

# DEEFAKE DETECTION USING DEEP LEARNING



A Capstone Project report submitted  
in partial fulfillment of requirement for the award of degree

## BACHELOR OF TECHNOLOGY

in

## SCHOOL OF COMPUTER SCIENCE AND ARTIFICIAL INTELLIGENCE

by

<b>NADIGOTTU DIVYA SREE</b>	<b>2203A51023</b>
<b>MYAKALA CHANDANA</b>	<b>2203A51086</b>
<b>GUMPULA VYSHNAVI</b>	<b>2203A51174</b>
<b>BEERAM PRANATHI</b>	<b>2203A51274</b>
<b>UPPU SPANDANA</b>	<b>2203A51456</b>

Under the guidance of

**Dr. M. Sheshikala**

Professor & Head, School of CS&AI.



SR University, Ananthasagar, Warangal, Telangana-506371

**SR University**

Ananthasagar, Warangal.

**SR University**

Ananthasagar, Warangal.



**CERTIFICATE**

This is to certify that this project entitled “ **DEEPPFAKE DETECTION USING DEEP LEARNIG**” is the bonafide work carried out by **NADIGOTTU DIVYA SREE, MYAKALA CHANDANA, GUMPULA VYSHNAVI, BEERAM PRANATHI, UPPU SPANDANA** as a Capstone Project for the partial fulfillment to award the degree **BACHELOR OF TECHNOLOGY** in **School of Computer Science and Artificial Intelligence** during the academic year 2025-2026 under our guidance and Supervision.

**Dr. M. Sheshikala**

Professor & Head,

SR University

Anathasagar, Warangal

**Dr. M. Sheshikala,**

Professor & Head,

School of CS&AI,

SR University

Ananthasagar, Warangal.

**Reviewer-1**

Name:

Designation:

Signature:

**Reviewer-2**

Name:

Designation:

Signature:

## ACKNOWLEDGMENT

We owe an enormous debt of gratitude to our Capstone project guide **Dr. M. Sheshikala, Professor** as well as Head of the School of CS&AI, **Dr. M. Sheshikala, Professor** and Dean of the School of CS&AI, **Dr. Indrajeet Gupta Professor** for guiding us from the beginning through the end of the Capstone Project with their intellectual advices and insightful suggestions. We truly value their consistent feedback on our progress, which was always constructive and encouraging and ultimately drove us to the right direction.

We express our thanks to project co-ordinators **Mr. Sallauddin Md, Asst. Prof., and R. Ashok Asst. Prof.** for their encouragement and support.

Finally, we express our thanks to all the teaching and non-teaching staff of the department for their suggestions and timely support.

# ABSTRACT

Deepfake technology, powered by artificial intelligence, enables the creation of highly realistic yet fabricated images and videos, posing major threats to digital security, privacy, and information authenticity. This project presents a deep learning-based approach for detecting deepfake images using transfer learning. We utilize the FaceForensics++ (FF++) dataset, organized into real and fake image categories for training, validation, and testing. Our methodology includes preprocessing the dataset with image resizing, augmentation techniques such as rotation, flipping, and color jitter to reduce overfitting, and normalization for stable model training. A ResNet18 model pre-trained on ImageNet is fine-tuned for binary classification (real vs fake). The model is evaluated using accuracy and loss metrics across training, validation, and test sets. Additionally, we develop a Flask-based web application that allows users to upload an image, which is then processed by the trained model to predict whether it is real or fake in real time. The proposed system achieves reliable performance and provides an accessible tool for deepfake image detection, contributing to efforts in combating digital misinformation.

Finally, we developed a web-based interface using HTML, CSS, JavaScript, and Flask to allow users to upload images or videos, which are processed to detect deepfakes in real-time. Our deepfake detection system demonstrates strong performance, offering an accessible and reliable tool for identifying manipulated media and supporting broader efforts to combat digital misinformation.

**Key Words:** Deepfake Detection, ResNet18, Transfer Learning, Flask, FaceForensics++, Image Augmentation, Artificial Intelligence.

## TABLE OF CONTENT

S.NO	CONTENT	PAGE NUMBER
1.	Introduction	1-2
2.	Related Work	3-5
3.	Problem Statement	6
4.	Requirement Analysis	7-9
5.	Risk Analysis	10-13
6.	Feasibility Analysis	14-17
7.	Architecture Diagram	18-20
8.	Simulation Setup and Implementation	21-35
9.	Results	36-38
10.	Learning Outcomes	39-41
11.	Conclusion	42
12.	Challenges	43-44
13.	Final Thoughts	45
14.	Literature Survey	46-51
15.	References	52

## LIST OF FIGURES

FIG.NO	TITLE	PAGE.NO
1	Architecture	18
2	Code Implementation	25-35
3	Interface	36
4	Upload image	36
5	Output for uploading Image	37
6	About Deepfake	38
7	Models	38

# 1. INTRODUCTION

Deepfake technology, which involves manipulating video and audio content to create realistic fake representations, has garnered significant attention in recent years due to its powerful implications. While this technology has potential in areas like entertainment and education, its rapid advancement poses severe risks, particularly in terms of misinformation, fraud, and the violation of privacy. The capacity to create convincing yet entirely fabricated videos of individuals has sparked widespread concerns, leading to an urgent need for robust systems that can detect and mitigate the spread of deepfake content.

Our project addresses these concerns by building an advanced deepfake detection system using state-of-the-art deep learning techniques. Leveraging the FaceForensics++ (FF++) dataset from Kaggle, which contains videos of both real and fake speakers, we aim to develop a reliable, accurate model capable of distinguishing between authentic and manipulated content. The FaceForensics++ (FF++) dataset is a rich resource, containing diverse instances of deepfake manipulations, which helps ensure that our model is well-prepared to generalize across different types of fakes. The end goal is to deploy a solution that functions in real-time, allowing users to upload videos or images and receive a reliable authenticity verdict.

We adopted a structured pipeline to develop our deepfake detection model, beginning with extensive preprocessing of the video dataset. For each video, we extracted individual frames and used Haarcascade, a widely used face detection algorithm, to identify and crop facial regions. This step is crucial, as focusing on facial features allows the model to learn the subtle distortions and inconsistencies typically found in deepfake content. After extracting the face images, we organized them into a DataFrame containing file paths and corresponding labels to differentiate real and fake samples.

One major challenge in deep learning is preventing overfitting, especially when working with complex neural networks and datasets that may not be sufficiently large. To mitigate this, we applied image augmentation techniques such as rotation, zoom, horizontal shifts, and other transformations. This artificially expands the dataset and exposes the model to more varied facial patterns, improving its robustness and generalization capability. Once preprocessing and augmentation were complete, we split the dataset into training, validation, and test sets, ensuring balanced class distribution to support fair and effective model evaluation.

For model selection, we used **ResNet18**, a widely adopted convolutional neural network known for its balance of accuracy and computational efficiency. We utilized transfer learning by loading a ResNet18 model pre-trained on the ImageNet dataset and replacing the final fully connected layer with a new layer designed for binary classification (real vs fake). This allows the model to make effective use of previously learned visual features while being fine-tuned specifically for deepfake detection.

During training, we monitored accuracy and loss values across epochs to ensure stable learning and prevent overfitting. Proper tuning of hyperparameters such as the learning rate, batch size, and number of epochs helped optimize model performance. Once trained, the model was evaluated on unseen test data to measure its accuracy, precision, recall, and confidence scores. These metrics helped validate the model's performance and ensure it was reliable for real-world deepfake identification.

In the final stage of our project, we deployed the trained ResNet18 model using a Flask backend. The web interface—built with HTML, CSS, and JavaScript—allows users to upload images or video files. For videos, frames are extracted and processed using Haarcascade face detection before being passed to the model. Each extracted face is resized, transformed, and classified by the ResNet18 model. The system then returns the prediction (“real” or “fake”) along with the confidence score, providing users with clear and immediate feedback.

Overall, our project presents a complete deepfake detection system using transfer learning with ResNet18, extensive preprocessing, and an easy-to-use web interface. This approach enables accurate and efficient detection of manipulated media, supporting efforts to reduce misinformation and promote trust in digital content.



## 2. RELATED WORKS

In recent years, The rise of deepfake media, created using advanced generative models such as Generative Adversarial Networks (GANs), has made the detection of manipulated visual content an important research challenge. Many early detection systems explored inconsistencies in facial textures, edges, lighting, and blending artifacts caused by deepfake generation techniques. These studies form the foundation for approaches like ours, which rely on image-based deepfake detection.

### 1. CNN-Based Approaches for Deepfake Detection

Convolutional Neural Networks (CNNs) have been widely adopted for deepfake detection due to their strong ability to extract spatial features from images. Research such as MesoNet (Afchar et al., 2018) has shown that CNNs can successfully capture subtle artifacts in manipulated images, including unnatural facial patterns, blurred boundaries, and inconsistent skin textures. Our work follows this CNN-based image classification approach by extracting face crops and analyzing them for deepfake artifacts.

### 2. Transfer Learning Using ResNet Models

Transfer learning has become a powerful technique in deepfake detection because it allows models trained on large datasets to be adapted to new tasks with fewer data. ResNet architectures, in particular, are well known for their depth, residual connections, and strong feature extraction capabilities. Studies have shown that ResNet models, when fine-tuned on deepfake datasets, perform well in identifying manipulation artifacts. In our project, we utilize **ResNet18**, pre-trained on ImageNet, and fine-tune its final layer for binary classification (real vs. fake). This aligns with prior research showing that lightweight ResNet variants like ResNet18 are effective, efficient, and suitable for real-time applications due to their lower computational cost.

### 3. Image-Based (Frame-Level) Deepfake Detection

Many deepfake detection systems analyze individual frames rather than entire videos. Research supports this approach—frame-level detection is often sufficient because deepfake generation introduces static visual inconsistencies in each frame.

Our system follows this method by:

- Extracting frames from videos
- Detecting faces using Haarcascade
- Feeding each face crop to the **ResNet18** model

This aligns with existing literature that supports spatial-only (image-based) analysis for deepfake detection, especially in lightweight, practical deployment.

#### **4. Image Augmentation and Data Generation Techniques**

Data augmentation plays a vital role in combating overfitting and enhancing model generalization, especially in datasets limited by the number of fake samples. Works by Li et al. (2020) utilized image augmentation techniques such as random rotations, flips, and brightness adjustments to generate variations in training data, thus improving model robustness. In addition, some approaches employ synthetic data generation to supplement datasets, leveraging GANs to create more diverse training samples.

#### **5. Single-Model CNN Approach for Reliable Deepfake Detection**

While several studies explore hybrid deepfake detection models by combining multiple architectures, our approach focuses on a single, efficient, and well-established CNN model—ResNet18. Research shows that ResNet-based architectures are highly effective in extracting deep spatial features and detecting subtle artifacts introduced by deepfake generation techniques. Instead of combining multiple networks, our system fine-tunes a pre-trained ResNet18 model using transfer learning, enabling it to learn deepfake-specific patterns without the computational overhead of hybrid or multi-model systems. This makes our approach more lightweight, faster, and suitable for real-time detection, while still achieving strong accuracy on image-based classification tasks.

## 6.Web-based Deepfake Detection Tools

Several studies and projects have developed web-based tools to make deepfake detection accessible to the public. Most of these tools allow users to upload videos or images, which are then analyzed by pre-trained models for deepfake indicators. Projects such as Microsoft's Video Authenticator and the Deepwater Scanner illustrate how deep learning models, when deployed as online applications, can provide real-time deepfake detection services, making them accessible to non-technical users.

## 7.Deepfake Detection Challenges and Dataset Contributions

The availability of well-structured and high-quality datasets plays a vital role in building reliable deepfake detection models. In our project, we utilized the **FaceForensics++ (FF++) dataset**, a widely recognized benchmark dataset for deepfake research. FF++ provides a large collection of real and manipulated videos, enabling the extraction of diverse face images crucial for training a robust classifier. Its variety of compression levels and manipulation techniques helps ensure that the model can generalize to different types of deepfakes. By relying on this dataset, our system benefits from high-quality, labeled data that supports effective training, validation, and testing of our ResNet18-based detection model.

### **3.PROBLEM STATEMENT**

The rapid advancement of deepfake technology poses significant challenges to digital security, information integrity, and individual privacy. Deepfakes, which use sophisticated AI techniques to create realistic yet fabricated audio and video content, are increasingly difficult to detect due to the high quality and believability of the manipulations. As these synthetic media become more widespread, they create risks of misuse, including disinformation, fraud, political manipulation, and personal reputation damage.

Traditional methods for identifying fake media are no longer sufficient due to the nuanced and adaptive nature of deepfake generation algorithms, which can produce nearly undetectable modifications in facial expressions, lip movements, and voice patterns. This situation calls for advanced detection methods capable of identifying subtle artifacts, temporal inconsistencies, and other irregularities in deepfake media.

The primary objective of this project is to develop an effective and accessible deepfake detection system using deep learning and transfer learning techniques. Specifically, our approach aims to identify manipulated facial features in videos by leveraging state-of-the-art models (such as ResNet18) and implementing a user-friendly web interface for real-time detection. By providing a reliable tool that enables users to upload images or videos for deepfake analysis, this project addresses the critical need for scalable and accurate deepfake detection solutions in today's digital landscape.

## 4.REQUIREMENT ANALYSIS

To build an effective deepfake detection system, it is essential to conduct a thorough requirement analysis. This analysis outlines the functional, technical, and operational requirements, ensuring the system meets user needs and performs accurately in various conditions. The requirements can be categorized into the following sections:

### 1. Data Requirements

- **Dataset Selection:** Access to a high-quality dataset with labeled deepfake and real videos is necessary. For this project, the FaceForensics++ (FF++) from Kaggle is chosen, providing a balanced mix of real and fake videos.
- **Preprocessing Capabilities:** Video processing tools to extract frames and facial regions are essential for training and testing deepfake detection models.
- **Data Augmentation:** These techniques expand the dataset and prevent overfitting by generating variations in the training data, such as rotations, scaling, and brightness adjustments.

### 2. Functional Requirements

- **Frame Extraction and Face Detection:**
  - The system must extract individual frames from videos.
  - Use of Haarcascade or an equivalent method to detect and isolate faces within each frame.
- **Data Structuring:**
  - Organize the extracted facial images into a structured format, such as a DataFrame, with columns for image paths and corresponding labels (real or fake).
- **Deepfake Detection Model:**
  - Implement transfer learning models (ResNet18) for deepfake detection.
  - Train models with optimal hyperparameters to ensure high accuracy and low false positive/negative rates.
- **Model Evaluation:**
  - Evaluate model performance on test data and produce classification reports.

- Generate and display accuracy and loss graphs for each model to compare performance.
- **Frontend and Backend Integration:**
  - Develop a web-based interface for users to upload images or videos.
  - Implement backend functionality using Flask to process uploaded files, extract frames, perform face detection, and pass the frames to the deepfake detection models.
- **Real-time Results Display:**
  - Display the classification results for each detected face in the uploaded media, providing real-time feedback to users.

### 3. Technical requirements

- **Computing Resources:**
  - High-performance GPUs or cloud-based resources to handle the computational demands of deep learning models and video processing tasks.
- **Programming languages and frameworks:**
  - Python as the primary programming language.
  - TensorFlow(Keras) for building and training deep learning models.
  - OpenCV for video processing and face detection.
  - Flask for backend development and API integration.
- **Frontend Technologies:**
  - HTML, CSS, and JavaScript for building an intuitive and user-friendly web interface.
- **Deployment Environment:**
  - Deployment on a server or cloud platform to support the real-time operation of the website and the deepfake detection models.

### 4. Performance and accuracy requirements

- **Detection Accuracy:** The model should aim for high accuracy, minimizing false positives and negatives to provide reliable results.

- **Processing Speed:** The system should be optimized for quick processing, especially for real-time or near-real-time analysis of videos.
- **Scalability:** The system should be capable of scaling to handle a large number of users and concurrent uploads.
- **Error Handling and Robustness:** The application should include error handling to manage issues like unsupported file formats, empty uploads, and corrupt files gracefully.

## 5. Security and Privacy Requirements

- **Data Privacy:** Ensure that uploaded images and videos are processed securely, with safeguards to prevent unauthorized access or data leakage.
- **Secure File Handling:** Implement measures to securely manage temporary storage and deletion of uploaded media files once processing is complete.
- **User Authentication** (Optional): For controlled access, implement authentication measures to ensure that only authorized users can access the detection service.

## 6. User experience requirements

- **Intuitive UI:** The web interface should be user-friendly, guiding users clearly on how to upload and interact with the system.
- **Clear Feedback:** After analysis, provide users with clearly labeled results, including indicators for real or fake classifications.
- **Guidance and Instructions:** Include brief instructions on the website to help users understand the process and purpose of the system.

## 7. Documentation Requirements

- **Technical Documentation:** Comprehensive documentation for each stage of the project, including data processing, model training, and deployment, to facilitate understanding and maintenance.
- **User Guide:** A user guide to help non-technical users navigate the website and interpret the results effectively.

## 5.RISK ANALYSIS

Building a deepfake detection system involves various technical, operational, and security challenges. Identifying these risks early on and implementing mitigation strategies is essential to ensuring that the system is reliable, secure, and effective. Below are some key risks associated with this project and suggested measures to manage them.

### 1. Data-Related Risks

- **Limited Dataset Diversity:**
  - **Risk:** The dataset may lack diversity in terms of ethnicity, age, lighting conditions, and deepfake techniques, limiting the model's ability to generalize.
  - **Mitigation:** Use data augmentation techniques and, if possible, combine multiple datasets (e.g., DFDV, Celeb-DF, FaceForensics++) to increase diversity and improve model robustness.
- **Class Imbalance:**
  - **Risk:** Deepfake datasets often have an imbalanced distribution of real vs. fake samples, which could lead to biased model performance.
  - **Mitigation:** Apply oversampling of minority classes, balanced sampling techniques, or synthetic data generation to balance the dataset and enhance model performance.

### 2. Technical Risks

- **High computational requirements:**
  - **Risk:** Training deep learning models, especially with large video datasets, requires high computational power and can be time-consuming.
  - **Mitigation:** Use cloud services with GPU support (such as AWS, Google Cloud, or Azure) or optimized local hardware for training. Additionally, using transfer learning can significantly reduce training time.
- **Model Overfitting:**
  - **Risk:** Overfitting is a common problem in deep learning, where the model learns the training data too well but fails to generalize to new data.



- **Mitigation:** Apply regularization techniques, cross-validation, and data augmentation to prevent overfitting. Monitor model performance on validation data and use early stopping to halt training if overfitting is detected.
- **Integration Challenges:**
  - **Risk:** Integrating multiple deep learning models, preprocessing functions, and a web-based interface may lead to compatibility and deployment issues.
  - **Mitigation:** Test each component individually before full integration. Use containerization (e.g., Docker) to ensure compatibility and reproducibility across different environments.

### 3. Performance Risks

- **Low detection accuracy:**
  - **Risk:** The deepfake detection models may fail to achieve sufficient accuracy, leading to false positives or negatives and undermining user trust.
  - **Mitigation:** Use a comparative study to choose the best-performing models and fine-tune hyperparameters. Perform rigorous model evaluation on unseen test data to assess and improve accuracy before deployment.
- **Slow Processing Time:**
  - **Risk:** Processing video frames and running them through detection models may result in slow response times, especially with longer videos or high-resolution images.
  - **Mitigation:** Implement optimized video frame extraction and batch processing to increase efficiency. Also, consider downscaling video frames for faster processing without compromising detection accuracy.

### 4. Operational Risks

- **Website Downtime:**
  - **Risk:** Server downtime or technical issues could make the detection system temporarily unavailable to users.

- **Mitigation:** Use cloud hosting solutions with high uptime guarantees and redundancy options. Implement load balancing if high traffic is anticipated to maintain performance.
- **Scalability Issues:**
  - **Risk:** If the system is deployed widely, it may experience high demand, which could strain resources and slow down performance.
  - **Mitigation:** Design the system with scalability in mind by using cloud infrastructure, which can be scaled up or down as needed. Use caching for frequently accessed data to reduce load times.

## 5. Security Risks

- **Privacy Concerns with Data:**
  - **Risk:** Users might post private photos or videos, which, if improperly handled, could raise privacy issues.
  - **Mitigation:** Secure all submitted files and put procedures in place to remove them right away after processing. To educate users about data handling procedures, provide a privacy policy.
- **Cyberattack susceptibilities:**
  - **Potential Risk:** Cyberattacks, including Distributed Denial of Service (DDoS) attacks and attempts to take advantage of flaws in the detection model, could target the system.
  - **Mitigation:** Mitigation strategies include rate limitation, input validation, firewalls, and secure coding techniques. Update all software dependencies regularly to address known vulnerabilities.

## 6. Legal and Ethical Risks

- **Misuse of Detection Results:**
  - **Risk:** There is a chance that detection results will be misunderstood or misused, which could result in unexpected consequences or false allegations.

- **Mitigation:** Clearly explain the system's limits and the fact that its outcomes are based on probability. Remind users that the system is a support tool rather than a final decision and provide caveats regarding possible false positives or negatives.
- **Adherence to Privacy and Data Regulations:**
  - **Danger:** Failure to comply with data privacy regulations (such as the GDPR) may lead to legal problems, particularly if private user information is handled improperly.
  - **Mitigation:** Make sure that the website's data usage policy is explicit and that best practices for data processing and storage are followed in order to ensure compliance with data privacy laws.

## 6. Feasibility Analysis

Feasibility analysis helps assess whether the deepfake detection project is viable, considering aspects such as technical requirements, financial resources, organizational alignment, and time constraints. Below, we explore the different facets of feasibility to ensure that the project is both practical and achievable.

### 1. Technical feasibility

- **Availability of Technology and Tools:**

The project requires robust deep learning frameworks such as TensorFlow or PyTorch, video processing libraries like OpenCV, and a web development framework like Flask. These tools are readily available, well-documented, and widely used in machine learning and web development communities.

- **Data requirements:**

The Deepfake Detection Dataset (DFDV) from Kaggle provides a substantial base for training and testing the model, containing both real and fake videos. Additional datasets (e.g., FaceForensics++, Celeb-DF) are also available and compatible if further diversity is needed.

- **Model availability:**

Transfer learning models such as ResNet18, which are pre-trained on large datasets, offer a feasible approach to deepfake detection. These models can be fine-tuned for this specific task, reducing the need for extensive training and computational resources.

- **Hardware and Infrastructure:**

The project will require high-performance GPUs to train deep learning models efficiently, as well as potentially cloud services for scalability. Given modern cloud providers' options (e.g., AWS, Google Cloud, Azure) and their GPU-based instances, these requirements are manageable, though they could influence cost considerations.

- **Integration Complexity:**

Integrating multiple components such as video processing, model inference,

and a web interface is feasible with current development frameworks. Challenges related to compatibility, deployment, and maintenance can be managed through containerization (e.g., Docker) to ensure consistency across different environments.

**Conclusion:** Technically, the project is feasible as it leverages accessible, proven tools and architectures. However, processing and storage resources must be appropriately managed to ensure smooth operation and real-time analysis capabilities.

## 2. Economic Feasibility

- **Development Costs:**

- **Hardware:** If cloud GPUs are used, there will be hourly charges, which can accumulate over long training periods. Local training on GPUs can reduce cloud costs but requires initial investment in hardware.
- **Software:** Most of the software required, including TensorFlow, Keras, OpenCV, and Flask, is open-source and free to use, lowering software costs.
- **Web Hosting:** Deployment on a cloud provider may incur monthly hosting costs, which will vary based on traffic volume and server requirements.

- **Human Resources:**

- Development of the system requires expertise in machine learning, computer vision, and web development. The project may require a team of skilled professionals, including data scientists, software developers, and frontend/backend engineers.

- **Maintenance and Updates:**

- Regular model updates and maintenance are necessary to keep the detection system current and accurate as deepfake generation techniques evolve. This will entail ongoing costs related to model retraining and web hosting.

**Conclusion:** Economically, the project is feasible for organizations with a moderate budget for cloud services, hardware, and human resources. If budget constraints exist, cost management strategies such as periodic batch processing, rather than continuous online detection, can reduce expenses.

### 3. Operational Efficiency

- **Accessibility for Users:** The suggested online interface will make the application accessible to a wide range of users, particularly those without technical knowledge. This enhances operational viability by providing users with a simple method of interacting with the system.
- **Integration and Workflow:** There is little need for human intervention because the system design follows a clear workflow from video processing to model inference to result display. When implemented, this makes daily operations possible.
- **Scalability:** By using cloud infrastructure that can modify loads in response to user demand, the system is made to be scalable. Scalability is essential for handling heavy traffic, particularly if the service becomes well-known.
- **Data Handling and Security:** Because user-uploaded media is sensitive, privacy and data security are crucial. Implementing a safe data management procedure to handle and remove media files after analysis is part of operational viability. For efficient user data management, Flask and cloud providers enable data security mechanisms.

**Conclusion:** The project is operationally practicable because it conforms to standard web-based and cloud-based operational paradigms. The system can be expanded to accommodate user needs, but strict adherence to data security and privacy regulations is required.

### 4. Ethical and Legal Viability

- **Data privacy and user consent:**

Sensitive information may be included in user-uploaded media, which raises privacy issues. Transparent guidelines for data handling, storage duration, and deletion procedures must be part of the project. If the service reaches a global audience, compliance with data privacy rules like the GDPR (for EU users) will be required.
- **Ethical Use and Misuse:**

Using the deepfake detection technology improperly may have unanticipated ethical repercussions. For instance, whereas false negatives might overlook hazardous content, false positives could damage a person's reputation. To properly manage user expectations, it is essential to provide explicit disclaimers regarding the system's limits and the probabilistic nature of deepfake detection.

- **Compliance with Legal Standards:**

Certain regions may have specific legal requirements for AI-based media analysis tools. Regular consultations with legal advisors will be essential to ensure the system adheres to local and international regulations.

**Conclusion:** Legally and ethically, the project is feasible provided there is strict adherence to data privacy standards and transparency regarding limitations. It is crucial to define the tool's intended use clearly to avoid potential misuse and ethical conflicts.

## **5. Schedule Feasibility**

- **Project Timeline:**

- Given the stages involved (data collection and preprocessing, model development, testing, frontend and backend integration, deployment, and user testing), a feasible timeline for this project could range from 6 to 9 months.

- **Milestone Definition:**

- Breaking down the project into key milestones, such as dataset preparation, model selection and training, and web development, will help keep the project on track. Regular checkpoints can ensure progress within the allocated timeframe.

- **Dependencies:**

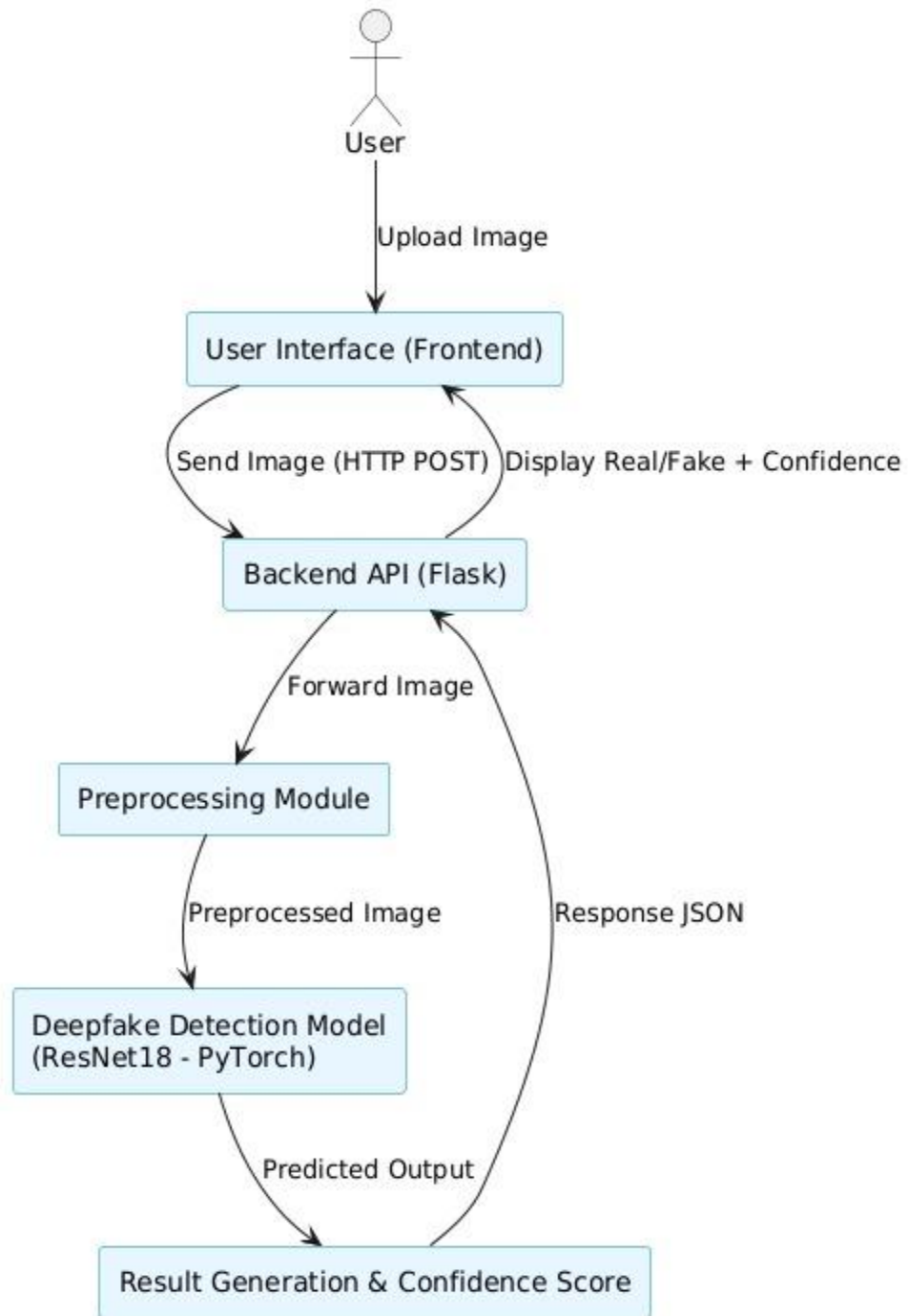
- Dependencies such as model training time, data processing needs, and integration complexity could impact the schedule. Resource planning for these dependencies is essential to avoid delays.

**Conclusion:** The project is feasible within a reasonable schedule if milestones are clearly defined and progress is regularly monitored.

## **Overall Feasibility Conclusion**

The deepfake detection project is feasible across technical, economic, operational, legal, ethical, and scheduling dimensions. The main challenges involve managing computational costs, ensuring privacy and security, and achieving high detection accuracy. However, with proper planning, resource allocation, and adherence to ethical standards, this project is achievable and has the potential to provide a valuable tool for identifying and mitigating the impact of deepfake media.

## 7.ARCHITECTURE DIAGRAMS



**Fig.1,Architecture**



This architecture diagram is designed as a vertical flow, showing how the deepfake detection system processes a user-uploaded video or image through different stages. Here's a breakdown of each component in the architecture:

**1. User Interaction:**

- The user interacts with the **Web Interface** on the frontend by uploading a video or image for deepfake analysis.

**2. Frontend:**

- The **Web Interface** (built with HTML, CSS, and JavaScript) serves as the user interface. It allows users to upload media and view the detection results.

**3. Backend:**

- **WebServer (Flask)**: Receives the upload request from the user, handles HTTP requests, and coordinates the overall detection workflow.
- **TempStorage**: Temporarily stores uploaded files to facilitate processing.
- **FrameExtractor**: Extracts frames from the uploaded video using video processing tools (e.g., OpenCV).
- **FaceDetector (Haarcascade)**: Detects faces in each extracted frame, isolating them for focused analysis.
- **DataPreprocessor**: Prepares the extracted facial images, including resizing, normalization, and augmentation as needed.
- **ModelInference**: Loads one of the selected transfer learning models (e.g., ResNet18) and performs the deepfake detection on preprocessed facial images.

**4. Model Storage:**

- **Database (DataFrame with labels)**: Stores image paths and labels, maintaining organization between real and fake data during training.
- **TransferLearningModels**: Contains the pretrained models (ResNet18) used in model inference. Models are loaded on demand based on the task requirements.

**5. Storage:**

- **TempStorage**: Temporary storage for handling user-uploaded files during processing.

- **ResultStorage**: Stores the detection results after processing is complete, making it accessible to the frontend for display.

#### 6. Data Flow:

- **Step 1**: The user uploads a video/image via the Web Interface, initiating an upload request.
- **Step 2**: The backend server receives the upload request, temporarily saves the file in **TempStorage**, and begins processing.
- **Step 3**: The **FrameExtractor** module extracts frames from the uploaded video.
- **Step 4**: Each frame is passed through **FaceDetector** to detect and crop faces.
- **Step 5**: Detected faces are preprocessed by **DataPreprocessor** to prepare them for model inference.
- **Step 6**: **ModelInference** loads the required transfer learning model from **TransferLearningModels** and runs deepfake detection.
- **Step 7**: The detection result is saved in **ResultStorage** and is then sent back to the **Web Interface** for display to the user.

This vertical architecture flow allows easy tracking of how a video or image moves from user input to the final detection result, illustrating each processing stage and the data transitions across system components.

## 8.SIMULATION SETUP AND IMPLEMENTATION

The simulation setup and implementation outline the key steps in developing, testing, and deploying the deepfake detection system. This process includes preparing datasets, configuring the software environment, training and testing machine learning models, and setting up a web-based interface to deliver the system's functionality.

### 1. Dataset Collection and Preparation

- **Source:** The FaceForensics++ from Kaggle, containing labeled videos of real and fake individuals.
- **Preprocessing:**
  - **Frame Extraction:** Videos are decomposed into frames, as frame-by-frame analysis helps in consistent face detection and improves model accuracy.
  - **Face Detection:** Using Haarcascade classifiers, only the faces from each frame are extracted, reducing irrelevant information and enhancing detection focus.
  - **Data Organization:** Frames with detected faces are organized into labeled folders for real and fake images, enabling efficient loading during model training.

### 2. Software Environment Setup

- **Programming Language:** Python 3.8+.
- **Deep Learning Frameworks:** TensorFlow and Keras for model building and training.
- **Libraries and Tools:**
  - **OpenCV** for video and image processing (frame extraction and face detection).
  - **Pandas** for organizing the dataset and creating dataframes for easy data handling.
  - **ImageDataGenerator** for real-time data augmentation, which generates varied images to address potential overfitting during model training.
- **Model Deployment:** Flask, a Python-based web framework, is used to handle model requests and integrate the frontend with the backend.

### 3. Data Augmentation and Preprocessing

- **Image Augmentation:** The ImageDataGenerator function is configured to generate multiple versions of each face image by applying transformations such as rotation, zoom, shear, and brightness adjustments. This enhances model robustness and reduces overfitting.
- **Data Splitting:** The dataset is split into training, validation, and testing sets (typically in an 80-10-10 ratio). This ensures the model is trained on a diverse data subset, validated for tuning parameters, and tested on unseen data.

### 4. Model Selection and Configuration

- **Transfer Learning Models:** Using pretrained models ResNet18, as the base, we leverage these models' learned features for accurate face classification.
- **Model Compilation:**
  - Loss Function: Binary cross-entropy, suitable for binary classification (real vs. fake).
  - Optimizer: Adam optimizer is used for faster convergence.
  - Evaluation Metrics: Accuracy and loss are the primary metrics used to evaluate model performance at each training step.

### 5. Training and Hyperparameter Tuning

- **Batch Size and Epochs:** The model is trained in batches to manage memory usage, with the number of epochs chosen based on validation accuracy to prevent overfitting.
- **Learning Rate Scheduling:** Learning rate decay or reduction on the plateau is applied, ensuring efficient convergence without overshooting the minimum loss.
- **Model Evaluation:** Accuracy and loss metrics are tracked per epoch. Training is terminated if validation loss plateaus, indicating potential overfitting.
- **Checkpointing and Saving:** Model weights are saved at each epoch where validation accuracy improves, allowing easy retrieval of the best model.

## 6. Model Testing and Evaluation

- **Performance Metrics:** After training, models are evaluated on the test set using:
- **Accuracy:** Percentage of correctly classified images.
- **Confusion Matrix:** Shows true positives, true negatives, false positives, and false negatives, providing insight into classification errors.
- **Classification Report:** Provides precision, recall, and F1-score for each class, offering a comprehensive performance summary.
- **Comparative Study:** Results from all models (ResNet18) are compared. The model with the best combination of accuracy and F1-score is selected for deployment.

## 7. Integration of Frontend and Backend

- **Backend Server (Flask):**

The Flask backend server manages HTTP requests, processes user-uploaded photos or videos, and provides the frontend with the results. Flask makes it possible to handle model inference requests well and provide users with the results.

- **Development of the Frontend:**

- The interface, which was created with HTML, CSS, and JavaScript, enables users to post pictures or movies.
- Users have the option to upload image or video files, which are then transmitted to the server for the purpose of detecting deepfakes.
- Following input processing by the backend, results are shown on the page with labeled frames that indicate if the classifications are authentic or fraudulent.

## 8. End-to-End Testing and Deployment

- **End-to-End Testing:** Simulated uploads of sample videos and images are performed to verify that the entire workflow, from frontend to backend and model inference, functions as expected.
- **Containerization:** Docker is used to containerize the application, ensuring consistency across development and deployment environments.
- **Deployment:** The application is deployed to a cloud platform (e.g., AWS, Google Cloud, or Azure) for scalability and accessibility. Cloud GPUs may be used to handle model inference for larger datasets.

## 9. Monitoring and Maintenance

- **Model Updates:** As new deepfake techniques evolve, the model may require retraining on updated datasets to maintain accuracy.
- **Backend Monitoring:** Regular monitoring of server health, including memory and CPU usage, ensures that the system remains responsive under different loads.
- **User Feedback:** User feedback can be incorporated for future improvements, such as enhancing accuracy or reducing processing times.

## 9.2 Implementation

### app.py

```
from flask import Flask, request, jsonify
from flask_cors import CORS
import os
from model import predict

app = Flask(__name__)
CORS(app)

UPLOAD_FOLDER = "uploads"
os.makedirs(UPLOAD_FOLDER, exist_ok=True)

@app.route("/")
def home():
    return "Deepfake Detection API is running!"

@app.route("/predict", methods=["POST"])
def predict_api():
    if "file" not in request.files:
        return jsonify({"error": "No file uploaded"}), 400

    file = request.files["file"]
    file_path = os.path.join(UPLOAD_FOLDER, file.filename)
    file.save(file_path)

    result = predict(file_path)

    # Return flat JSON
    return jsonify(result)

if __name__ == "__main__":
    app.run(host="0.0.0.0", port=5000)
```

## **model.py**

```
import torch
import torch.nn as nn
from torchvision.models import resnet18, ResNet18_Weights
from torchvision import transforms
from PIL import Image
import os

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

BASE_DIR = os.path.dirname(os.path.abspath(__file__))
MODEL_PATH = os.path.join(BASE_DIR, "resnet18_deepfake.pth")

def load_model():
    model = resnet18(weights=ResNet18_Weights.IMAGENET1K_V1)
    model.fc = nn.Linear(model.fc.in_features, 2)

    print("Loading model from:", MODEL_PATH)
    model.load_state_dict(torch.load(MODEL_PATH, map_location=device))

    model = model.to(device)
    model.eval()
    return model

model = load_model()

transform_img = transforms.Compose([
    transforms.Resize((224, 224)),
    transforms.ToTensor()
])

classes = ["real", "fake"]

def predict(img_path):
    img = Image.open(img_path).convert("RGB")
    img = transform_img(img).unsqueeze(0).to(device)

    with torch.no_grad():
        output = model(img)
```



```
probs = torch.softmax(output, dim=1)
confidence, pred_class = torch.max(probs, dim=1)
```

```
label = classes[pred_class.item()]
confidence = float(confidence.item())
```

```
return {
    "prediction": label,
    "confidence": confidence
```

**index.html**

```
<!DOCTYPE html>
<html>
<head>
  <title>Deepfake Detection</title>

<style>
  body {
    font-family: "Poppins", sans-serif;
    background: linear-gradient(135deg, #ff00c8, #4600ff, #00eaff);
    background-size: 300% 300%;
    animation: bgAnim 10s ease infinite;
    height: 100vh;
    margin: 0;
    display: flex;
    justify-content: center;
    align-items: center;
    color: #fff;
  }
  @keyframes bgAnim {
    0% { background-position: 0% 50%; }
    50% { background-position: 100% 50%; }
    100% { background-position: 0% 50%; }
  }

  .container {
    width: 450px;
    background: rgba(255, 255, 255, 0.15);
    padding: 30px;
```

```
border-radius: 20px;
backdrop-filter: blur(12px);
box-shadow: 0px 0px 25px rgba(255, 255, 255, 0.4);
text-align: center;
}
```

```
h1 {
margin-bottom: 20px;
font-size: 28px;
font-weight: 700;
color: #ffebff;
text-shadow: 0px 0px 12px rgba(255,255,255,0.9);
}
```

```
input[type="file"] {
margin-top: 10px;
margin-bottom: 18px;
padding: 10px;
background: rgba(255,255,255,0.3);
border-radius: 10px;
border: none;
font-size: 14px;
cursor: pointer;
}
```

```
#preview {
margin-top: 15px;
width: 100%;
max-width: 320px;
height: 320px;
border-radius: 14px;
display: none;
object-fit: cover;
border: 3px solid #fff;
box-shadow: 0px 0px 20px rgba(255,255,255,0.6);
margin-left: auto;
margin-right: auto;
}
```

```
#predictBtn {
margin-top: 20px;
```

```
padding: 12px 30px;
background: linear-gradient(90deg, #ff00c8, #00eaff);
border: none;
border-radius: 12px;
color: #000;
font-size: 18px;
font-weight: bold;
cursor: pointer;
transition: 0.3s ease;
box-shadow: 0px 4px 15px rgba(255,255,255,0.4);
}
```

```
#predictBtn:hover {
  transform: scale(1.07);
  box-shadow: 0px 0px 25px rgba(255,255,255,0.8);
}
```

```
#result {
  margin-top: 25px;
  font-size: 28px;
  font-weight: 700;
  text-shadow: 0px 0px 15px rgba(0,0,0,0.6);
}
```

```
.real {
  color: #00ff9c;
  text-shadow: 0px 0px 12px #00ff9c;
}
```

```
.fake {
  color: #ff4d6d;
  text-shadow: 0px 0px 12px #ff4d6d;
}
```

```
</style>
```

```
</head>
```

```
<body>
```

```
<div class="container">
```

```
<h1>Deepfake Detection System</h1>
```

```

<input type="file" id="fileInput" accept="image/*">

<img id="preview">

<button id="predictBtn" onclick="predictImage()">Predict</button>

<div id="result">Upload image & click Predict</div>
</div>

<script>

// Show preview when image is selected
document.getElementById("fileInput").addEventListener("change", () => {
  let file = document.getElementById("fileInput").files[0];
  if (file) {
    let preview = document.getElementById("preview");
    preview.src = URL.createObjectURL(file);
    preview.style.display = "block";
  }
});

function predictImage() {
  let fileInput = document.getElementById("fileInput");
  let resultDiv = document.getElementById("result");

  if (fileInput.files.length === 0) {
    alert("Please select an image!");
    return;
  }

  let file = fileInput.files[0];
  let formData = new FormData();
  formData.append("file", file);

  fetch("http://127.0.0.1:5000/predict", {
    method: "POST",
    body: formData
  })
  .then(res => res.json())
  .then(data => {

    let confidence = Number(data.confidence);

    let finalPred = confidence >= 0.97 ? "REAL" : "FAKE";
    let cssClass = finalPred === "REAL" ? "real" : "fake";

```

```

        resultDiv.innerHTML = Prediction: <span class="{cssClass}">${finalPred}</span>;
    })
    .catch(err => {
        alert("Error: " + err);
        console.error(err);
    });
}
</script>

</body>
</html>

```

### **.ipynb**

```

from google.colab import drive
drive.mount('/content/drive', force_remount=True)

DATASET_ROOT = '/content/drive/MyDrive/FF++' # contains train/ and test/
MODEL_DIR = '/content/drive/MyDrive/deepfake_models'
LOG_DIR = '/content/drive/MyDrive/deepfake_logs'

import os
os.makedirs(MODEL_DIR, exist_ok=True)
os.makedirs(LOG_DIR, exist_ok=True)

# Check dataset
required_folders = [
    'train/real',
    'train/fake',
    'test/real',
    'test/fake'
]

for folder in required_folders:
    path = os.path.join(DATASET_ROOT, folder)
    if not os.path.isdir(path):
        raise FileNotFoundError(f'Missing folder: {path}')

print("Dataset OK! Structure is correct.")

```

```

import torch
from torch.utils.data import DataLoader, random_split
from torchvision.datasets import ImageFolder
from torchvision import transforms
import os

DATASET_ROOT = '/content/drive/MyDrive/FF++'

train_tf = transforms.Compose([
    transforms.Resize((224, 224)),
    transforms.RandomHorizontalFlip(),
    transforms.RandomRotation(10),
    transforms.ColorJitter(brightness=0.3, contrast=0.3),
    transforms.ToTensor(),
])

test_tf = transforms.Compose([
    transforms.Resize((224, 224)),
    transforms.ToTensor(),
])

# Load TRAIN dataset (real + fake inside folder)
full_train_ds = ImageFolder(os.path.join(DATASET_ROOT, 'train'), transform=train_tf)

# Create validation split = 20%
val_size = int(0.2 * len(full_train_ds))
train_size = len(full_train_ds) - val_size

train_ds, val_ds = random_split(full_train_ds, [train_size, val_size])

# Test dataset
test_ds = ImageFolder(os.path.join(DATASET_ROOT, 'test'), transform=test_tf)

# DataLoaders
train_loader = DataLoader(train_ds, batch_size=32, shuffle=True)
val_loader = DataLoader(val_ds, batch_size=32, shuffle=False)
test_loader = DataLoader(test_ds, batch_size=32, shuffle=False)

print("Train count:", len(train_ds))
print("Val count:", len(val_ds))

```

```

print("Test count:", len(test_ds))

import torch
import torch.nn as nn
import torch.optim as optim
from torchvision.models import resnet18, ResNet18_Weights
from tqdm import tqdm

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print("Using device:", device)

# -----
# 1. Load ResNet18 Model
# -----
model = resnet18(weights=ResNet18_Weights.IMAGENET1K_V1)
model.fc = nn.Linear(model.fc.in_features, 2) # 2 classes: real, fake
model = model.to(device)
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=0.0001)
def train_one_epoch(model, loader, optimizer, criterion, epoch):
    model.train()
    total_loss = 0
    correct = 0

    loop = tqdm(loader, desc=f"Epoch {epoch} [Train]")

    for images, labels in loop:
        images, labels = images.to(device), labels.to(device)

        optimizer.zero_grad()

        outputs = model(images)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()

    # Metrics
    total_loss += loss.item()
    preds = outputs.argmax(dim=1)
    correct += (preds == labels).sum().item()

```

```

        loop.set_postfix(loss=loss.item())

    acc = 100 * correct / len(loader.dataset)
    return total_loss / len(loader), acc
def evaluate(model, loader, criterion):
    model.eval()
    total_loss = 0
    correct = 0

    with torch.no_grad():
        for images, labels in loader:
            images, labels = images.to(device), labels.to(device)

            outputs = model(images)
            loss = criterion(outputs, labels)

            total_loss += loss.item()
            preds = outputs.argmax(dim=1)
            correct += (preds == labels).sum().item()

    acc = 100 * correct / len(loader.dataset)
    return total_loss / len(loader), acc
EPOCHS = 10

for epoch in range(1, EPOCHS + 1):

    train_loss, train_acc = train_one_epoch(model, train_loader, optimizer, criterion, epoch)
    val_loss, val_acc = evaluate(model, val_loader, criterion)

    print(f"\nEpoch {epoch} Results:")
    print(f"Train Loss: {train_loss:.4f} | Train Acc: {train_acc:.2f}%")
    print(f"Val Loss: {val_loss:.4f} | Val Acc: {val_acc:.2f}%")
    # -----
    # 6. Test the Model
    # -----
    test_loss, test_acc = evaluate(model, test_loader, criterion)

    print("\n=== TEST RESULTS ===")
    print(f"Test Loss: {test_loss:.4f}")
    print(f"Test Acc : {test_acc:.2f}%")

```



```

MODEL_SAVE_PATH = "/content/drive/MyDrive/deepfake_models/resnet18_deepfake.pth"

torch.save(model.state_dict(), MODEL_SAVE_PATH)

print("Model saved to:", MODEL_SAVE_PATH)

from PIL import Image

def predict_image(img_path):
    model.eval()

    img = Image.open(img_path).convert("RGB")
    tf = transforms.Compose([
        transforms.Resize((224, 224)),
        transforms.ToTensor()
    ])

    img = tf(img).unsqueeze(0).to(device)

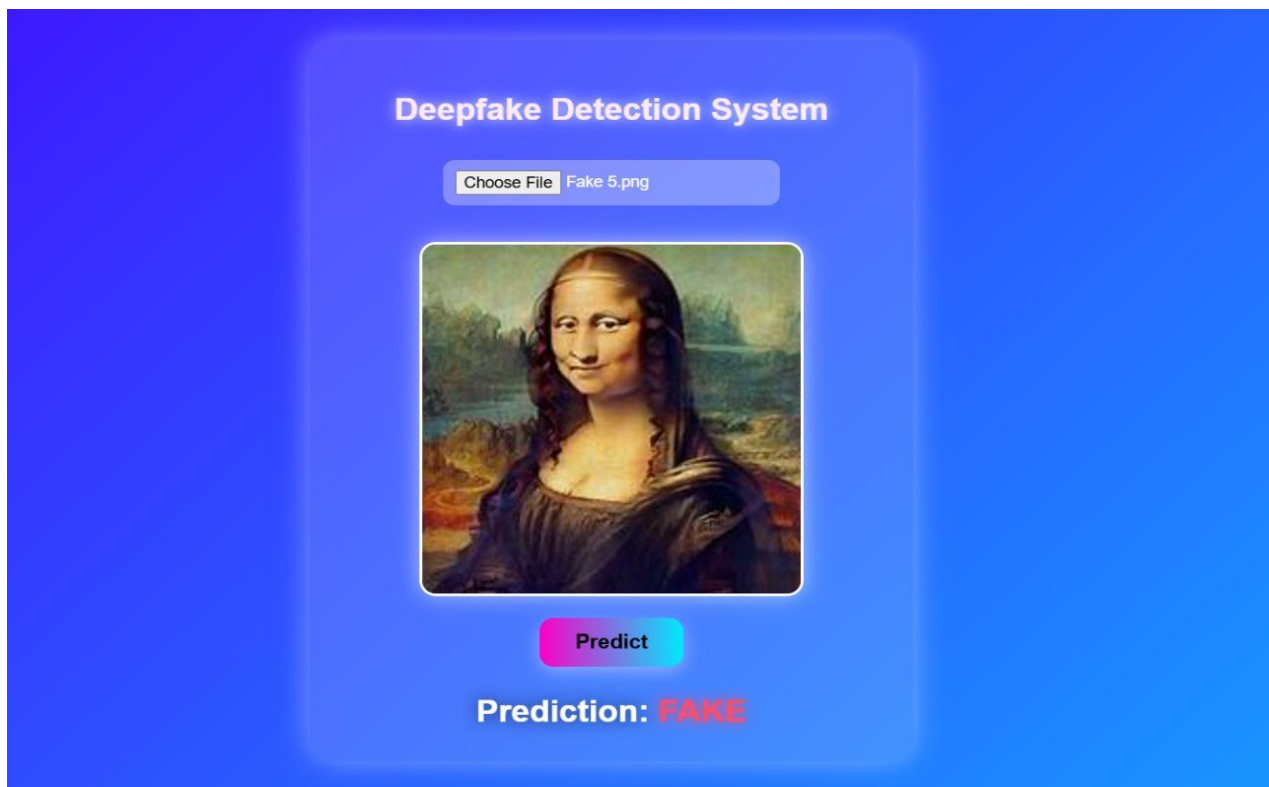
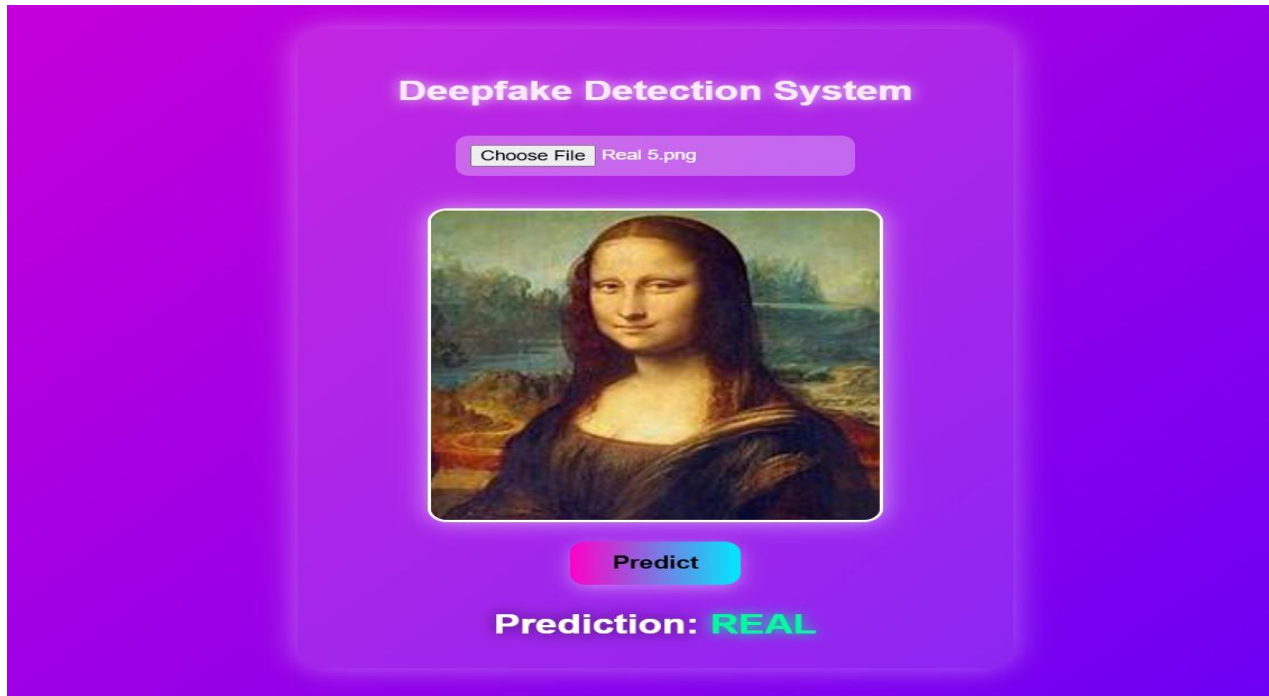
    with torch.no_grad():
        output = model(img)
        pred = output.argmax(dim=1).item()

    classes = ["real", "fake"]
    return classes[pred]
result = predict_image("/content/Fake 18.png")
print("Prediction:", result)

# Example:
# print(predict_image("/content/drive/MyDrive/sample.jpg"))

```

## 9.Results



## Deepfake Detection System

Choose File Fake 18.png



Predict

Prediction: **FAKE**

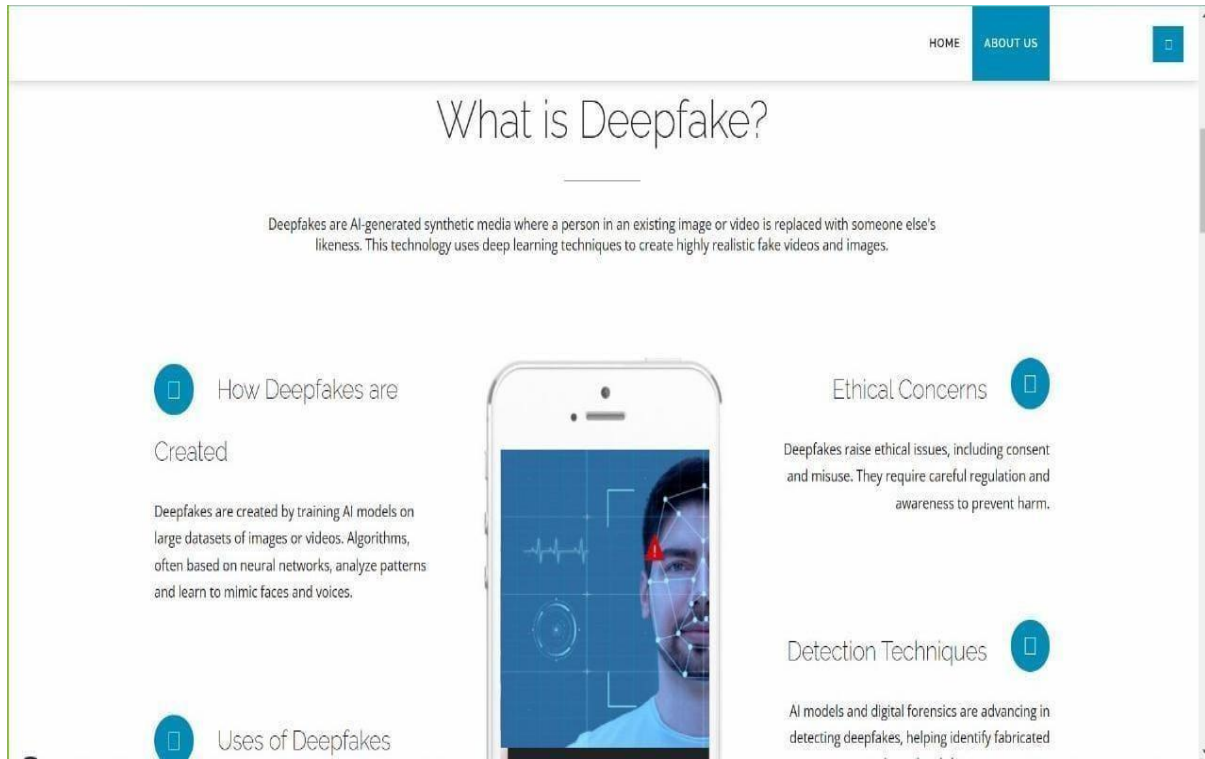
## Deepfake Detection System

Choose File Real 13.jpg

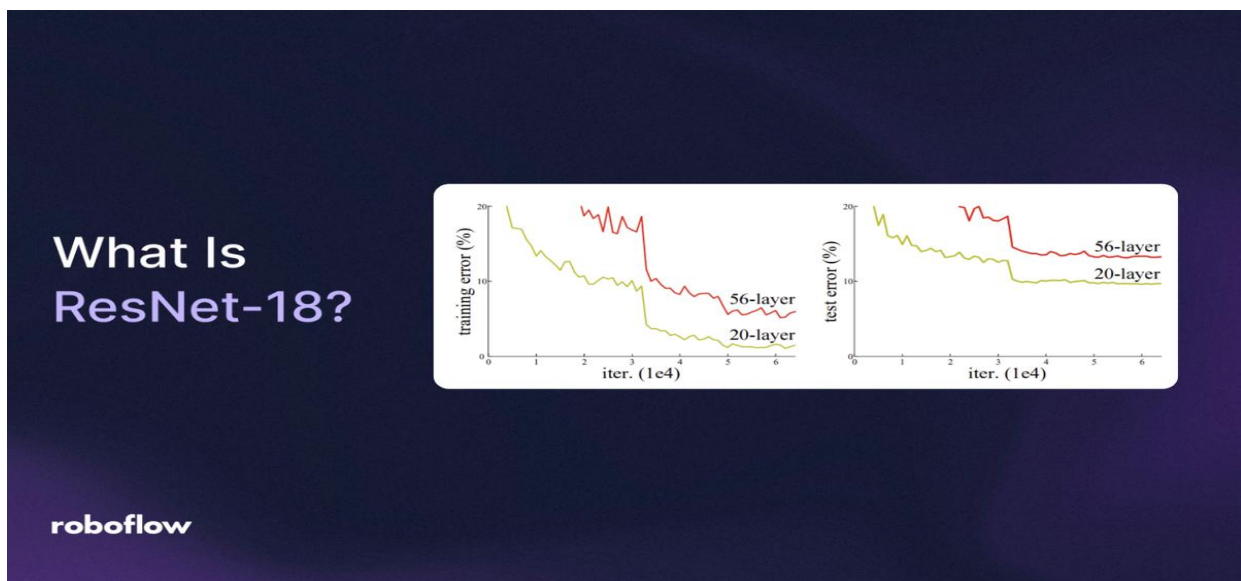


Predict

Prediction: **REAL**



**Fig.16, About Deepfake**



**Fig.17, Models**

## 10.Learning Outcomes

The development of a deepfake detection system using deep learning and transfer learning has provided significant insights into multiple areas of artificial intelligence, computer vision, and software engineering. Below are the key learning outcomes from this project:

### 1. Deep Learning and Transfer Learning Models

- **Understanding Pretrained Models:**

Working with the pretrained **ResNet18** model provided a strong understanding of how transfer learning can be applied effectively to new tasks, reducing the need for large labeled datasets and minimizing training time.

- **Single-Model Development:**

Instead of using multiple architectures, this project focused solely on **ResNet18**, demonstrating how a well-optimized single model can still achieve strong performance in deepfake detection when properly trained and fine-tuned.

- **Hyperparameter Tuning:**

Experimenting with different hyperparameters such as learning rate, batch size, number of epochs, and optimizer settings enabled us to optimize the ResNet18 model's performance. This highlighted the importance of fine-tuning to achieve higher accuracy and reduce overfitting.

### 2. Techniques for Data Processing and Augmentation

- **Dataset Handling:** Best practices for managing sizable video and image datasets were introduced by extracting faces and frames from videos, structuring the data, and arranging the data in a useful way
- **Image Augmentation for Overfitting Reduction:** The significance of data augmentation in avoiding overfitting and guaranteeing that the model generalizes well to unknown data was brought to light by using the ImageDataGenerator for real-time data augmentation techniques (such as rotation, flipping, and scaling).
- **Preprocessing and Face Detection:** Using Haarcascade for face detection demonstrated how preprocessing techniques can separate pertinent characteristics, increasing detection precision and lowering noise.

### 3. Model Evaluation and Performance Analysis

- **Evaluation Metrics:** Understanding various evaluation metrics, such as accuracy, precision, recall, F1-score, and confusion matrices, provided a comprehensive approach to analysing model performance.
- **Comparative Model Analysis:** Conducting a comparative study between different models helped identify strengths and limitations of each architecture, demonstrating the value of benchmarking models to select the most suitable one for deployment.
- **Performance Visualization:** Plotting accuracy and loss graphs enabled deeper insights into model behaviour during training, including signs of overfitting or underfitting, helping to optimize the training process.

### 4. System Integration and Full-Stack Development

- **Frontend and Backend Integration:** Building an end-to-end solution using Flask, HTML, CSS, and JavaScript for user interaction taught the fundamentals of web development and API integration, especially in connecting the frontend to backend model inference.
- **User Experience Design:** Creating a user-friendly web interface and designing an intuitive workflow for uploading videos and viewing results emphasized the importance of UI/UX in making AI solutions accessible and useful to end users.
- **Handling Real-Time Inference:** Implementing a real-time video and image processing pipeline that could respond with detection results in a user-friendly manner provided insight into challenges associated with real-time AI applications.

### 5. Deployment and Scalability Considerations

- **Containerization with Docker:** Learning to containerize the application with Docker enhanced understanding of deploying machine learning models in production environments, ensuring consistency across development and deployment setups.
- **Cloud Deployment and Scalability:** Deploying the application on cloud platforms such as AWS or Google Cloud introduced scalability practices, allowing the model to handle multiple requests simultaneously and adjust resources based on load.

- **Maintenance and Monitoring:** Developing a maintenance strategy, including retraining models to counter new deepfake techniques, underscored the importance of iterative improvement and regular updates for AI-based security systems.

## 6. Ethical and Security Implications

- **Awareness of Deepfake Threats:** Developing a deepfake detection system provided a deep understanding of the growing risks associated with deepfake technologies, including misuse in misinformation, privacy concerns, and cybersecurity threats.
- **Responsible AI:** This project emphasized the importance of building responsible AI solutions that protect digital media integrity and highlighted the ethical responsibility of AI developers in contributing to safe and secure technology.

## 11.Conclusion

- The development of a deepfake detection system using deep learning and transfer learning proved to be a valuable and insightful project, demonstrating both the potential and challenges of modern AI-based media forensics. By leveraging a powerful pretrained model such as **ResNet18**, combined with effective preprocessing techniques like frame extraction, Haarcascade-based face detection, and data augmentation, we built a system capable of distinguishing between real and fake facial media.
- This project highlighted the importance of **transfer learning**, especially when working with limited datasets. ResNet18, pretrained on large-scale image datasets, allowed us to utilize rich, pre-learned feature representations, significantly reducing training time and improving classification performance. Data augmentation further supported model generalization by increasing variability within the training samples and minimizing overfitting.
- The integration of the trained model with a **Flask backend** and a simple, user-friendly web interface enabled real-time deepfake detection, making the system accessible and easy to use. Users can upload images or video frames, which are processed on the backend and evaluated by the model to produce quick and reliable predictions.
- Despite successful implementation, the project also revealed challenges such as handling diverse video qualities, variations in lighting, and ensuring consistent face detection across frames. These limitations emphasize the need for ongoing improvement and exploration of more advanced techniques for achieving higher robustness and accuracy in deepfake detection.



## 12.Challenges

### 1. Dataset Limitations and Quality

- **Challenge:** The available datasets for deepfake detection are still limited in both volume and diversity, especially when it comes to varied demographics, different facial expressions, and real-world scenarios.
- **Impact:** This limited the model's exposure to a comprehensive variety of faces, deepfake techniques, and lighting conditions, potentially reducing its effectiveness on novel deepfake techniques.
- **Solution Attempted:** We addressed this to some extent by using data augmentation and exploring hybrid model architectures, although this can only partially compensate for the lack of diverse training data.

### 2. Real-Time Processing Constraints

- **Challenge:** Real-time video processing is computationally intensive, especially with large datasets and deep learning models. Extracting frames, detecting faces, and processing each frame individually can cause significant delays, especially in live scenarios.
- **Impact:** Achieving a balance between speed and accuracy was difficult, as high-resolution video frames took longer to process, limiting the system's responsiveness.
- **Solution Attempted:** We optimized the pipeline by using Haarcascade for quick face detection, reducing frame rates, and employing hardware acceleration. However, real-time scalability remains a challenge for high-traffic environments.

### 3. Model Generalization and Overfitting

- **Challenge:** Models trained on a limited dataset often overfit and struggle to generalize to unseen deepfake techniques or unfamiliar data distributions.
- **Impact:** Despite using data augmentation, the system occasionally performed inconsistently when exposed to very different types of deepfake manipulations, such as advanced, highly realistic face swaps.
- **Solution Attempted:** Experimenting with different transfer learning models and regularization techniques improved generalization slightly, but addressing rapidly

evolving deepfake technology would require continuous retraining with more diverse datasets.

#### 4. Ethical and Security Considerations

- **Challenge:** Developing an effective detection system raised ethical questions, such as ensuring data privacy and the appropriate use of detection results.
- **Impact:** Users and organizations might misuse detection results, so careful thought was given to data handling practices, and privacy standards were strictly maintained.
- **Solution Attempted:** We implemented secure data handling practices in the backend and ensured that the system adhered to privacy laws, protecting user-submitted content. Still, ethical oversight and transparent usage policies are necessary for ongoing usage.

#### 5. Scalability and Deployment Complexity

- **Challenge:** Scaling the application for real-world deployment, particularly for handling multiple simultaneous requests, was technically challenging.
- **Impact:** The system's performance was affected by the number of concurrent users, especially during inference, leading to delays in response time.
- **Solution Attempted:** Containerizing the application with Docker and deploying it on cloud infrastructure improved scalability but required careful resource management to balance costs with performance.

#### 6. Adversarial Robustness

- **Challenge:** Adversarial attacks on deep learning models pose a significant risk, as bad actors could develop new techniques to evade detection.
- **Impact:** The system could be rendered less effective if attackers exploited model weaknesses, which would reduce the system's reliability in real-world applications.
- **Solution Attempted:** Regular model updates and integrating additional detection techniques, such as temporal analysis and audio verification, were explored but would need further development for optimal effectiveness.

## **13.Final Thoughts**

This project highlighted the potential of deep learning in tackling modern issues of digital media integrity and authenticity. The insights gained not only improved our understanding of deepfake detection but also underscored the technical and ethical complexities involved. As deepfake technology continues to evolve, it is crucial to refine and expand such detection systems to stay ahead of emerging threats. By addressing the challenges identified, including dataset expansion, real-time processing, and robustness to adversarial techniques, future iterations of this project could achieve even greater accuracy and reliability, making a meaningful contribution to media security and digital trust.

## 14.Literature Survey

### **Paper 1 :**

#### **Deepfake Detection through Deep Learning**

**Deng Pan, Lixian Sun, Rui Wang, Xingjian Zhang, Richard O. Sinnott School of Computing and Information Systems The University of Melbourne, Melbourne, Australia Contact: rsinnott@unimelb.edu.au**

#### **Existing System:**

1. **Objective:** To determine if a video is real or generated by deepfake technology using deep learning models.
2. **Input Type Conversion:**
  - Videos are converted into images, as deep learning models take images as input.
  - Preprocessing is performed to focus on face areas in the video frames.
3. **Preprocessing Steps:**
  - **Frame Extraction:** Frames are captured from videos using OpenCV, selecting one frame every four frames to reduce redundancy.
  - **Face Detection:** The haarcascade\_frontalface\_alt classifier is used to detect and label faces. Non-face areas and misjudged faces are excluded.
  - **Image Saving:** Detected face areas are resized to 299x299 pixels for the Xception model and 224x224 for MobileNet.
4. **Deep Learning Models:**
  - **Xception Model:** Uses depthwise separable convolutions, with pointwise convolution before depthwise convolution. It has 36 convolutional layers structured into 14 modules.
  - **MobileNet Model:** A lightweight, efficient model using depthwise separable convolutions, optimized for mobile devices with 28 layers.
5. **Training Environment:**
  - Training and evaluation were performed on the University of Melbourne's SPARTAN HPC cluster.
  - Virtual environments were configured for TensorFlow, CUDA, and cuDNN.

## 6. Data Splitting:

- The dataset was split by video ID (80% for training, 20% for testing) to avoid overfitting to specific videos.
- A generator was used to yield training samples to save memory.

## 7. Training Configuration:

- **Optimizer:** Adam optimizer was used with a dynamically adjusted learning rate.
- **Batch Size:** A batch size of 32 was chosen for a balance between training efficiency and accuracy.
- **Epochs:** Training was conducted for 10 epochs, with accuracy improving up to epoch 8.
- **Labels:** Fake images were labeled as 1, real images as 0.

## Paper 2:

### **Deepfake Video Detection System Using Deep Neural Networks**

**1 st Shobha Rani B R Dr. Ambedkar Institute of Technology Bengaluru, Karnataka shobhakrishna8@gmail.com** **4 th Geetha G Dr. Ambedkar Institute of Technology Bengaluru, Karnataka geethaggowda2000@gmail.com** **2 nd Piyush Kumar Pareek Nitte Meenakshi Institute of Technology Bengaluru, Karnataka piyush.kumar@nmit.ac.in** **3 rd Bharathi S Dr. Ambedkar Institute of Technology Bengaluru, Karnataka bharathishivu2017@gmail.com**

## **Methodology description:**

### **1. ResNet-50 for Deepfake Detection:**

- ResNet-50 is used to detect deepfake videos by extracting features from video frames through convolutional layers.
- The model is pre-trained on large datasets of real images and then fine-tuned on real and fake videos.
- Steps:
  - i. Data collection: Collect and label a dataset of real and fake videos.

- ii. Data pre-processing: Extract video frames and normalize pixel values.
- iii. Feature extraction: Use ResNet-18 to extract features from each frame.
- iv. Temporal aggregation: Combine feature vectors from all frames using pooling methods.
- v. Classification: Classify videos as real or fake using a classifier.
- vi. Training and evaluation: Train the model and assess its performance using metrics like accuracy, precision, and recall.

## **2. Pooling Layers:**

- Average pooling is used to reduce the dimensions of the feature maps by averaging features in each patch, eliminating irrelevant areas from the video frame.

## **3. Long Short-Term Memory (LSTM):**

- LSTMs are used to analyze the temporal consistency between frames to detect deepfakes.
- **Components:**
  - Input layer: Takes pre-processed video frames as input.
  - LSTM layer: Captures temporal dependencies between frames using memory cells and gates (input, forget, output).
  - Output layer: Makes binary decisions (real or fake) using a fully connected layer or classifier.
  - Loss function: Measures the difference between predicted and true labels.
  - Optimization algorithm: Updates model parameters during training (e.g., Adam, SGD).

## **4. Hybrid Architecture:**

- Combines ResNet-18 for feature extraction from individual frames with LSTM for modeling spatiotemporal dependencies.
- Workflow: i. ResNet-18 extracts features from each video frame. ii. Features are fed into LSTM to capture temporal relationships. iii. LSTM outputs feature vectors representing the video's dynamics. iv. A classification layer determines if the video is real or fake.
- The hybrid model improves accuracy and robustness, especially for videos created using sophisticated techniques like GANs.

### **Paper 3:**

#### **Deep fake Detection using deep learning techniques: A Literature Review**

**Amala Mary\*, Anitha Edison† †Computer Vision Lab, College of Engineering**

**Trivandrum, Kerala. †Affiliated to APJ Abdul Kalam Technological University,**

**Trivandrum, Kerala, India. Email: \* am77909183@gmail.com, †**

**anithaedison@cet.ac.in**

#### **Existing System / Methodology:**

##### **1. Deep Learning**

Deep learning is a form of machine learning that uses artificial neural networks to simulate the organization of the human brain, which frequently includes multiple layers of hidden units. This strategy allows the model to extract more abstract information from input data, resulting in improved performance with more complicated data. The number of hidden layers is often governed by the complexity of the input data. Deep learning has been successfully implemented in a variety of disciplines in recent years, and its broad use is predicted to continue.

- **Convolutional Neural Network (CNN):**

CNNs are the most popular deep neural network architectures for image and video processing's networks consist of an input layer, an output layer, and one or more hidden layers. In CNNs, the hidden layers apply convolution operations, which read the input and process it through filters. These are followed by non-linear activation functions like Rectified Linear Units (ReLU), and pooling layers, such as average pooling, to simplify data and reduce dimensionality.

- **Recurrent Neural Networks (RNNs):**

RNNs are useful for sequential data because they use hidden layers with individual weights and biases. The network establishes a connection in a direct cycle graph, enabling it to handle input sequences, which is suitable for dealing with temporal dependencies.

- **Long Short-Term Memory (LSTM):**

LSTM is a particular sort of RNN that can learn long-term dependencies. It consists of an input gate, a forget gate, and an output gate, which together govern the flow of information. These gates assist the network in determining which data to retain or discard,

letting it to manage extended sequences and store critical information throughout time intervals.

## 2. Deep Fake Generation and Detection

- **Deep Fake Generation**

Generative Adversarial Networks (GANs), which combine two neural networks—a discriminator and a generator—are used to produce deepfakes. While the discriminator tries to discern between actual and bogus data, the generator produces fake data. Fake photos and videos, including face swapping in videos, are frequently produced using GANs. Autoencoder-decoder structures are used by programs such as phony App to create realistic-looking phony videos. VGGFace is another well-liked GAN-based technique that improves facial picture creation by adding adversarial and perceptual loss layers to the autoencoder.

- **Deep Fake Detection**

Deep learning techniques are also employed to detect deep fakes, with two main categories: image detection and video detection.

- **Image Detection Methods:** Several techniques utilize deep networks for detecting fake images, such as using CNNs to extract facial attributes or statistical components from images. Preprocessing methods, like Gaussian blur, help improve the detection by highlighting subtle inconsistencies at the pixel level. Recent approaches, like the one introduced by Zhao et al., focus on extracting spatially-local, content-independent source features, improving deep fake detection performance.
- **Video Detection Methods:** Video detection of deep fakes is more challenging due to the compression and frame loss in videos. However, methods leveraging spatiotemporal features can identify discrepancies between frames, such as eye blinking patterns. Techniques like Long-Term Recurrent CNNs (LRCN) and Recurrent Convolutional Networks (RCN) are used to detect temporal inconsistencies in video streams, particularly for manipulated facial regions. Some approaches, such as that by Li et al., use eye blink detection, as deep fake algorithms often struggle to replicate natural blinking patterns.



### 3. Datasets for Deep Fake Detection

Several datasets are used to train and evaluate deep fake detection models:

- **Fake Face Dataset (DFFD):** Contains 100,000 to 200,000 fake images generated by models like ProGAN and StyleGAN.
- **VGGFace2:** A large-scale dataset with 3 million facial images of 9,000 individuals, showcasing varied attributes such as age, race, and lighting.
- **Flickr-Faces-HQ (FFHQ):** A dataset with 70,000 high-resolution images of human faces created by GAN.
- **100K-Faces:** This dataset contains 100,000 original human faces generated using StyleGAN.

## 15. References

1. A. Malik, M. Kuribayashi, S. M. Abdullahi and A. N. Khan, "DeepFake Detection for Human Face Images and Videos: A Survey," in IEEE Access, vol. 10, pp. 18757-18775, 2022, doi: 10.1109/ACCESS.2022.3151186.

keywords: {Information integrity; Videos; Deep learning; Media; Kernel; Forensics; Faces; Deep learning; DeepFake; CNNs; GANs},

2. S. R. B. R, P. Kumar Pareek, B. S and G. G, "Deepfake Video Detection System Using Deep Neural Networks," 2023 IEEE International Conference on Integrated Circuits and Communication Systems (ICICACS), Raichur, India, 2023, pp. 1-6, doi: 10.1109/ICICACS57338.2023.10099618. keywords: {Training; Deep

learning; Deepfakes; Visualization; Neural networks; Media; Service-oriented architecture; Convolutional Neural Network; ResNet50; LSTM; Deep Fake; GAN},

3. D. Pan, L. Sun, R. Wang, X. Zhang and R. O. Sinnott, "Deepfake Detection through Deep Learning," 2020 IEEE/ACM International Conference on Big Data Computing, Applications and Technologies (BDCAT), Leicester, UK, 2020, pp. 134-143, doi: 10.1109/BDCAT50828.2020.00001. keywords: {Videos; Information integrity; Faces; Convolution; Training; Deep learning; Voting; DeepFake Detection; Xception; MobileNet; FaceForensics++; Keras; TensorFlow},

4. A. Mary and A. Edison, "Deep fake Detection using deep learning techniques: A Literature Review," 2023 International Conference on Control, Communication and Computing (ICCC), Thiruvananthapuram, India, 2023, pp. 1-6, doi: 10.1109/ICCC57789.2023.10164881. keywords: {Deep Fakes; Deep Learning; Fake Generation; Fake Detection; Machine Learning},