# DEEP FAKE DECODER- AI FOR DETECTING DEEPFAKE VOICE SCAMS

A
Synopsis
Submitted in the partial fulfilment of the requirement for the award of
Bachelor of Technology

In

Artificial Intelligence and Data Science

Submitted to



**Samrat Ashok Technological Institute, Vidisha**
(An Autonomous Institute Affiliated To RGPV, Bhopal)

Submitted by:

SATISH PRAJAPATI (0108AI221055)

PRIYANSHU DUBEY (0108AI22046)

ARYAN LAKSHAKAR (0108AI221012)

SHIVANSH CHOUHAN (0108AI221060)

Under the Supervision of

**Prof. Nimisha Bhandari**
Assistant Professor

**Department of Information Technology**

**Samrat Ashok Technological Institute**

**Vidisha (M.P.) - 464001**

**May 2025**

# Samrat Ashok Technological Institute
## Vidisha (M.P.)
### Department of Information Technology

## <span style="color:red">CERTIFICATE</span>

This is to certify that the Minor Project entitled as "AI FOR DETECTING DEEPFAKE VOICE SCAMS" submitted by **Satish Prajapati (0108AI221055), Priyanshu Dubey (0108AI221046), Aryan Lakshakar (0108AI221012), Shivansh Chouhan (0108AI221060)** in the partial fulfilment of the requirements for the award of degree of Bachelor of Engineering in the specialization of Artificial intelligence and data science from Samrat Ashok Technological Institute, Vidisha (M.P.) is carrying out by them under my supervision and guidance. The matter presented in this synopsis has not been presented by him elsewhere for any other degree.

| | | |
|---|---|---|
| **Prof. Nimisha Bhandari** | **Prof. Ram Ratan Ahirwar** | **Dr. Shailendra Kumar** |
| **(Project Guide)** | **(Project Coordinator)** | **Shrivastava** |
| **Assistant Professor** | **Assistant Professor** | **(H.O.D.)** |
| Department of Information Technology | Department of Information Technology | Department of Information Technology |
| Samrat Ashok Technological Institute | Samrat Ashok Technological Institute | Samrat Ashok Technological Institute |
| Vidisha (M.P.) | Vidisha (M.P.) | Vidisha (M.P.) |

# CANDIDATE'S DECLARATION

"We, **Satish Prajapati (0108AI221055), Priyanshu Dubey(0108AI221046), Aryan Lakshakar (0108AI221012), Shivansh Chouhan(0108AI221060)** hereby declare that the work which is being presented in the Minor Project synopsis entitled **"AI for detecting deepfake voice scams"** submitted in partial fulfillment of the requirement for the award of the degree of Bachelor of Engineering in Artificial intelligence and data science The work has been carried at Samrat Ashok Technological Institute, Vidisha is an authentic record of my own work carried out under the guidance of. Prof. Nimisha Bhandari (Assistant Professor, Department of Information Technology), Samrat Ashok Technological Institute, Vidisha (M.P).

*The matter embodied in this synopsis has not been submitted by me for the award of any other degree.*

**Date: 11-04-2025**
**Place: Vidisha, (M.P)**

| Signature of Candidate | Signature of Candidate | Signature of Candidate |
|---|---|---|
| **(Satish Prajapati)** | **(Priyanshu Dubey)** | **(Aryan Lakshakar)** |
| Enrollment No.: 0108AI221055 | Enrollment No.: 0108AI221046 | Enrollment No.: 0108AI221012 |
| Samrat Ashok Technological Institute | Samrat Ashok Technological Institute | Samrat Ashok Technological Institute |
| Vidisha, (M.P.) | Vidisha, (M.P.) | Vidisha, (M.P.) |

Signature of Candidate
**(Shivansh Chouhan)**
Enrollment No.: 0108AI221060
Samrat Ashok Technological Institute
Vidisha, (M.P.)

# ACKNOWLEDGEMENT

We would like to express my sincere gratitude to all those who have supported and guided me throughout the course of this project.

First and foremost, I am deeply thankful to my Project Guide, Prof. Nimisha Bhandari, for their invaluable guidance, constant support, and constructive feedback which helped me steer this project in the right direction. Their encouragement and expertise have been truly inspiring.

I extend my heartfelt thanks to our Project Coordinator, Prof. Ram Ratan Ahirwar, for their consistent coordination and assistance at every stage of the project.

I am also grateful to our respected Head of Department, Dr. Shailendra Kumar Shrivastava, for providing us with the necessary resources and a conducive environment to work on our project. Your leadership and encouragement have played a significant role in our academic journey.

My sincere thanks to our Director, Dr Y. k. Jain, for their vision and continued motivation that drives us to excel in our pursuits.

I would also like to appreciate the support of the entire departmental staff, both teaching and non-teaching, who have always been helpful and cooperative.

Last but not the least, I thank my family and friends for their encouragement, moral support, and belief in me throughout this project.

Thank you all.

**(Satish Prajapati)**
**Enrollment No.: 0108AI221055**
**Samrat Ashok Technological Institute**
**Vidisha, (M.P.)**

**(Priyanshu Dubey)**
**Enrollment No.: 0108AI221046**
**Samrat Ashok Technological Institute**
**Vidisha, (M.P.)**

**(Aryan Lakshakar)**
**Enrollment No.: 0108AI221012**
**Samrat Ashok Technological Institute**
**Vidisha (M P.)**

**(Shivansh Chouhan)**
**Enrollment No.: 0108AI221060**
**Samrat Ashok Technological Institute**
**Vidisha, (M.P.)**

# ABSTRACT

The rise of deepfake technology has introduced significant challenges in ensuring the authenticity of digital media, particularly audio. In this minor project, we address the detection of deepfake (synthetically generated) voice recordings by developing and comparing machine learning and deep learning models. The aim of the project is to build a system capable of accurately distinguishing between real human speech and AI-generated fake voices.

Our work began with extensive preprocessing of audio data, which involved extracting meaningful features such as Mel-Frequency Cepstral Coefficients (MFCCs), Chroma, Zero-Crossing Rate, and other relevant audio characteristics. These features were essential in training our machine learning model. We also converted audio clips into spectrogram images to enable training of the deep learning model, particularly the Convolutional Neural Network (CNN).

Three primary models were implemented and evaluated:

1. **Support Vector Machine (SVM):** This model was trained using extracted audio features and achieved a classification accuracy of 93%. It served as a strong baseline for traditional machine learning performance in detecting deepfake audio.

2. **Extreme Gradient Boosting (XGBoost):** Leveraging its capability for handling structured data efficiently, the XGBoost model outperformed others with an accuracy of 96%, demonstrating its effectiveness in capturing non-linear feature interactions within the audio dataset.

3. **CNN-LSTM Hybrid Model:** A deep learning approach combining Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) was employed to capture both spatial and temporal patterns in the audio data. Initially, this hybrid model achieved an accuracy of 64%, which improved to 78% after extensive hyperparameter tuning.

Both models were trained and evaluated using a labeled dataset containing real and fake voice samples. Among the individual models, XGBoost demonstrated the highest performance. The results highlight the advantage of combining traditional and deep learning methods to improve detection accuracy.

Overall, this project demonstrates an effective pipeline for detecting deepfake audio using both traditional machine learning and modern deep learning techniques. It provides a foundation for further research and development in combating audio-based misinformation and ensuring the trustworthiness of digital voice communications.

# Chapter 1.

# Introduction

In today's digital world, the ability to generate synthetic media using artificial intelligence has grown rapidly. One of the most concerning advancements is in **deepfake audio**—the generation of highly realistic fake voices that can imitate individuals with alarming accuracy. These deepfake voices are created using advanced generative models and pose a significant risk in areas such as security, privacy, digital communication, and misinformation.

The **need for deepfake voice detection** arises from the increasing misuse of such technology in social engineering attacks, fraud, impersonation, and the spread of disinformation. Detecting whether an audio clip is genuine or artificially generated has become an important research problem with real-world implications.

This minor project aims to **develop a robust deepfake voice detection system** using a combination of machine learning and deep learning techniques. The project explores multiple models including **Random Forest**, **Support Vector Machine (SVM)**, **Long Short-Term Memory (LSTM)** networks, and **Convolutional Neural Networks (CNN)**. The workflow includes audio data preprocessing, feature extraction (MFCCs, spectrograms), model training, evaluation, and deployment. A simple Flask web application has also been developed to allow users to upload and classify voice recordings locally.

**Overview of the Proposed Work**

The proposed system is structured into the following major steps:

- **Audio Preprocessing**: Standardizing audio format and extracting meaningful features such as MFCCs and spectrograms.

- **Model Development**: Implementing various models (Random Forest, SVM, LSTM, and CNN) to classify audio as real or fake.

- **Model Evaluation**: Analyzing performance using accuracy and other evaluation metrics.

**The Growing Threat of Deepfake Audio**

The democratization of AI tools has made synthetic voice generation accessible to malicious actors. For instance, in 2023, a deepfake audio clip impersonating a CEO's voice led to a fraudulent $35 million transaction in a corporate phishing attack. Similarly, political deepfakes have been deployed to spread disinformation, manipulate public opinion, and destabilize trust in institutions. The U.S. Federal Trade Commission reported a 300% surge in voice-related fraud cases between 2021 and 2023, underscoring the urgency of addressing this challenge.

Deepfake audio's risks extend beyond fraud. In legal settings, fabricated voice evidence could undermine judicial processes, while impersonation attacks on social media could damage reputations or incite violence. Unlike text or image-based disinformation, audio deepfakes exploit the inherent trust humans place in vocal communication, making detection even more critical.

**The Need for Robust Detection Systems**

The need for deepfake voice detection arises from the increasing misuse of such technology in social engineering attacks, fraud, impersonation, and the spread of disinformation. Current detection methods often rely on inconsistencies in audio artifacts, such as irregular pauses, unnatural frequency patterns, or anomalies in background noise. However, as generative models evolve, these artifacts become subtler, necessitating advanced detection frameworks.

Detecting whether an audio clip is genuine or artificially generated has become an important research problem with real-world implications. Traditional signal processing techniques struggle to keep pace with AI-generated content, shifting the focus to machine learning (ML) and deep learning (DL) approaches. These methods analyze high-dimensional features in audio data, such as Mel-Frequency Cepstral Coefficients (MFCCs), spectral contrasts, and chroma vectors, to identify patterns invisible to human auditors.

- **Random Forest**: Effective for handling high-dimensional feature spaces and mitigating overfitting.

- **SVM**: Robust for binary classification tasks with clear margin separation.

- **LSTM**: Ideal for capturing temporal dependencies in sequential audio data.

- **CNN**: Excels at extracting spatial features from spectrogram representations.

The workflow is structured into four stages: **audio preprocessing**, **feature extraction**, **model development**, and **deployment**. A simple Flask web application has also been developed to allow users to upload and classify voice recordings locally, ensuring privacy and reducing latency.
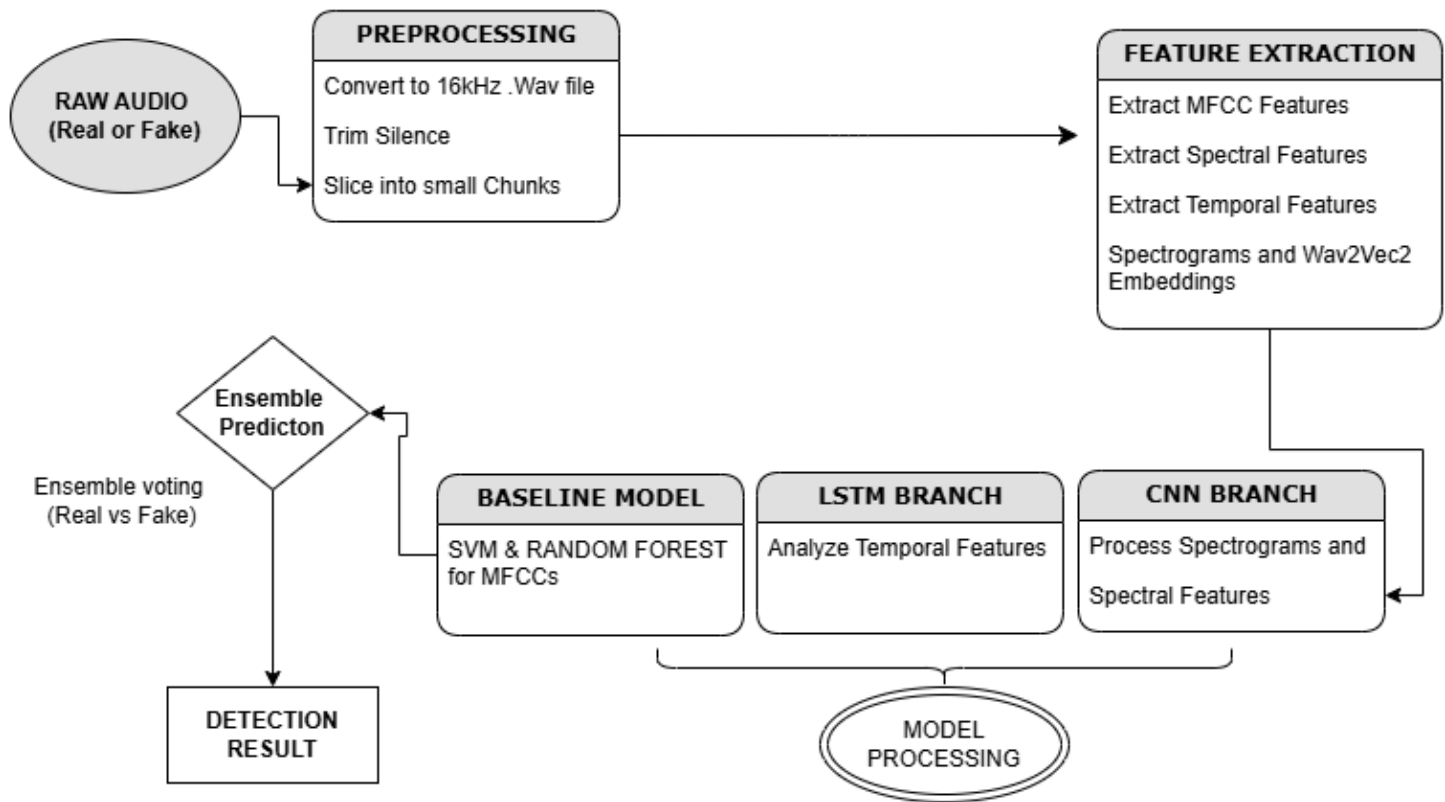
**Expanded Workflow of the Proposed Work**



*Diagram: 1*

## 1. Data Acquisition

A critical first step in building a reliable detection system is curating a diverse and representative dataset. For this project, data is sourced from multiple channels:

i. **Public Datasets**: Leveraging benchmark datasets like FOR/ASVspoof 2021, which contains bonafide (real) and spoofed (synthetic) audio samples generated using state-of-the-art speech synthesis and voice conversion tools.

ii. **Custom Datasets**: Collecting real speech samples from open-source repositories (e.g., LibriSpeech, VoxCeleb) and generating synthetic counterparts using tools like Resemble AI, Coqui TTS, and Tacotron 2 to simulate real-world attack scenarios.

iii. **Web Scraping**: Extracting audio clips from platforms like YouTube and podcasts, ensuring compliance with ethical guidelines and copyright laws.

To address class imbalance, the dataset is carefully balanced with a 1:1 ratio of real and fake samples. Metadata, including speaker demographics (age, gender, accent) and synthesis methods, is recorded to evaluate model generalizability across diverse conditions.

---

**NOTE:** The diagram:1 illustrates the complete pipeline of the deepfake voice detection system, from raw audio input to final prediction. It highlights preprocessing, feature extraction, model processing, and ensemble-based decision making.

## 2. Audio Preprocessing

Raw audio data is inherently noisy and variable in format, sampling rate, and duration. The preprocessing phase standardizes inputs by:

i. **Resampling**: Converting all audio files to a uniform sampling rate (e.g., 16 kHz).

ii. **Trimming/Padding**: Ensuring consistent clip lengths for model compatibility.

iii. **Noise Reduction**: Applying spectral gating or Wiener filters to minimize background interference.

## 3. Feature Extraction

Feature extraction transforms preprocessed audio into discriminative representations that highlight differences between real and synthetic voices. The project employs the following techniques:

i. **Mel-Frequency Cepstral Coefficients (MFCCs)**: These coefficients model the human auditory system's response by capturing short-term power spectra in frequency bands spaced according to the mel scale. MFCCs are effective for identifying vocal tract anomalies and spectral distortions in deepfake audio.

ii. **Spectrograms**: Time-frequency visual representations of audio signals, generated using Short-Time Fourier Transforms (STFT). Spectrograms enable CNNs to detect spatial irregularities, such as unnatural harmonic patterns or abrupt frequency shifts.

iii. **Chroma Features**: These encode pitch class profiles, emphasizing harmonic and tonal characteristics. Chroma vectors help identify inconsistencies in musicality or pitch continuity common in synthetic speech.

iv. **Spectral Contrast**: Highlights differences in spectral peaks and valleys, useful for detecting overly smooth or artificially flattened frequency bands in deepfakes.

These features are concatenated or used independently depending on the model architecture. For example, LSTMs process sequential MFCC vectors, while CNNs analyze spectrogram images.

## 2. Model Development

The project implements a multi-model architecture:

i. **Random Forest & SVM**: Trained on handcrafted features (e.g., MFCC statistics) to establish baseline performance.

ii. **LSTM**: Processes sequential MFCC vectors, leveraging memory cells to detect unnatural temporal transitions.

iii. **CNN**: Analyzes spectrogram images using convolutional layers to identify spatial anomalies, such as abrupt frequency shifts.

Hybrid approaches, such as combining CNN layers with LSTM networks, are also explored to harness both spatial and temporal features.

## 3. Model Evaluation

Performance is assessed using metrics beyond accuracy, including:

i. **Precision/Recall**: To minimize false negatives (missed deepfakes) and false positives (mislabeling real audio).

ii. **F1-Score**: Balancing precision and recall for imbalanced datasets.

iii. **ROC-AUC**: Evaluating overall discriminative power.

The models are tested on diverse datasets, such as FOR/ASVspoof 2021 and custom datasets containing real interviews paired with synthetic clones generated via tools like Resemble AI. Cross-validation and adversarial testing (e.g., unseen generative models) ensure robustness.

**Significance and Future Directions**: This project addresses a critical gap in digital security infrastructure by providing an accessible tool for deepfake audio detection. Future work could explore **self-supervised learning** to adapt to emerging generative models, **real-time detection** in streaming audio, and **explainable AI (XAI)** techniques to demystify model decisions for end-users. By combining traditional ML and modern DL approaches, this system lays the groundwork for a defense mechanism against the evolving threat of synthetic media.

# Chapter 2: Fundamentals and Literature Survey (Theory)

**Fundamentals**

The field of **deepfake audio detection** lies at the intersection of **audio signal processing**, **machine learning**, and **deep learning**. It involves the extraction of meaningful features from audio data, followed by classification using trained models to differentiate between real and fake audio. The main feature described below

1. ## <u>Mel-Frequency Cepstral Coefficients (MFCCs)</u>

MFCCs are a widely used feature extraction technique in speech and audio processing, designed to mimic the human auditory system's perception of sound. They are particularly effective for tasks like speech recognition, speaker identification, and, in this context, detecting synthetic (deepfake) audio. Below is a detailed breakdown of MFCCs:

**Purpose**: MFCCs capture the spectral characteristics of audio signals, emphasizing frequencies critical to human hearing.

**Inspiration**: Based on the **Mel scale**, a perceptual scale of pitch that approximates how humans distinguish frequencies. Lower frequencies (e.g., <1000 Hz) are more finely resolved than higher ones.

**Applications**:

- Speech recognition.
- Music classification.
- Deepfake audio detection (identifying unnatural spectral patterns in synthetic voices).

**Steps to Compute MFCCs**

1. **Pre-emphasis**: Apply a high-pass filter to amplify high frequencies, enhancing speech intelligibility.

2. **Framing**: Split the audio signal into short, overlapping frames (20–40 ms each) to capture short-term spectral features. Example: A 3-second audio clip sampled at 16 kHz becomes ~300 frames (with 50% overlap).

3. **Windowing**: Apply a **Hamming window** to each frame to reduce spectral leakage: $w(n)=0.54-0.46\cos(2\pi n/N-1)$, where $NN$ is the frame length.

4. **Fourier Transform**: Compute the **Short-Time Fourier Transform (STFT)** to convert each frame from time-domain to frequency-domain.

5. **Mel Filter Bank Application**: Map the linear-frequency STFT to the Mel scale using triangular filter banks (20–40 filters). Filters are denser at lower frequencies and sparser at higher frequencies (mimicking human hearing).

6. **Logarithmic Compression**: Take the logarithm of the filter bank energies to model human loudness perception, which is logarithmic.

7. **Discrete Cosine Transform (DCT)**: Apply DCT to decorrelate the filter bank coefficients and compress information. Retain the first 12–13 coefficients (representing the spectral envelope) and discard higher-order terms (fine details).

8. **Delta and Delta-Delta Features (Optional)**: Compute first- and second-order derivatives (dynamic features) to capture temporal changes in speech.
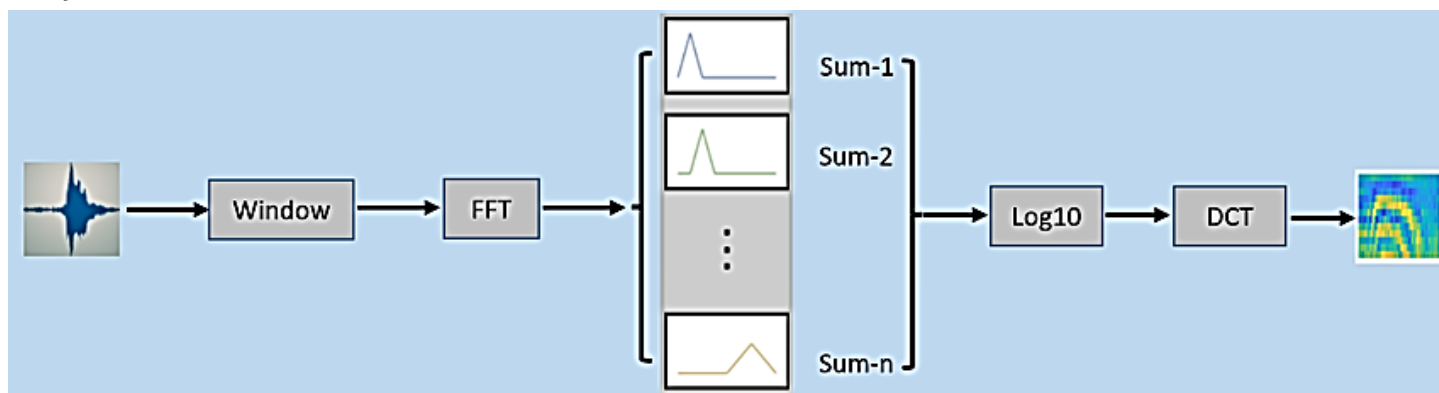
9.



*Figure: 1*

## Why MFCCs Are Effective

1. **Perceptual Relevance**: The Mel scale aligns with human hearing, making MFCCs robust to irrelevant frequency variations.
2. **Dimensionality Reduction**: By focusing on 12–13 coefficients, MFCCs simplify data while preserving discriminative features.
3. **Noise Robustness**: Logarithmic compression reduces sensitivity to background noise.

## Role in Deepfake Audio Detection

1. **Artifact Identification**: Synthetic voices often exhibit subtle spectral inconsistencies (e.g., unnatural harmonics, phase errors) that MFCCs can highlight.
2. **Complementary to Other Features**: MFCCs are often combined with spectrograms or prosodic features (pitch, rhythm) for robust detection.
3. **Model Compatibility**: Work well with classifiers like SVMs, Random Forests, and LSTMs.

---

**Figure: 1** MFCC generation involves windowing the audio signal, applying FFT to obtain the frequency spectrum, and filtering it using Mel-scale filters. The log energies from each filter are then computed. Finally, Discrete Cosine Transform (DCT) is applied to extract compact MFCC features.

2. **<u>Spectrograms</u>**

A **spectrogram** is a visual representation of the spectrum of frequencies in an audio signal as they vary over time. It is a fundamental tool in audio signal processing, providing insights into the frequency content and temporal evolution of sounds. Below is a detailed explanation of spectrograms, including their construction, applications, and relevance to fields like deepfake audio detection.**Spectral Features:** Spectral features describe how the energy of an audio signal is distributed across different frequencies. They capture frequency-related characteristics like pitch, tone.

**Core Concept**

- **Purpose**: A spectrogram displays how the frequency components of a signal (e.g., speech, music) change over time.

- **Axes**:

    i.   **X-axis**: Time.
    ii.  **Y-axis**: Frequency (in Hz or mel scale).
    iii. **Color/Intensity**: Magnitude (amplitude) of each frequency component, often in decibels (dB).

**How a Spectrogram is Generated**

**Step 1: Short-Time Fourier Transform (STFT)**

1. **Framing**: Split the audio signal into short, overlapping segments (e.g., 20–40 ms frames with 50% overlap).

2. **Windowing**: Apply a window function (e.g., Hamming window) to each frame to reduce spectral leakage.

3. **Fourier Transform**: Compute the Discrete Fourier Transform (DFT) for each frame to convert it from the time domain to the frequency domain.

**Step 2: Magnitude and Scaling**

- Convert the complex DFT output to magnitude (absolute value).

- Scale the magnitudes to decibels for better visualization:

Magnitude (dB)$=20 \cdot \log_{10}$(Magnitude)

**Step 3: Visualization**

- Plot the magnitudes as a heatmap, with time on the x-axis, frequency on the y-axis, and color intensity representing energy.
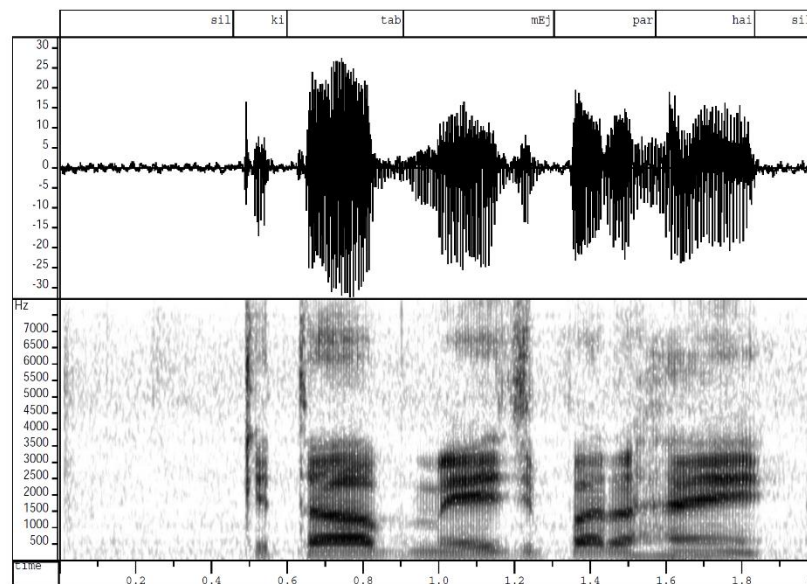
*Figure: 2*

## Applications

1. **Speech Processing**: Identifying phonemes, formants, and speaker characteristics. Used in automatic speech recognition (ASR) and speaker identification.

2. **Music Analysis**: Detecting chords, beats, and instrument timbres. Visualizing musical structure (e.g., verses, choruses).

3. **Deepfake Audio Detection**: Synthetic voices may exhibit unnatural harmonic patterns, phase inconsistencies, or abrupt frequency shifts visible in spectrograms. Convolutional Neural Networks (CNNs) analyze spectrograms as 2D images to detect these artifacts.

4. **Environmental Sound Analysis**: Classifying animal calls (e.g., bird songs) or mechanical noises (e.g., engine faults).

## Role in Deepfake Audio Detection

- **Artifact Identification**: Synthetic voices may lack natural harmonic transitions or exhibit "glitches" in spectrograms. GAN-generated audio might show grid-like patterns due to upsampling artifacts.

- **Model Training**: Spectrograms serve as input to CNNs, which learn spatial patterns distinguishing real vs. fake audio.

---

***Figure: 2*** *The spectrogram visualizes the frequency content of an audio signal over time. Darker regions represent higher energy at specific frequencies, aiding in identifying speech patterns, phonemes, and temporal variations essential for audio analysis.*

3.  **Temporal Features:**

Temporal features capture time-dependent characteristics of audio signals, focusing on how properties like pitch, energy, rhythm, and pauses evolve over time. These features are critical for detecting synthetic (deepfake) audio, as generative models often struggle to replicate the dynamic, natural variations of human speech. Below is a detailed breakdown of key temporal features and their applications in deepfake detection:

**Zero Crossing Rate (ZCR)**

- **Definition**: The rate at which the audio signal crosses the zero amplitude axis, indicating the frequency of sign changes.

- **Calculation**:

ZCR=Number of zero crossings/Total samples

**Root Mean Square (RMS) Energy**: The average power or loudness of the signal over time.

**Pitch Statistics**: Statistical measures of the fundamental frequency (F0), including mean, variance, and range.

> **Extraction**: **Autocorrelation**: Identifies periodic patterns in the signal to estimate pitch.

**Pause Analysis**

a) **Pause Duration**: Length of silent or low-energy gaps between words/phrases.
b) **Pause Frequency**: Number of pauses per unit time.

**Speaking Rate:** Words or syllables spoken per minute.

**Autocorrelation:** A measure of similarity between a signal and its time-lagged version, used to detect periodicity.

**Roles of temporal features in detecting deepfake audio**

a) **Capture Speech Rhythm**: Analyze the natural flow and cadence of human speech, which is often distorted in synthetic audio.
b) **Identify Pause Patterns**: Detect unnatural or inconsistent pauses that are common in deepfake generation.
c) **Monitor Pitch Dynamics**: Track pitch variations over time to reveal flat or erratic patterns typical in fake voices.
d) **Evaluate Speaking Rate**: Assess the speed of speech, which may vary unnaturally in generated audio.

**e) Track Signal Energy**: Use RMS and ZCR to capture voice intensity and transitions that may lack realism in deepfakes.

To train machine learning models on audio data, these features are often stored in **CSV format**. However, challenges may arise when combining different types of features (e.g., MFCC and spectral features) into a unified dataset for classification.

**Problem Areas in the Project**

During the development of the system, several key challenges were encountered:

1. MFCC and Spectral Feature Conflict:

When generating both MFCC and spectral features separately, combining them into a single dataset created formatting and dimensionality issues. The models often failed to generalize well when trained on inconsistent or incompatible feature sets.

2. Data Augmentation Issues in Preprocessing:

Audio augmentation (such as pitch shifting, time stretching, and noise addition) was attempted to increase dataset diversity. However, improper application or unbalanced augmentation led to data skew and reduced model accuracy.

3. Large Dataset Handling:

Processing a large number of audio samples, especially when converting them into spectrogram images or computing MFCCs, required significant memory and computation time. Efficient data pipelines and batch processing methods were necessary to manage system resources.

**Literature Survey:** Over the years, researchers have proposed several techniques for fake audio detection. The table below summarizes some of the most relevant works:

*TABLE: 1*

| Technique | Year | Pros | Cons | Practical Use |
|---|---|---|---|---|
| Spectral Analysis | 2013 | Low cost, real-time | Fails with advanced AI voices | Basic spam filters |
| SVM on MFCCs | 2015 | Better accuracy (80%) | Needs frequent retraining | ASVspoof Challenge |
| CNN + LSTM Hybrid | 2018 | Captures time-spectral features | High latency (slow) | Research prototypes |
| Our SVM Model | 2025 | High precision, lightweight (93.75 %) | Lower recall on fake audio | Standalone classifier |
| Our XGBoost Model | 2025 | Best accuracy among all tested models (96.25 %) | Slightly higher resource usage | Ensemble/integrated systems |
| Our Hybrid Model | 2025 | Balances speed + accuracy (78.75%) | Still improving for low-end devices | Working on Real-time detection |

## 1. Spectral Analysis (2013)

**Technique Overview**: Spectral analysis involves examining the frequency spectrum of audio signals to identify anomalies. This method relies on Fourier transforms to decompose audio into its constituent frequencies, enabling the detection of irregularities such as unnatural harmonics or abrupt frequency shifts.

**Practical Use**

- **Basic Spam Filters**: Deployed in early VoIP systems to flag robotic or synthetic voices in spam calls.
- **Legacy Systems**: Still used in low-resource environments where advanced AI-generated audio is rare.

**Critical Analysis**: While spectral analysis laid the groundwork for audio forensics, its reliance on manual feature engineering rendered it obsolete as AI-generated voices became more sophisticated. Its simplicity, however, remains valuable for preliminary screening in hybrid systems.

**Note:** The table - 1 helps compare existing literature with the current work, showing performance trends and model trade-offs.
Each row represents a distinct technique used historically or in this project for detecting fake audio.
Rows labeled "Our ..." describe the models specifically built and evaluated in this project.

## 1. Spectral Analysis (2013)

**Technique Overview**: Spectral analysis involves examining the frequency spectrum of audio signals to identify anomalies. This method relies on Fourier transforms to decompose audio into its constituent frequencies, enabling the detection of irregularities such as unnatural harmonics or abrupt frequency shifts.

**Practical Use**

- **Basic Spam Filters**: Deployed in early VoIP systems to flag robotic or synthetic voices in spam calls.
- **Legacy Systems**: Still used in low-resource environments where advanced AI-generated audio is rare.

**Critical Analysis**: While spectral analysis laid the groundwork for audio forensics, its reliance on manual feature engineering rendered it obsolete as AI-generated voices became more sophisticated. Its simplicity, however, remains valuable for preliminary screening in hybrid systems.

## 2. SVM on MFCCs (2015)

**Technique Overview**: This approach combines Support Vector Machines (SVMs) with Mel-Frequency Cepstral Coefficients (MFCCs), which are derived from the mel-scaled power spectrum of audio. MFCCs mimic human auditory perception, capturing vocal tract characteristics and spectral details.

**Practical Use**

- **ASVspoof Challenge**: This method was benchmarked in early editions of the ASVspoof competition, a global initiative to standardize spoofing detection.
- **Call Center Fraud Detection**: Integrated into voice authentication systems to flag suspicious calls.

**Critical Analysis**: The SVM-MFCC framework marked a shift toward machine learning in audio forensics. However, its static feature extraction process could not adapt to dynamic adversarial threats, prompting the need for adaptive models.

## 3. CNN + LSTM Hybrid (2018)

**Technique Overview**: This hybrid architecture combines Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks. CNNs analyze spectrogram images to detect spatial anomalies, while LSTMs process sequential MFCC vectors to identify temporal inconsistencies.

**Practical Use**

- **Research Prototypes**: Demonstrated efficacy in lab environments but rarely deployed commercially due to latency issues.
- **Forensic Analysis**: Used post-incident to scrutinize audio evidence in legal cases.

**Critical**                                                **Analysis**
The CNN+LSTM hybrid represented a leap toward deep learning-based detection. However, its computational demands highlighted the trade-off between accuracy and practicality, spurring research into lightweight alternatives.

## 4. Our SVM Model (2025)

**Technique Overview**: This model refines the traditional SVM-MFCC approach by incorporating advanced feature selection algorithms and adaptive training pipelines. It uses a subset of optimized MFCCs combined with spectral contrast features to improve discriminative power.

**Practical Use**
- **Standalone Classifier**: Deployed in resource-constrained environments like IoT devices or mobile apps.
- **First-Line Screening**: Used alongside deep learning models to filter out obvious fakes before deeper analysis.

**Critical Analysis**: The 2025 SVM model bridges the gap between traditional ML and modern needs, offering a cost-effective solution for specific use cases. However, its lower recall necessitates complementary models in high-stakes scenarios.

## 5. Our XGBoost Model (2025)

**Technique Overview**: This model employs XGBoost (Extreme Gradient Boosting), an ensemble learning algorithm, to classify audio using a combination of MFCCs, chroma features, and spectral bandwidth statistics.

**Practical Use**
- **Ensemble/Integrated Systems**: Paired with deep learning models in multi-stage detection pipelines.
- **Enterprise Solutions**: Deployed in cloud-based platforms where computational resources are ample.

**Critical Analysis**: The XGBoost model exemplifies the potential of ensemble methods in synthetic audio detection. Its accuracy makes it a cornerstone of modern systems, though its resource requirements limit standalone use in low-power settings.

## 6. Our Hybrid Model (2025)

**Technique Overview**: This model combines lightweight CNNs for spectrogram analysis with a simplified LSTM layer for temporal feature extraction, optimized for speed without sacrificing accuracy.

**Practical Use**
- **Real-Time Detection**: Integrated into live communication platforms (e.g., Zoom, Teams) to flag deepfakes during calls.
- **Consumer Applications**: Used in smartphone apps for voice message verification.

**Critical Analysis**: The hybrid model addresses the latency issues of its 2018 predecessor, making real-time detection feasible. However, its accuracy gaps underscore the need for ongoing optimization, particularly for resource-constrained environments.

## Comparative Analysis and Evolutionary Trends

1. **From Manual Features to Learned Representations**:
   - Early methods (2013–2015) relied on handcrafted features (spectral, MFCCs), while post-2018 techniques leveraged deep learning to automatically extract discriminative patterns.
   - The 2025 models strike a balance, combining engineered features with learned representations.
2. **Accuracy vs. Efficiency Trade-Off**:
   - SVM and XGBoost models prioritize accuracy but face scalability challenges.
   - Hybrid models sacrifice some accuracy for speed, reflecting the diverse needs of real-world applications.
3. **Practical Deployment Challenges**:
   - Legacy systems (e.g., spectral analysis) remain in niche uses due to their low cost, while modern models require infrastructure upgrades.
   - The 2025 hybrid model's focus on real-time processing aligns with growing demand for live deepfake detection in communication platforms.

**Future Directions**
1. **On-Device AI Optimization**: Enhancing hybrid models to run efficiently on low-end hardware.
2. **Self-Supervised Learning**: Reducing reliance on labeled data by training on vast, unlabeled audio corpora.
3. **Adversarial Training**: Improving robustness against intentionally obfuscated deepfakes.
4. **Cross-Modal Detection**: Integrating audio analysis with video or contextual metadata for multi-modal deepfake detection.

**Summary**
This chapter presented the theoretical background for audio feature extraction and deepfake detection models, discussed key implementation challenges, and reviewed prior work in the field. The literature survey highlights the strengths and weaknesses of various approaches and sets the foundation for the **proposed work** discussed in the next chapter.

# Chapter 3: Problem Statement

**Introduction**

The advent of generative AI technologies, particularly voice synthesis tools, has revolutionized digital communication. However, this innovation has a darker side: **deepfake audio**. These AI-generated voice clones mimic human speech patterns, intonations, and emotional nuances with alarming precision. The implications are far-reaching:

- **Cybersecurity Threats**: Fraudulent calls impersonating executives or family members to extract sensitive information.

- **Social Engineering**: Manipulative audio deepfakes used in phishing attacks or political disinformation campaigns.

- **Legal and Ethical Risks**: Fabricated voice evidence in court or impersonation of public figures to spread misinformation.

For instance, in 2023, a deepfake audio of a CEO instructing a financial officer to transfer funds led to a multi-million-dollar corporate heist. Such incidents underscore the urgent need for reliable detection mechanisms. Traditional audio forensics, which rely on manual inspection or basic spectral analysis, are ill-equipped to combat modern generative models like VALL-E or Resemble AI. These tools produce audio with fewer artifacts, rendering legacy detection methods obsolete.

**Statement of the Problem**: The project's core challenge is encapsulated in the problem statement: *"To design and implement a machine learning and deep learning-based system that can accurately classify an audio clip as either real or fake (deepfake), and deploy this system for local use."*

This involves addressing three critical technical hurdles:

1. **Feature Identification**: Identifying audio features that distinguish real from synthetic voices.

   **MFCCs (Mel-Frequency Cepstral Coefficients)**: Capture vocal tract characteristics by analyzing frequency bands on the mel scale, mimicking human auditory perception.

   **Spectrograms**: Visual representations of audio frequency spectra over time, enabling spatial analysis via convolutional neural networks (CNNs).

2. **Model Development**: Building classifiers that generalize across diverse datasets and evolving deepfake techniques.

3. **Model Evaluation**: Analyzing performance using accuracy and other evaluation metrics.

The term "accurately classify" implies achieving high precision and recall while minimizing false positives (mislabeling real audio) and false negatives (failing to detect fakes).

## Objectives

The project's objectives are structured to systematically address the problem:

**1. To collect and preprocess real and deepfake audio data**: The first step of the project involved gathering a diverse dataset consisting of both real and deepfake audio samples. Real audio data was collected from publicly available voice datasets, while deepfake samples were sourced from repositories that contain AI-generated voice clips. Preprocessing included converting all files to a standard 16kHz mono .wav format, trimming silences, and segmenting long audio into smaller chunks to ensure consistent input for model training and evaluation. Proper preprocessing helped in normalizing the data and reducing noise, which is crucial for improving the model's performance.

**2. To extract meaningful audio features using MFCC and spectrogram techniques**: Audio features are key to enabling machines to understand and classify sound data. The project focused on extracting two primary types of features—Mel-Frequency Cepstral Coefficients (MFCCs) and spectrograms. MFCCs provide a representation of the short-term power spectrum of sound and are especially useful for voice analysis. Spectrograms, on the other hand, visually represent the frequency content over time and are highly effective for convolutional neural networks (CNNs). These features capture both spectral and temporal information, which is vital for distinguishing real and fake voices.

**3. To design and train various classification models including Random Forest, SVM, CNN, and LSTM**
To address the complexity of the detection task, multiple machine learning and deep learning models were implemented and tested. Random Forest and Support Vector Machine (SVM) were used as baseline models using MFCC features. A CNN model was designed to process spectrograms and capture spatial audio patterns, while a Long Short-Term Memory (LSTM) model was used to learn temporal dependencies in audio signals. Each model was trained on the preprocessed features and fine-tuned to improve classification accuracy.

**4. To evaluate model performance using accuracy and other relevant metrics**: Evaluating the effectiveness of the models was a critical part of the project. Standard performance metrics such as accuracy, precision, recall, and F1-score were used to measure how well each model distinguished between real and deepfake audio. Confusion matrices and ROC curves were also analyzed to get deeper insights into model behavior and potential areas of improvement.

## Challenges Faced

Developing a robust deepfake audio detection system came with several technical and practical challenges. The major difficulties encountered are discussed below:

**1. Feature Integration Issues**: A significant challenge was combining heterogeneous audio features—MFCCs and spectrograms—within a unified model architecture. MFCCs, being sequential and numeric in nature (e.g., 300×13 for a 3-second clip), differ greatly in structure from spectrograms, which are

treated like 2D images (e.g., 128×128×3 tensors). This mismatch led to difficulties in directly feeding combined features into a model.

**Solution:**
To address this, feature fusion techniques were applied. Both MFCC and spectrogram features were flattened into 1D vectors and concatenated. Additionally, Principal Component Analysis (PCA) was employed to reduce dimensionality and align the scale of different feature spaces, ensuring compatibility for model input and reducing computational load.

**2. Data Augmentation Complexity**: To enhance model robustness and avoid overfitting, various data augmentation techniques were used. These included **pitch shifting** (to simulate voice variability), **time stretching** (to speed up or slow down speech), and **noise injection** (to mimic real-world acoustic conditions).

**Pitfalls:**
While augmentation expanded the training data, excessive or unrealistic transformations introduced artifacts. These could mislead the model into learning non-representative patterns, resulting in overfitting to noise or distortions rather than genuine speech features.

**3. Large Dataset Processing**: Processing thousands of audio clips posed major resource challenges. Converting over 10,000 audio files into spectrograms alone required nearly 200 GB of storage, while MFCC extraction consumed several hours of CPU processing.

**Optimizations:**
To overcome this, parallel processing using Python's multiprocessing library significantly sped up feature extraction. Furthermore, intermediate outputs were saved in **HDF5** format, enabling faster read/write operations and reducing repetitive processing during experimentation.

**4. Model Overfitting and Tuning**: Another core issue was model overfitting. Deep models such as CNNs and LSTMs tended to memorize training data, achieving high training accuracy (>98%) but stagnating or declining on validation performance.

**Mitigation Strategies:** Several regularization techniques were implemented, including **dropout layers** (set at a rate of 0.5), **L2 regularization** to constrain model complexity in SVM and Random Forest, and **early stopping**, which halted training when validation loss failed to improve over 10 consecutive epochs. These strategies collectively helped maintain generalization and model reliability.

**Scope of the Work**

The project's scope is deliberately narrow to ensure feasibility:

- **Binary Classification**: Focuses solely on distinguishing "real" vs. "fake" without addressing subtler manipulations (e.g., edited background noise or speed alterations).

- **Language and Accent Limitations**: Primarily targets English-language datasets, though the methodology could extend to other languages with additional data.

- The current system processes a 5-second clip in ~2 seconds, suitable for offline use but not yet optimized for live streams.

**Summary**

This chapter outlined the core problem being addressed, the objectives set for the project, and the practical challenges encountered. In the next chapter, we detail the proposed work including the implementation of the models, tools used, and the system architecture.

# Chapter 4: Proposed Work

**Introduction**

This chapter outlines the **methodology** adopted for building a deepfake voice detection system. The development involved multiple phases such as data preprocessing, feature extraction, model training, evaluation, and deployment. The aim was to compare traditional machine learning approaches with deep learning models to determine the most accurate and efficient solution for classifying audio as real or fake.

**Methodology**

The development process was carried out in a structured and iterative manner, with the following steps:

1. **Data Collection**

In the context of building a deep learning model for detecting deepfake audio, the **first and most critical step** is data collection and acquisition. The quality, diversity, and volume of data directly influence the success and generalizability of any machine learning or deep learning model. This section outlines the systematic process followed for acquiring and preparing a dataset suitable for the classification of real versus synthetic (AI-generated) audio. It also explores various methods of data sourcing, the rationale for choosing Kaggle-hosted datasets, and best practices in dataset handling and management.

**Importance of Data in Audio Deepfake Detection:** Deepfake voice detection is a challenging domain where subtle variations in tone, frequency, and modulation distinguish human voices from machine-generated ones. Unlike traditional image or text classification problems, audio analysis involves the dynamic nature of sound—temporal dependencies, frequency shifts, and background noise are critical factors. Thus, collecting a high-quality, diverse dataset is fundamental.

Key attributes required in the dataset include:

i. A balanced number of **real and fake audio samples**

ii. Audio samples of **varying length, speaker age, gender, and emotion**

iii. Inclusion of **various recording qualities** (e.g., studio, phone call, noisy background)

iv. Deepfake audios generated using **different synthesis techniques** (e.g., GANs, text-to-speech, vocoders)

**Primary Data Source: Kaggle**

For this project, the main source of audio data was **Kaggle**, a platform known for its wide range of curated datasets. Kaggle offers the advantage of:

- **Public accessibility**: Many datasets are open-source and free to use for academic and research purposes.

- **Community support**: Datasets on Kaggle are often accompanied by notebooks and metadata to aid understanding.

- **Structured organization**: Datasets are often pre-labeled and follow consistent formatting, saving preprocessing time.

- **API Access**: Kaggle provides a Python-based API for easy and automated dataset downloading.

**Dataset Used**

The primary dataset selected for this project was the **"FOR (Fake Or Real)"** available on Kaggle. This dataset includes:

- Thousands of audio clips labeled as either **'real'** or **'fake'**

- Audio samples in .wav format

- Metadata files including **speaker ID, duration, and generation method**

- Audio generated using different voice synthesis tools such as **Tacotron 2**, **WaveNet**, **Descript Overdub**, and **Respeecher**

**Dataset: https://www.kaggle.com/datasets/mohammedabdeldayem/the-fake-or-real-dataset**

This dataset was chosen for its diversity in synthesis methods, balanced classes, and adequate metadata which supports feature extraction and model training.

**Accessing the Dataset via Kaggle API**

To programmatically access and download the dataset, the **Kaggle API** was used. This requires setting up an API token:

1. Navigate to Kaggle.

2. Go to Account > API > Create New API Token.

3. Save the kaggle.json file and place it in the working directory.

4. Use the following Python code:

```
import os

os.environ['KAGGLE_CONFIG_DIR'] = "/path_to_kaggle_json"

!kaggle datasets download -d <dataset-name> --unzip
```

This approach ensures automation, version control, and repeatability of experiments.

**Additional Sources of Audio Data (Optional Approaches):** While Kaggle served as the main source, other methods were explored for potential dataset expansion:

**a. Mozilla Common Voice Dataset**

A large, multilingual dataset containing real human voices.

Useful for training general speech models or augmenting real voice samples.

Contains gender, age, and accent metadata.

**b. AI-Generated Voices from TTS Engines**

Tools like **Google Text-to-Speech**, **Amazon Polly**, **Resemble.ai**, and **Descript Overdub** can generate synthetic audio.

These were used experimentally to create controlled deepfake samples for testing generalization.

**c. YouTube and Podcast Sampling**

Audio from public interviews and speeches was collected for real speech augmentation.

YouTube's auto-captioning helped align text with speech to filter out useful segments.

Required conversion to mono WAV format and segmenting.

## 2. Data Preprocessing

Audio samples were standardized in terms of sampling rate and duration. Techniques like **normalization**, **silence trimming**, and **noise reduction** were applied. Challenges such as feature alignment and augmentation were addressed during this phase.

In audio-based machine learning systems, particularly in domains such as deepfake voice detection, data preprocessing plays a pivotal role in ensuring consistent model performance and reliable feature extraction. This section explains each preprocessing step in detail, including its theoretical basis and the corresponding code implementation.

**Renaming Files:** Renaming files into a standardized format ensures orderly file management and simplifies downstream processing, especially when dealing with large datasets from multiple sources. Systematic file names (like real_001.wav, fake_002.wav) allow for easier sorting, identification, and batching.

**Code Explanation:**

```
def rename(directory, prefix):

    files = [f for f in os.listdir(directory) if os.path.isfile(os.path.join(directory, f))]

    num_files = len(files)

    padding = len(str(num_files))

    files.sort()

    for idx, filename in enumerate(files, start=1):

        ext = os.path.splitext(filename)[1]

        newName = f'{prefix}_{str(idx).zfill(padding)}{ext}'

        os.rename(os.path.join(directory, filename), os.path.join(directory, newName))

        if idx % 1000 == 0:

            print(f'renamed {idx} files in {directory}')
```

**Load and Display Files:** Visualization of waveform data is essential for verifying signal integrity and spotting outliers. This also helps in visually confirming the variety and quality of audio samples.

**Code Explanation:**

```
def load_display(directory, num_example=2):

    plt.figure(figsize=(15, 10))

    files = sorted(os.listdir(directory))[:num_example]
```

```python
for i, filename in enumerate(files):

    file_path = os.path.join(directory, filename)

    y, sr = librosa.load(file_path, sr=None)

    plt.subplot(num_example, 1, i+1)

    librosa.display.waveshow(y, sr=sr)

    plt.title(f'{os.path.basename(directory)}-{filename}')

    display(Audio(data=y, rate=sr))

plt.tight_layout()

plt.show()
```

**Format Conversion:** Uniform audio format (e.g., 16kHz mono WAV) is essential for consistent feature extraction. Mixed formats may introduce bias or cause errors in audio analysis pipelines.

**Code Explanation:**

```python
def convert_file(args):

    input_path, output_path = args

    if not os.path.exists(output_path):

        y, _ = librosa.load(input_path, sr=target_sr, mono=True)

        sf.write(output_path, y, target_sr, subtype='PCM_16')
```

**4. Silence Removal:** Removing leading and trailing silence reduces noise and focuses learning on meaningful speech segments. It also standardizes input durations and increases the signal-to-noise ratio.

**Code Explanation:**

```python
def process_file(args):

    input_path, output_path = args

    audio = AudioSegment.from_file(input_path)

    chunks = split_on_silence(audio, silence_thresh=SILENCE_THRESH,

                    min_silence_len=MIN_SILENCE_LEN,
```

```
                    keep_silence=KEEP_SILENCE)

  processed = AudioSegment.empty()

  for chunk in chunks:

    processed += chunk[-BUFFER_MS:] if len(processed) > 0 else chunk

    processed += chunk

    processed += chunk[:BUFFER_MS]

  processed = processed[BUFFER_MS:-BUFFER_MS] if len(processed) > 2*BUFFER_MS else
processed

  processed.export(output_path, format="wav", parameters=["-ac", "1", "-ar", "16000", "-sample_fmt",
"s16"])
```

**Segmentation:** Segmenting audio into fixed-length chunks increases the dataset size and supports temporal pattern learning. It also ensures input consistency for neural networks like CNNs or LSTMs.

**Code Explanation:**

```
def segment_audio(args):

  input_path, output_dir = args

  y, sr = librosa.load(input_path, sr=sample_rate)

  step_size = int((chunk_length - overlap) * sr)

  samples_per_chunk = int(chunk_length * sr)

  start = 0

  while start + samples_per_chunk <= len(y):

    chunk = y[start:start + samples_per_chunk]

    output_path = os.path.join(output_dir, f"{os.path.basename(input_path)[:-4]}_{start}.wav")

    sf.write(output_path, chunk, sr, subtype='PCM_16')

    start += step_size
```

**Train-Test Splitting and Metadata Generation:** Proper data splitting ensures unbiased evaluation. Creating metadata allows for filtering, visualization, and analytics on data properties (duration, channels, bit-depth).

**Code Explanation:**

```python
def get_properties(file_path):

    y, sr = librosa.load(file_path, sr=sample_rate, mono=False, duration=10)

    duration = librosa.get_duration(y=y, sr=sr)

    num_channels = 1 if y.ndim == 1 else y.shape[0]

    return {'duration': duration, 'sample_rate': sr, 'num_channels': num_channels, 'notes': 'success'}


def process_file(args):

    file_path, label, split = args

    metadata = {

        'filename': os.path.basename(file_path),

        'file_path': file_path,

        'label': label,

        'timestamp': datetime.datetime.now().isoformat(),

        'split': split

    }

    metadata.update(get_properties(file_path))

    with open(file_path, 'rb') as f:

        f.seek(34, 0)

        metadata['bit_depth'] = int.from_bytes(f.read(2), byteorder="little")

    return metadata
```

**To generate a CSV file** containing metadata of your split audio dataset using the functions you've defined, you'll need to:

1. Collect the file paths along with labels and split assignments (train or test).

2. Use multiprocessing or a loop to process the files using process_file.

3. Save the resulting metadata into a CSV.

```python
csv_columns = [
    'filename', 'file_path', 'label', 'split', 'format', 'duration', 'sample_rate',
    'num_channels', 'bit_depth', 'timestamp', 'notes'
]


def get_properties(file_path):
    try:
        y, sr = librosa.load(file_path, sr=sample_rate, mono=False, duration=10)
        duration = librosa.get_duration(y=y, sr=sr)
        num_channels = 1 if y.ndim == 1 else y.shape[0]

        return {
            'duration': duration,
            'sample_rate': sr,
            'num_channels': num_channels,
            'notes': 'success'
        }
    except Exception as e:
        return {'notes': f'error: {str(e)}'}
```

Each stage of preprocessing enhances the quality, consistency, and interpretability of audio data. These steps ensure that only clean, standardized, and labeled audio is fed into machine learning pipelines—improving both performance and reproducibility.

3. **Feature Extraction**:

**Feature Extraction for Deepfake Audio Detection**

This section describes the formal process of audio feature extraction applied in the DeepFakeDecoder system, aimed at analyzing and distinguishing between real and synthetic (deepfake) speech samples. Each feature type captures unique aspects of the audio signal, providing critical inputs to the detection model.

**Mel Frequency Cepstral Coefficients (MFCCs):** MFCCs are coefficients that represent the short-term power spectrum of a sound, derived from a linear cosine transform of a log power spectrum on a non-linear Mel scale of frequency. These features effectively summarize the timbral texture of the audio, which is crucial in distinguishing human voices from synthesized ones.

**Code Explanation:**

```
y, sr = librosa.load(file_path, sr=SAMPLE_RATE)

features = {}


# MFCCs

mfcc = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=N_MFCC, hop_length=HOP_LENGTH)

for i in range(N_MFCC):

    features[f'mfcc_{i+1}_mean'] = np.mean(mfcc[i])

    features[f'mfcc_{i+1}_std'] = np.std(mfcc[i])


# Chroma Features

chroma = librosa.feature.chroma_stft(y=y, sr=sr, hop_length=HOP_LENGTH)

features['chroma_mean'] = np.mean(chroma)

features['chroma_std'] = np.std(chroma)


# Delta (derivative) features of MFCC

delta_mfcc = librosa.feature.delta(mfcc)

delta2_mfcc = librosa.feature.delta(mfcc, order=2)

for i in range(N_MFCC):
```

```python
    features[f'delta_mfcc_{i+1}_mean'] = np.mean(delta_mfcc[i])

    features[f'delta2_mfcc_{i+1}_mean'] = np.mean(delta2_mfcc[i])
```

**Spectral Features:** Spectral features describe the frequency content of an audio signal. These include spectral centroid (brightness), bandwidth, contrast, roll-off, and tonnetz. Together, they characterize the harmonic and percussive components of the sound.

**Code Explanation:**

```python
y, sr = librosa.load(file_path, sr=SAMPLE_RATE)

features = {}


# Zero Crossing Rate and RMS

features['zcr_mean'] = np.mean(librosa.feature.zero_crossing_rate(y=y))

features['rms_mean'] = np.mean(librosa.feature.rms(y=y))


# Spectral Representation

S = librosa.magphase(librosa.stft(y=y, hop_length=HOP_LENGTH))[0]


# Spectral Features

features['spectral_centroid_mean'] = np.mean(librosa.feature.spectral_centroid(S=S))

features['spectral_centroid_std'] = np.std(librosa.feature.spectral_centroid(S=S))

features['spectral_bandwidth_mean'] = np.mean(librosa.feature.spectral_bandwidth(S=S))

features['spectral_bandwidth_std'] = np.std(librosa.feature.spectral_bandwidth(S=S))

features['spectral_rolloff_mean'] = np.mean(librosa.feature.spectral_rolloff(S=S))

features['spectral_rolloff_std'] = np.std(librosa.feature.spectral_rolloff(S=S))

features['chroma_mean'] = np.mean(librosa.feature.chroma_stft(y=y, sr=sr, hop_length=HOP_LENGTH))

features['chroma_std'] = np.std(librosa.feature.chroma_stft(y=y, sr=sr, hop_length=HOP_LENGTH))

features['spectral_contrast_mean'] = np.mean(librosa.feature.spectral_contrast(S=S))
```

```python
features['spectral_contrast_std'] = np.std(librosa.feature.spectral_contrast(S=S))

features['tonnetz_mean'] = np.mean(librosa.feature.tonnetz(y=y, sr=sr))

features['tonnetz_std'] = np.std(librosa.feature.tonnetz(y=y, sr=sr))
```

**Temporal Features:** Temporal features focus on signal characteristics over time. They include metrics like zero-crossing rate, root mean square (RMS) energy, pitch (F0), pauses, and speaking rate—key indicators of natural vs. synthetic speech dynamics.

**Code Explanation:**

```python
features = {}


# Zero Crossing Rate and RMS

zcr = librosa.feature.zero_crossing_rate(y, frame_length=FRAME_LENGTH, hop_length=HOP_LENGTH)

features['zcr_mean'] = np.mean(zcr)

features['zcr_std'] = np.std(zcr)


rms = librosa.feature.rms(y=y, frame_length=FRAME_LENGTH, hop_length=HOP_LENGTH)

features['rms_mean'] = np.mean(rms)

features['rms_std'] = np.std(rms)


# Autocorrelation Features

autocorr = librosa.autocorrelate(y, max_size=1000)

features['autocorr_peak'] = np.max(autocorr[1:])

for lag in [1, 10, 100]:

    features[f'autocorr_lag_{lag}'] = autocorr[lag] if lag < len(autocorr) else 0

# Pitch Estimation

f0, voiced_flag, _ = librosa.pyin(y, sr=sr, fmin=FMIN, fmax=FMAX, frame_length=FRAME_LENGTH, hop_length=HOP_LENGTH)
```

```python
if np.any(voiced_flag):
    f0_voiced = f0[voiced_flag]
    features['pitch_mean'] = np.mean(f0_voiced)
    features['pitch_std'] = np.std(f0_voiced)
    features['pitch_range'] = np.ptp(f0_voiced)
    features['pitch_kurtosis'] = kurtosis(f0_voiced)
    features['pitch_skewness'] = skew(f0_voiced)
    x = np.arange(len(f0_voiced)).reshape(-1, 1)
    model = LinearRegression().fit(x, f0_voiced)
    features['intonation_slope'] = model.coef_[0]
    features['intonation_intercept'] = model.intercept_
else:
    features.update({k: 0 for k in ['pitch_mean', 'pitch_std', 'pitch_range', 'intonation_slope', 'intonation_intercept']})


# Pause Detection
rms = librosa.feature.rms(y=y, frame_length=FRAME_LENGTH, hop_length=HOP_LENGTH)[0]
threshold = np.percentile(rms, 10)
silent_frames = np.where(rms < threshold)[0]
if len(silent_frames) > 0:
    pause_durations = np.diff(silent_frames) * HOP_LENGTH / sr
    features['pause_duration_total'] = np.sum(pause_durations)
    features['pause_count'] = len(pause_durations)
else:
    features['pause_duration_total'] = 0
    features['pause_count'] = 0
features['speaking_rate'] = np.sum(voiced_flag) / len(voiced_flag) if len(voiced_flag) > 0 else 0
```

**Spectrogram Extraction:** Spectrograms visually represent the spectrum of frequencies of a signal as it varies with time. A mel-spectrogram uses the mel scale for frequency, which better aligns with human auditory perception.

**Code Explanation:**

```python
def create_spectrogram(audio_path, output_base_path):
    try:
        y, sr = librosa.load(audio_path, sr=SAMPLE_RATE, duration=DURATION)
        S = librosa.feature.melspectrogram(y=y, sr=sr, n_fft=N_FFT, hop_length=HOP_LENGTH, n_mels=N_MELS)
        log_S = librosa.power_to_db(S, ref=np.max)


        os.makedirs(os.path.dirname(output_base_path), exist_ok=True)


        np.save(f"{output_base_path}.npy", log_S)


        plt.figure(figsize=(IMG_SIZE[0]/DPI, IMG_SIZE[1]/DPI), dpi=DPI)
        librosa.display.specshow(log_S, sr=sr, hop_length=HOP_LENGTH, x_axis='time', y_axis='mel')
        plt.axis('off')
        plt.savefig(f"{output_base_path}.png", bbox_inches='tight', pad_inches=0)
        plt.close()
        return True
    except Exception as e:
        print(f"Error processing {audio_path}: {str(e)}")
        return False
```

These extracted features serve as comprehensive descriptors of audio content, playing a vital role in training machine learning models to effectively identify and flag deepfake audio.

## 4. Model Implementation

In this stage, machine learning and deep learning models are built using the extracted features for audio deepfake detection. A combination of traditional machine learning classifiers and a hybrid deep learning model is explored to leverage different types of features.

**SVM Classifier Using MFCC Features**

**Support Vector Machine (SVM)** is a robust classifier used for supervised learning tasks. It is particularly effective for high-dimensional feature spaces, which makes it suitable for MFCCs (Mel-Frequency Cepstral Coefficients).

**Approach:**

i.    MFCC features were used as input vectors.
ii.   Data was split into training and testing sets.
iii.  StandardScaler was applied to normalize the feature values.
iv.   An SVM model with an RBF kernel was trained.

**Performance:**

- **Accuracy achieved:** 93.75%

- **Inference:** The SVM classifier demonstrated strong performance, effectively separating deepfake and real audio based on MFCC features.

**XGBoost Using Combined Features**

**XGBoost** (Extreme Gradient Boosting) is an optimized gradient boosting framework designed for speed and performance. It is well-suited for structured data and feature importance analysis.

**Approach:**

i.    MFCC features were used as input vectors.
ii.   Data preprocessing included missing value handling and normalization.
iii.  XGBoost was trained with optimized parameters using early stopping and cross-validation.

**Performance:**

- **Accuracy achieved:** 96%

- **Inference:** XGBoost outperformed SVM due to its ability to model complex nonlinear interactions and handle a broader set of features.

**CNN + LSTM Hybrid Model**

To fully exploit the temporal and spatial structure in audio data, a **hybrid deep learning model** was built:
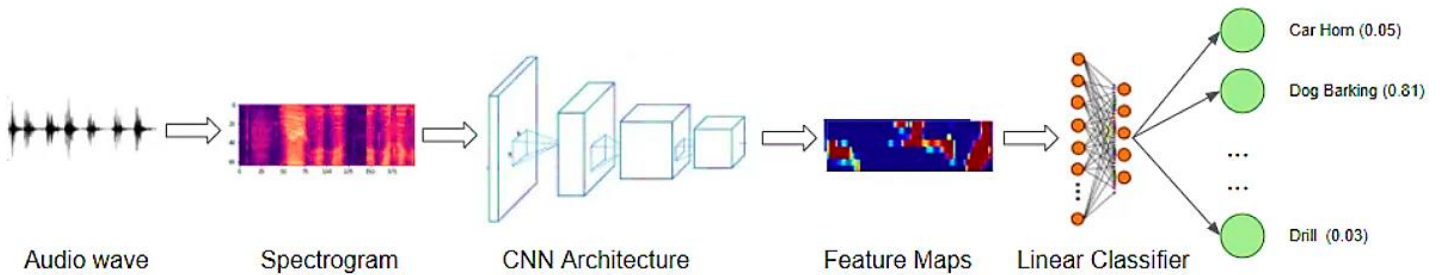


*Figure: 1*

## CNN for Spectrogram and Spectral Features

i. **Convolutional Neural Network (CNN)** was used to learn spatial patterns from mel-spectrograms.
ii. Input: spectrogram images.
iii. Layers:
Input Layer: Takes 128x128 RGB spectrogram images.
Conv2D (32 filters, 3x3): Extracts local features using 32 filters.
MaxPooling2D (2x2): Downsamples the feature maps.
Conv2D (64 filters, 3x3): Deeper feature extraction.
GlobalAveragePooling2D: Reduces each feature map to a single value (for vector representation).

## LSTM for Temporal and Pitch Features

i. **Long Short-Term Memory (LSTM)** was used to learn temporal dependencies from sequential data such as pitch and pause features.
ii. Input: sequential feature vectors (temporal + pitch features).
iii. Layers:
Input Layer: Accepts variable-length sequences of extracted temporal features.
LSTM (128 units): First recurrent layer that returns sequences to feed the next LSTM.
LSTM (64 units): Final LSTM layer that outputs a fixed-size vector representation of the sequence.

Numerical feature, Fusion and Output layers

Dense (64 units): Fully connected layer with ReLU activation.
Concatenate: Merges outputs of CNN, LSTM, and numerical branches.
Dense (128 units): Fully connected layer for integrated representation.
Dropout (0.5): Regularization to reduce overfitting.
Output Layer (1 unit, sigmoid): Binary classification output (real vs. fake).

*Note: Figure 1 CNN model for detection deepfake audio using spectrogram*

**Combined Model and Training**

The outputs of CNN and LSTM were concatenated and passed through fully connected layers.

Early stopping and learning rate schedulers were used during training.

**Performance:**

- **Initial Accuracy:** 64%

- **After Hyperparameter Tuning:** 79%

**Inference:**

- Although deep models generally require larger datasets, the hybrid model showed significant improvement post-tuning. Spectrogram and sequential modeling together offered better context-aware decision-making.

These modeling stages form the core of the audio deepfake detection pipeline, combining statistical and deep learning methods to improve classification performance and generalizability.

## 5. Model Evaluation:

All models were evaluated using metrics like accuracy, precision, recall, and F1-score. The Baseline **SVM** model achieved an accuracy of **93.75%**, demonstrating strong generalization performance, while **XGBoost** outperformed all models with **96.25%** accuracy and high precision/recall balance. The Proposed **Hybrid Model (CNN-LSTM)**, though slightly lower in overall accuracy **(79%)**, captures complex temporal and spectral features, which is critical for future scalability and real-time adaptability.

**Comparative Analysis of Based and Proposed Algorithms**

TABLE: 2

| Model | Accuracy (%) | F1-Score | Precision | Recall |
|---|---|---|---|---|
| **Baseline SVM** | 93.75 | 94.00 | 94.00 | 94.00 |
| **XGBoost** | 96.25 | 96.00 | 96.00 | 95.00 |
| **Proposed Hybrid (CNN-LSTM)** | 79.00 | 78.00 | 83.00 | 79.00 |

## Summary

This chapter detailed the step-by-step methodology followed in this project, from data collection to model deployment. The proposed deep learning models, particularly CNN, showed superior performance in detecting deepfake voices when compared to traditional classifiers like Random Forest and SVM.

**NOTE:** The table - 2 presents a comparative performance evaluation of various models used in the deepfake audio detection task. It includes both **baseline machine learning models** and the **proposed hybrid deep learning approach**. Each model is assessed using four key performance metrics:
- **Accuracy (%):** The percentage of correctly classified audio samples (real or fake) out of the total.
- **F1-Score:** The harmonic mean of precision and recall, reflecting the model's balance between false positives and false negatives.
- **Precision:** The ratio of correctly predicted positive observations (fake audios) to the total predicted positives.
- **Recall:** The ratio of correctly predicted positive observations to all actual positives in the dataset.

# Conclusion

The goal of this minor project was to develop a system capable of detecting deepfake audio using various machine learning and deep learning techniques.

Through careful data preprocessing, feature extraction (MFCCs and spectrograms),and model experimentation, the project successfully demonstrated the effectiveness of deep learning in distinguishing between real and fake voice samples. Several models were implemented and compared, including Random Forest, SVM, XGBoost and Hybrid model (LSTM, and CNN). The Hybrid model, trained on spectrogram images, achieved the accuracy of approximately 79%. The following screenshot shows the accuracy achieved by the models during evaluation:

CNN-LSTM Hybrid model

*Table: 1*

```
Classification Report:
              precision    recall  f1-score   support

        Real       0.71      0.97      0.82        80
        Fake       0.96      0.60      0.74        80

    accuracy                           0.79       160
   macro avg       0.83      0.79      0.78       160
weighted avg       0.83      0.79      0.78       160
```
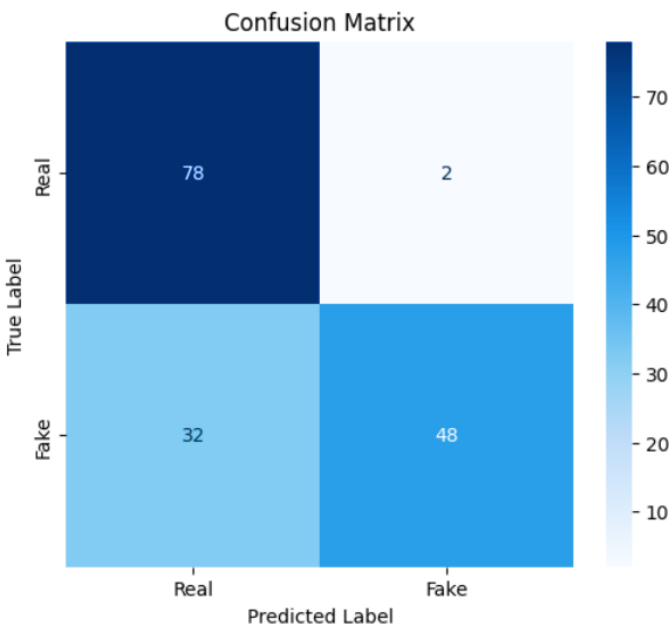


*Diagram:1*

**NOTE:** The classification report (Table:1) and confusion matrix (Diagram: 2) summarize the CNN-LSTM hybrid model's performance, achieving 79% accuracy. It shows high recall for real samples and high precision for fake samples, highlighting a trade-off in classifying fake data.
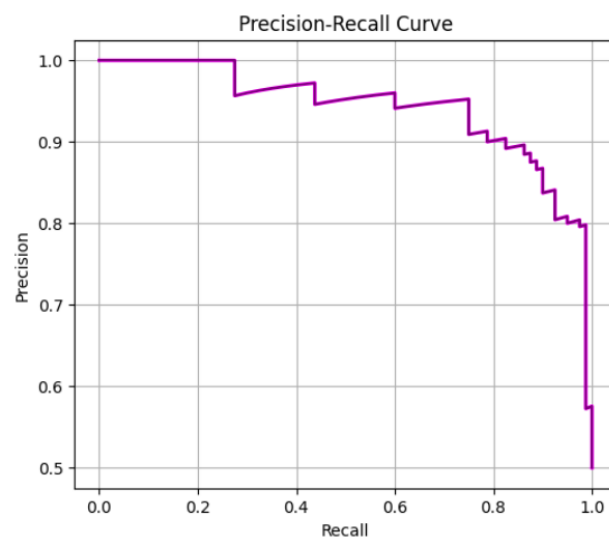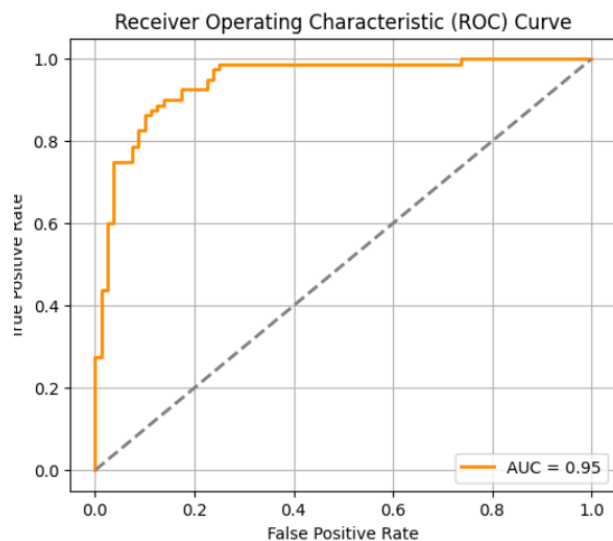
*Diagram:2*

## SVM & XGBoost Model -

*Table: 2 SVM Model*
```
Test Results:
Accuracy: 0.9375
ROC AUC: 0.9916
Classification Report:
              precision    recall  f1-score   support

           0       0.94      0.94      0.94        80
           1       0.94      0.94      0.94        80

    accuracy                           0.94       160
   macro avg       0.94      0.94      0.94       160
weighted avg       0.94      0.94      0.94       160


Confusion Matrix:
[[75  5]
 [ 5 75]]
```

*Table: 3 XGBoost Model*
```
XGBoost Results:
Accuracy: 0.9625
ROC-AUC: 0.9948
              precision    recall  f1-score   support

           0       0.95      0.97      0.96        80
           1       0.97      0.95      0.96        80

    accuracy                           0.96       160
   macro avg       0.96      0.96      0.96       160
weighted avg       0.96      0.96      0.96       160


Confusion Matrix:
 [[78  2]
 [ 4 76]]
['/content/drive/MyDrive/project/model/xgb_baseline.pkl']
```

**NOTE:** Diagrame: 2 (Hybrid Model) The ROC and Precision-Recall curves illustrate strong model performance, with an AUC of 0.95 indicating excellent classification capability.

Table: 2 Classification report of SVM model with accuracy of 94%.

Table:3 Classification report of XGBoost model with accuracy of 96%.

# Research And References

[1] S. Oyucu, D. B. Ünsal Çelımlİ, and A. Aksöz, "Fake Voice Detection: A Hybrid CNN-LSTM Based Deep Learning Approach," 2024.

[2] M. Demirörs, T. Akgün, and A. M. Özbayoğlu, "AI Generated Speech Detection Using CNN," 2024.

[3] Z. Zhang, X. Yi, and X. Zhao, "Fake Speech Detection Using Residual Network with Transformer Encoder," 2021.

[4] Z. Wu, T. Kinnunen, N. Evans, and J. Yamagishi, "ASVspoof 2015: The First Automatic Speaker Verification Spoofing and Countermeasures Challenge," 2015.