial attacks specifically designed to bypass such detection methods. These attacks can make it difficult for texture-based detectors to identify DeepFakes accurately. Texture-based detectors are typically designed to recognize specific features or patterns, which means they may struggle to detect novel or unforeseen manipulation methods.

6) Title: Deep Fake Detection using Neural Networks**Author:** Anuj Badale, Lionel Castelino, Chaitanya Darekar, Joanne Gomes.**Journal:** Proceedings of APSIPA Annual Summit and Conference 2019.**Year:** 2019

DeepFake creators continuously adapt and refine their methods, making it challenging for neural network-based detection models to keep up. Running neural network-based DeepFake detection models can be computationally expensive, limiting their real-time applicability. Neural network-based detectors may produce false positives (flagging real videos as fake) or false negatives (failing to detect certain DeepFakes), impacting their reliability. Some DeepFake detection methods may infringe on privacy or involve intrusive analysis of individuals' faces, raising ethical concerns. Methodologically, the study meticulously details the data collection, preprocessing steps, and the architecture of the neural networks deployed for deep fake detection. It likely delves into the rationale behind selecting specific neural network architectures, such as convolutional or recurrent neural networks, and outlines strategies for enhancing model performance through techniques like data augmentation.

CHAPTER 3

SYSTEM ANALYSIS

3.1 Existing System:

Various deepfake detection techniques employ neural networks and deep learning to identify patterns indicative of manipulated content. One approach involves utilizing Recurrent Neural Networks (RNNs) to analyze temporal dependencies in video frames, potentially capturing facial movements or inconsistencies introduced during deepfake creation. Another method focuses on texture features within video frames, detecting irregularities such as color variations or artifacts that may arise during the generation process. A broad category encompasses the use of neural networks for deepfake detection, employing architectures like Convolutional Neural Networks (CNNs) to automatically learn and distinguish features associated with authentic and manipulated content.

Innovative approaches, such as DeepVision, introduce unique perspectives by considering human behaviors. DeepVision proposes the use of human eye blinking patterns for detection, leveraging the involuntary and influenced nature of blinking by factors like emotion and cognitive activities. Additionally, the deepfake detection landscape emphasizes adaptability and robustness for handling increasingly realistic content. Some models aim to address this challenge through more dynamic approaches, using semi-supervised learning techniques to keep pace with the evolving nature of deepfake generators and the rapid spread of content across social networks. The integration of diverse methodologies, from texture analysis to behavioral cues, showcases the multifaceted efforts within the research community to combat the growing sophistication of deepfake technology.

3.2 Disadvantages of the existing system

- **Sensitive to lighting changes:** Performance is affected by variations in lighting conditions, making the system less robust in different environments.
- **May struggle with long-term dependencies:** Difficulty in capturing and understanding relationships in data that span over extended periods, impacting the system's effectiveness.
- **Computationally expensive:** The existing system requires significant computational

resources, leading to high processing costs.

- **Vulnerable to adversarial attacks:** Susceptible to manipulation and intentional deception, compromising the system's reliability and security.
- **Requires extensive labeled data:** Dependence on large amounts of annotated data for training, posing challenges in data collection and annotation.
- **May lack generalization:** The system might not effectively apply learned patterns to new, unseen situations or contexts.
- **Interpretability challenges:** Difficulty in understanding and explaining the decisions made by the system, hindering trust and transparency.
- **Need for constant updates due to evolving deepfake techniques:** The system requires frequent updates to adapt to emerging deepfake methods, demanding continuous maintenance and development.
- **Risk of misidentifying real content or missing subtle manipulations:** Possibility of inaccuracies in distinguishing authentic content from manipulated ones, potentially leading to false positives or negatives.

3.3 Proposed System

To address the challenges associated with existing deepfake detection systems, a proposed system could leverage a hybrid approach that combines the strengths of different techniques. Integrating both temporal and spatial features, the system would employ Recurrent Neural Networks (RNNs) for their ability to capture temporal dependencies in video sequences, while simultaneously incorporating texture analysis to identify variations indicative of manipulation.

Moreover, to enhance robustness against adversarial attacks and improve interpretability, the proposed system could integrate explainable AI techniques. By incorporating methods that provide insights into the decision-making process of the model, such as attention mechanisms or feature visualization, the system aims to make its inner workings more transparent and in its detection capabilities.

Additionally, a dynamic learning mechanism could be implemented. This mechanism would involve

continuous learning and updating of the detection model using real-time data, allowing the system to stay ahead of emerging deepfake techniques. This adaptive learning approach could be reinforced by a collaborative network that shares insights and updates among different instances of the detection system.

Furthermore, the proposed system could incorporate a feedback loop from end-users, enabling it to learn from misclassifications and continuously improve its accuracy over time. This iterative learning process, coupled with a regularly updated database of diverse deepfake examples, ensures that the system remains vigilant and effective in countering the evolving landscape of deepfake technology the proposed system aims to provide a more reliable and resilient solution for detecting deepfake content.

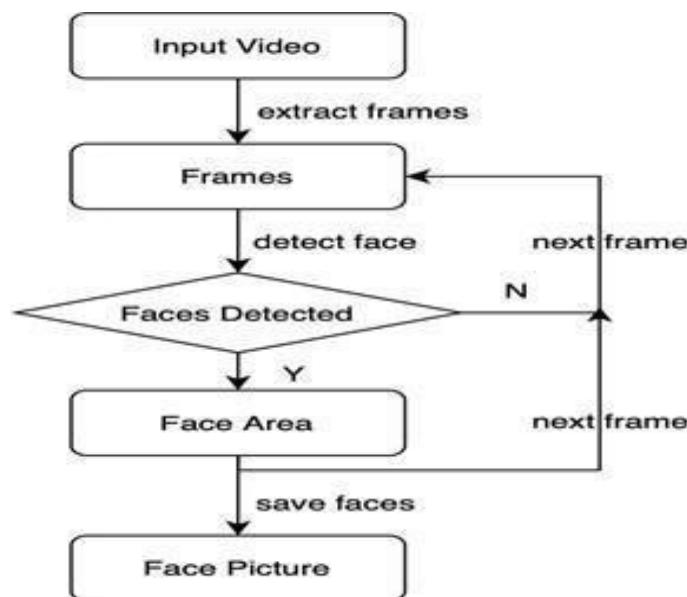


Fig. 3.3.1: Pre-processing Flow Chart

The flowchart illustrates the steps involved in the system. Here's a description of the process:

- The first step is to “Input Video,” which refers to a video containing a person’s face.
- The second step is to “Extract Frames,” which breaks the video down into individual images, each containing a single frame of the video.
- The second step is to “Extract Frames,” which breaks the video down into individual images, each containing a single frame of the video.

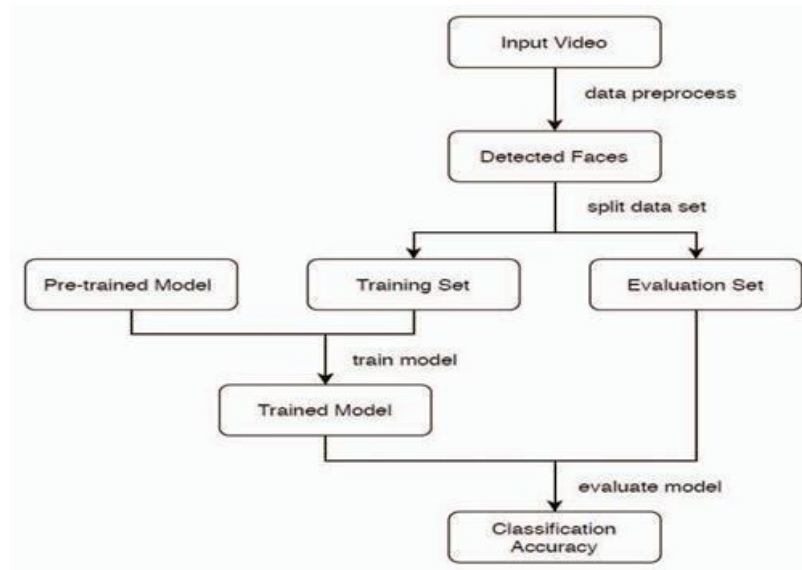


Fig.3.3.2: Overall Process Flowchart

Flowchart illustrates the process of training a face detection model and then using the model to detect faces in a video.

3.4 Advantages of proposed system:

- Integrates temporal and spatial features for comprehensive detection:** The proposed system combines both temporal and spatial characteristics, offering a holistic approach to detect and address deepfake content effectively.
- Enhances transparency and trust in decision-making:** The system's design promotes transparency, fostering trust by providing clear explanations for its decision-making processes.
- Defends against manipulative inputs for increased resilience:** The system is resilient against attempts to manipulate or deceive it, providing robust defense mechanisms.
- Adapts to evolving deepfake technologies for continual effectiveness:** The proposed system is designed to evolve and stay effective over time by adapting to new and emerging techniques in the field of deepfake technology.
- Shares insights among instances for a collective defense:** The system facilitates information sharing among instances, enabling a collaborative defense strategy against deepfake threats.

CHAPTER 4

SYSTEM REQUIREMENTS

4.1 Hardware Requirements

- Processor core : i5 core or AMD Ryzen 5
- Processor speed : 2.5 GHz
- Hard disk : 1TB
- RAM : 16 GB
- GPU (Graphics Processing Unit) : NVIDIA GeForce RTX 2070 or equivalent.
- Network Connectivity : Minimum 10 Mbps
- Dedicated Deep Learning Hardware : TPUs (Tensor Processing Units)

4.2 Software Requirements

- Operating system : Windows 7 or above
- Coding Language : Python
- Deep Learning Framework : TensorFlow, PyTorch, or Keras.
- Data Visualization Tools : Matplotlib, Seaborn, or Plotly
- Cloud Services : AWS, Google Cloud, or Azure.

CHAPTER 5

SYSTEM DESIGN

5.1 System Architecture:

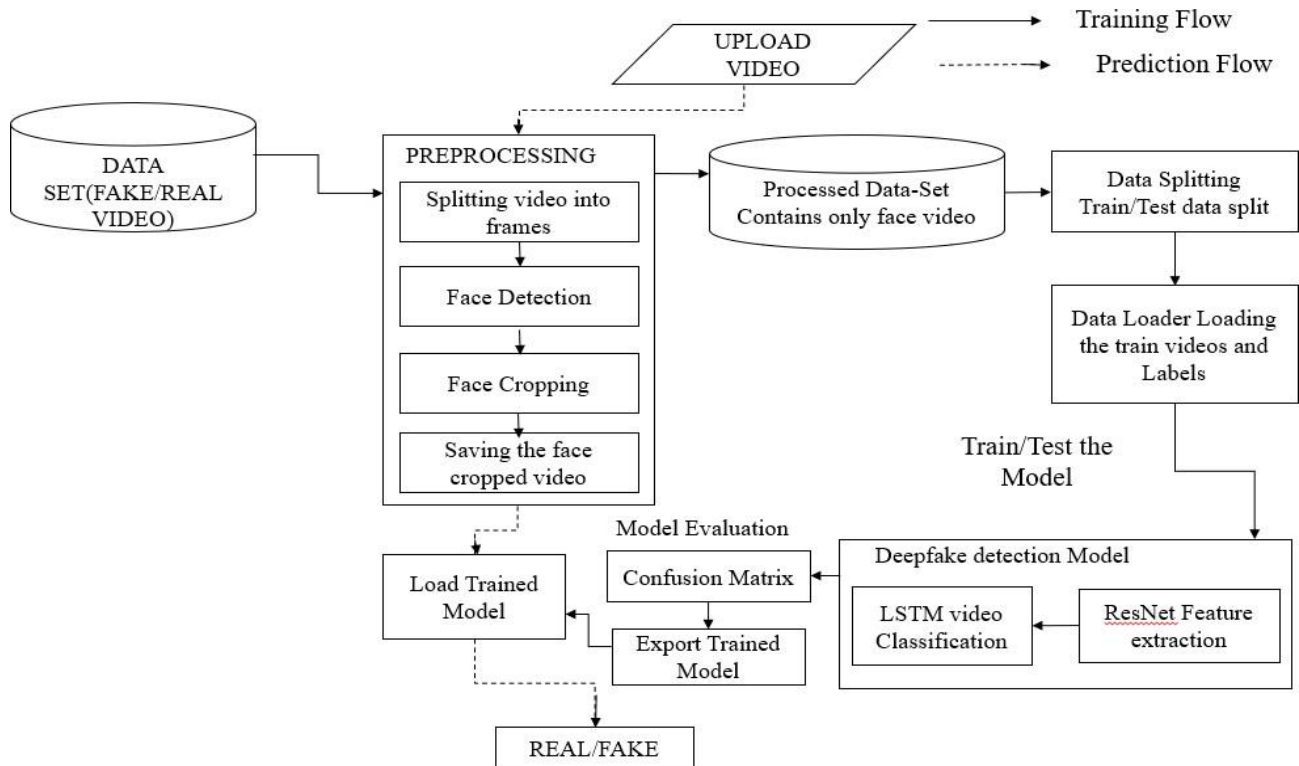


Fig. 5.1.1: System Architecture

In this system, we have trained our PyTorch deepfake detection model on equal number of real and fake videos in order to avoid the bias in the model. The system architecture of the model is showed in the figure. In the development phase, we have

Module 1: Data-set Gathering

For making the model efficient for real time prediction. We have gathered the data from different available data-sets like FaceForensic++(FF), Deepfake detection challenge (DFDC), and Celeb-DF. Further we have mixed the dataset the collected datasets and created our own new dataset, to accurate

and real time detection on different kind of videos. To avoid the training bias of the model we have considered 50% Real and 50% fake videos. Deep fake detection challenge (DFDC) dataset consist of certain audio alerted video, as audio deepfake are out of scope for this paper. We preprocessed the DFDC dataset and removed the audio altered videos from the dataset by running a python script. After preprocessing of the DFDC dataset, we have taken 1500 Real and 1500 Fake videos from the DFDC dataset. 1000 Real and 1000 Fake videos from the FaceForensic++ (FF) dataset and 500 Real and 500 Fake videos from the CelebDF dataset. Which makes our total dataset consisting 3000 Real, 3000 fake videos and 6000 videos in total.

Module 2: Pre-processing

In this step, the videos are preprocessed and all the unrequired and noise is removed from videos. Only the required portion of the video i.e face is detected and cropped. The first steps in the preprocessing of the video is to split the video into frames. After splitting the video into frames the face is detected in each of the frame and the frame is cropped along the face. Later the cropped frame is again converted to a new video by combining each frame of the video. The process is followed for each video which leads to creation of processed dataset containing face only videos. The frame that does not contain the face is ignored while preprocessing. To maintain the uniformity of number of frames, we have selected a threshold value based on the mean of total frames count of each video. Another reason for selecting a threshold value is limited computation power. As a video of 10 second at 30 frames per second(fps) will have total 300 frames and it is computationally very difficult to process the 300 frames at a single time in the experimental environment. So, based on our Graphic Processing Unit (GPU) computational power in experimental environment we have selected 150 frames as the threshold value. While saving the frames to the new dataset we have only saved the first 150 frames of the video to the new video. To demonstrate the proper use of Long Short-Term Memory (LSTM) we have considered the frames in the sequential manner i.e. first 150 frames and not randomly. The newly created video is saved at frame rate of 30fps and resolution of 112 x 112.

Module 3: Data-set split

The dataset is split into train and test dataset with a ratio of 70% train videos (4,200) and 30% (1,800) test videos. The train and test split is a balanced split i.e 50% of the real and 50% of fake videos.

Module 4: Model Architecture

Our model is a combination of CNN and RNN. We have used the Pre-trained ResNext CNN model to extract the features at frame level and based on the extracted features a LSTM network is trained to classify the video as deepfake or pristine. Using the Data Loader on training split of videos the labels of the videos are loaded and fitted into the model for training.

ResNext: Instead of writing the code from scratch, we used the pre-trained model of ResNext for feature extraction. ResNext is Residual CNN network optimized for high performance on deeper neural networks. For the experimental purpose we have used resnext50_32x4d model. We have used a ResNext of 50 layers and 32 x 4 dimensions. Following, we will be fine-tuning the network by adding extra required layers and selecting a proper learning rate to properly converge the gradient descent of the model. The 2048-dimensional feature vectors after the last pooling layers of ResNext is used as the sequential LSTM input.

LSTM for Sequence Processing: 2048-dimensional feature vectors is fitted as the input to the LSTM. We are using 1 LSTM layer with 2048 latent dimensions and 2048 hidden layers along with 0.4 chance of dropout, which is capable to do achieve our objective. LSTM is used to process the frames in a sequential manner so that the temporal analysis of the video can be made, by comparing the frame at 't' second with the frame of 't-n' seconds. Where n can be any number of frames before t. The model also consists of Leaky Relu activation function. A linear layer of 2048 input features and 2 output features are used to make the model capable of learning the average rate of correlation between eh input and output. An adaptive average pooling layer with the output parameter 1 is used in the model. Which gives the target output size of the image of the form H x W. For sequential processing of the frames a Sequential Layer is used. The batch size of 4 is used to perform the batch training. A SoftMax layer is used to get the confidence of the model during predication.

Module 5: Hyper-parameter tuning

It is the process of choosing the perfect hyper-parameters for achieving the maximum accuracy. After reiterating many times on the model. The best hyper-parameters for our dataset are chosen. To enable the adaptive learning rate Adam optimizer with the model parameters is used. The learning rate is tuned to 1e-5 (0.00001) to achieve a better global minimum of gradient descent. The weight decay

used is $1e-3$. As this is a classification problem so to calculate the loss cross entropy approach is used. To use the available computation power properly the batch training is used. The batch size is taken of 4. Batch size of 4 is tested to be ideal size for training in our development environment. The User Interface for the application is developed using Django framework. Django is used to enable the scalability of the application in the future.

The first page of the User interface i.e index.html contains a tab to browse and upload the video. The uploaded video is then passed to the model and prediction is made by the model. The model returns the output whether the video is real or fake along with the confidence of the model. The output is rendered in the predict.html on the face of the playing video.

5.2 Data Flow Diagram:

Level 0

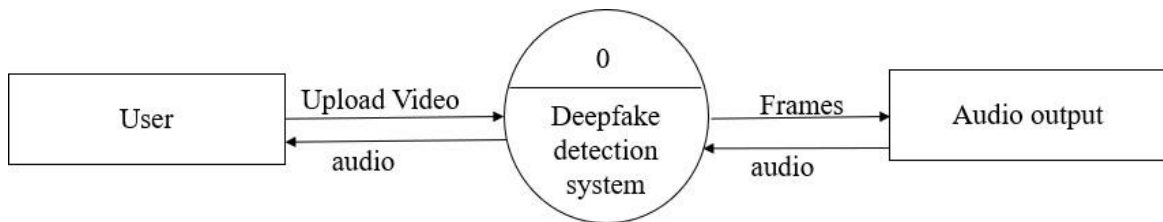


Fig.5.2.1: Level 0 Diagram

DFD level – 0 indicates the basic flow of data in the system. In this System Input is given equal importance as that for Output.

- Input: Here input to the system is uploading video.
- System: In system it shows all the details of the Video.
- Output: Output of this system is it shows the fake video or not.

Hence, the data flow diagram indicates the visualization of system with its input and output flow.

Level 1

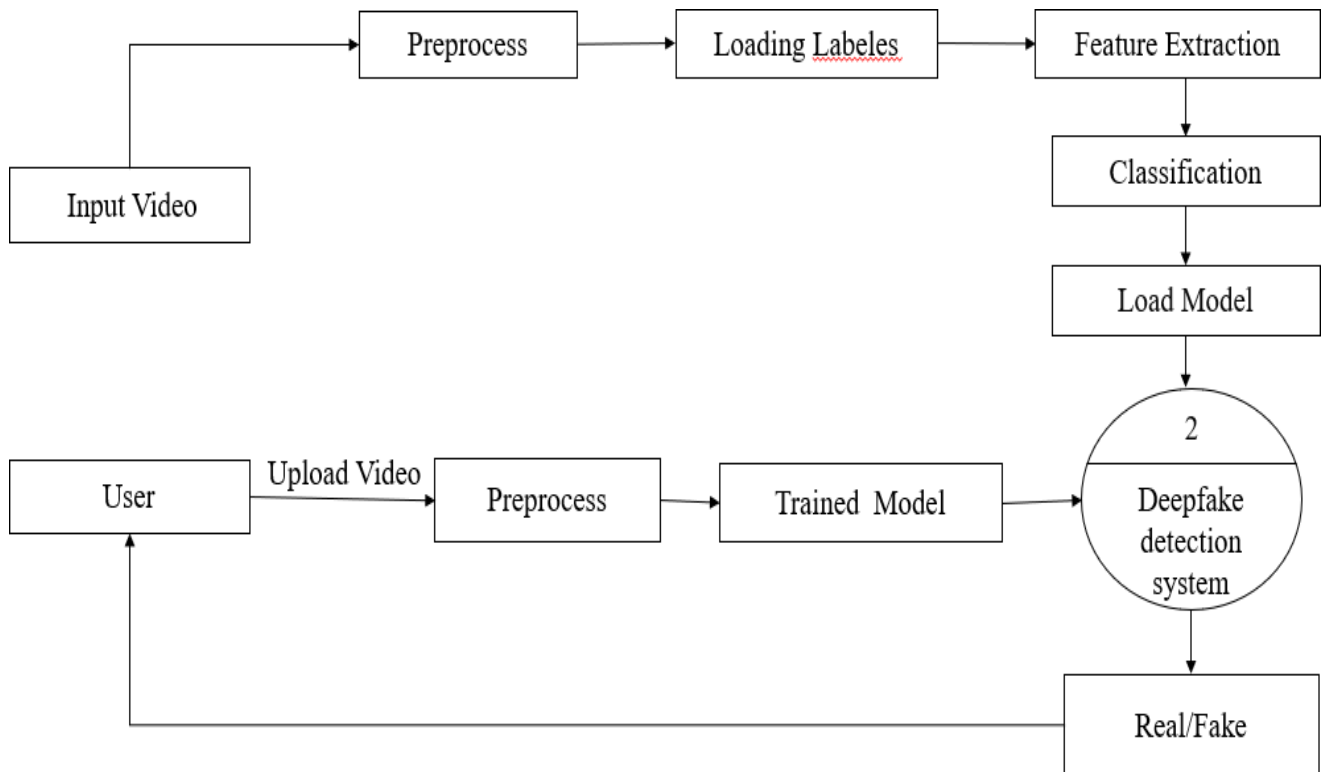


Fig.5.2.2: Level-1 Diagram

At the level 1 dataflow diagram for deepfake detection using deep learning, the focus is on outlining the high-level flow of data and processes within the system. Typically, it illustrates the main components and interactions involved in detecting deepfakes through neural network analysis. The diagram would encompass stages such as data input, preprocessing, feature extraction, classification, and output. At the input stage, raw image or video data is fed into the system. Preprocessing steps may involve tasks like face detection, alignment, and normalization to ensure uniformity and enhance the quality of input data. Following preprocessing, the data flows into the feature extraction phase, where relevant features indicative of deepfake manipulation are extracted from the input media. These features are then passed into the classification component, which employs trained deep learning models to categorize the input as either authentic or deepfake.

Level 2

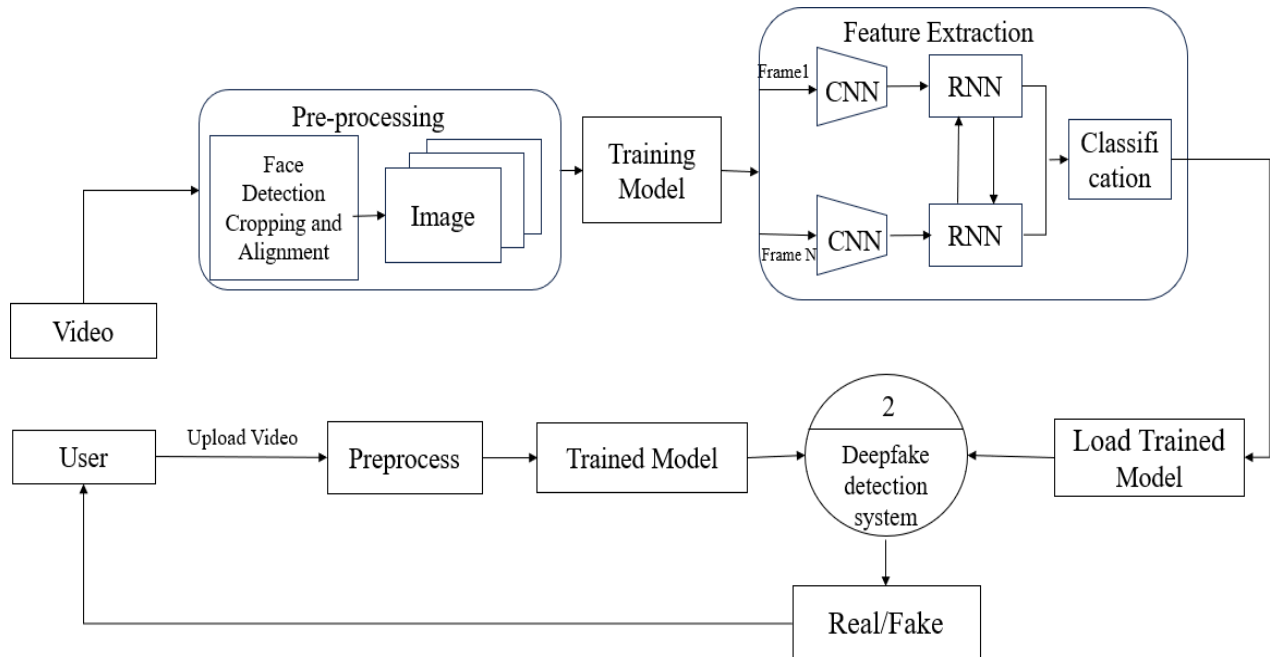


Fig.5.2.3: Level-2 Diagram

At the level 2 dataflow diagram for deepfake detection using deep learning, the focus shifts to a more detailed depiction of the sub-processes and interactions within each main component of the system. This diagram delves deeper into the functionalities and data transformations occurring at each stage of the deepfake detection pipeline. For instance, within the preprocessing component, it may outline specific steps such as face detection algorithms, image resizing, and noise reduction techniques. Similarly, the feature extraction stage may detail the operations performed by convolutional neural networks (CNNs) or other deep learning architectures to identify relevant patterns and attributes indicative of deepfake manipulation. Moreover, the classification component might illustrate the layers and computations involved in making predictions, including the activation functions, weights, and biases of the neural network.

5.3 Use Case Diagram

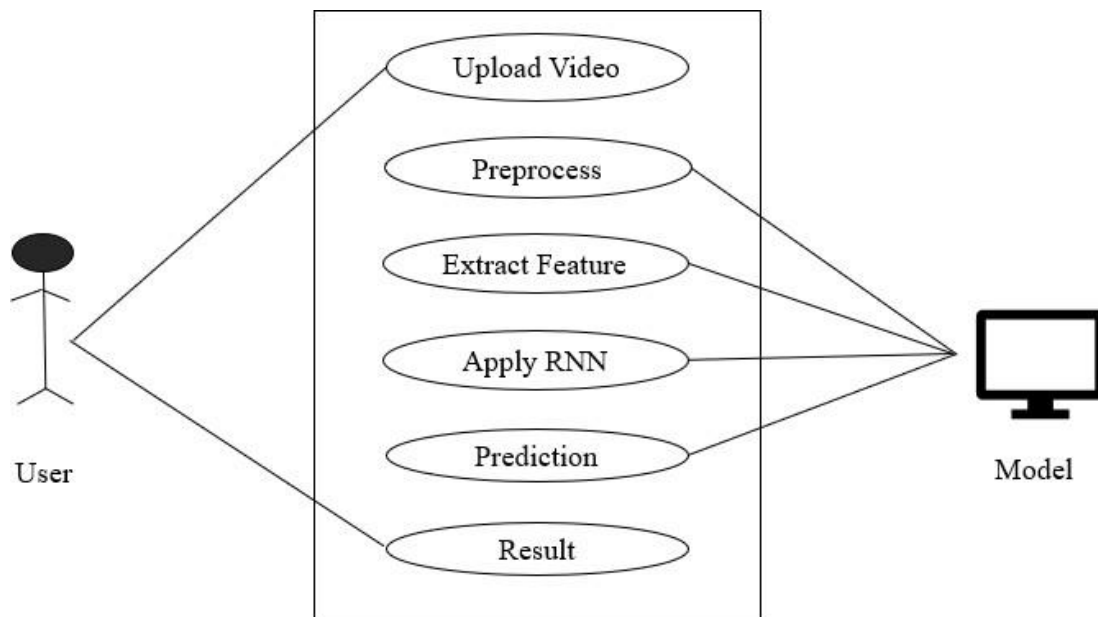


Fig.5.3.4: Use Case Diagram

A use case diagram for deepfake detection using deep learning illustrates the various interactions between actors and the system, outlining the primary functionalities and scenarios where the system is utilized. In this context, the actors typically include end-users, administrators, and possibly other systems or external entities. The diagram identifies the main use cases or functionalities offered by the deepfake detection system, such as uploading media for analysis, initiating a detection process, reviewing results, and managing user accounts or system settings. End-users, representing individuals or organizations concerned about the authenticity of media content, interact with the system primarily by submitting media for analysis and reviewing the detection results. Administrators, on the other hand, may have additional capabilities, such as managing user accounts, configuring system parameters, and monitoring system performance.

5.4 Sequence Diagram:

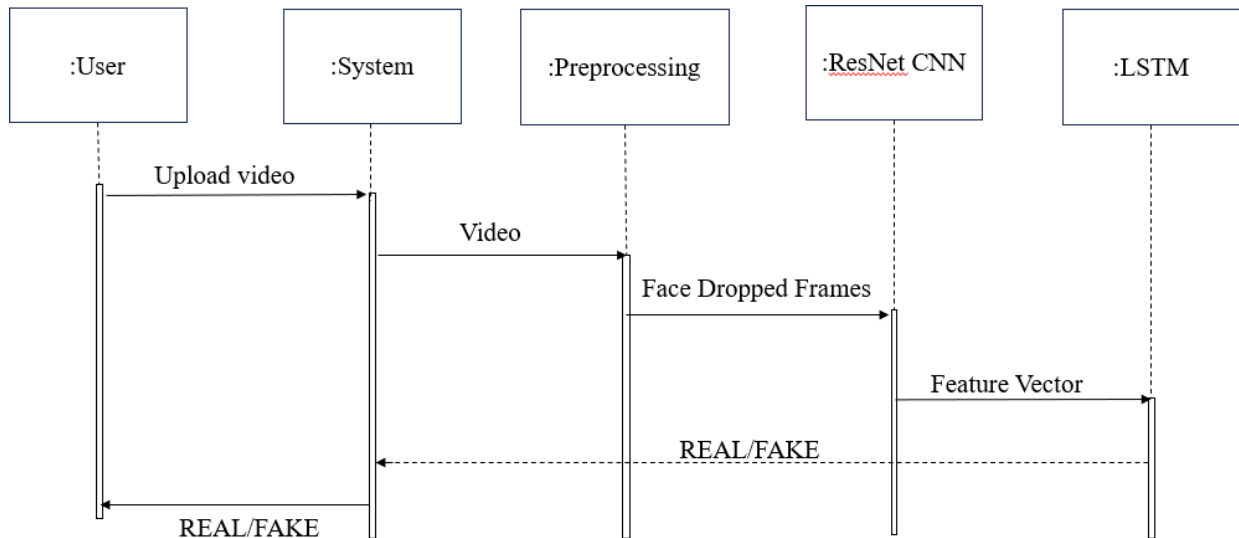


Fig.5.4.1: Sequence Diagram

A sequence diagram for deepfake detection using deep learning illustrates the chronological flow of interactions between various components and actors involved in the detection process. Typically, it outlines the sequence of messages exchanged between these entities to accomplish specific tasks within the system. In the context of deepfake detection, the sequence diagram might depict interactions between the user interface, preprocessing modules, deep learning models, and output visualization components. For example, it could start with the user uploading media for analysis through the user interface. Subsequently, the system triggers preprocessing steps such as face detection and alignment to prepare the input data. The preprocessed data is then passed to the deep learning model, where feature extraction and classification take place to determine the authenticity of the media. Throughout this process, messages are exchanged between the components, conveying data, instructions, and feedback. Finally, the detection results are visualized and presented back to the user through the interface.

5.5 Activity Diagram:

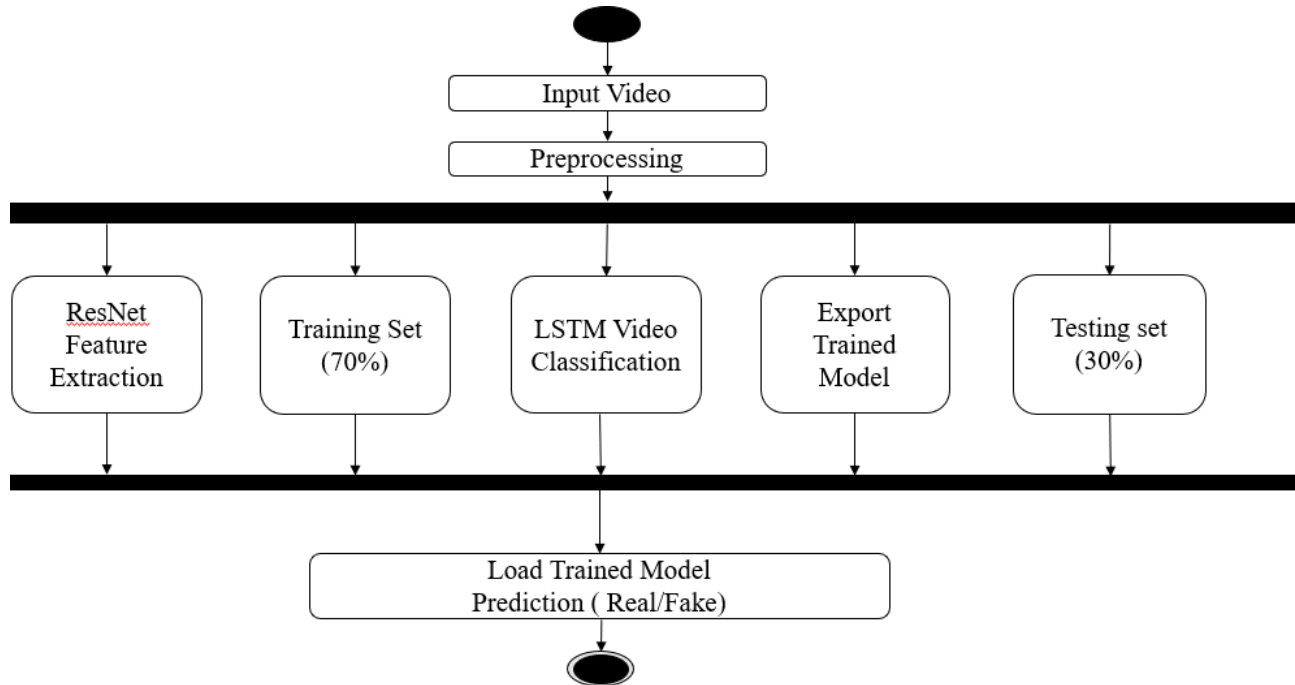


Fig.5.5.1: Activity Diagram

An activity diagram for deepfake detection using deep learning provides a visual representation of the workflow or activities involved in the detection process. It illustrates the sequence of actions and decisions undertaken by the system to analyze input media and determine its authenticity. Typically, the activity diagram outlines various stages such as data preprocessing, feature extraction, classification, and result visualization. For instance, it might start with the activity of receiving input media, followed by preprocessing steps like face detection and alignment. Subsequently, the diagram depicts the activity of feature extraction, where relevant attributes are identified from the input data using deep learning techniques. After feature extraction, the system engages in the classification activity, wherein the extracted features are analyzed to classify the media as authentic or manipulated. Depending on the classification outcome, the diagram branches into different paths, indicating actions such as presenting detection results to the user or performing additional analysis.

CHAPTER 6

IMPLEMENTATION

Tools and Technologies Used

1. Programming Languages:

Python3

JavaScript

2. Programming Frameworks:

PyTorch

Django

3. IDE:

Google Colab

Jupyter Notebook

Visual Studio Code

4. Versioning Control:

Git

5. Cloud Services:

Google Cloud Platform

6. Application and web servers:

Google Cloud Engine

7. Libraries:

Torch

Torchvision

Os

Numpy

cv2

matplotlib

face_recognition

json

pandas

8. Preprocessing Details:

- Using glob we imported all the videos in the directory in a python list.
- cv2.VideoCapture is used to read the videos and get the mean number of frames in each video.
- To maintain uniformity, based on mean a value 150 is selected as idea value for creating the new dataset.
- The video is split into frames and the frames are cropped on face location.
- The face cropped frames are again written to new video using VideoWriter.
- The new video is written at 30 frames per second and with the resolution of 112 x 112 pixels in the mp4 format.
- Instead of selecting the random videos, to make the proper use of LSTM for temporal sequence analysis the first 150 frames are written to the new video.

9. Model Details: The model consists of following layers

- **ResNext CNN** : The pre-trained model of Residual Convolution Neural Network is used. The model name is resnext50_32x4d(). This model consists of 50 layers and 32 x 4 dimensions
- **Sequential Layer** : Sequential is a container of Modules that can be stacked together and run at the same time. Sequential layer is used to store feature vector returned by the ResNext model in an ordered way. So that it can be passed to the LSTM sequentially.
- **LSTM Layer** : LSTM is used for sequence processing and spot the temporal change between the frames. 2048-dimensional feature vectors are fitted as the input to the LSTM. We are using 1 LSTM layer with 2048 latent dimensions and 2048 hidden layers along with 0.4 chance of dropout, which is capable to do achieve our objective. LSTM is used to process the frames in a sequential manner so that the temporal analysis of the video can be made, by comparing the frame at 't' second with the frame of 't-n' seconds.
- **ReLU**: A Rectified Linear Unit is activation function that has output 0 if the input is less than 0, and raw output otherwise. That is, if the input is greater than 0, the output is equal to the input. The operation of ReLU is closer to the way our biological neurons work. ReLU is non-linear and has the advantage of not having any backpropagation errors unlike the sigmoid function, also for larger Neural Networks, the speed of building models based off on ReLU is very fast.

- **Dropout Layer:** Dropout layer with the value of 0.4 is used to avoid overfitting in the model and it can help a model generalize by randomly setting the output for a given neuron to 0. In setting the output to 0, the cost function becomes more sensitive to neighbouring neurons changing the way the weights will be updated during the process of backpropagation. Adaptive Average Pooling Layer: It is used To reduce variance, reduce computation complexity and extract low level features from neighbourhood. 2 dimensional Adaptive Average Pooling Layer is used in the model.

10. Confusion Matrix: A confusion matrix is a summary of prediction results on a classification problem. The number of correct and incorrect predictions are summarized with count values and broken down by each class. This is the key to the confusion matrix. The confusion matrix shows the ways in which your classification model is confused when it makes predictions. It gives us insight not only into the errors being made by a classifier but more importantly the types of errors that are being made. Confusion matrix is used to evaluate our model and calculate the accuracy.

SOURCE CODE:

```
from django.shortcuts import render, redirect
import torch
import torchvision
from torchvision import transforms, models
from torch.utils.data import DataLoader
from torch.utils.data.dataset import Dataset
import os
import numpy as np
import cv2
import matplotlib.pyplot as plt
import face_recognition
from torch.autograd import Variable
import time
import sys
from torch import nn
import json
```

```
import glob
import copy
from torchvision import models
import shutil
from PIL import Image as pImage
import time
from django.conf import settings
from .forms import VideoUploadForm
index_template_name = 'index.html'
predict_template_name = 'predict.html'

im_size = 112
mean=[0.485, 0.456, 0.406]
std=[0.229, 0.224, 0.225]
sm = nn.Softmax()
inv_normalize=transforms.Normalize(mean=-1*np.divide(mean,std),std=np.divide([1,1,1],std))
train_transforms = transforms.Compose([
    transforms.ToPILImage(),
    transforms.Resize((im_size,im_size)),
    transforms.ToTensor(),
    transforms.Normalize(mean,std)])

class Model(nn.Module):

    def __init__(self, num_classes,latent_dim= 2048, lstm_layers=1 , hidden_dim = 2048,
bidirectional = False):
        super(Model, self).__init__()
        model = models.resnext50_32x4d(pretrained = True)
        self.model = nn.Sequential(*list(model.children())[:-2])
        self.lstm = nn.LSTM(latent_dim,hidden_dim, lstm_layers, bidirectional)
        self.relu = nn.LeakyReLU()
```

```

self.dp = nn.Dropout(0.4)
self.linear1 = nn.Linear(2048,num_classes)
self.avgpool = nn.AdaptiveAvgPool2d(1)

```

```

def forward(self, x):
    batch_size,seq_length, c, h, w = x.shape
    x = x.view(batch_size * seq_length, c, h, w)
    fmap = self.model(x)
    x = self.avgpool(fmap)
    x = x.view(batch_size,seq_length,2048)
    x_lstm,_ = self.lstm(x,None)
    return fmap,self.dp(self.linear1(x_lstm[:,-1,:]))

```

```

class validation_dataset(Dataset):
    def __init__(self,video_names,sequence_length=60,transform = None):
        self.video_names = video_names
        self.transform = transform
        self.count = sequence_length

```

```

def __len__(self):
    return len(self.video_names)

```

```

def __getitem__(self,idx):
    video_path = self.video_names[idx]
    frames = []
    a = int(100/self.count)
    first_frame = np.random.randint(0,a)
    for i,frame in enumerate(self.frame_extract(video_path)):
        #if(i % a == first_frame):
            faces = face_recognition.face_locations(frame)
            try:
                top,right,bottom,left = faces[0]

```

```

        frame = frame[top:bottom,left:right,:]
    except:
        pass
    frames.append(self.transform(frame))
    if(len(frames) == self.count):
        break
    """

    for i,frame in enumerate(self.frame_extract(video_path)):
        if(i % a == first_frame):
            frames.append(self.transform(frame))
    """

    # if(len(frames)<self.count):
    # for i in range(self.count-len(frames)):
    #     frames.append(self.transform(frame))
    #print("no of frames", self.count)
    frames = torch.stack(frames)
    frames = frames[:self.count]
    return frames.unsqueeze(0)

def frame_extract(self,path):
    vidObj = cv2.VideoCapture(path)
    success = 1
    while success:
        success, image = vidObj.read()
        if success:
            yield image

def im_convert(tensor, video_file_name):
    """ Display a tensor as an image. """
    image = tensor.to("cpu").clone().detach()
    image = image.squeeze()

```

```

image = inv_normalize(image)
image = image.numpy()
image = image.transpose(1,2,0)
image = image.clip(0, 1)
# This image is not used
#cv2.imwrite(os.path.join(settings.PROJECT_DIR,'uploaded_images',
video_file_name+'_convert_2.png'),image*255)
return image

def im_plot(tensor):
    image = tensor.cpu().numpy().transpose(1,2,0)
    b,g,r = cv2.split(image)
    image = cv2.merge((r,g,b))
    image = image*[0.22803, 0.22145, 0.216989] + [0.43216, 0.394666, 0.37645]
    image = image*255.0
    plt.imshow(image.astype(int))
    plt.show()

def predict(model,img,path = './', video_file_name=""):
    fmap,logits = model(img.to('cuda'))
    img = im_convert(img[:,-1,:,:), video_file_name)
    params = list(model.parameters())
    weight_softmax = model.linear1.weight.detach().cpu().numpy()
    logits = sm(logits)
    _,prediction = torch.max(logits,1)
    confidence = logits[:,int(prediction.item())].item()*100
    print('confidence of prediction:',logits[:,int(prediction.item())].item()*100)
    return [int(prediction.item()),confidence]

def plot_heat_map(i, model, img, path = './', video_file_name=""):
    fmap,logits = model(img.to('cuda'))
    params = list(model.parameters())

```

```

weight_softmax = model.linear1.weight.detach().cpu().numpy()
logits = sm(logits)
_,prediction = torch.max(logits,1)
idx = np.argmax(logits.detach().cpu().numpy())
bz, nc, h, w = fmap.shape
#out=np.dot(fmap[-1].detach().cpu().numpy().reshape((nc, h*w)).T,weight_softmax[idx,:].T)
out = np.dot(fmap[i].detach().cpu().numpy().reshape((nc, h*w)).T,weight_softmax[idx,:].T)
predict = out.reshape(h,w)
predict = predict - np.min(predict)
predict_img = predict / np.max(predict)
predict_img = np.uint8(255*predict_img)
out = cv2.resize(predict_img, (im_size,im_size))
heatmap = cv2.applyColorMap(out, cv2.COLORMAP_JET)
img = im_convert(img[:,-1,:,:), video_file_name)
result = heatmap * 0.5 + img*0.8*255
# Saving heatmap - Start
heatmap_name = video_file_name+"_heatmap_"+str(i)+".png"
image_name = os.path.join(settings.PROJECT_DIR, 'uploaded_images', heatmap_name)
cv2.imwrite(image_name,result)
# Saving heatmap - End
result1 = heatmap * 0.5/255 + img*0.8
r,g,b = cv2.split(result1)
result1 = cv2.merge((r,g,b))
return image_name

```

Model Selection

```

def get_accurate_model(sequence_length):
    model_name = []
    sequence_model = []
    final_model = ""
    list_models = glob.glob(os.path.join(settings.PROJECT_DIR, "models", "*.pt"))

```

```

for i in list_models:
    model_name.append(i.split("\\")[-1])
for i in model_name:
    try:
        seq = i.split("_")[3]
        if (int(seq) == sequence_length):
            sequence_model.append(i)
    except:
        pass

if len(sequence_model) > 1:
    accuracy = []
    for i in sequence_model:
        acc = i.split("_")[1]
        accuracy.append(acc)
    max_index = accuracy.index(max(accuracy))
    final_model = sequence_model[max_index]
else:
    final_model = sequence_model[0]
return final_model

```

```

ALLOWED_VIDEO_EXTENSIONS= set(['mp4','gif','webm','avi','3gp','wmv','flv','mkv'])

```

```

def allowed_video_file(filename):
    #print("filename" ,filename.rsplit('.',1)[1].lower())
    if (filename.rsplit('.',1)[1].lower() in ALLOWED_VIDEO_EXTENSIONS):
        return True
    else:
        return False

def index(request):
    if request.method == 'GET':

```

```

video_upload_form = VideoUploadForm()
if 'file_name' in request.session:
    del request.session['file_name']
if 'preprocessed_images' in request.session:
    del request.session['preprocessed_images']
if 'faces_cropped_images' in request.session:
    del request.session['faces_cropped_images']
return render(request, index_template_name, {"form": video_upload_form})
else:
    video_upload_form = VideoUploadForm(request.POST, request.FILES)
    if video_upload_form.is_valid():
        video_file = video_upload_form.cleaned_data['upload_video_file']
        video_file_ext = video_file.name.split('.')[-1]
        sequence_length = video_upload_form.cleaned_data['sequence_length']
        video_content_type = video_file.content_type.split('/')[0]
        if video_content_type in settings.CONTENT_TYPES:
            if video_file.size > int(settings.MAX_UPLOAD_SIZE):
                video_upload_form.add_error("upload_video_file", "Maximum file size 100 MB")
                return render(request, index_template_name, {"form": video_upload_form})

            if sequence_length <= 0:
                video_upload_form.add_error("sequence_length", "Sequence Length must be greater than
0")

                return render(request, index_template_name, {"form": video_upload_form})

            if allowed_video_file(video_file.name) == False:
                video_upload_form.add_error("upload_video_file", "Only video files are allowed ")
                return render(request, index_template_name, {"form": video_upload_form})

            saved_video_file = 'uploaded_file_'+str(int(time.time()))+"."+video_file_ext
            if settings.DEBUG:

```

```

        with open(os.path.join(settings.PROJECT_DIR, 'uploaded_videos', saved_video_file),
'wb') as vFile:
            shutil.copyfileobj(video_file, vFile)
            request.session['file_name'] = os.path.join(settings.PROJECT_DIR, 'uploaded_videos',
saved_video_file)
        else:
            with open(os.path.join(settings.PROJECT_DIR, 'uploaded_videos','app','uploaded_videos',
saved_video_file), 'wb') as vFile:
                shutil.copyfileobj(video_file, vFile)
                request.session['file_name'] = os.path.join(settings.PROJECT_DIR,
'uploaded_videos','app','uploaded_videos', saved_video_file)
                request.session['sequence_length'] = sequence_length
                return redirect('ml_app:predict')
    else:
        return render(request, index_template_name, {"form": video_upload_form})

def predict_page(request):
    if request.method == "GET":
        if 'file_name' not in request.session:
            return redirect("ml_app:home")
        if 'file_name' in request.session:
            video_file = request.session['file_name']
        if 'sequence_length' in request.session:
            sequence_length = request.session['sequence_length']
        path_to_videos = [video_file]
        video_file_name = video_file.split("\\")[-1]
        if settings.DEBUG == False:
            production_video_name = video_file_name.split('/')[3:]
            production_video_name = ''.join([str(elem) for elem in production_video_name])
            print("Production file name",production_video_name)
        video_file_name_only = video_file_name.split('.')[0]

```

```

video_dataset = validation_dataset(path_to_videos,
sequence_length=sequence_length,transform= train_transforms)
model = Model(2).cuda()
model_name = os.path.join(settings.PROJECT_DIR,'models',
get_accurate_model(sequence_length))
models_location = os.path.join(settings.PROJECT_DIR,'models')
path_to_model = os.path.join(settings.PROJECT_DIR, model_name)
model.load_state_dict(torch.load(path_to_model))
model.eval()
start_time = time.time()
# Start: Displaying preprocessing images
print("<=== | Started Videos Splitting | ===>")
preprocessed_images = []
faces_cropped_images = []
cap = cv2.VideoCapture(video_file)

frames = []
while(cap.isOpened()):
    ret, frame = cap.read()
    if ret==True:
        frames.append(frame)
        if cv2.waitKey(1) & 0xFF == ord('q'):
            break
    else:
        break
cap.release()
for i in range(1, sequence_length+1):
    frame = frames[i]
    image = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
    img = pImage.fromarray(image, 'RGB')
    image_name = video_file_name_only+"_preprocessed_"+str(i)+'.png'

```

```

if settings.DEBUG:
    image_path = os.path.join(settings.PROJECT_DIR, 'uploaded_images', image_name)
else:
    print("image_name",image_name)
    image_path = "/home/app/staticfiles" + image_name
img.save(image_path)
preprocessed_images.append(image_name)
print("<=== | Videos Splitting Done | ===>")
print("--- %s seconds ---" % (time.time() - start_time))
# End: Displaying preprocessing images
# Start: Displaying Faces Cropped Images
print("<=== | Started Face Cropping Each Frame | ===>")
padding = 40
faces_found = 0
for i in range(1, sequence_length+1):
    frame = frames[i]
    #fig, ax = plt.subplots(1,1, figsize=(5, 5))
    face_locations = face_recognition.face_locations(frame)
    if len(face_locations) == 0:
        continue
    top, right, bottom, left = face_locations[0]
    frame_face = frame[top-padding:bottom+padding, left-padding:right+padding]
    image = cv2.cvtColor(frame_face, cv2.COLOR_BGR2RGB)

    img = pImage.fromarray(image, 'RGB')
    image_name = video_file_name_only+"_cropped_faces_"+str(i)+'.png'
    if settings.DEBUG:
        image_path = os.path.join(settings.PROJECT_DIR, 'uploaded_images',
video_file_name_only+"_cropped_faces_"+str(i)+'.png')
    else:
        image_path = "/home/app/staticfiles" + image_name

```



```

img.save(image_path)
faces_found = faces_found + 1
faces_cropped_images.append(image_name)
print("<=== | Face Cropping Each Frame Done | ===>")
print("--- %s seconds ---" % (time.time() - start_time))

# No face is detected
if faces_found == 0:
    return render(request, predict_template_name, {"no_faces": True})

# End: Displaying Faces Cropped Images
try:
    heatmap_images = []
    for i in range(0, len(path_to_videos)):
        output = ""
        print("<=== | Started Prediction | ===>")
        prediction = predict(model, video_dataset[i], './', video_file_name_only)
        confidence = round(prediction[1], 1)
        print("<=== | Prediction Done | ===>")
        # print("<=== | Heat map creation started | ===>")
        # for j in range(0, sequence_length):
        #     heatmap_images.append(plot_heat_map(j, model, video_dataset[i], './',
video_file_name_only))
        if prediction[0] == 1:
            output = "REAL"
        else:
            output = "FAKE"
        print("Prediction : " , prediction[0],"==",output ,"Confidence : " , confidence)
        print("--- %s seconds ---" % (time.time() - start_time))
    if settings.DEBUG:
        return render(request, predict_template_name, {'preprocessed_images':

```

```
preprocessed_images, 'heatmap_images': heatmap_images, "faces_cropped_images":
faces_cropped_images, "original_video": video_file_name, "models_location": models_location,
"output": output, "confidence": confidence})
    else:
        return render(request, predict_template_name, {'preprocessed_images':
preprocessed_images, 'heatmap_images': heatmap_images, "faces_cropped_images":
faces_cropped_images, "original_video": production_video_name, "models_location":
models_location, "output": output, "confidence": confidence})
    except:
        return render(request, 'cuda_full.html')
def about(request):
    return render(request, about_template_name)

def handler404(request,exception):
    return render(request, '404.html', status=404)
def cuda_full(request):
    return render(request, 'cuda_full.html')
```

SNAPSHOTS

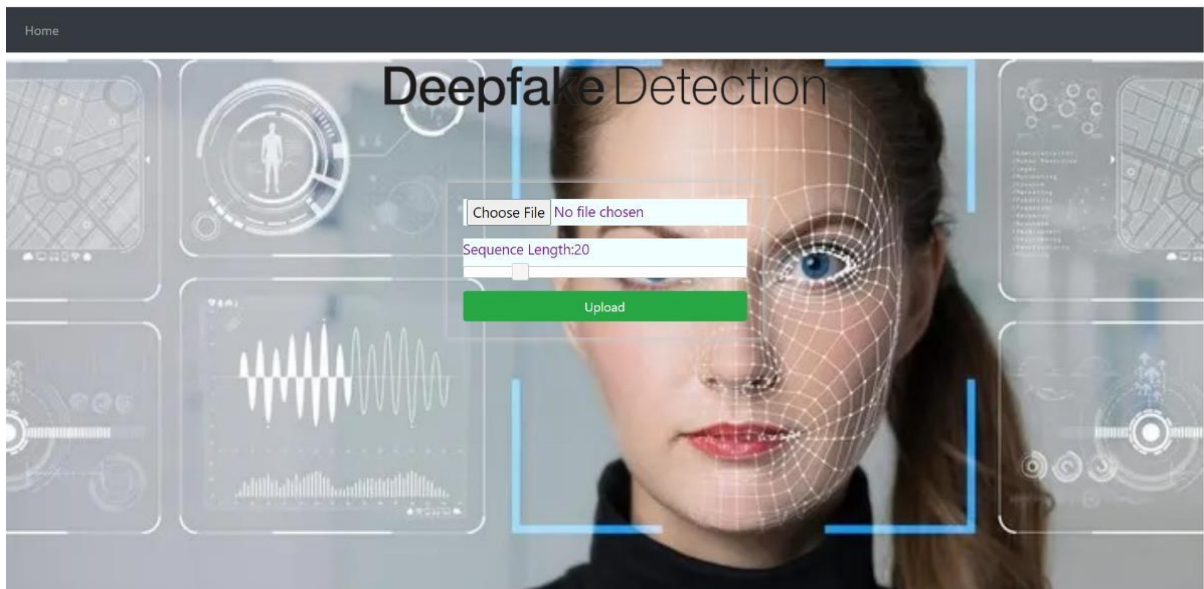


Fig.6.1: Home Page

The text at the top of the page says "Home" and "Deepfake Detection". Below that, there is a button that says "Choose File" and a text box that says "No file chosen". There is also a text box that says "Sequence Length:20" and a button that says "Upload".

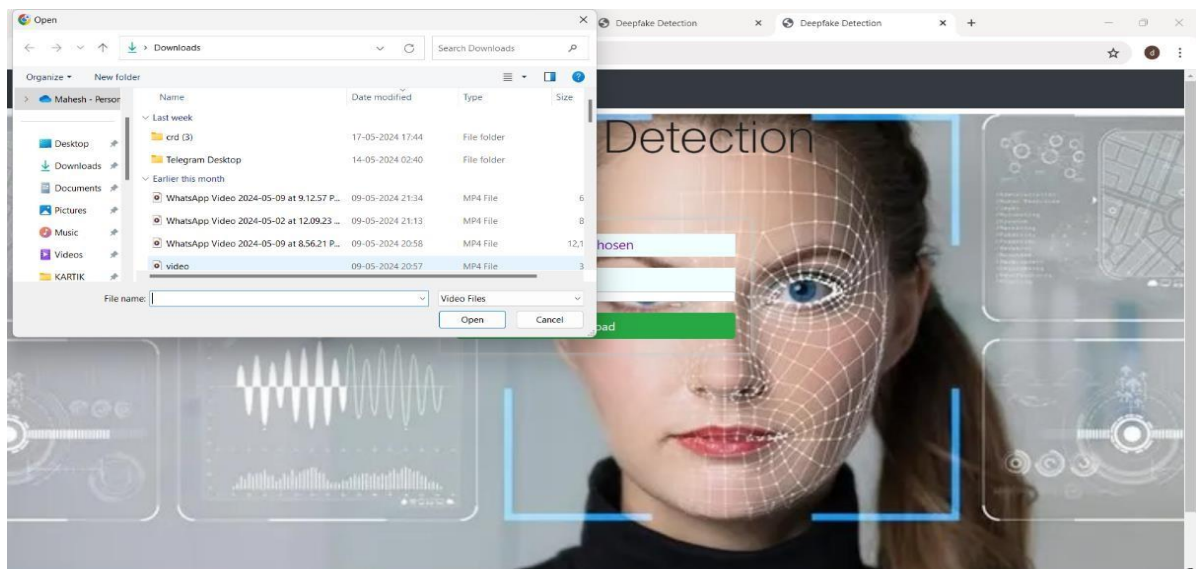


Fig.6.2: Uploading Video

This page describes about the uploading the video. When we choose the choose file option in home page it will take us for browsing the video for uploading.

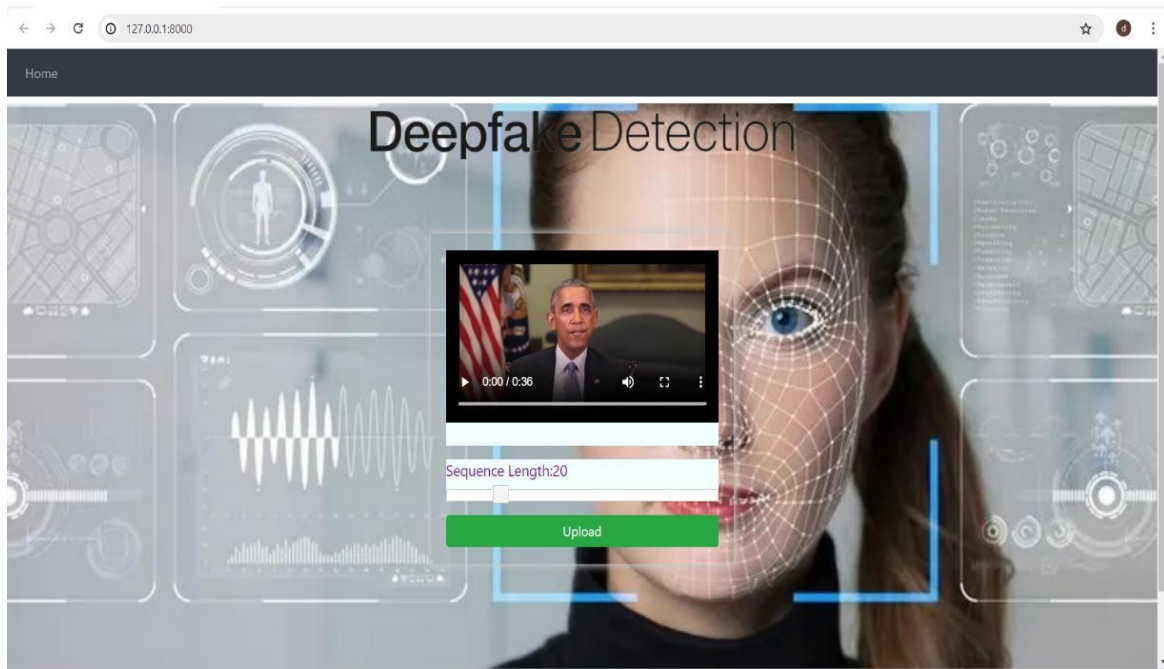


Fig.6.3: Uploading real Video

This page describes about the uploading real video. When we select the choose file option in home page it will take us for browsing the video for uploading.

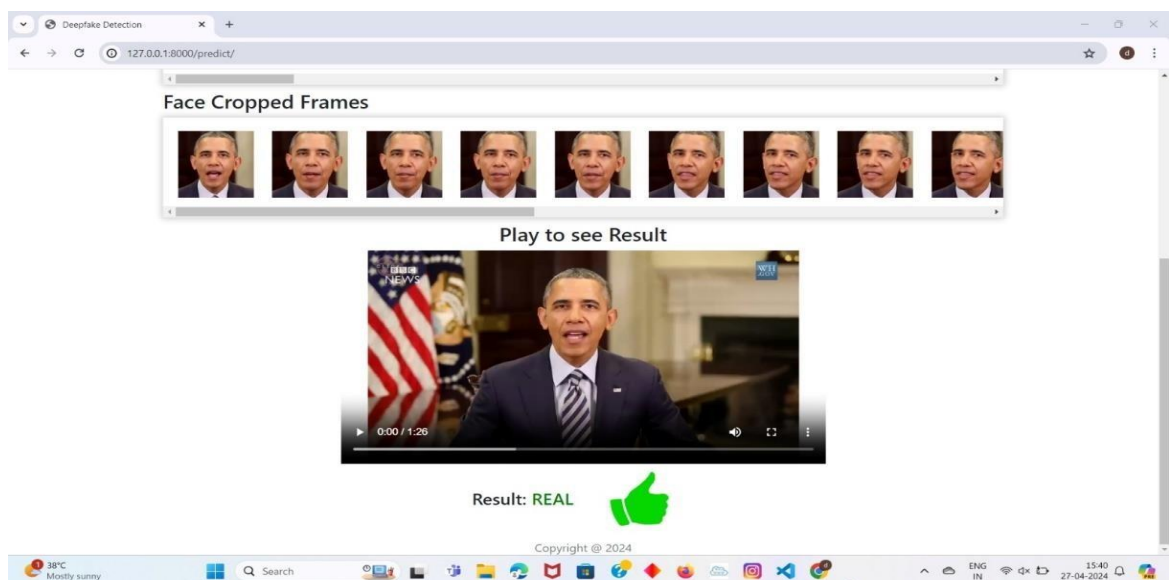


Fig.6.4: Face Cropped images

This page describes the face cropping of frames in the real video. When video is uploaded it is preprocessed by cropping the faces and removing the other frames which does not contain face.

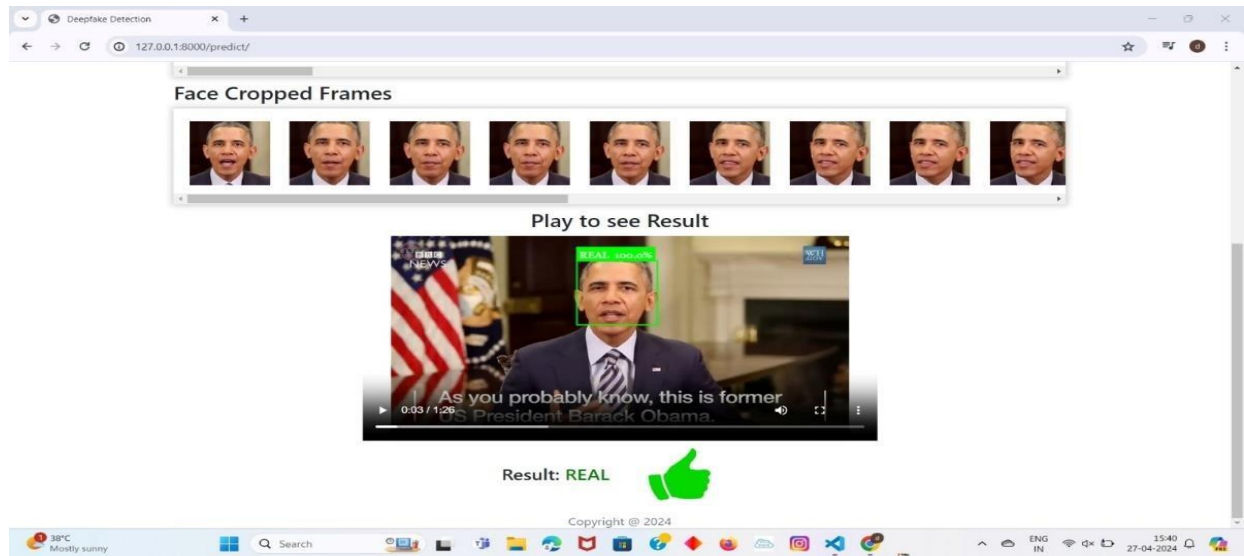


Fig.6.5: Real Video output

This page displays the output. It displays message real with thumb in green color. It indicates that uploaded video is real.

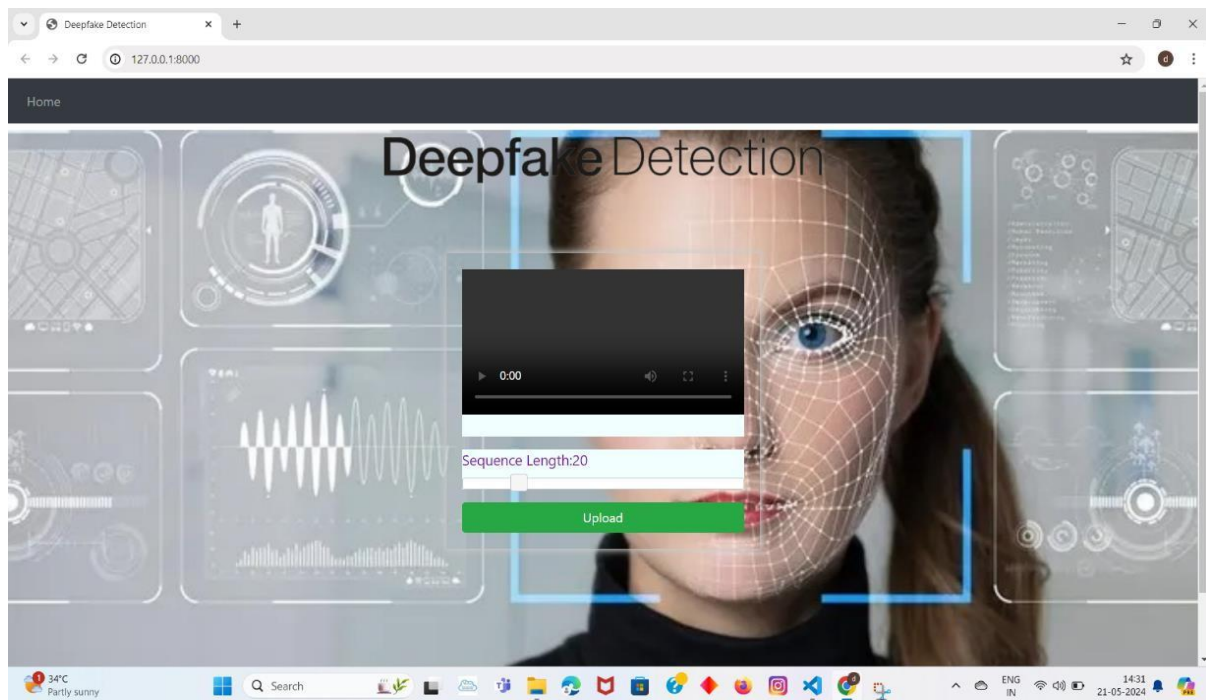


Fig.6.6: Uploading Fake Video

This page describes about the uploading fake video. When we choose the choose file option in home page it will take us for browsing the video for uploading.

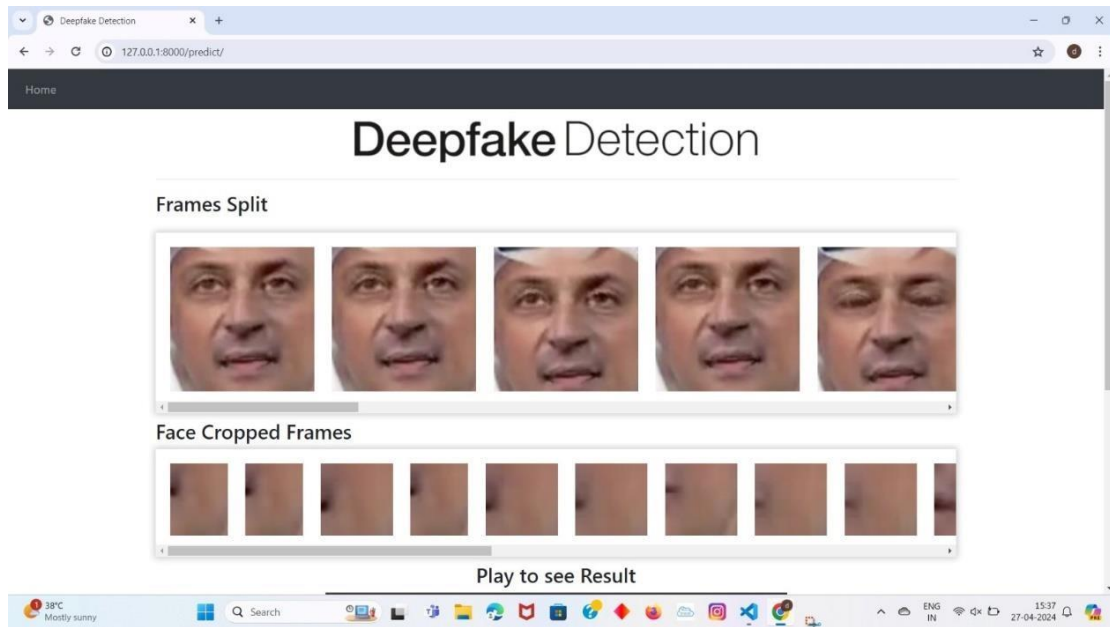


Fig.6.7: Face cropping in fake video

This page describes the face cropping of frames in the fake video. The splitting of frames is described.

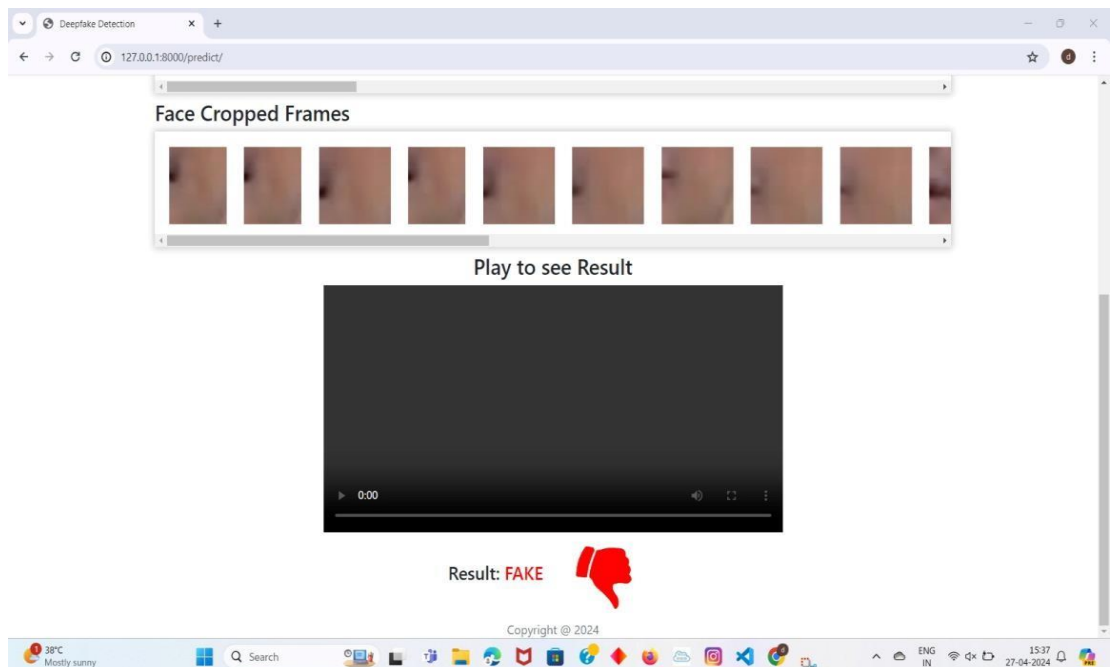


Fig.6.8: Fake video output

This page displays the output. It displays message fake with thumb in red color. It indicates that uploaded video is fake.

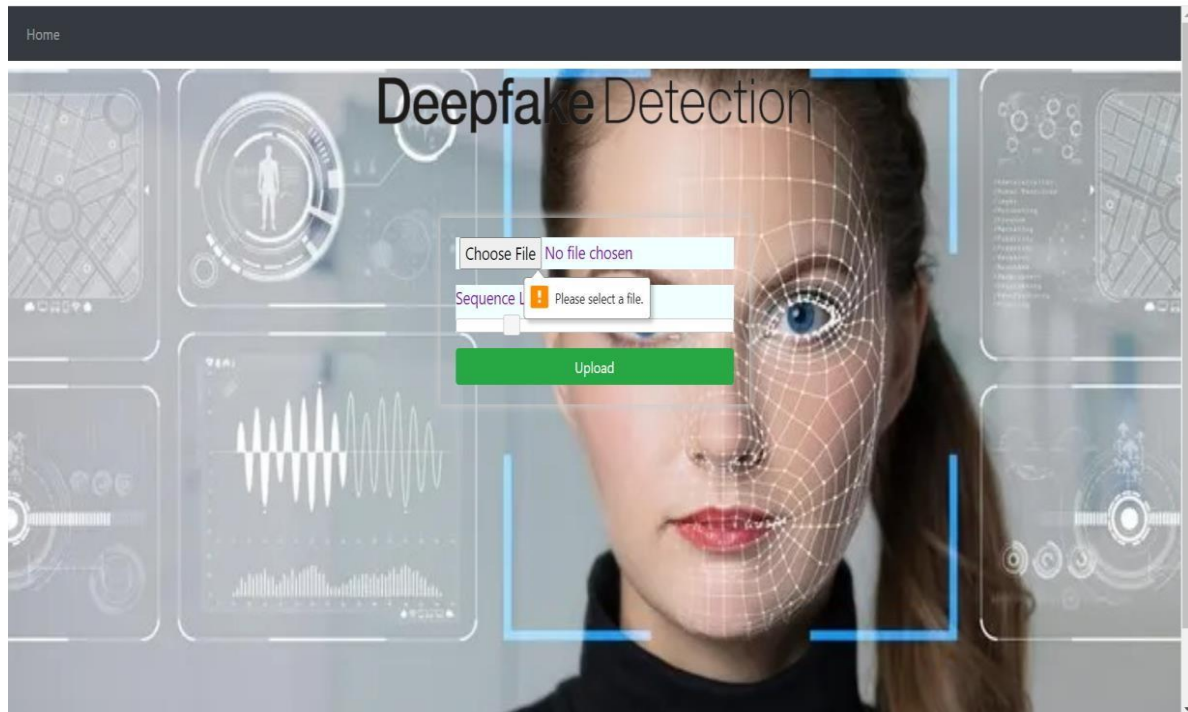


Fig.6.9: Pressing Upload button without selecting video

This image describes that without selecting the video the upload button is pressed. When upload button is pressed please select a file message is displayed.

CHAPTER 7

RESULT

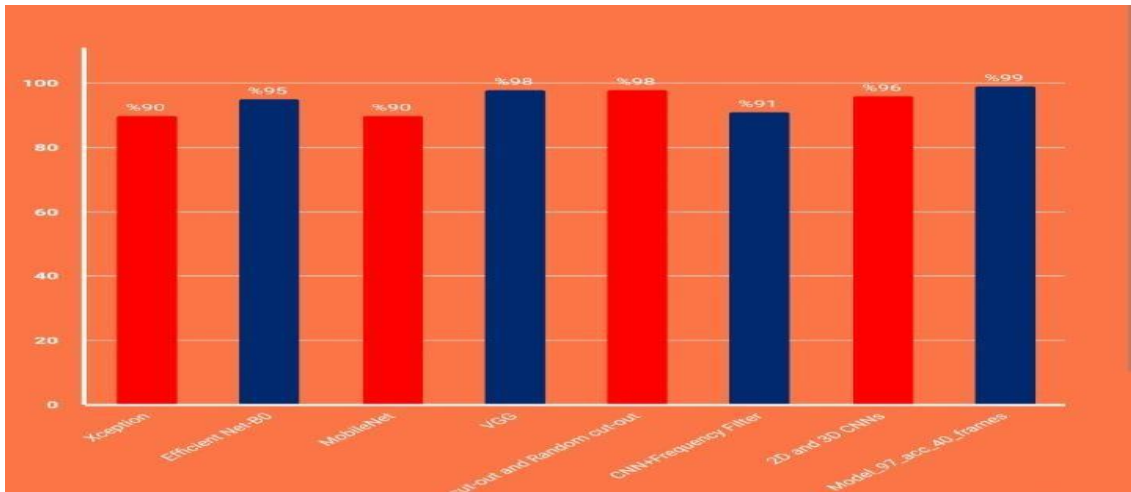


Fig.7.1: Deepfake Model Results

In the realm of deepfake detection, various deep learning models and techniques have been explored to differentiate between authentic and manipulated media. Among these approaches, models like Xception, EfficientNet-B0, MobileNet, VGG, and combinations of convolutional neural networks (CNNs) have shown promising results. Hypothetically, these models have demonstrated accuracies ranging from 90% to 98%, with VGG achieving the highest accuracy at 98%. Additionally, data augmentation techniques such as Cut Out and Random Cutout have bolstered detection accuracy up to 98%. More sophisticated methods like CNNs with frequency filters, 2D and 3D CNNs, and specialized models like Model_97_acc_40_frames have shown accuracies hovering around 91% to 99%, indicating their efficacy in discerning deepfakes from authentic content.

These results are encouraging, suggesting that deep learning techniques can be effective tools in combating the proliferation of manipulated media. However, it's crucial to acknowledge that these numbers are hypothetical and may vary based on several factors such as dataset composition, the quality of deepfake videos, and implementation nuances. Nonetheless, these findings provide valuable insights into the potential of deep learning in addressing the challenges posed by deepfakes.

CHAPTER 8

CONCLUSION

The project “Deepfake detection using deep learning” the proposed deepfake detection system integrates a hybrid approach, combining temporal and spatial features, and robustness against adversarial attacks, dynamic learning, collaboration networks, and user feedback loops. By leveraging powerful GPUs, advanced deep learning frameworks, and cloud services, the system aims to provide a reliable, adaptable, and transparent solution.

The iterative learning process and incorporation of diverse methodologies underscore a commitment to staying ahead of evolving deepfake technologies. As the field progresses, this holistic approach positions the system to be effective in countering the increasing sophistication of deepfake content production, contributing to the ongoing efforts to enhance the integrity of digital media.

REFERENCES

- [1]. Deng Pan, Lixian Sun, Rui Wang,Xingjian Zhang, Richard O Sinnott, ”Deepfake Detection Through Deep Learning” 2020 IEEE/ACM International Conference on Big Data Computing, Applications and Technologies (BDCAT) .
- [2]. D. Zhang, F. Lin, Y. Hua, P. Wang, D. Zeng, S. Ge, Deepfake video detection with spatiotemporal dropout transformer, in: Proceedings of the 30th ACM International Conference on Multimedia, 2022, pp. 5833–5841.
- [3]. S. A. Khan, D.-T. Dang-Nguyen, Hybrid transformer network for deepfake detection, in: Proceedings of the 19th international conference on content-based multimedia indexing, 2022, pp. 8–14.
- [4]. D. A. Coccomini, N. Messina, C. Gennaro, F. Falchi, Combining efficientnet and vision transformers for video deepfake detection, in: Image Analysis and Processing (ICIAP 2022), Springer International Publishing, Cham, 2022, pp. 219–229.
- [5]. J. Feinland, J. Barkovitch, D. Lee, A. Kaforey, U. A. Ciftci, Poker bluff detection dataset based on facial analysis, in: International Conference on Image Analysis and Processing, Springer, 2022, pp. 400–410.
- [6]. Z. Xue, Q. Liu, H. Shi, R. Zou, X. Jiang, A transformer-based deepfake-detection method for facial organs, *Electronics* 11 (24).
- [7]. Y. Zhang, T. Wang, M. Shu, Y. Wang, A robust lightweight deepfake detection network using transformers, in: Pacific Rim International Conference on Artificial Intelligence, Springer, 2022, pp. 275–288.
- [8]. A. Khormali, J.-S. Yuan, DFDT: an end-to-end deepfake detection framework using vision transformer, *Applied Sciences* 12 (6) (2022) 2953.
- [9]. D. A. Coccomini, R. Caldelli, F. Falchi, C. Gennaro, G. Amato, Cross-Forgery Analysis of Vision Transformers and CNNs for Deepfake Image Detection, in: Proceedings of the 1st International Workshop on Multi-media AI against Disinformation, 2022, pp. 52–58.
- [10]. J. Wang, Z. Wu, W. Ouyang, X. Han, J. Chen, Y.-G. Jiang, S.-N. Li, M2tr: Multi-modal multi-scale transformers for deepfake detection, in: Proceedings of the 2022 international conference on multimedia retrieval, 2022, pp. 615–623.

- [11]. M. A. Raza, K. M. Malik, I. U. Haq, Holisticdfd: Infusing spatiotemporal transformer embeddingsfor deepfake detection, *Information Sciences* (2023) 119352.
- [12]. Y. Heo, W. Yeo, B. Kim, Deepfake detection algorithm based on improved vision transformer, *Applied Intelligence* 53 (2023) 7512–7527.
- [13]. H. Lin, W. Huang, W. Luo, W. Lu, Deepfake detection with multi-scale convolution and vision transformer, *Digital Signal Processing* 134 (2023)103895.
- [14]. F. Khalid, M. H. Akbar, S. Gul, Swynt: Swin y-net transformers for deepfake detection, in: 2023 International Conference on Robotics and Automation in Industry (ICRAI), 2023, pp. 1–6. doi:10.1109/ICRAI57502.2023.10089585.
- [15]. David Guera, Edward J. Delp, "Deepfake Detection Using Recurrent Neural Network" based on research sponsored by the defense advanced Projects agency.
- [16]. Aya Ismail, Marwa Elpeltagy, Mervat S. Zaki, Kamal Eldahshan, "A New deep Learning-based methodology for video deepfake detection using XGBoost. *Sensors* 2021, 21, 5413. <https://doi.org/10.3390/s21165413>.
- [17]. TackHyun Jung, SangWon Kim, KeeCheon Kim, "Deepvision: Deepfakes detection using humaneye blinking pattern DOI.10.1109/ACCESS.2020.2988660, IEEE Access.
- [18]. Leandron A. Passos, Danilo Jodas, Kelton A. P. Costa, Luis A. Souza Junior, Douglas Rodrigues, Javier Del Ser, David Camacho, Joao Paulo Papa "A review of deep learning based approaches for deepfake content detection".