

# **REVEALING AND CLASSIFICATION OF DEEP FAKE IMAGES WITH VIDEOS USING CUSTOMIZED DEEP LEARNING MODELS**

**A PROJECT REPORT**

*Submitted by*

**PRAVEEN KUMAR K**

**REGISTER NO: 717822Z135**

*in partial fulfillment for the award of the degree of*

**MASTER OF COMPUTER APPLICATIONS**

**SCHOOL OF COMPUTER APPLICATIONS  
KARPAGAM COLLEGE OF ENGINEERING  
COIMBATORE**

**ANNA UNIVERSITY, CHENNAI**

**MAY 2024**



**KARPAGAM COLLEGE OF ENGINEERING**

**COIMBATORE – 641 032**

**ANNA UNIVERSITY, CHENNAI**

**BONAFIDE CERTIFICATE**

Certified that this Report titled “**REVEALING AND CLASSIFICATION OF DEEP FAKE IMAGES WITH VIDEOS USING CUSTOMIZED DEEP LEARNING MODELS**” is the bonafide work of **PRAVEEN KUMAR K (Reg. No.717822z135)** who carried out the work under my supervision. Certified further that to the best of my knowledge the work reported herein does not form part of any other thesis or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

Signature of the HOD with date

Dr. K. Anuradha

Director

School of Computer Applications

Karpagam College of Engineering

Coimbatore – 641032

Signature of the Supervisor with date

Mr. R. Ramprashath

Assistant Professor

School of Computer Applications

Karpagam College of Engineering

Coimbatore – 641032

---

Certified that the candidate was examined during the viva voce examinations held on \_\_\_\_\_

Signature of the Internal Examiner with date

Signature of External Examiner with date



**NoviTech**  
the innovation partner

02<sup>nd</sup> January 2024

**Mr. Praveen Kumar K**

**Reg no: 717822z135**

**Karpagam College of Engineering  
Coimbatore.**

**Dear Praveen Kumar K,**

**Sub : Internship Confirmation**

On behalf of, **NoviTech R&D Pvt Ltd.** We would like to inform that you are being accepted as one of our interns. We are pleased to inform you that you have been qualified as per the requirements for the internship. You will be working with our Technical team. Your internship will begin with effective from **January 2024 to April 2024**. You will be assigned to various tasks which relates to the project assigned to you after whom your performance will be assessed and appraised.

**For NoviTech R&D Pvt Ltd**

  
Authorized Signatory

NoviTech R&D Pvt Ltd.,  
2<sup>nd</sup> Floor, Sai Sruthi Complex, Ramar Koil Street, Ram Nagar,  
Coimbatore - 641 009, Tamilnadu, INDIA.

☎ 0422 - 4645101  
🌐 [www.novitechrd.com](http://www.novitechrd.com)  
✉ [support@novitechrandd.com](mailto:support@novitechrandd.com)



## **DECLARATION**

I hereby declare that this project report entitled **“REVEALING AND CLASSIFICATION OF DEEP FAKE IMAGES WITH VIDEOS USING CUSTOMIZED DEEP LEARNING MODELS”** submitted by me for the degree of **“MASTER OF COMPUTER APPLICATIONS”** at **Karpagam College of Engineering, Coimbatore** is the record of original work done by me under the guidance and supervision of **Mr. R. Ramprashath, Assistant Professor** at the School of Computer Applications, Karpagam College of Engineering, Coimbatore – 641032 and has not formed the basis for the award of any degree, or diploma or titles in this institution or any other Institution of higher learning.

Date:

**Name and Signature of the Candidate**

Place: Coimbatore

## ACKNOWLEDGMENT

**First and foremost praises and thanks to the almighty for her showers and blessings throughout my project work to complete it successfully.**

I extend my gratitude to the Management of Karpagam College of Engineering, Coimbatore for the excellent infrastructure and support facilities to undergo the project work.

I am very grateful to **Dr. V. Kumar Chinnaiyan, the Principal** and **Dr. K. Anuradha, Director**, School of Computer Applications for provided the facilities, support and permission to carried out my project work at our esteemed institution.

I record my sincere gratitude to my Project Coordinator **Mr. R. Ramprashath** for giving inputs, encouragement for the continuous improvement during the progress and to complete this project work.

I would like to express my sincere gratitude to my Supervisor **Mr. R. Ramprashath** for the continuous support for my PG study, for his motivation and adequate guidance which helped me to achieve success in all my accomplishments and to complete this project work.

I also thank all the teaching faculty members and non-teaching Staff members of the **School of Computer Applications**, Karpagam College of Engineering, Coimbatore for their kindness and support.

I would like to thank **my parents, family members and friends** who sacrificed their time and energy to complete the project work successfully.

**PRAVEEN KUMAR K**

## **ABSTRACT**

Deep fakes are becoming more common; they include editing previously published films and photos to produce content that appears authentic but is wholly fake. The development process has been considerably expedited by the widespread availability of deep learning techniques, such as autoencoders, Generative Adversarial Networks (GANs), and user-friendly software. These sophisticated algorithms adeptly fuse and modify visual and audio elements, facilitating the production of content that closely mimics genuine footage, even for those without specialized knowledge. The malicious manipulation of images and videos poses significant security and societal concerns. With an emphasis on facial alteration, the goal of this research is to create a deep learning perfect for the detection and classification of deepfake images and videos. The dataset used for the project is either Face Forensics++, Celeb-DF, or the Deepfake Detection Challenge Dataset (DFDC), available on Kaggle, consisting of real and deepfake images and videos. By utilising Recurrent and Convolutional Neural Networks, we have made development in DF detection. Commencing with preprocessing the data, extracting frames from the videos, and separating the dataset into training and validation sets. For the detection and classification of deepfake images and videos, OpenCV, and Face Recognition for facial detection, Convolutional neural networks (CNNs) are used by the system to extract features at the frame level. A recurrent neural network is trained using these features (RNN). Various techniques such as data augmentation, learning rate scheduling, and early stopping enhance model performance. This comprehensive approach ensures accurate discrimination between authentic and deep fake content, addressing concerns regarding the integrity of digital media.

## TABLE OF CONTENTS

CHAPTER NO	TITLE	PAGE NO
	<b>ABSTRACT</b>	<b>iii</b>
	<b>LIST OF FIGURES</b>	<b>vii</b>
	<b>LIST OF TABLES</b>	<b>viii</b>
<b>1</b>	<b>SYNOPSIS</b>	<b>1</b>
<b>2</b>	<b>TECHNICAL KEYWORDS</b>	<b>2</b>
	2.1 AREA OF PROJECT	2
	2.2 TECHNICAL KEYWORDS	2
<b>3</b>	<b>INTRODUCTION</b>	<b>3</b>
	3.1 PROJECT IDEA	3
	3.2 MOTIVATION OF THE PROJECT	3
	3.3 LITERATURE SURVEY	4
<b>4</b>	<b>PROBLEM DEFINITION AND SCOPE</b>	<b>6</b>
	4.1 PROBLEM STATEMENT	6
	4.1.1 GOALS AND OBJECTIVES	6
	4.1.2 STATEMENT OF SCOPE	6
	4.2 MAJOR CONSTRAINTS	7
	4.3 METHODOLOGIES OF PROBLEM SOLVING	7
	4.3.1 ANALYSIS	7
	4.3.2 DESIGN	8
	4.3.3 DEVELOPMENT	8
	4.3.4 EVALUATION	9
	4.4 OUTCOME	9
	4.5 APPLICATIONS	9
	4.6 HARDWARE RESOURCES REQUIRED	9
	4.7 SOFTWARE RESOURCES REQUIRED	9
<b>5</b>	<b>PROJECT PLAN</b>	<b>11</b>
	5.1 PROJECT MODEL ANALYSIS	11
	5.1.1 RECONCILED ESTIMATES	12
	5.1.2 COST ESTIMATION USING	12

	COCOMO(CONSTRUCTIVE COST) MODEL	13
	5.2 RISK MANAGEMENT W.R.T. NP HARD ANALYSIS	13
	5.2.1 RISK IDENTIFICATION	13
	5.2.2 RISK ANALYSIS	14
	5.3 PROJECT SCHEDULE	14
	5.3.1 PROJECT TASK SET	15
	5.3.2 TIMELINE CHART	
<b>6</b>	<b>SYSTEM ANALYSIS</b>	<b>16</b>
	6.1 INTRODUCTION	16
	6.1.1 PURPOSE AND SCOPE OF DOCUMENT	16
	6.1.2 USE CASE VIEW	16
	6.2 FUNCTIONAL MODEL AND DESCRIPTION	17
	6.2.1 DATA FLOW DIAGRAM	17
	6.2.2 ACTIVITY DIAGRAM	19
	6.2.3 NON FUNCTIONAL REQUIREMENTS	20
	6.2.4 SEQUENCE DIAGRAM	21
<b>7</b>	<b>DETAILED DESIGN DOCUMENT</b>	<b>22</b>
	7.1 INTRODUCTION	22
	7.1.1 SYSTEM ARCHITECTURE	22
	7.2 ARCHITECTURAL DESIGN	24
	7.2.1 MODULE 1 : DATA-SET GATHERING	24
	7.2.2 MODULE 2 : PRE-PROCESSING	25
	7.2.3 MODULE 3 : DATA-SET SPLIT	26
	7.2.4 MODULE 4 : MODEL ARCHITECTURE	27
	7.2.5 MODULE 5 : HYPER-PARAMETER TUNING	29
<b>8</b>	<b>PROJECT IMPLEMENTATION</b>	<b>30</b>
	8.1 INTRODUCTION	30
	8.2 TOOLS AND TECHNOLOGIES USED	31
	8.2.1 PLANNING	31
	8.2.2 UML TOOLS	31
	8.2.3 PROGRAMMING LANGUAGES	31
	8.2.4 PROGRAMMING FRAMEWORKS	31



	8.2.5 IDE	31
	8.2.6 VERSIONING CONTROL	31
	8.2.7 CLOUD SERVICES	31
	8.2.8 APPLICATION AND WEB SERVERS	31
	8.2.9 LIBRARIES	32
	8.3 ALGORITHM DETAILS	32
	8.3.1 DATASET DETAILS	32
	8.3.2 PREPROCESSING DETAILS	32
	8.3.3 MODEL DETAILS	33
	8.3.4 MODEL TRAINING DETAILS	36
	8.3.5 MODEL PREDICTION DETAILS	38
<b>9</b>	<b>SOFTWARE TESTING</b>	<b>39</b>
	9.1 TYPE OF TESTING USED	39
	9.2 TEST CASES AND TEST RESULTS	40
<b>10</b>	<b>CONCLUSION AND FUTURE ENHANCEMENTS</b>	<b>41</b>
	10.1 CONCLUSION	41
	10.2 FUTURE ENHANCEMENTS	41
	10.3 PROJECT COMPLETION CERTIFICATE	
<b>11</b>	<b>APPENDICES</b>	<b>43</b>
	11.1 SOURCE CODE	43
	11.2 SCREENSHOTS	49
	11.3 PAPER PRESENTATION CERTIFICATE	
	11.4 PAPER PROCEEDING	
<b>12</b>	<b>REFERENCES</b>	<b>58</b>
	12.1 REFERENCES	58

## LIST OF FIGURES

S.NO	FIGURE NO	FIGURE CAPTION	PAGE NO
1.	5.1	SPIRAL METHODOLOGY SDLC	11
2.	6.1	USE CASE DIAGRAM	16
3.	6.2	DFD LEVEL 0	17
4.	6.3	DFD LEVEL 1	18
5.	6.4	DFD LEVEL 2	18
6.	6.5	TRAINING WORKFLOW	19
7.	6.6	TESTING WORKFLOW	20
8.	6.7	SEQUENCE DIAGRAM	21
9.	7.1	SYSTEM ARCHITECTURE	22
10.	7.2	DEEPFAKE GENERATION	23
11.	7.3	FACE SWAPPED DEEPFAKE GENERATION	23
12.	7.4	DATASET	25
13.	7.5	PRE-PROCESSING OF VIDEO	26
14.	7.6	TRAIN TEST SPLIT	27
15.	7.7	OVERVIEW OF OUR MODEL	28
16.	8.1	RESNEXT ARCHITECTURE	33
17.	8.2	RESNEXT WORKING	34
18.	8.3	OVERVIEW OF RESNEXT ARCHITECTURE	34
19.	8.4	OVERVIEW OF LSTM ARCHITECTURE	35
20.	8.5	INTERNAL LSTM ARCHITECTURE	35
21.	8.6	RELU ACTIVATION FUNCTION	35
22.	8.7	DROPOUT LAYER OVERVIEW	36
23.	8.8	SOFTMAX LAYER	37

## LIST OF TABLES

S.NO	TABLE NO	NAME OF THE TABLE	PAGE NO
1.	4.6.1	HARDWARE REQUIREMENTS	9
2.	5.1.1.1	COST ESTIMATION	12
3.	5.2.2.1	RISK DESCRIPTION	14
4.	5.2.2.2	RISK PROBABILITY DEFINITIONS	14
5.	9.2.1	TEST CASE REPORT	40

## CHAPTER 1

### SYNOPSIS

Deep fake is a technique for human image synthesis based on neural network tools like GAN(Generative Adversarial Network) or Auto Encoders etc. These tools super impose target images onto source videos using a deep learning techniques and create a realistic looking deep fake video. These deep-fake video are so real that it becomes impossible to spot difference by the naked eyes. In this work, we describe a new deep learning-based method that can effectively distinguish AI-generated fake videos from real videos. We are using the limitation of the deep fake creation tools as a powerful way to distinguish between the pristine and deep fake videos. During the creation of the deep fake the current deep fake creation tools leaves some distinguishable artifacts in the frames which may not be visible to the human being but the trained neural networks can spot the changes. Deepfake creation tools leave distinctive artefacts in the resulting Deep Fake videos, and we show that they can be effectively captured by Res-Next Convolution Neural Networks.

Our system uses a Res-Next Convolution Neural Networks to extract frame-level features. These features are then used to train a Long Short Term Memory(LSTM) based Recurrent Neural Network(RNN) to classify whether the video is subject to any kind of manipulation or not, i.e whether the video is deep fake or real video. We proposed to evaluate our method against a large set of deep fake videos collected from multiple video websites. We are tried to make the deep fake detection model perform better on real time data. To achieve this we trained our model on combination of available data-sets. So that our model can learn the features from different kind of images. We extracted a adequate amount of videos from Face-Forensic++[1], Deepfake detection challenge[2], and Celeb-DF[3] data-sets. We also evaluated our model against the large amount of real time data like YouTube data-set to achieve competitive results in the real time scenario's.

## **CHAPTER 2**

### **TECHNICAL KEYWORD**

#### **2.1 Area of Project**

Our project is a Deep learning project which is a sub branch of Artificial Intelligence and deals with the human brain inspired neural network technology. Computer vision plays an important role in our project. It helps in processing the video and frames with the help of Open-CV. A PyTorch trained model is a classifier to classify the source video as deepfake or pristine.

#### **2.2 Technical Keywords**

- Deep learning
- Computer vision
- Res-Next Convolution Neural Network
- Long short-term memory (LSTM)
- OpenCV
- Face Recognition
- GAN (Generative Adversarial Network)
- PyTorch.

## **CHAPTER 3**

### **INTRODUCTION**

#### **3.1 Project Idea**

In the world of ever growing Social media platforms, Deepfakes are considered as the major threat of the AI. There are many Scenarios where these realistic face swapped deepfakes are used to create political distress, fake terrorism events, revenge porn, blackmail peoples are easily envisioned. Some of the examples are Brad Pitt, Angelina Jolie nude videos.

It becomes very important to spot the difference between the deepfake and pristine video. We are using AI to fight AI. Deepfakes are created using tools like FaceApp[11] and Face Swap [12], which using pre-trained neural networks like GAN or Auto encoders for these deepfakes creation. Our method uses a LSTM-based artificial neural network to process the sequential temporal analysis of the video frames and pre-trained ResNext CNN to extract the frame-level features. ResNext Convolution neural network extracts the frame-level features and these features are further used to train the Long Short Term Memory based artificial Recurrent Neural Network to classify the video as Deepfake or real. To emulate the real-time scenarios and make the model perform better on real-time data, we trained our method with large amount of balanced and combination of various available dataset like FaceForensic++[1], Deepfake detection challenge[2], and Celeb-DF[3].

Further to make the ready to use for the customers, we have developed a front end application where the user the user will upload the video. The video will be processed by the model and the output will be rendered back to the user with the classification of the video as deepfake or real and confidence of the model.

#### **3.2 Motivation of the Project**

The increasing sophistication of mobile camera technology and the ever-growing reach of social media and media sharing portals have made the creation and propagation of digital videos more convenient than ever before. Deep learning has given rise to technologies that would have been thought impossible only a handful of years ago. Modern generative models are one example of these, capable of synthesizing hyper realistic images, speech, music, and even video. These models have found use in a wide variety of applications, including making the world more accessible through text-to-speech, and helping generate training data for medical imaging.

Like any transformative technology, this has created new challenges. So-called "deep fakes" produced by deep generative models that can manipulate video and audio clips. Since their first appearance in late 2017, many open-source deep fake generation methods and tools have emerged now, leading to a growing number of synthesized media clips. While many are likely intended to be humorous, others could be harmful to individuals and society. Until recently, the number of fake videos and their degrees of realism has been increasing due to availability of the editing tools, the high demand on domain expertise.

Spreading of the Deep fakes over the social media platforms have become very common leading to spamming and peculating wrong information over the platform. Just imagine a deep fake of our prime minister declaring war against neighboring countries, or a Deep fake of reputed celebrity abusing the fans. These types of the deep fakes will be terrible, and lead to threatening, misleading of common people.

To overcome such a situation, Deep fake detection is very important. So, we describe a new deep learning-based method that can effectively distinguish AI- generated fake videos (Deep Fake Videos) from real videos. It's incredibly important to develop technology that can spot fakes, so that the deep fakes can be identified and prevented from spreading over the internet.

### **3.3 Literature Survey**

Face Warping Artifacts [15] used the approach to detect artifacts by comparing the generated face areas and their surrounding regions with a dedicated Convolutional Neural Network model. In this work there were two-fold of Face Artifacts.

Their method is based on the observations that current deepfake algorithm can only generate images of limited resolutions, which are then needed to be further transformed to match the faces to be replaced in the source video. Their method has not considered the temporal analysis of the frames.

Detection by Eye Blinking [16] describes a new method for detecting the deep- fakes by the eye blinking as a crucial parameter leading to classification of the videos as deepfake or pristine. The Long-term Recurrent Convolution Network (LRCN) was used for temporal analysis of the cropped frames of eye blinking. As today the deepfake generation algorithms have become so powerful that lack of eye blinking can not be the only clue for detection of the deepfakes. There must be certain other parameters must be considered for the detection of deep- fakes like teeth enchantment, wrinkles on faces, wrong placement of eyebrows etc.

Capsule networks to detect forged images and videos [17] uses a method that uses a capsule network to detect forged, manipulated images and videos in different scenarios, like replay attack detection and computer-generated video detection.

In their method, they have used random noise in the training phase which is not a good option. Still the model performed beneficial in their dataset but may fail on real time data due to noise in training. Our method is proposed to be trained on noiseless and real time datasets.

Recurrent Neural Network [18] (RNN) for deepfake detection used the approach of using RNN for sequential processing of the frames along with ImageNet pre-trained model. Their process used the HOHO [19] dataset consisting of just 600 videos.

Their dataset consists small number of videos and same type of videos, which may not perform very well on the real time data. We will be training our model on large number of Realtime data.

Synthetic Portrait Videos using Biological Signals [20] approach extract biological signals from facial regions on pristine and deepfake portrait video pairs. Applied transformations to compute the spatial coherence and temporal consistency, capture the signal characteristics in feature vector and photoplethysmography (PPG) maps, and further train a probabilistic Support Vector Machine (SVM) and a Convolutional Neural Network (CNN). Then, the average of authenticity probabilities is used to classify whether the video is a deepfake or a pristine.

Fake Catcher detects fake content with high accuracy, independent of the generation, content, resolution, and quality of the video. Due to lack of a discriminator leading to the loss in their findings to preserve biological signals, formulating a differentiable loss function that follows the proposed signal processing steps is not straight forward process.



## CHAPTER 4

### PROBLEM DEFINITION AND SCOPE

#### 4.1 Problem Statement

Convincing manipulations of digital images and videos have been demonstrated for several decades through the use of visual effects, recent advances in deep learning have led to a dramatic increase in the realism of fake content and the accessibility in which it can be created. These so-called AI-synthesized media (popularly referred to as deep fakes). Creating the Deep Fakes using the Artificially intelligent tools are simple task. But, when it comes to detection of these Deep Fakes, it is major challenge. Already in the history there are many examples where the deepfakes are used as powerful way to create political tension[14], fake terrorism events, revenge porn, blackmail peoples etc. So it becomes very important to detect these deepfake and avoid the percolation of deepfake through social media platforms. We have taken a step forward in detecting the deep fakes using LSTM-based artificial Neural network.

##### 4.1.1 Goals and objectives

Goal and Objectives:

- Our project aims at discovering the distorted truth of the deep fakes.
- Our project will reduce the Abuses' and misleading of the common people on the world wide web.
- Our project will distinguish and classify the video as deepfake or pristine.
- Provide a easy to use system for used to upload the video and distinguish whether the video is real or fake.

##### 4.1.2 Statement of scope

There are many tools available for creating the deep fakes, but for deep fake detection there is hardly any tool available. Our approach for detecting the deep fakes will be great contribution in avoiding the percolation of the deep fakes over the world wide web. We will be providing a web-based platform for the user to upload the video and classify it as fake or real. This project can be scaled up from developing a web-based platform to a browser plugin for automatic deep fake detection's. Even big application like WhatsApp, Facebook can integrate this project with their application for easy pre-detection of deep fakes before sending to another user. A description of the software with Size of input,

bounds on input, input validation, input dependency, i/o state diagram, Major inputs, and outputs are described without regard to implementation detail.

## 4.2 Major Constraints

- **User:** User of the application will be able to detect whether the uploaded video is fake or real, Along with the model confidence of the prediction.
- **Prediction:** The User will be able to see the playing video with the output on the face along with the confidence of the model.
- **Easy and User-friendly User-Interface:** Users seem to prefer a more simplified process of Deep Fake video detection. Hence, a straight forward and user-friendly interface is implemented. The UI contains a browse tab to select the video for processing. It reduces the complications and at the same time enriches the user experience.
- **Cross-platform compatibility:** with an ever-increasing target market, accessibility should be your main priority. By enabling a cross-platform compatibility feature, you can increase your reach across different platforms. Being a server side application it will run on any device that has a web browser installed in it.

## 4.3 Methodologies of Problem Solving

### 4.3.1 Analysis

- **Solution Requirement**

We analysed the problem statement and found the feasibility of the solution of the problem. We read different research papers as mentioned in 3.3. After checking the feasibility of the problem statement. The next step is the data-set gathering and analysis. We analysed the data set in different approaches of training like negatively or positively trained i.e. training the model with only fake or real videos but found that it may lead to addition of extra bias in the model leading to inaccurate predictions. So after doing a lot of research we found that the balanced training of the algorithm is the best way to avoid the bias and variance in the algorithm and get a good accuracy.

- **Solution Constraints**

We analysed the solution in terms of cost, speed of processing, requirements, level of expertise, availability of equipment's.

- **Parameter Identified**

1. Blinking of eyes
2. Teeth enchantment
3. Bigger distance for eyes
4. Moustaches
5. Double edges, eyes, ears, nose
6. Iris segmentation
7. Wrinkles on face
8. Inconsistent head pose
9. Face angle
10. Skin tone
11. Facial Expressions
12. Lighting
13. Different Pose
14. Double chins
15. Hairstyle
16. Higher cheek bones

### **4.3.2 Design**

After research and analysis we developed the system architecture of the solution as mentioned in the Chapter 6. We decided the baseline architecture of the Model which includes the different layers and their numbers.

### **4.3.3 Development**

After analysis we decided to use the PyTorch framework along with python3 language for programming. PyTorch is chosen as it has good support to CUDA i.e Graphic Processing Unit (GPU) and it is customizable. Google Cloud Platform for training the final model on large number of data-set.

#### 4.3.4 Evaluation

We evaluated our model with a large number of real time dataset which include YouTube videos dataset. Confusion Matrix approach is used to evaluate the accuracy of the trained model.

#### 4.4 Outcome

The outcome of the solution is trained deepfake detection models that will help the users to check if the new video is deepfake or real.

#### 4.5 Applications

Web based application will be used by the user to upload the video and submit the video for processing. The model will pre-process the video and predict whether the uploaded video is a deepfake or real video.

#### 4.6 Hardware Resources Required

In this project, a computer with sufficient processing power is needed. This project requires too much processing power, due to the image and video batch processing.

- **Client-side Requirements:** Browser: Any Compatible browser device

**Table 4.1:** Hardware Requirements

Sr. No.	Parameter	Minimum Requirement
1	Intel Xeon E5 2637	3.5 GHz
2	RAM	16 GB
3	Hard Disk	100 GB
4	Graphic card	NVIDIA GeForce GTX Titan (12 GB RAM)

## **4.7 Software Resources Required**

Platform :

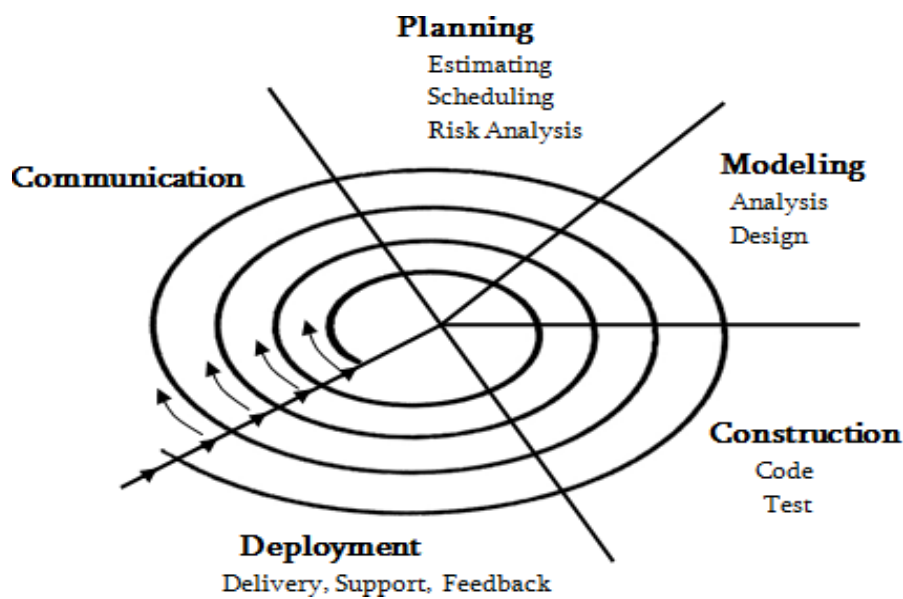
1. Operating System: Windows 7+
2. Programming Language : Python 3.0
3. Framework: PyTorch 1.4 , Django 3.0
4. Cloud platform: Google Cloud Platform
5. Libraries : OpenCV, Face-recognition

## CHAPTER 5

### PROJECT PLAN

#### 5.1 Project Model Analysis

We Use Spiral model As the Software development model focuses on the people doing the work, how they work together and risk handling. We are using Spiral because, It ensures changes can be made quicker and throughout the development process by having consistent evaluations to assess the product with the expected outcomes re-requested. As we developed the application in the various modules, spiral model is best suited for this type of application. An Spiral approach provides a unique opportunity for clients to be involved throughout the project, from prioritizing features to iteration planning and review sessions to frequent algorithms containing new features. However, this also requires clients to understand that they are seeing a work in progress in exchange for this added benefit of transparency. As our model consistsof lot of risk and spiral model is capable of handling the risks that's the reason weare using spiral model for product development.



**Figure 5.1: Spiral Methodology SDLC**

### 5.1.1 Reconciled Estimates

#### 1. Cost Estimate : Rs 11,600

**Table 5.1:** Cost Estimation

Cost(in Rs)	Description
5260	Pre-processing the dataset on GCP
2578	Training models on on GCP
761	Google Colab Pro subscription
3000	Deploying project to GCP using Cloud engine

#### 2. Time Estimates : 12 Months (refer Appendix B)

### 5.1.2 Cost Estimation using COCOMO(Constructive Cost) Model

Since we have small team , less-rigid requirements, long deadline we are using theorganic COCOMO[23] model.

1. **Efforts Applied:** It defines the Amount of labor that will be required to complete a task. It is measured in person-months units.

$$EffortApplied(E) = a_b(KLOC)^{b_b}E = 2.4(20.5)^{1.05}$$

$$E = 57.2206PM$$

2. **Development Time:** Simply means the amount of time required for the completion of the job, which is, of course, proportional to the effort put. It is measured in the units of time such as weeks, months.

$$DevelopmentTime(D) = c_b(E)^{d_b}D = 2.5(57.2206)^{0.38}$$

$$D = 11.6M$$

3. **People Required:** The number of developed needed to complete the project.

$$PeopleRequired(P) = \frac{E}{D}$$

$$P = \frac{57.2206}{11.6}$$

$$P = 4.93$$

## 5.2 Risk Management w.r.t. NP Hard analysis

### 5.2.1 Risk Identification

Before the training, we need to prepare thousands of images for both persons. We can take a shortcut and use a face detection library to scrape facial pictures from their videos. Spend significant time to improve the quality of your facial pictures. It impacts your final result significantly.

1. Remove any picture frames that contain more than one person.
2. Make sure you have an abundance of video footage. Extract facial pictures contain different pose, face angle, and facial expressions.
3. Some resembling of both persons may help, like similar face shape.

### 5.2.2 Risk Analysis

In Deepfakes, it creates a mask on the created face so it can blend in with the target video. To further eliminate the artifacts

1. Apply a Gaussian filter to further diffuse the mask boundary area.
2. Configure the application to expand or contract the mask further.
3. Control the shape of the mask.



**Table 5.2:** Risk Description

ID	Risk Description	Probability	Impact		
			Schedule	Quality	Overall
1	Does it over blur comparing with other non-facial areas of the video?	Low	Low	High	High
2	Does it flick?	High	Low	High	High
3	Does it have a change of skin tone near the edge of face?	Low	High	High	Low
4	Does it have a double chin, double eyebrows, double edges on the face?	High	Low	High	Low
5	When the face is partially blocked by hands or other things, does it flick or get blurry?	High	High	High	High

**Table 5.3:** Risk Probability definitions

Probability	Value	Description
High	Probability of occurrence is	> 75%
Medium	Probability of occurrence is	26 – 75%
Low	Probability of occurrence is	< 25%

## 5.3 Project Schedule

### 5.3.1 Project task set

Major Tasks in the Project stages are

- Task 1: Data-set gathering and analysis

This task consists of downloading the dataset. Analysing the dataset and making the dataset ready for the preprocessing.

- Task 2 : Module 1 implementation

Module 1 implementation consists of splitting the video to frames and cropping each frame consisting of face.

- Task 3: Pre-processing

Pre-processing includes the creation of the new dataset which includes only face cropped videos.

- Task 4: Module 2 implementation

Module 2 implementation consists of implementation of DataLoader for loading the video and labels. Training a base line model on small amount of data.

- Task 5 : Hyper parameter tuning

This task includes the changing of the Learning rate, batch size, weight decay and model architecture until the maximum accuracy is achieved.

- Task 6 : Training the final model

The final model on large dataset is trained based on the best hyper parameter identified in the Task 5.

- Task 7 : Front end Development

This task includes the front end development and integration of the back-end and front-end.

- Task 8 : Testing

The complete application is tested using unit testing,

### **5.3.2 Timeline chart**

Please refer Annex C for the planner

## CHAPTER 6

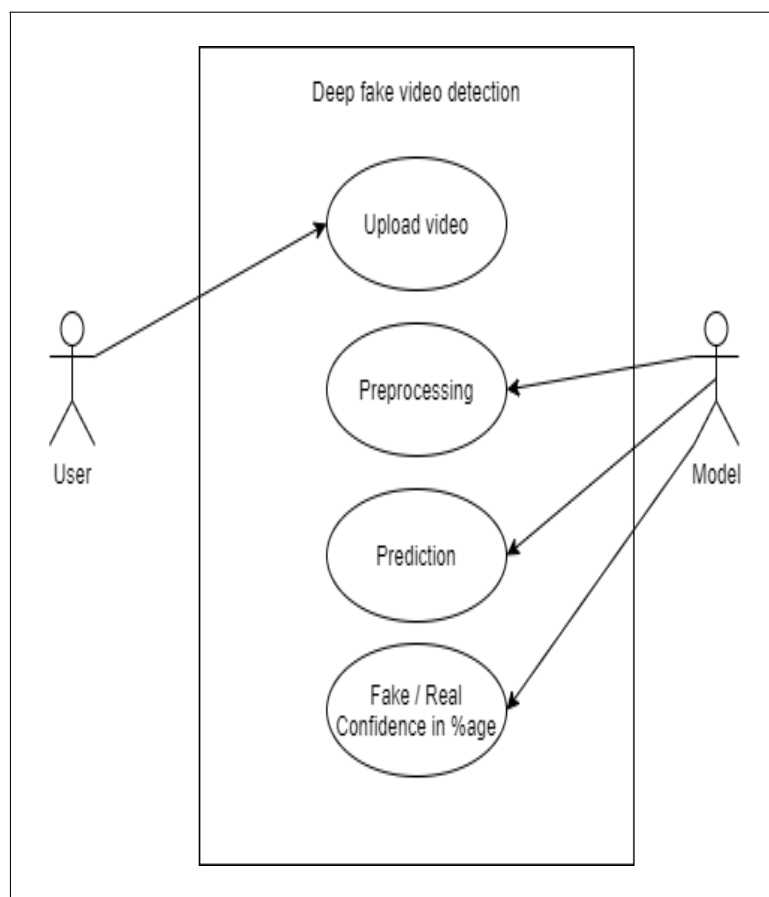
### SYSTEM ANALYSIS

#### 6.1 Introduction

##### 6.1.1 Purpose and Scope of Document

This document lays out a project plan for the development of Deepfake video de- detection using neural network. The intended readers of this document are current and future developers working on Deepfake video detection using neural network and the sponsors of the project. The plan will include, but is not restricted to, a summary of the system functionality, the scope of the project from the perspective of the “Deepfake video detection” team (me and my mentors), use case diagram, Data flow diagram, activity diagram, functional and non-functional requirements, project risks and how those risks will be mitigated, the process by which we will develop the project, and metrics and measurements that will be recorded throughout the project.

##### 6.1.2 Use Case View



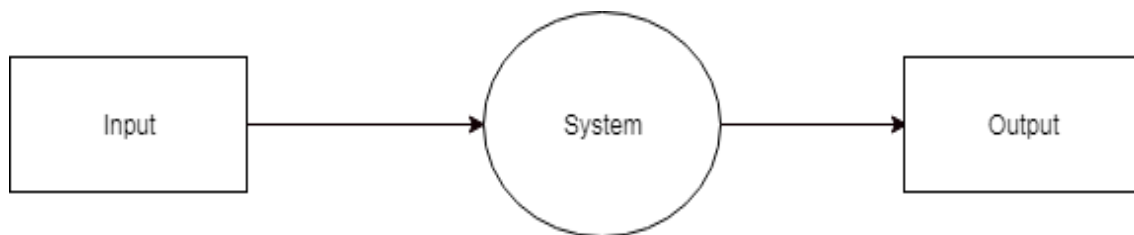
**Figure 6.1: Use case diagram**

## 6.2 Functional Model and Description

A description of each major software function, along with data flow (structured analysis) or class hierarchy (Analysis Class diagram with class description for object oriented system) is presented.

### 6.2.1 Data Flow Diagram

#### DFD Level-0



**Figure 6.2: DFD Level 0**

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.

### DFD Level-1

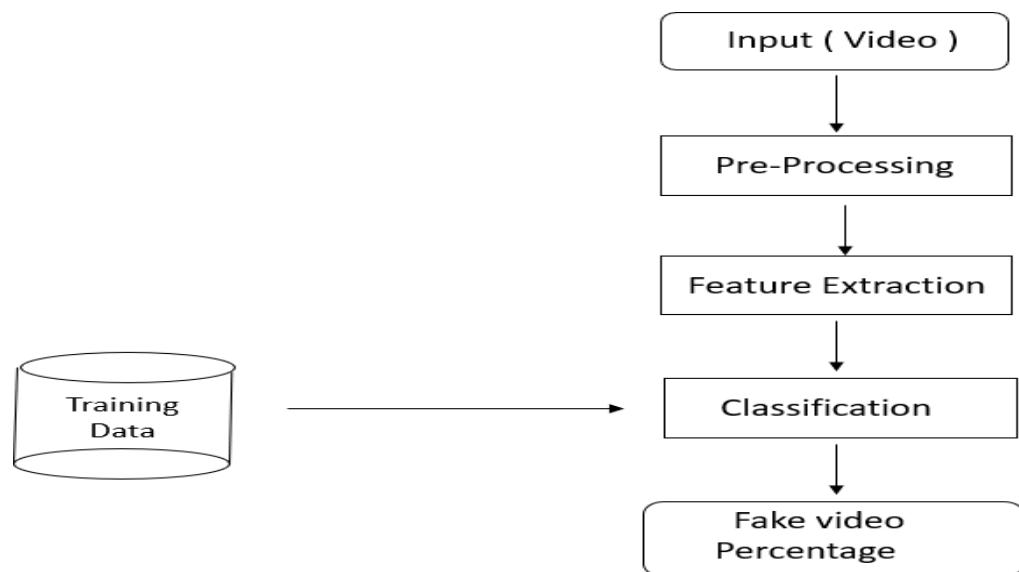


Figure 6.3: DFD Level 1

### DFD Level-2

[1] DFD level-2 enhances the functionality used by user etc.

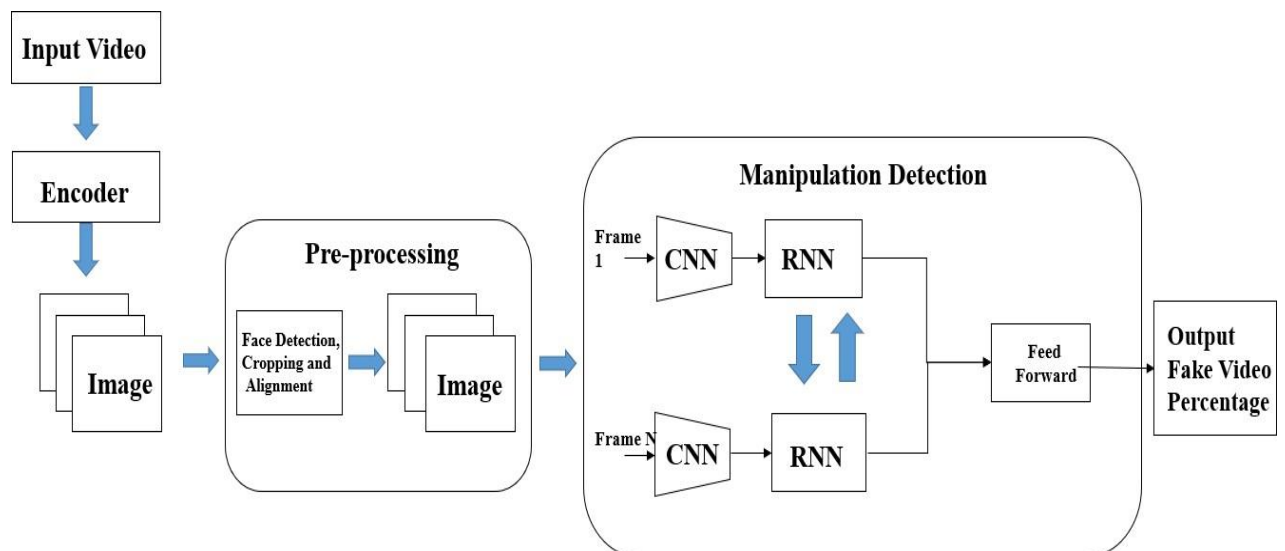


Figure 6.4: DFD Level 2

### 6.2.2 Activity Diagram:

#### Training Workflow:

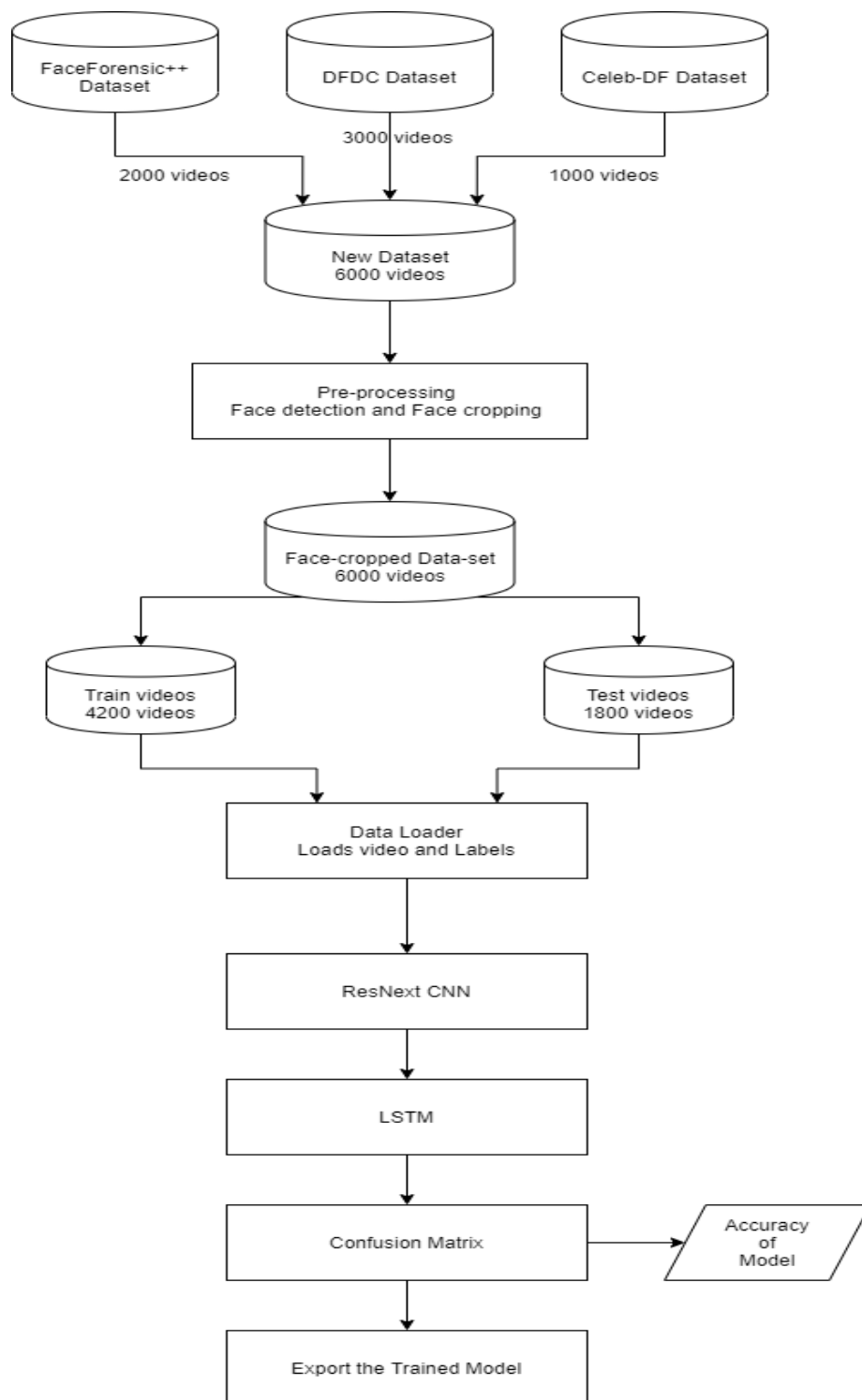
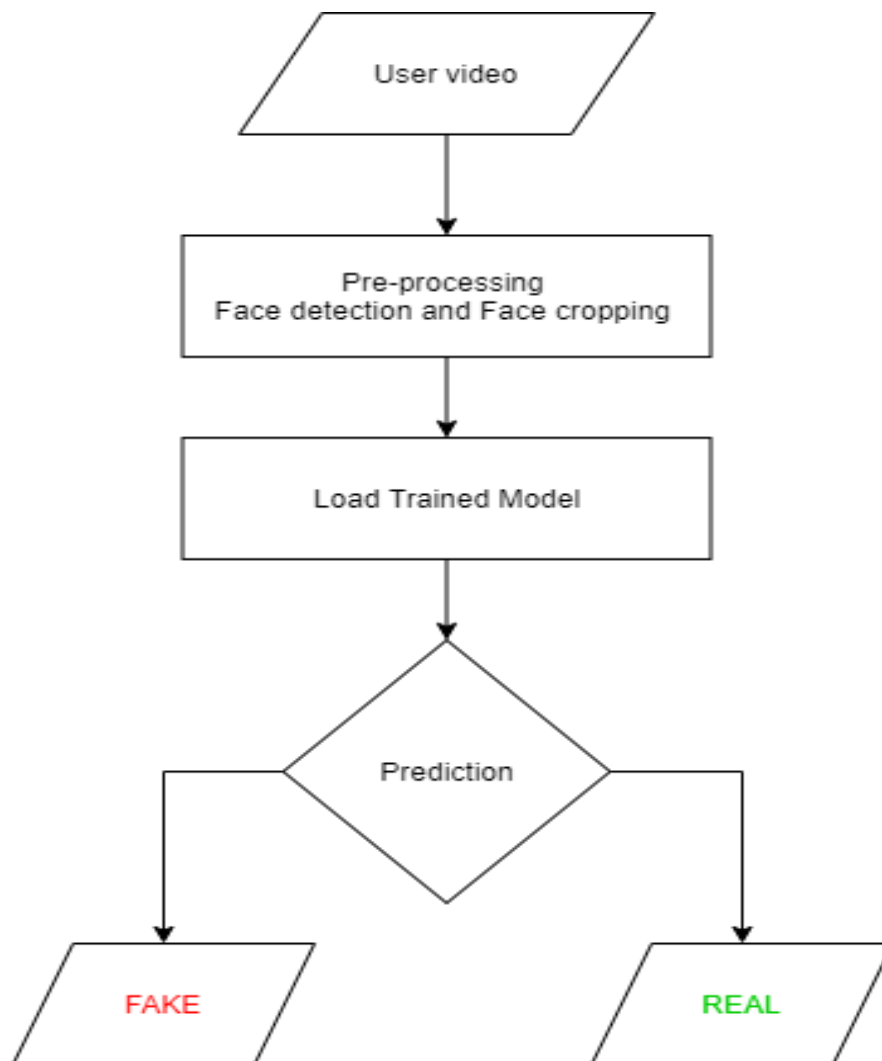


Figure 6.5: Training Workflow

**Testing Workflow:****Figure 6.6: Testing Workflow****6.2.3 Non Functional Requirements:****Performance Requirement**

- The software should be efficiently designed so as to give reliable recognition of fake videos and so that it can be used for more pragmatic purpose.
- The design is versatile and user friendly.
- The application is fast, reliable and time saving.

- The system have universal adaptations.
- The system is compatible with future upgradation and easy integration.

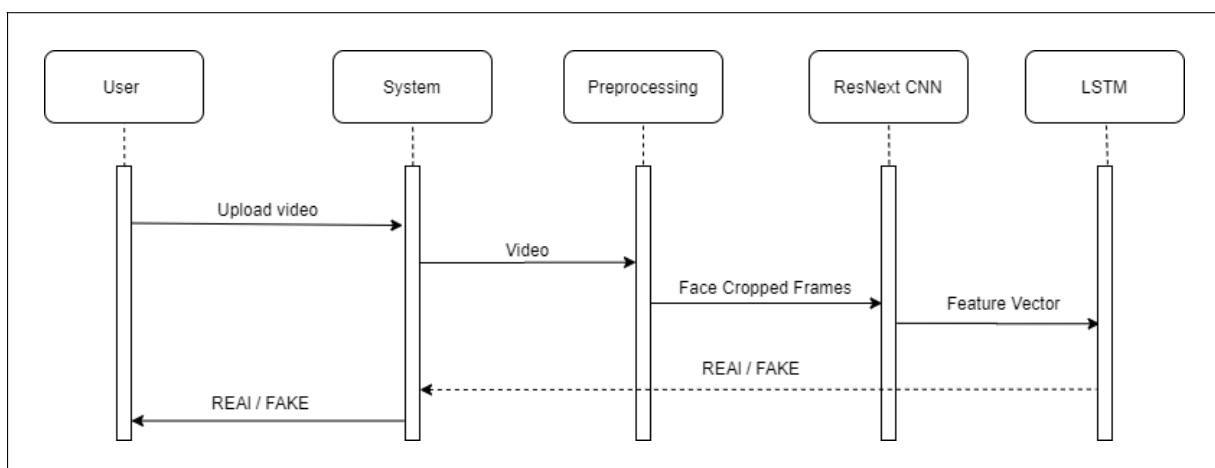
#### Safety Requirement

- The Data integrity is preserved. Once the video is uploaded to the system. It is only processed by the algorithm. The videos are kept secured from the human interventions, as the uploaded video is not are not able for human manipulation.
- To extent the safety of the videos uploaded by the user will be deleted after 30 min from the server.

#### Security Requirement

- While uploading the video, the video will be encrypted using a certain symmet- ric encryption algorithm. On server also the video is in encrypted format only. The video is only decrypted from preprocessing till we get the output. After getting the output the video is again encrypted.
- This cryptography will help in maintain the security and integrity of the video.
- SSL certification is made mandatory for Data security.

#### 6.2.4 Sequence Diagram



**Figure 6.7: Sequence Diagram**

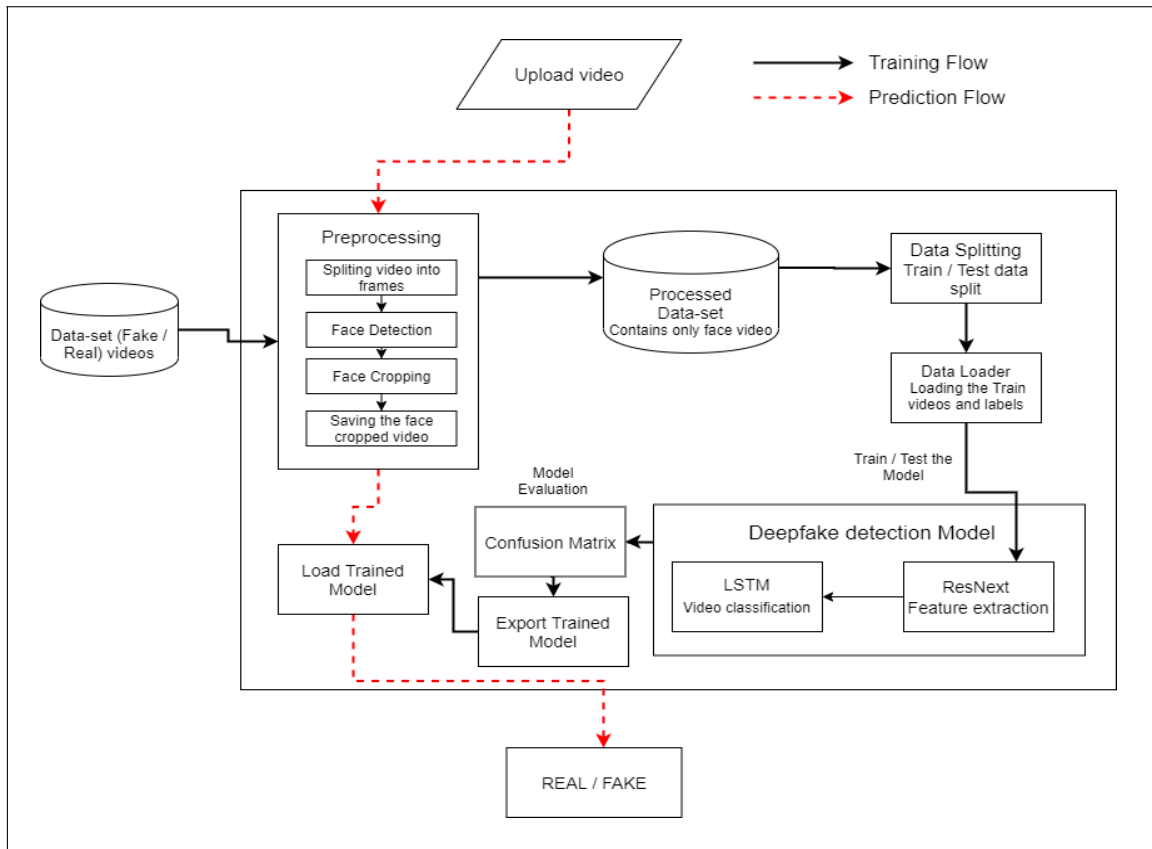


## CHAPTER 7

### DETAILED DESIGN DOCUMENT

#### 7.1 Introduction

##### 7.1.1 System Architecture



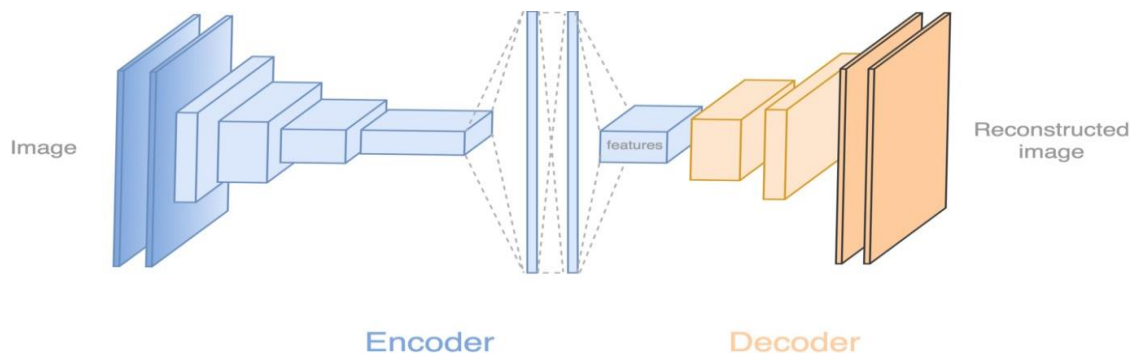
**Figure 7.1: System Architecture**

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

taken a dataset, preprocessed the dataset and created a new processed dataset which only includes the face cropped videos.

## • Creating deepfake videos

To detect the deepfake videos it is very important to understand the creation process of the deepfake. Majority of the tools including the GAN and autoencoders takes a source image and target video as input. These tools split the video into frames, detect the face in the video and replace the source face with target face on each frame. Then the replaced frames are then combined using different pre-trained models. These models also enhance the quality of video by removing the left-over traces by the deepfake creation model. Which result in creation of a deepfake looks realistic in nature. We have also used the same approach to detect the deepfakes. Deepfakes created using the pretrained neural networks models are very realistic that it is almost impossible to spot the difference by the naked eyes. But in reality, the deepfakes creation tools leaves some of the traces or artifacts in the video which may not be noticeable by the naked eyes. The motive of this paper to identify these unnoticeable traces and distinguishable artifacts of these videos and classified it as deepfake or real video.



**Figure 7.2: Deepfake generation**



**Figure 7.3: Face Swapped deepfake generation**

### **Tools for deep fake creation.**

1. Faceswap
2. Faceit
3. Deep Face Lab
4. Deepfake Capsule GAN
5. Large resolution face masked

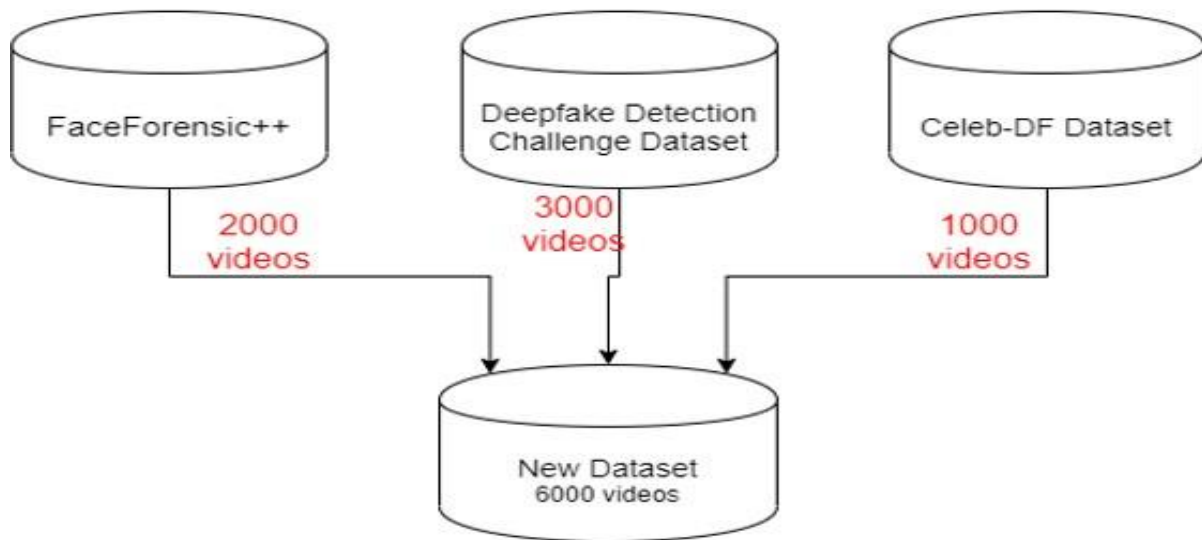
## **7.2 Architectural Design**

### **7.2.1 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)[1], Deepfake detection challenge(DFDC)[2], and Celeb-DF[3]. 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 [3] consist of certain audioaltered 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)[1] dataset and 500 Real and 500 Fake videos from the Celeb-DF[3] dataset. Which makes our total dataset consisting 3000 Real, 3000 fake videos and 6000 videos in total. Figure 2 depicts the distribution of the data-sets.

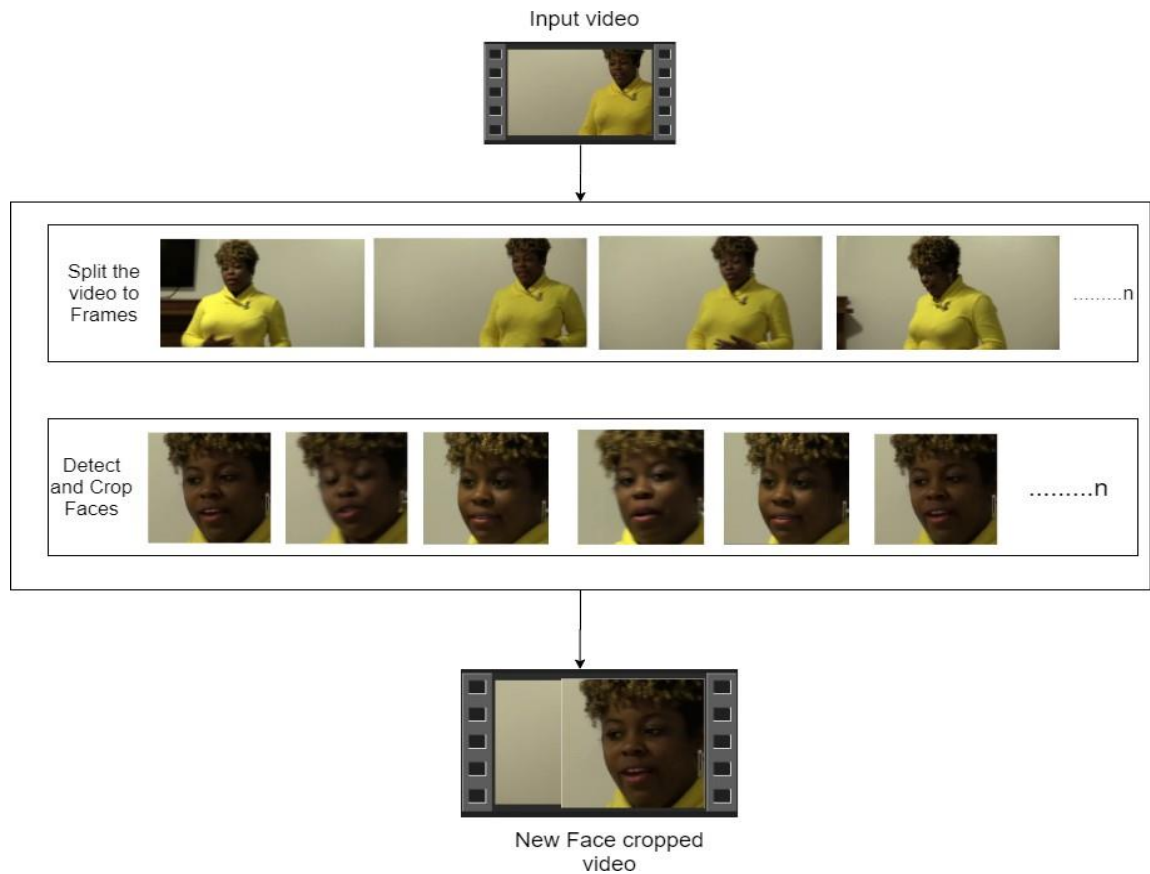


**Figure 7.4: Dataset**

### 7.2.2 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.

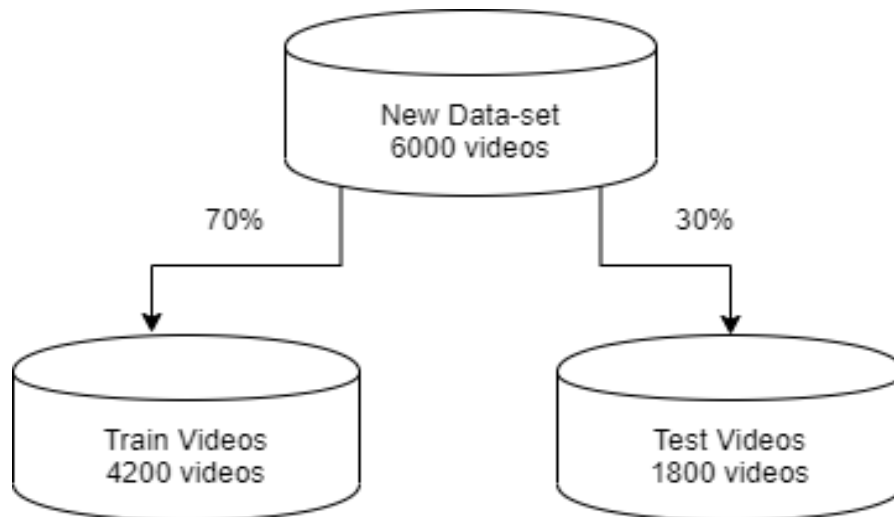
To maintain the uniformity of number of frames, we have selected a threshold value based on the mean of total frames count of each video. Another reason for selecting a threshold value is limited computation power. As a video of 10 second at 30 frames per second (fps) will have total 300 frames and it 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 30 fps and resolution of 112 x 112.



**Figure 7.5: Pre-processing of video**

### 7.2.3 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 in each split.



**Figure 7.6: Train test split**

#### 7.2.4 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.

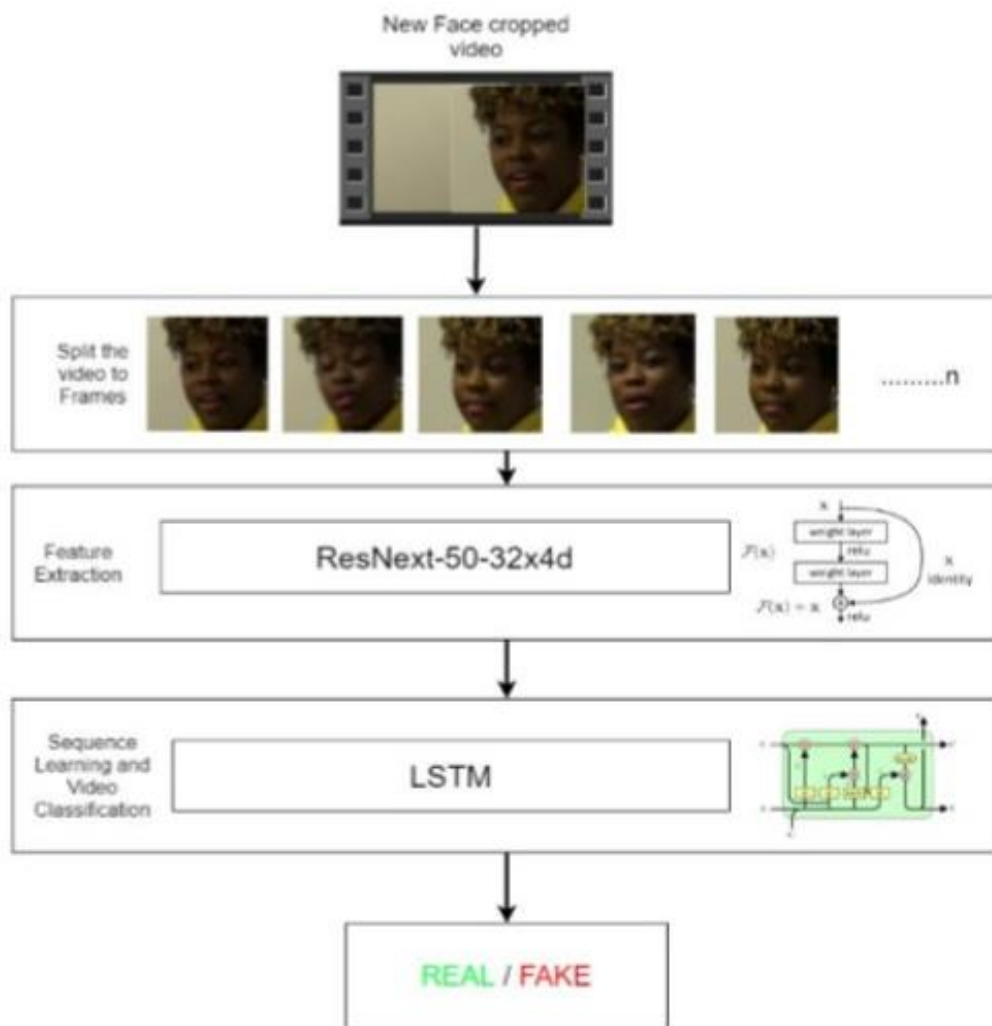
##### **LSTM for Sequence Processing:**

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.



**Figure 7.7: Overview of our model**

### 7.2.5 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[21] 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.



## CHAPTER 8

### PROJECT IMPLEMENTATION

#### 8.1 Introduction

There are many examples where deepfake creation technology is used to mis-lead the people on social media platform by sharing the false deepfake videos of the famous personalities like Mark Zuckerberg Eve of House A.I. Hearing, Don-ald Trump's Breaking Bad series where he was introduces as James McGill, Barack Obama's public service announcement and many more [5]. These types of deepfakes creates a huge panic among the normal people, which arises the need to spot these deepfakes accurately so that they can be distinguished from the real videos.

Latest advances in the technology have changed the field of video manipulation. The advances in the modern open source deep learning frameworks like TensorFlow, Keras, PyTorch along with cheap access to the high computation power has driven the paradigm shift. The Conventional autoencoders[10] and Generative Adversarial Network (GAN) pretrained models have made the tampering of the realistic videos and images very easy. Moreover, access to these pretrained models through the smartphones and desktop applications like FaceApp and Face Swap has made the deepfake creation a childish thing. These applications generate a highly realistic synthesized transformation of faces in real videos. These apps also provide the user with more functionalities like changing the face hair style, gender, age and other attributes. These apps also allow the user to create a very high quality and indistin- guishable deepfakes. Although some malignant deepfake videos exist, but till now they remain a minority. So far, the released tools [11,12] that generate deepfake videos are being extensively used to create fake celebrity pornographic videos or revenge porn [13]. Some of the examples are Brad Pitt, Angelina Jolie nude videos. The real looking nature of the deepfake videos makes the celebrities and other fa- mous personalities the target of pornographic material, fake surveillance videos, fake

news and malicious hoaxes. The Deepfakes are very much popular in creating the political tension [14]. Due to which it becomes very important to detect the deepfake videos and avoid the percolation of the deepfakes on the social media platforms.

## **8.2 Tools and Technologies Used**

### **8.2.1 Planning**

1. OpenProject

### **8.2.2 UML Tools**

1. draw.io

### **8.2.3 Programming Languages**

1. Python3
2. JavaScript

### **8.2.4 Programming Frameworks**

1. PyTorch
2. Django

### **8.2.5 IDE**

1. Google Colab
2. Jupyter Notebook
3. Visual Studio Code

### **8.2.6 Versioning Control**

1. Git

### **8.2.7 Cloud Services**

1. Google Cloud Platform

### **8.2.8 Application and web servers:**

1. Google Cloud Engine

### **8.2.9 Libraries**

1. torch
2. torchvision
3. os
4. numpy
5. cv2
6. matplotlib
7. face\_recognition
8. json
9. pandas
10. copy
11. glob
12. random
13. sklearn

## **8.3 Algorithm Details**

### **8.3.1 Dataset Details**

Refer 7.2.1

### **8.3.2 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.

### 8.3.3 Model Details

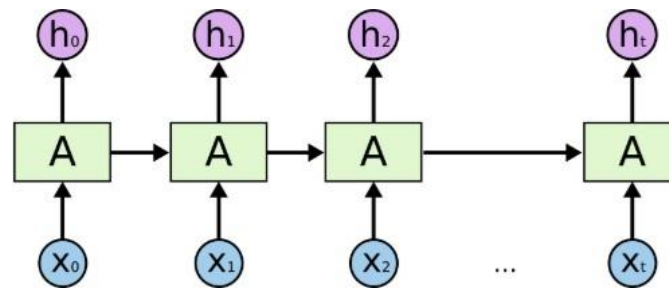
The model consists of following layers:

- **ResNext CNN** : The pre-trained model of Residual Convolution Neural Network is used. The model name is resnext50\_32x4d()[22]. This model consists of 50 layers and 32 x 4 dimensions. Figure shows the detailed implementation of model.

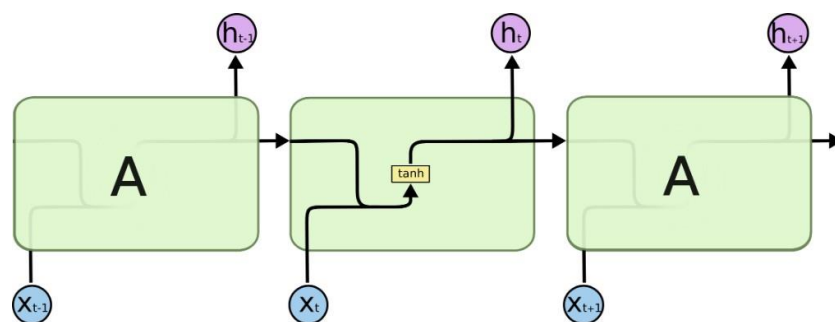
stage	output	<b>ResNeXt-50 (32×4d)</b>
conv1	112×112	7×7, 64, stride 2
		3×3 max pool, stride 2
conv2	56×56	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128, C=32 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3	28×28	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256, C=32 \\ 1 \times 1, 512 \end{bmatrix} \times 4$
conv4	14×14	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512, C=32 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$
conv5	7×7	$\begin{bmatrix} 1 \times 1, 1024 \\ 3 \times 3, 1024, C=32 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1	global average pool 1000-d fc, softmax
# params.		<b>25.0×10<sup>6</sup></b>

**Figure 8.1: ResNext Architecture**



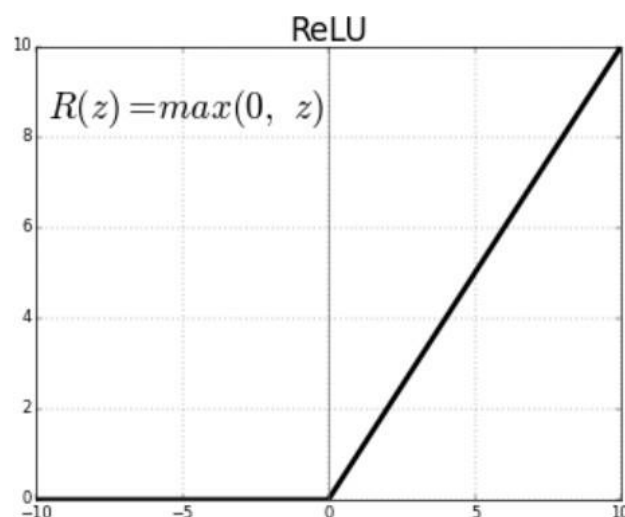


**Figure 8.4: Overview of LSTM Architecture**



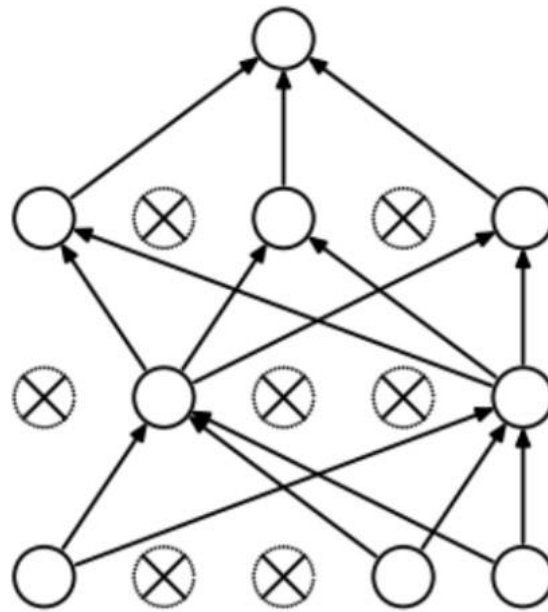
**Figure 8.5: Internal LSTM Architecture**

- **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.



**Figure 8.6: Relu Activation function**

- **Dropout Layer** :Dropout layer with the value of 0.4 is used to avoid over-fitting 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.



**Figure 8.7: Dropout layer overview**

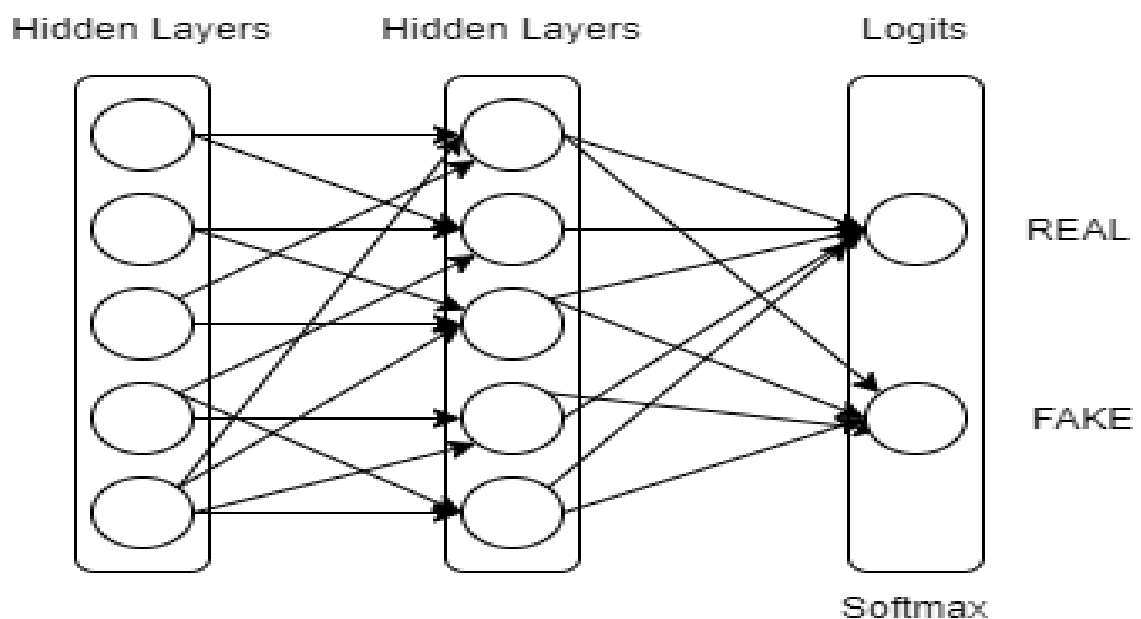
- **Adaptive Average Pooling Layer** : It is used To reduce variance, reduce computation complexity and extract low level features from neighbourhood.2 dimensional Adaptive Average Pooling Layer is used in the model.

#### 8.3.4 Model Training Details

- **Train Test 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 in each split. Refer figure 7.6
- **Data Loader**: It is used to load the videos and their labels with a batch size of 4.

- **Training:** The training is done for 20 epochs with a learning rate of  $1e-5$  (0.00001), weight decay of  $1e-3$  (0.001) using the Adam optimizer.
- **Adam optimizer[21]:** To enable the adaptive learning rate Adam optimizer with the model parameters is used.
- **Cross Entropy:** To calculate the loss function Cross Entropy approach is used because we are training a classification problem.
- **Softmax Layer:** A Softmax function is a type of squashing function. Squashing functions limit the output of the function into the range 0 to 1. This allows the output to be interpreted directly as a probability. Similarly, softmax functions are multi-class sigmoids, meaning they are used in determining probability of multiple classes at once. Since the outputs of a softmax function can be interpreted as a probability (i.e. they must sum to 1), a softmax layer is typically the final layer used in neural network functions. It is important to note that a softmax layer must have the same number of nodes as the output layer.

In our case softmax layer has two output nodes i.e REAL or FAKE, also Softmax layer provide us the confidence(probability) of prediction.



**Figure 8.8: Softmax Layer**



- **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.

- **Export Model:** After the model is trained, we have exported the model. So that it can be used for prediction on real time data.

### 8.3.5 Model Prediction Details

- The model is loaded in the application
- The new video for prediction is preprocessed (refer 8.3.2, 7.2.2) and passed to the loaded model for prediction
- The trained model performs the prediction and return if the video is a real or fake along with the confidence of the prediction.

## **CHAPTER 9**

### **SOFTWARE TESTING**

#### **9.1 Type of Testing Used**

##### **Functional Testing**

1. Unit Testing
2. Integration Testing
3. System Testing
4. Interface Testing

##### **Non-functional Testing**

1. Performance Testing
2. Load Testing
3. Compatibility Testing

## 9.2 Test Cases and Test Results

### Test Cases

**Table 9.1:** Test Case Report

Case id	Test Case Description	Expected Result	Actual Result	Status
1	Upload a word file in-stead of video	Error message: Onlyvideo files allowed	Error message: Onlyvideo files allowed	Pass
2	Upload a 200MB video file	Error message: Maxlimit 100MB	Error message: Maxlimit 100MB	Pass
3	Upload a file without any faces	Error message: No faces detected. Cannot process the video.	Errormessage:No facesdetected. Cannot process video.	Pass
4	Videos with many faces	Fake / Real	Fake	Pass
5	Deepfake video	Fake	Fake	Pass
6	Enter /predict in URL	Redirect to /upload	Redirect to /upload	Pass
7	Press upload buttonwithout selecting video	Alert message: Please select video	Alert message: Pleaseselect video	Pass
8	Upload a Real video	Real	Real	Pass
9	Upload a face croppedreal video	Real	Real	Pass
10	Upload a face cropped fake video	Fake	Fake	Pass

## CHAPTER 10

### CONCLUSION AND FUTURE ENHANCEMENTS

#### 10.1 Conclusion

We presented a neural network-based approach to classify the video as deep fake or real, along with the confidence of proposed model. Our method is capable of predicting the output by processing 1 second of video (10 frames per second) with a good accuracy. We implemented the model by using pre-trained ResNext CNN model to extract the frame level features and LSTM for temporal sequence processing to spot the changes between the  $t$  and  $t-1$  frame. Our model can process the video in the frame sequence of 10,20,40,60,80,100.

#### 10.2 Future Enhancements

There is always a scope for enhancements in any developed system, especially when the project build using latest trending technology and has a good scope in future.

- Web based platform can be upscaled to a browser plugin for ease of access to the user.
- Currently only Face Deep Fakes are being detected by the algorithm, but the algorithm can be enhanced in detecting full body deep fakes.



**NoviTech**  
the innovation partner

Date: 18.05.2024

### PROJECT COMPLETION CERTIFICATE

This is to inform you that **Mr. Praveen Kumar K (Reg. No. 717822Z135)** doing **MASTER OF COMPUTER APPLICATIONS** in Karpagam College of Engineering, Coimbatore is successfully completed his project work entitled **“REVEALING AND CLASSIFICATION OF DEEP FAKE IMAGES WITH VIDEOS USING CUSTOMIZED DEEP LEARNING MODELS”** in our concern from January 2024 to April 2024.

During this period, we found him to be sincere in work and regular in his attendance 100%. We are certifying that his conduct and character has been found to be good.

FOR NOVITECH R&D PVT LTD



AUTHORITY SIGNATURE

## CHAPTER 11

### APPENDIX

#### 11.1 SOURCE CODE

```
from google.colab import drive
drive.mount('/content/drive')
#import libraries
!pip3 install face_recognition

import torch
import torchvision
from torchvision import transforms
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

#import libraries
import torch
from torch.autograd import Variable
import time
import os
import sys
import os
from torch import nn
from torchvision import models

#Model with feature visualization
from torch import nn
from torchvision import models
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, :]))
```

```
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))
```

```
def im_convert(tensor):
```

```
    """ 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)
    cv2.imwrite('./2.png', image*255)
    return image
```

```
def predict(model, img, path = './'):
```

```
    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)
confidence = logits[:,int(prediction.item())].item()*100
print('confidence of prediction:',logits[:,int(prediction.item())].item()*100)
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)
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,:,:])
result = heatmap * 0.5 + img*0.8*255
cv2.imwrite('/content/1.png',result)
result1 = heatmap * 0.5/255 + img*0.8
r,g,b = cv2.split(result1)
result1 = cv2.merge((r,g,b))
plt.imshow(result1)
plt.show()
return [int(prediction.item()),confidence]
#img = train_data[100][0].unsqueeze(0)
#predict(model,img)

#!pip3 install face_recognition
import torch
import torchvision
from torchvision import transforms
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

```



```

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
        #print("no of frames",len(frames))
        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_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()

#Code for making prediction
im_size = 112
mean=[0.485, 0.456, 0.406]
std=[0.229, 0.224, 0.225]

train_transforms = transforms.Compose([
    transforms.ToPILImage(),
    transforms.Resize((im_size,im_size)),
    transforms.ToTensor(),
    transforms.Normalize(mean,std)])

path_to_videos = ['/content/drive/My Drive/Balanced_Face_only_data/aagfhgtpmv.mp4',
                  '/content/drive/My Drive/Balanced_Face_only_data/aczrgyricp.mp4',
                  '/content/drive/My Drive/Balanced_Face_only_data/agdkmztvby.mp4',
                  '/content/drive/My Drive/Balanced_Face_only_data/abarnvbtwb.mp4']

path_to_videos = ['/content/drive/My Drive/Youtube_Face_only_data/000_00ert3.mp4',
                  '/content/drive/My Drive/Youtube_Face_only_data/000.mp4',
                  '/content/drive/My Drive/Youtube_Face_only_data/002_006.mp4',
                  '/content/drive/My Drive/Youtube_Face_only_data/002.mp4'

]

path_to_videos= ["/content/drive/My Drive/DFDC_REAL_Face_only_data/aabqyygbaa.mp4"]

video_dataset = validation_dataset(path_to_videos,sequence_length = 20,transform =
train_transforms)
model = Model(2).cuda()
path_to_model = '/content/drive/My Drive/Models/model_87_acc_20_frames_final_data.pt'
model.load_state_dict(torch.load(path_to_model))
model.eval()
for i in range(0,len(path_to_videos)):
    print(path_to_videos[i])

```

```

prediction = predict(model,video_dataset[i],'./')
if prediction[0] == 1:
    print("REAL")
else:
    print("FAKE")

```

#Optional : If you want to pass full frame for prediction instead of face cropped frame

#code for full frame processing

```

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)):
            frames.append(self.transform(frame))
            if(len(frames) == self.count):
                break
        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

```

## 11.2 Screen shots

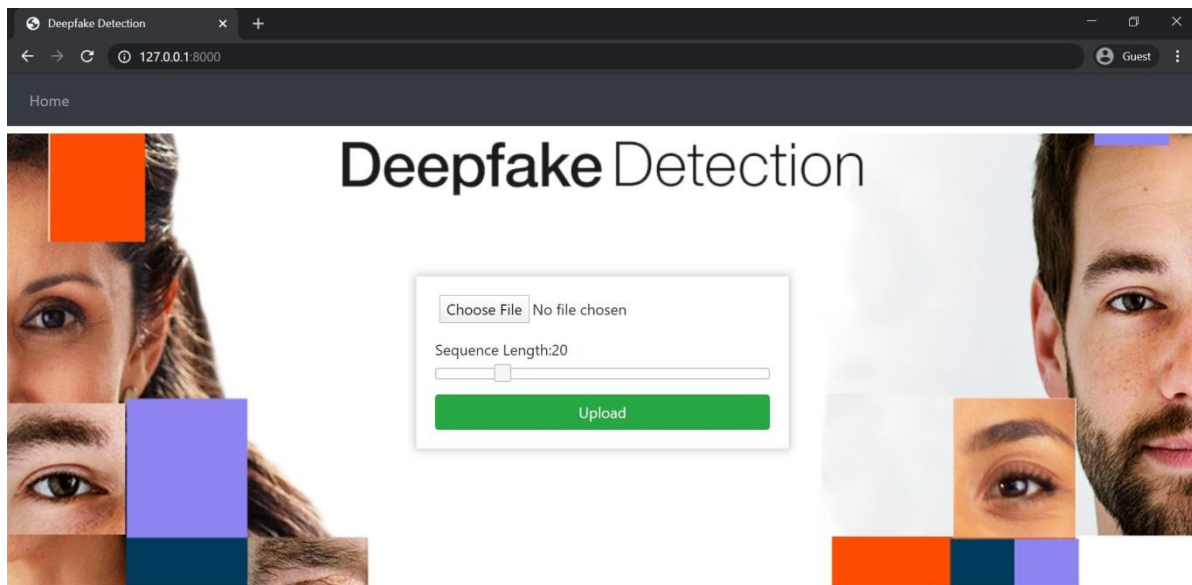


Figure 11.1: Home Page

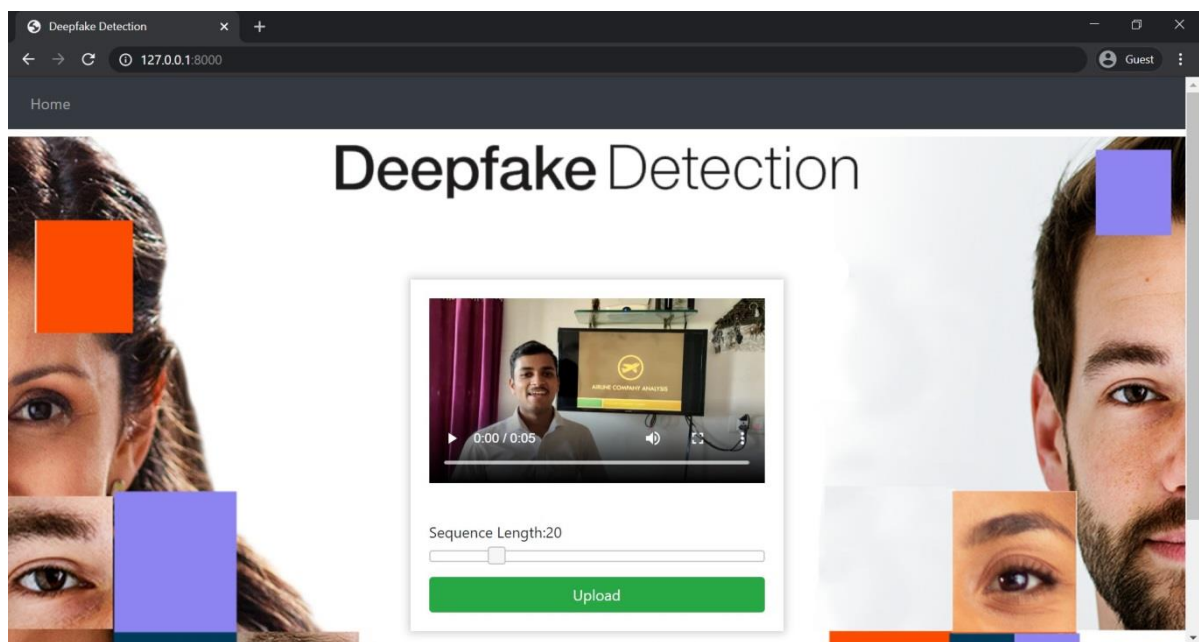


Figure 11.2: Uploading Real Video

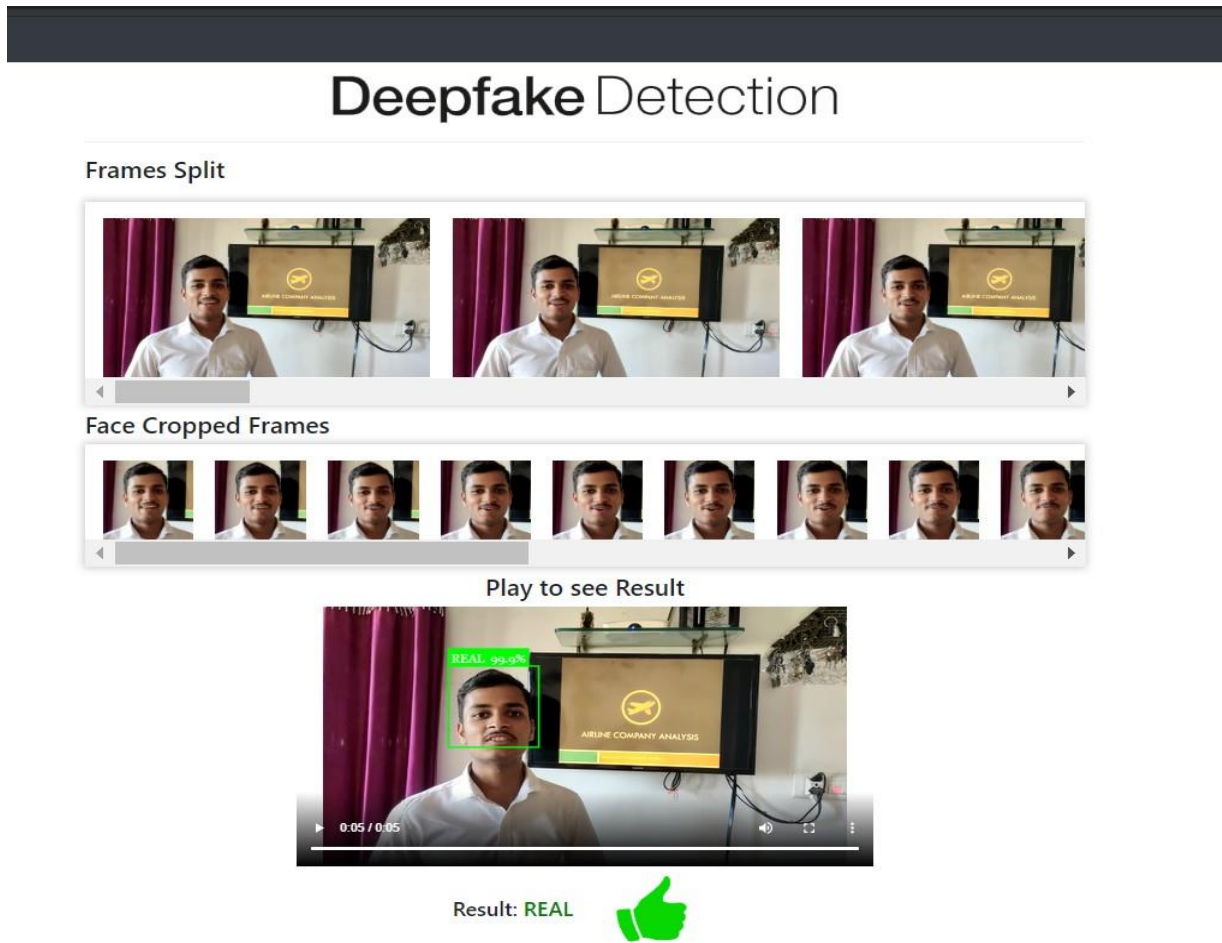


Figure 11.3: Real Video Output

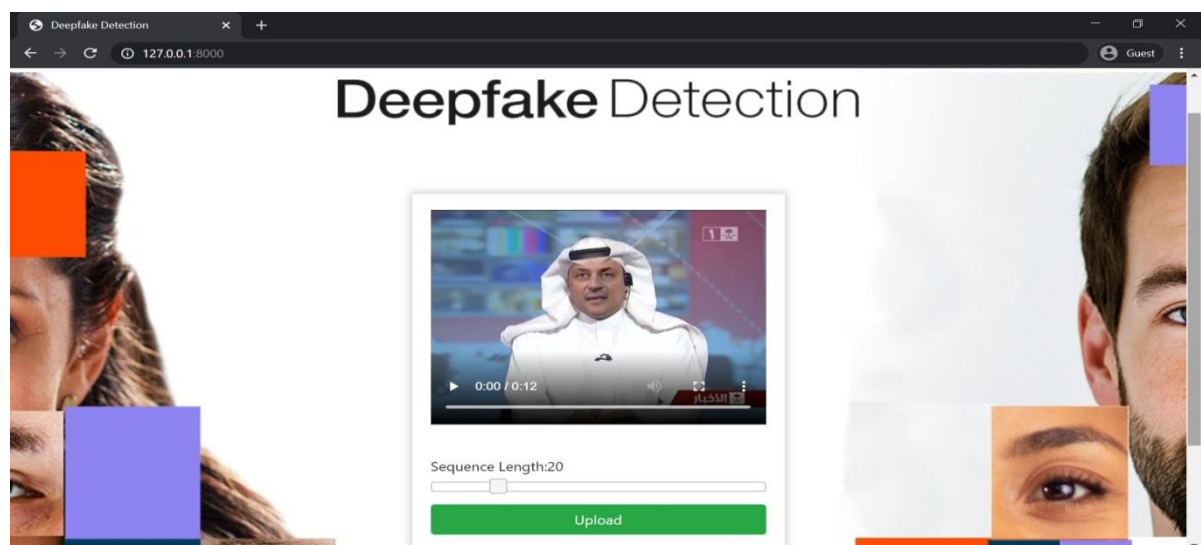


Figure 11.4: Uploading Fake Video

# Deepfake Detection

Frames Split



Face Cropped Frames



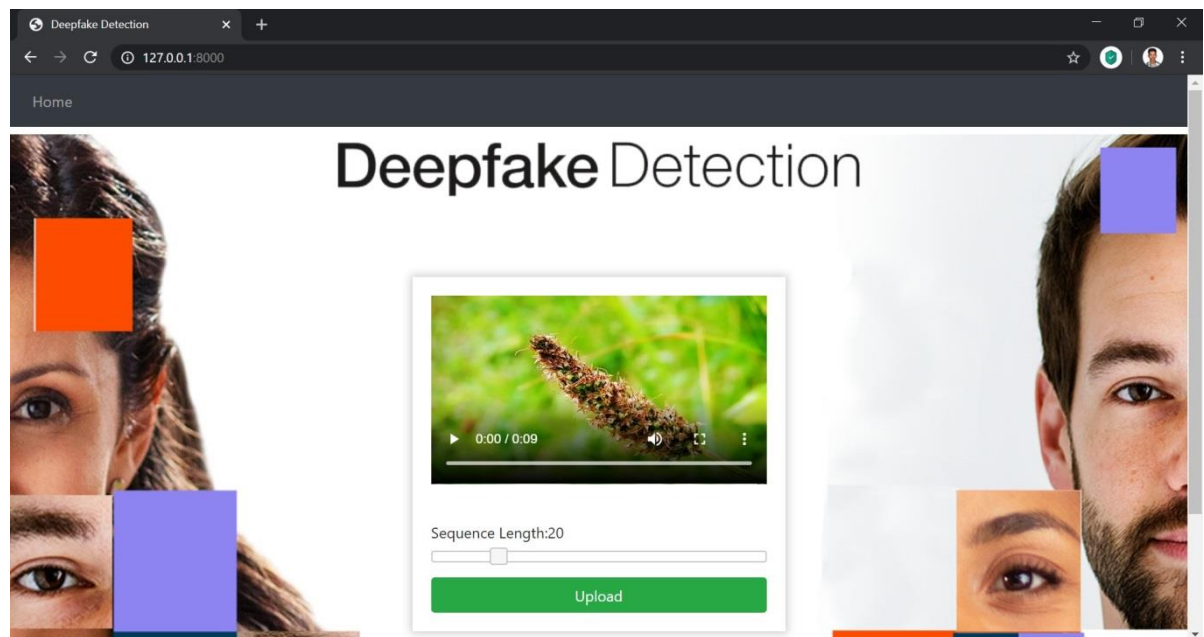
Play to see Result



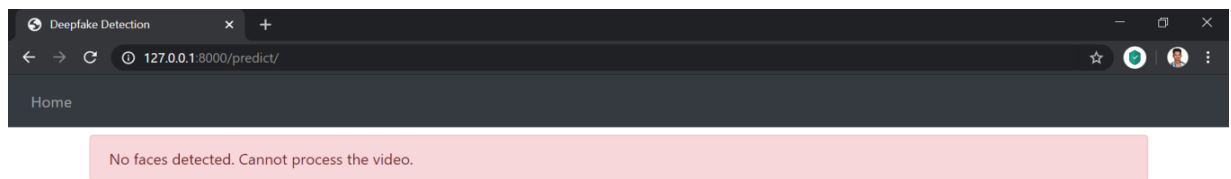
Result: **FAKE**



**Figure 11.5: Fake video Output**

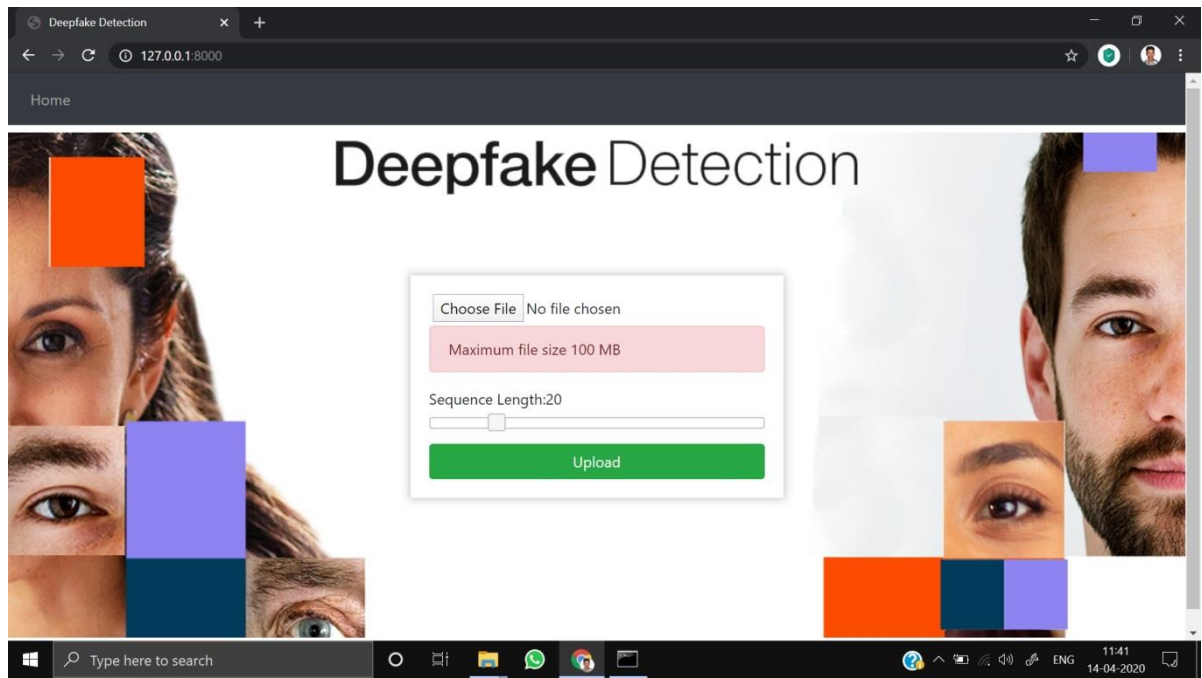


**Figure 11.6: Uploading Video with no faces**

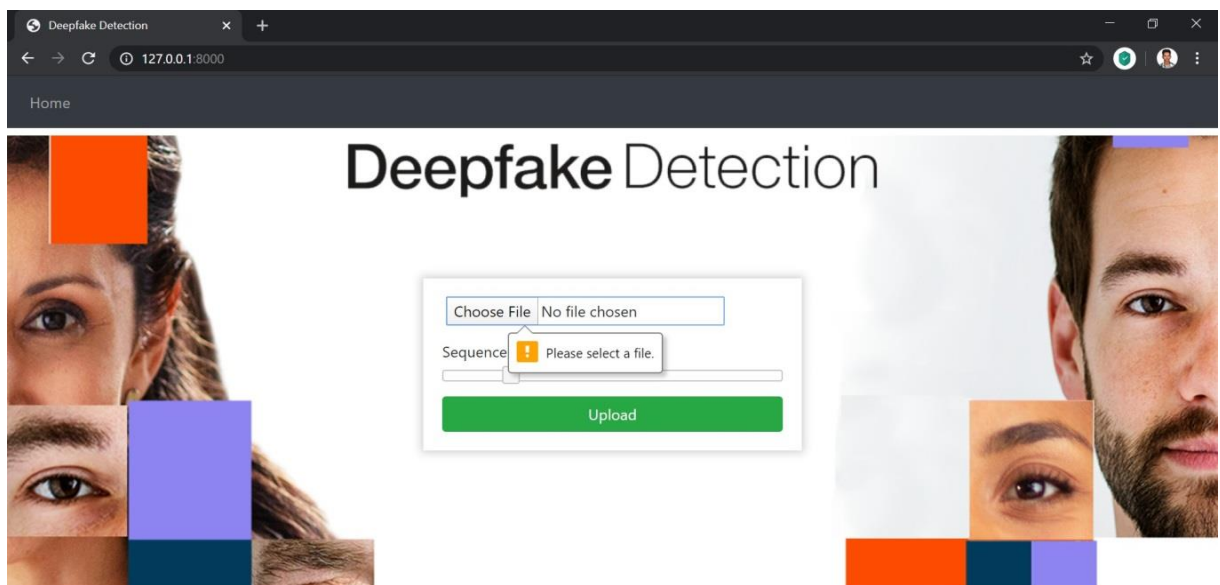


**Figure 11.7: Output of Uploaded video with no faces**





**Figure 11.8: Uploading file greater than 100MB**



**Figure 11.9: Pressing Upload button without selecting video**





12<sup>th</sup> ICCET 2024

*Proceedings of International Conference  
on Contemporary  
Engineering and Technology 2024*

*Chennai, India*

March 23rd - 24th, 2024

ISBN 978-81-965908-5-7

579	ICCET241428	ROBUST IRIS RECOGNITION ALGORITHM FOR IDENTIFICATION
580	ICCET241658	AN ENSEMBLE MODEL FOR DIABETIC PREDICTION AND ANALYSIS
581	ICCET241381	DYNAMIC AD-HOC NETWORK FOR TRAFFIC DETECTION IN 5G NETWORK WITH EFFICIENT LOCALIZATION USING A CLUSTERING ALGORITHM
582	ICCET241719	GRAPHICAL PASSWORD AUTHENTICATION SYSTEM
583	ICCET241408	DETECTION AND CLASSIFICATION OF SKIN DISEASE FROM IMAGES USING CNN WITH LSTM
584	ICCET241741	A DECENTRALIZED ESCROW PROTOCOL THAT FACILITATES SECURE P2P PAYMENTS BETWEEN TRUSTLESS PARTIES
585	ICCET240999	BRAIN TUMOR DETECTION USING DEEP LEARNING (CNN)
586	ICCET241565	ENERGY EFFICIENT SRAM DESIGN WITH IMPROVED STATIC NOISE MARGIN
587	ICCET241539	OBSTACLE DETECTING ROBOT USING ARDUINO
588	ICCET241788	VEHICLE OVERLOAD PREVENTION AND ACCIDENT DETECTION USING IOT SENSORS
589	ICCET241400	EXPERIMENTAL INVESTIGATION OF INTERNAL GRINDING ATTACHMENT IN CENTRE LATHE- A REVIEW
590	ICCET241615	ENHANCED SECURE SYSTEM FOR GROUP SHARING AND EFFECTIVE TRACKING MODEL FOR UNAUTHORIZED ACCESS IN CLOUD COMPUTING
591	ICCET241758	ACTIVE MACHINE LEARNING ADVERSARIAL ATTACK DETECTION IN THE USER FEEDBACK PROCESS
592	ICCET241386	CLASSIFICATION OF SATELLITE IMAGES USING DEEP LEARNING FOR DISASTER MANAGEMENT
593	ICCET241122	MULTILEVEL HYBRID FIRE DETECTION AND PREVENTION USING IOT
594	ICCET241498	AGRINNOVATE – EMPOWERING AGRICULTURE THROUGH DATA DRIVEN PRECISION
595	ICCET241406	REVEALING AND CLASSIFICATION OF DEEP FAKE IMAGES WITH VIDEOS USING CUSTOMIZED DEEP LEARNING MODELS
596	ICCET241522	A COMPREHENSIVE LIFESPAN ANALYSIS OF RECHARGABLE BATTERIES THROUGH HYBRID CNN-LSTM-DNN METHOD
597	ICCET241483	EVALUATING CHATGPT FOR INDIRECT ASSESSMENTS IN HIGHER EDUCATION
598	ICCET241516	STRATEGIC INSIGHT ANALYTICS FOR MACHINE LEARNING DRIVEN CYBER THREAT FORECASTING IN SMART GRID POWER NETWORKS
599	ICCET241471	PREDICTING WIND TURBINE ENERGY FOR ENVIRONMENT RESILIENCE BY HARNESSING DATA FROM IOT THROUGH RNN AND VMD
600	ICCET241116	SAFEGUARDING USER DATA: BLOCKCHAIN AS AN ENABLER OF ADVANCED CONSENT MANAGEMENT SYSTEMS
601	ICCET241760	STUDENT ATTENTIVENESS FOR THE DEVELOPMENT OF PERSONALIZED LEARNING SYSTEMS
602	ICCET241368	A DEEP LEARNING-POWERED EMOTION-DRIVEN MUSIC RECOMMENDER USING AUDIO DATASET TESS
603	ICCET241562	GUARDIANEYE: A MULTIFACETED APPROACH TO REAL-TIME ONLINE PROCTORING WITH GAZE TRACKING AND FACIAL ASPECT RATIO ANALYSIS
604	ICCET241515	FASTEST ONLINE VERIFICATION SERVICE E-KYC USING BLOCKCHAIN

### **595. REVEALING AND CLASSIFICATION OF DEEP FAKE IMAGES WITH VIDEOS USING CUSTOMIZED DEEP LEARNING MODELS**

K. Praveen Kumar  
School of Computer Application  
Karpagam college of Engineering  
R. Ramprashath  
School of Computer Application  
Karpagam College of Engineering

Deep fakes are becoming more common; they include editing previously published films and photos to produce content that appears authentic but is wholly fake. The development process has been considerably expedited by the widespread availability of deep learning techniques, such as autoencoders, Generative Adversarial Networks (GANs), and user-friendly software. These sophisticated algorithms adeptly fuse and modify visual and audio elements, facilitating the production of content that closely mimics genuine footage, even for those without specialized knowledge. The malicious manipulation of images and videos poses significant security and societal concerns. With an emphasis on facial alteration, the goal of this research is to create a deep learning perfect for the detection and classification of deepfake images and videos. The dataset used for the project is either Face Forensics++, Celeb-DF, or the Deepfake Detection Challenge Dataset (DFDC), available on Kaggle, consisting of real and deepfake images and videos. By utilising Recurrent and Convolutional Neural Networks, we have made development in DF detection. Commencing with preprocessing the data, extracting frames from the videos, and separating the dataset into training and validation sets. For the detection and classification of deepfake images and videos, OpenCV, and Face Recognition for facial detection, Convolutional neural networks (CNNs) are used by the system to extract features at the frame level. A recurrent neural network is trained using these features (RNN). Various techniques such as data augmentation, learning rate scheduling, and early stopping enhance model performance. This comprehensive approach ensures accurate discrimination between authentic and deep fake content, addressing concerns regarding the integrity of digital media.

## CHAPTER 12

### REFERENCES

- [1] Yuezun Li, Ming-Ching Chang and Siwei Lyu “Exposing AI Created Fake Videos by Detecting Eye Blinking” in arxiv.
- [2] Huy H. Nguyen , Junichi Yamagishi, and Isao Echizen “ Using capsule networks to detect forged images and videos ”.
- [3] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari and Weipeng Xu “Deep Video Portraits” in arXiv:1901.02212v2.
- [4] Umur Aybars Ciftci, İlke Demir, Lijun Yin “Detection of Synthetic Portrait Videos using Biological Signals” in arXiv:1901.02212v2.
- [5] A. M. Almars, “Deepfakes Detection Techniques Using Deep Learning: A Survey,” *J. Comput. Commun.*, 2021, doi: 10.4236/jcc.2021.95003.
- [6] S. Y. Wang, O. Wang, R. Zhang, A. Owens, and A. A. Efros, “CNN-Generated Images Are Surprisingly Easy to Spot.. For Now,” 2020, doi: 10.1109/CVPR42600.2020.00872.
- [7] J. C. Dheeraj, K. Nandakumar, A. V. Aditya, B. S. Chethan, and G. C. R. Kartheek, “Detecting Deepfakes Using Deep Learning,” 2021, doi: 10.1109/RTEICT52294.2021.9573740.
- [8] M. Li, B. Liu, Y. Hu, and Y. Wang, “Exposing deepfake videos by tracking eye movements,” 2020, doi: 10.1109/ICPR48806.2021.9413139.
- [9] Guera, D., and Delp, E. J. (2018, November). Deepfake video detection using recurrent neural networks. In 2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS) IEEE.
- [10] Y. Al-Dhabi and S. Zhang, “Deepfake Video Detection by Combining Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN),” 2021, doi: 10.1109/CSAIEE54046.2021.9543264.
- [11] A. Badale, L. Castelino, and J. Gomes, “Deep Fake Detection using Neural Networks,” vol. 9, no. 3, pp. 349–354, 2021.
- [12] Y. Li, M. C. Chang, and S. Lyu, “In Ictu Oculi: Exposing AI created fake videos by detecting eye blinking,” 2019, doi: 10.1109/WIFS.2018.8630787.

- [13] S. Agarwal, H. Farid, Y. Gu, M. He, K. Nagano, and H. Li, “Protecting world leaders against deep fakes,” 2019.
- [14] M. Nagao, “Natural language processing and knowledge,” in *2005 International Conference on Natural Language Processing and Knowledge Engineering*, 2005, pp. 1-, doi: 10.1109/NLPKE.2005.1598694.
- [15] G. Jaiswal, “Hybrid Recurrent Deep Learning Model for DeepFake Video Detection,” in *2021 IEEE 8th Uttar Pradesh Section International Conference on Electrical, Electronics and Computer Engineering (UPCON)*, 2021, pp. 1–5, doi: 10.1109/UPCON52273.2021.9667632.
- [16] Long Short-Term Memory: From Zero to Hero with Pytorch: <https://blog.floydhub.com/long-short-term-memory-from-zero-to-hero-with-pytorch/>
- [17] Huy H. Nguyen , Junichi Yamagishi, and Isao Echizen “ Using capsule net works to detect forged images and videos ” in arXiv:1810.11215.
- [18]Tiago de Freitas Pereira, Andr e Anjos, Jos e Mario De Martino, and S ebastien Marcel, “Can face anti spoofing countermeasures work in a real world scenario?,”in ICB. IEEE, 2013.
- [19]Nicolas Rahmouni, Vincent Nozick, Junichi Yamagishi, and Isao Echizen, “Distinguishing computer graphics from natural images using convolution neural networks,” in WIFS. IEEE, 2017.
- [20] F. Song, X. Tan, X. Liu, and S. Chen, “Eyes closeness detection from still images with multi-scale histograms of principal oriented gradients,” *Pattern Recognition*, vol. 47, no. 9, pp. 2825–2838, 2014.