# DEEP FAKE AUDIO DETECTION USING MACHINE LEARNING

## A PROJECT REPORT

Submitted by

## MANOJKUMAR.B

## (P23DSC105)

In partial fulfilment for the award of the degree of

## MASTER OF SCIENCE

## in

## DATA SCIENCE



**POST GRADUATE DEPARTMENT OF DATA SCIENCE**

# NEHRU MEMORIAL COLLEGE

**(AUTONOMOUS)**

**[Nationally Reaccredited with 'A+' Grade by NAAC]**

**Affiliated with Bharathidasan University**

**Puthanampatti – 621 007.**

# NEHRU MEMORIAL COLLEGE

**(Autonomous)**

**[Nationally Reaccredited with 'A+' Grade by NAAC]**

**Affiliated with Bharathidasan University**

**Puthanampatti – 621 007.**

**POST GRADUATE DEPARTMENT OF DATA SCIENCE**

## BONAFIDE CERTIFICATE

Certified that this project report titled "DEEP FAKE AUDIO DETECTION USING MACHINE LEARNING", is a Bonafide work by,

<div style="text-align:center">

**NAME      : MANOJKUMAR.B**

**REG. NO  : P23DSC105**

</div>

submitted to Nehru Memorial College (Autonomous), Puthanampatti, Tiruchirappalli during the fourth semester in partial fulfilment of the requirement for the award of the degree of Master of Science in Data Science.

I further certify that the project work carried out by him is an independent work done under my supervision and guidance.

<table>
<tr>
<td><b>HEAD OF THE DEPARTMENT</b><br>(Department Of Data science)<br>(Mrs. N. KALPANA)</td>
<td><b>PROJECT SUPERVISOR</b><br>(Mrs. N. KALPANA)</td>
</tr>
</table>

Submitted on _____ for the Project Viva-voce to be held on _____

Examiners:

1. _____

2. _____

# ACKNOWLEDGEMENT

# ABSTRACT

This project presents a comprehensive system for detecting synthetic voice content using a multi-model machine learning framework, addressing the critical need for reliable detection as deepfake audio technology continues to evolve. The system features an extensive acoustic feature extraction pipeline that includes Mel-frequency cepstral coefficients (MFCCs), chromagrams, and mel spectrograms, which are vital for capturing subtle artifacts in artificially generated speech. Standardizing audio durations ensures consistent analysis across samples. The core innovation lies in training and evaluating multiple machine learning models— XGBoost, Random Forest, Support Vector Machine, and Gradient Boosting—to determine the most accurate classifier, while retaining all trained models for comparison and potential ensemble use. The backend, built with FastAPI, handles feature extraction, model training, and prediction, while a user-friendly Streamlit frontend enables intuitive audio analysis. Designed for flexibility, the system allows for both default and custom dataset training and provides real-time prediction with detailed confidence scores. By offering an open, modular, and extensible solution, this work significantly contributes to media forensics, with future scalability to accommodate advanced synthesis techniques and feature enhancements.

**Keywords**:

Deepfake Audio, Machine Learning, Feature Extraction, MFCC (Mel-Frequency Cepstral Coefficients), Chroma Features, Spectrogram, Gradient Boosting, FastAPI, Streamlit, Audio Classification

# TABLE OF CONTENT

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| Abbreviation | Full Form |
|---|---|
| AI | Artificial Intelligence |
| API | Application Programming Interface |
| CNN | Convolutional Neural Network |
| CPU | Central Processing Unit |
| GPU | Graphics Processing Unit |
| MFCC | Mel-Frequency Cepstral Coefficients |
| SDLC | System Development Life Cycle |
| SVM | Support Vector Machine |
| VAD | Voice Activity Detection |
| WAV | Waveform Audio File Format |
| MP3 | MPEG-1 Audio Layer III |
| FLAC | Free Lossless Audio Codec |
| SMOTE | Synthetic Minority Over-sampling Technique |
| FastAPI | Fast API Framework for Python |
| Streamlit | Python Library for Web App Development |
| XGBoost | Extreme Gradient Boosting |
| RBF | Radial Basis Function (Kernel used in SVM) |
| MFCCs | Mel-Frequency Cepstral Coefficients |

# CHAPTER - 1

## 1. INTRODUCTION

### 1.1 Background and Motivation:

The rapid advancement of artificial intelligence and deep learning technologies has revolutionized audio synthesis capabilities, enabling the creation of highly convincing synthetic voice replications known as "deepfakes." These technologies can now clone voices with minimal sample data, reproduce emotional inflections, and generate speech that is increasingly indistinguishable from authentic human recordings. What once required extensive technical expertise and significant computing resources has become accessible through user-friendly applications and open-source tools, allowing virtually anyone to create convincing voice imitations.

This technological democratization presents serious security and social challenges across multiple domains. Financial institutions face threats to voice-based authentication systems, where voice deepfakes could potentially bypass security measures for account access. Government and corporate entities confront risks of social engineering attacks, where synthetic voices impersonating executives or officials could be used to authorize fraudulent transactions or extract sensitive information. In the public information sphere, fabricated audio statements attributed to political figures or celebrities can spread misinformation rapidly, potentially influencing public opinion and undermining trust in legitimate communication channels.

Current detection methods face significant limitations in identifying these sophisticated voice deepfakes. Traditional audio forensics techniques often rely on specific artifacts produced by earlier generation algorithms, making them vulnerable as synthesis technologies evolve. Statistical approaches that worked against previous generations of voice synthesis fail against neural network-based methods that produce increasingly natural speech patterns. Additionally, the variety of synthesis methods in use—from GAN-based approaches to autoregressive models like WaveNet derivatives—means that detection systems trained on one type of deepfake often perform poorly against others.

The motivation for this research stems from these pressing challenges and the need for more robust, adaptable detection systems. By implementing a multi-model approach that leverages different machine learning classifiers (XGBoost, Random Forest, SVM, and Gradient Boosting), system aims to identify the subtle acoustic inconsistencies that persist even in highly

sophisticated deepfakes. The comprehensive feature extraction pipeline capturing MFCCs, chromagrams, and mel spectrograms provides a rich representation of audio characteristics that helps distinguish synthetic artifacts from natural speech variations.

Furthermore, the practical implementation through both a FastAPI backend and Streamlit frontend addresses the critical need for accessible detection tools that bridge the gap between technical capabilities and user-friendly applications. This approach enables both forensic experts and non-technical users to analyze suspicious audio content, democratizing detection capabilities to match the widespread availability of synthesis tools.

The significance of this work extends beyond technical innovation to the preservation of trust in audio communication. As synthetic voices become increasingly convincing, society needs reliable mechanisms to verify the authenticity of voice recordings. This research contributes to maintaining the integrity of audio evidence in legal proceedings, protecting individuals from voice-based impersonation attacks, and helping media organizations verify the authenticity of audio sources before publication.

## 1.2 Objectives and Scope of the Project:

This project aims to develop a practical system for detecting deepfake audio content through a combination of feature extraction techniques and machine learning classification. Based on the implemented code, the specific objectives of this project include:

1. **Implement Audio Feature Extraction**: Create a feature extraction pipeline that processes audio files to extract:

   - Mel-frequency cepstral coefficients (MFCCs)
   - Chroma features
   - Mel spectrogram features

2. **Train Multiple Classification Models**: Develop a system that trains and evaluates four different machine learning models:

   - XGBoost Classifier
   - Random Forest Classifier
   - Support Vector Machine (SVM)
   - Gradient Boosting Classifier

3.  **Compare Model Performance**: Automatically evaluate and compare the performance of each model to identify the most effective classifier for deepfake detection based on accuracy metrics.

4.  **Create a FAST API Backend**: Implement a Fast API-based service that provides:

    *   Endpoints for model training with custom datasets
    *   Prediction capabilities for analyzing uploaded audio files
    *   Performance reporting for trained models

5.  **Enable Flexible Dataset Management**: Design the system to work with:

    *   Default dataset paths for standard operation
    *   User-specified dataset paths through API parameters
    *   Standardized processing regardless of audio source

6.  **Develop a Stream-lit Frontend**: Build a user-friendly interface that complements the API backend (as mentioned in your description, though not shown in the provided code).

7.  **Provide Detailed Classification Results**: Return comprehensive prediction information including:

    *   Binary classification (REAL or FAKE)
    *   Confidence scores
    *   Probability distribution for both classes

The scope of this project focuses on binary classification of audio files as either authentic (REAL) or synthetic (FAKE). The system processes audio by standardizing the duration to 3 seconds and extracting acoustic features that help identify synthetic artifacts. It utilizes established machine learning algorithms rather than developing new classification methods, with emphasis on practical deployment through a web API service.

## 1.3 Problem Statement:

The increasing availability of voice synthesis technologies has made it possible for anyone to create fake audio that sounds like a real person's voice. These "deepfake" audio files

can be used to spread false information, impersonate others, or bypass voice authentication systems. This creates several key problems that this project aims to solve:

First, it is becoming difficult for humans to distinguish between real and fake audio by just listening. Even careful listeners can be fooled by well-made deepfakes, especially when the audio quality is not perfect or when they have no reason to suspect manipulation.

Second, there is a need for automated tools that can analyze audio files and determine whether they are authentic or synthetically generated. While some commercial solutions exist, they are often expensive and not accessible to everyone who needs this capability.

Third, different machine learning models may perform differently when detecting various types of fake audio. Some models might be better at catching certain synthesis techniques while missing others. Using just one model might not provide reliable results across all types of deepfakes.

Fourth, detection systems need to extract the right features from audio files to effectively identify manipulations. The subtle differences between real and fake audio might be found in specific acoustic properties that need to be properly captured and analyzed.

Finally, any detection solution needs to be user-friendly and accessible to people without technical expertise in audio processing or machine learning. The system should provide clear results that can be understood and acted upon by regular users.

**This project addresses these problems by creating a system that:**

- Extracts multiple audio features (MFCCs, chroma, and mel spectrograms)

- Trains and compares four different machine learning models

- Provides an API for integration with other systems

- Includes a user-friendly interface through a Streamlit application

- Gives clear predictions with confidence scores about whether an audio file is real or fake

## 1.4 About SRM Tech:

SRM Technologies Pvt. Ltd., established in 1999, is a globally recognized IT services and product engineering powerhouse, headquartered in Chennai, Tamil Nadu. With a strong global presence and a reputation for delivering cutting-edge technological solutions, SRM Tech has firmly positioned itself as a trusted partner for Fortune 500 companies. The company's expertise spans across digital transformation, cloud computing, embedded systems, automotive engineering, and next-generation technologies. Its Perungudi campus, situated in the prestigious RMZ Millenia Business Park, stands as a hub of innovation, housing high-performance teams engaged in software development, client success, and R&D initiatives.

As an integral part of the esteemed SRM Group, SRM Technologies is known for its commitment to excellence, continuous innovation, and a future-ready approach. The company has earned a distinguished reputation for partnering with industry leaders across critical sectors including healthcare, automotive, finance, and telecommunications. With a robust focus on emerging technologies such as artificial intelligence, machine learning, blockchain, and the Internet of Things (IoT), SRM Tech is shaping the digital future. Its dedication to building scalable, secure, and intelligent solutions enables enterprises worldwide to accelerate their digital journey and drive real-world impact.

### 1.4.1 Role and Responsibilities of SRM TECH:

- **Digital Transformation Services:** Empowering enterprises through end-to-end digital solutions, enhancing agility and competitiveness.
- **Product Engineering and Embedded Systems:** Delivering comprehensive product engineering services, from concept to manufacturing, across various industries.
- **Automotive Solutions:** Providing advanced automotive technologies, including in-vehicle software development and electric vehicle solutions.
- **Artificial Intelligence and Data Analytics:** Leveraging AI and data analytics to drive actionable insights and optimize business operations.
- **Cloud and Infrastructure Services:** Offering scalable and secure cloud solutions to ensure high availability and performance.
- **Cybersecurity and Quality Assurance:** Implementing robust security measures and quality assurance practices to protect digital assets.
- **Supply Chain Management Solutions:** Addressing complex supply chain challenges with transformative digital solutions.

- **Global Presence and Industry Collaboration:** Collaborating with leading industries worldwide to deliver tailored solutions that drive growth and innovation.

## 1.4.2 Global Presence

- **India**: Headquartered in Chennai, with additional offices in Bangalore, Hyderabad
- **United States**: Office located in Hamilton, New Jersey.
- **Japan**: Office situated in Tokyo
- **Germany**: Office at Am Seedamm
- **Canada**: Office in Windsor, Ontario
- **United Kingdom**: Office in London.
- **Middle East**: Office in Dubai.

## 1.4.3 Key Achievements and Milestones:

- **Industry Recognition**: SRM Technologies has been consistently recognized for its innovation, quality, and customer-centric approach. It has earned prestigious awards such as **Great Place to Work® Certification.**
- **Strategic Partnerships**: The company has forged strategic alliances with global technology giants like Microsoft, IBM, and Google Cloud, enabling it to offer advanced solutions and services to clients.
- **Successful Digital Transformation Projects**: SRM Technologies has successfully assisted multiple Fortune 500 companies in transforming their business processes through digital transformation, driving productivity, and operational efficiency.
- **Cutting-Edge Research and Development**: SRM Tech invests heavily in R&D, with several groundbreaking innovations in AI, IoT, blockchain, and automotive engineering that have revolutionized industries.
- **Sustainability Initiatives**: SRM Technologies has launched several sustainability-driven projects, focusing on reducing environmental impact through green technologies and contributing to global sustainability goals.
- **Global Clientele**: SRM Tech has expanded its client base to include multinational corporations across various sectors like healthcare, automotive, finance, and telecommunications, delivering tailored and scalable solutions.

# CHAPTER - 2

## 2. SYSTEM DEVELOPMENT LIFE CYCLE (SDLC)

The System Development Life Cycle (SDLC) approach was followed in developing the deepfake audio detection system, providing a structured methodology for planning, implementation, and deployment of this machine learning solution.

## 2.1 Phases of SDLC:

- Addressing the need for automated detection of deepfake audio content

- Establishing requirements for a system that can distinguish between real and synthetic voice recordings

- Defining success metrics based on classification accuracy and confidence scores

### Data Collection & Preprocessing:

- Utilizing organized directory structure with separate folders for real and fake audio samples

- Implementing standardized audio processing through librosa library

- Converting varied audio formats (WAV, MP3, FLAC, OGG) into consistent feature representations

- Normalizing audio duration to 3 seconds (trimming or padding as needed)

### User Interface Design & Data Visualization:

- Implementing visualization components in the Streamlit interface to display:

  - Audio waveforms for temporal pattern analysis
  - Spectrograms for frequency distribution examination

- Enabling user interaction with uploaded audio samples for inspection

### Model Selection & Development:

- Implementing four distinct classification models:

  - XGBoost Classifier

  - Random Forest Classifier

  - Support Vector Machine (SVM)

  - Gradient Boosting Classifier

- Developing feature extraction functions for MFCCs, chroma, and mel spectrograms

- Creating a REST API backend with FastAPI for model training and prediction

- Building a user-friendly interface with Streamlit for end-user interaction

### Model Training & Evaluation:

- Implementing train/test split methodology (80/20)

- Calculating comprehensive metrics including accuracy and classification reports

- Automating model comparison to identify the best-performing classifier

- Saving trained models for future use with appropriate versioning

## 2.2 Importance of SDLC Methodology:

### Ensures Data Quality:

- Standardized feature extraction process across all audio files

- Consistent audio duration handling improves feature consistency

- Enhances data quality by applying preprocessing steps to minimize background noise and irrelevant information.

### Improves Model Reliability:

- Comparative evaluation of multiple models enhances detection capability

- Best model selection mechanism ensures optimal performance

- Reduces the risk of overfitting by assessing the model's ability to generalize to unseen data.

### Facilitates Scalability:

- API-based architecture supports integration with other systems

- Model storage and versioning allows for system expansion

### Supports Continuous Improvement:

- Default model concept allows for updates as new training data becomes available

- Structured evaluation metrics provide baseline for future enhancements

## 2.3 SDLC Models Applied in the Project:

### Component-Based Development:

- Separation of frontend (Streamlit) and backend (FastAPI) components

- Modular code structure with specific functions for feature extraction, model training, and prediction

### Agile Methodology:

- Evidence of iterative development with multiple model implementations

- Flexible architecture allowing for rapid updates and improvements

### DevOps Integration:

- Implementation of automated model training and evaluation pipeline

- API endpoints supporting continuous deployment capabilities

# CHAPTER - 3

## 3. SYSTEM ANALYSIS

## 3.1 Requirement Analysis:

Reliable detection of deepfake audio is essential for safeguarding voice-based authentication, media forensics, and content integrity. In this phase identify and confirm all resources needed to build, train, and validate deepfake-voice detection models.

The key requirements for this project include:

- **Data Collection & Preprocessing:**
  - Gather genuine speech samples (various speakers, languages) and deepfake clips (multiple TTS/voice-conversion engines).
  - Normalize volume, remove silence, downsample/upsample as needed, and label each clip.
- **Feature Extraction & Engineering:**
  - Compute MFCCs, spectral contrast, chroma features and zero-crossing rate for every audio segment.
  - Optionally augment with noise or reverberation to improve robustness.
- **Model Training & Comparison:**
  - Train Random Forest, XGBoost, SVM and Gradient boosting methods on the extracted features.
  - Use grid or randomized search to tune hyperparameters for each model.

## 3.2 System Requirements:

To implement and validate the deepfake voice detection system, the project requires specific hardware, software, and data resources.

## 3.2.1 Hardware Requirements:

The system must have sufficient computational power to process audio files, train machine learning models, and perform real-time inference efficiently. The recommended hardware specifications include:

- **Processor:** Intel Core i7 / AMD Ryzen 7 (or higher) with multiple cores for handling real-time processing and model inference.

- **RAM:** Minimum 16GB (32GB recommended) for smooth audio data preprocessing and model training.

- **Storage:** At least 500GB SSD (preferably 1TB) to store datasets, extracted audio features, and trained models.

- **GPU (Optional):** NVIDIA RTX 3060 or higher (recommended for deep learning model training and faster inference).

- **Microphone (Optional):** For collecting real-time voice samples if needed.

- **Internet Connectivity:** Required for downloading pretrained models, datasets, and libraries.

## 3.2.2 Software Requirements:

The software environment should support audio processing, machine learning

- **Operating System:** Windows 10/11, or macOS
- **Programming Language:** Python 3.8+
- **Libraries:**
  - **Audio Processing:** Librosa, Pydub
  - **Deep Learning:** PyTorch or TensorFlow
  - **Machine Learning:** Scikit-learn
  - **Data Handling:** Pandas, NumPy
  - **Visualization:** Matplotlib, Seaborn
  - **Model Evaluation:** Scikit-learn (for metrics like accuracy, precision, recall, ROC curve)
- **API Development:** FastAPI (for building the REST API interface)
- **Development Environment:** VS Code, Jupyter Notebook, or PyCharm

## 3.3 Feasibility Study

The feasibility study assesses the practicality of implementing a deepfake voice detection system using machine learning and deep learning models across the following dimensions:

### 3.3.1 Technical Feasibility:

- **Data Availability:** The project uses publicly available datasets containing real and deepfake audio samples. These datasets are in WAV and MP3 formats, easily processed using Python libraries.

- **Software and Tools:** The implementation uses Python-based libraries like TensorFlow, PyTorch, Librosa, Scikit-learn, and Streamlit for data preprocessing, model training, evaluation, and frontend development.

- **Computational Requirements:** Machine learning models for voice detection require a strong CPU and optionally a GPU (NVIDIA RTX series recommended) for faster training and inference.

### 3.3.2 Operational Feasibility:

- **Ease of Implementation:** Deepfake detection methodologies are well-established, and Python libraries provide pre-built functions for audio processing, model training, and evaluation.

- **Skill Requirements:** The project requires skills in machine learning and audio signal processing. The development team must be proficient in Python and familiar with deepfake detection concepts..

- **Scalability:** The system is designed to process and detect fake voices in real-time and can be scaled to handle larger datasets or integrated into live applications if needed.

### 3.3.3 Economic Feasibility:

- **Cost of Computational Resources:** Since the project uses machine learning models it can run efficiently on a good CPU without requiring expensive GPUs. Cloud services (like Google Colab, AWS) may still be used for faster training but at minimal cost..

- **Software Licensing:** The project uses open-source libraries such as Scikit-learn, XGBoost, and Librosa, reducing software expenses.

### 3.3.4 Schedule Feasibility:

- **Project Timeline:** The project is divided into phases such as audio data collection, preprocessing, machine learning model training, evaluation, and reporting. Each phase is scheduled for smooth and timely completion.

- **Potential Delays:** Delays could occur due to issues like limited quality data, model retraining needs, or computational constraints. However, using pre-processed datasets and efficient coding practices can help minimize these risks.

## 3.4 Project Constraints:

The implementation of machine learning models for deepfake voice detection faces several constraints that may affect accuracy, scalability, and efficiency. These challenges include limited availability of labeled audio datasets, computational resource limitations for training models, potential variability in real-world audio quality, and dependencies on third-party libraries and tools for feature extraction and model evaluation.

### 3.4.1 Data Constraints:

- **Limited Availability of Authentic and Deepfake Audio Samples:** Access to large and balanced datasets containing both real and deepfake voice samples is limited, which can affect model training and evaluation.

- **Different Audio Formats and Sampling Rates:** Audio data collected may exist in various formats (e.g., WAV, MP3) and sampling rates, requiring standardization during preprocessing.

- **Quality Variations:** Some audio files may contain background noise, compression artifacts, or distortions, requiring careful preprocessing like denoising and normalization.

- **Dataset Imbalance:** There may be more real voice samples than deepfake samples, leading to class imbalance, which can impact machine learning model performance.

### 3.4.2 Computational Constraints:

- **Processing Power Requirement:** Training machine learning models (Random Forest, XGBoost) on audio features requires sufficient CPU and RAM resources for efficient processing.

- **Feature Extraction Overhead:** Extracting audio features like MFCCs, chroma, and spectral features from large numbers of files can be time-consuming and resource-intensive.

- **Model Training Time:** model training may take significant time, especially when using cross-validation for robust evaluation.

### 3.4.3 Methodological Constraints:

- **Detection Accuracy:** Machine learning models may struggle to detect highly sophisticated or subtle deepfake manipulations, leading to potential false positives or negatives.

- **Generalization Limitations**: Models trained on specific datasets may not perform equally well on unseen real-world data with different accents, noise conditions, or spoofing techniques.

- **Feature Limitation:** Relying only on hand-engineered audio features may limit model performance compared to more complex representations, but deep learning approaches are not considered in this project.

### 3.4.4 External Dependencies:

- **Dataset Availability:** The project depends on publicly available deepfake voice datasets, and limited availability or imbalance between real and fake samples can impact training.

- **Software and Library Dependencies**: Tools like Scikit-learn, Librosa, Pandas, and XGBoost are critical; any library updates, bugs, or compatibility issues could affect project continuity.

- **External Data Quality:** Variations in dataset quality (background noise, audio compression) may impact model performance and require additional preprocessing.

### 3.4.5 Time and Budget Constraints:

- **Project Timeline:** Efficient management of time is necessary for phases like data collection, feature extraction, model training, evaluation, and reporting.

- **Resource Constraints:** Limited access to high-performance machines may slow down training and experimentation phases.

- **Budget Limitations:** Since the project uses open-source tools and runs on local systems or free-tier cloud services, any requirement for additional computational resources could exceed the planned budget.

# CHAPTER – 4

## 4. DATA COLLECTION

## 4.1 Description of Data Sources

To effectively train and evaluate machine learning models for deepfake voice detection, high-quality audio datasets are essential. The data sources used in this project are carefully selected to include both real and manipulated (deepfake) audio recordings to provide a comprehensive dataset for training and testing. The following datasets are utilized:

- **Real Audio Data:** The real audio data comes from open-source datasets of human speech, such as the LibriSpeech corpus, which contains recordings of English audiobooks spoken by various speakers. This serves as the "real" class for the detection model.

- **Deepfake Audio Data:** Deepfake audio data is sourced from publicly available datasets like Fake Audio Detection (FAD) and Google's Deepfake Speech Detection Dataset, which include manipulated speech generated using voice cloning and synthetic speech synthesis technologies. These datasets are used as the "fake" class in model training.

- **Metadata:** Alongside audio files, relevant metadata (e.g., speaker demographics, recording environment) may be provided, which can help in evaluating the model's performance across different acoustic conditions.

## 4.2 Data Preprocessing

Data preprocessing is a crucial step in any machine learning project, particularly for speech-based tasks like deepfake voice detection. The goal is to prepare the raw audio data for input into the machine learning model by performing a series of steps to enhance quality and ensure consistency across the dataset. The following preprocessing steps will be applied to the audio data:

- **Audio Resampling:** Audio recordings may have different sampling rates. To standardize the dataset, all audio files will be resampled to a consistent rate (e.g., 16 kHz) to ensure uniformity in input data.

- **Silence Removal:** Background noise or silence periods in audio recordings may not provide useful information. Using techniques like Voice Activity Detection (VAD),

silence segments will be detected and removed to focus on the speech portion of the audio.

- **Noise Reduction:** Audio signals often contain background noise that can interfere with feature extraction. Preprocessing steps like noise reduction algorithms (e.g., spectral gating) will be applied to minimize non-speech noises.

- **Feature Extraction:** Raw audio data is converted into a format suitable for machine learning models by extracting meaningful features such as:

    - **Mel-Frequency Cepstral Coefficients (MFCCs):** These features represent the short-term power spectrum of audio and are commonly used in speech recognition tasks.

    - **Spectrograms:** Time-frequency representations of the audio signal will be computed to capture both temporal and spectral patterns.

    - **Chroma Features:** These features capture pitch and harmonic content, which are relevant for distinguishing human-like voice characteristics.

- **Labeling:** Each audio sample will be labeled as "real" (authentic) or "fake" (deepfake) based on its origin in the dataset. This is essential for supervised learning, where the model learns to differentiate between real and synthetic voices.

- **Normalization:** Audio features will be normalized (e.g., zero-mean, unit variance) to ensure consistent scale and improve the performance of the machine learning models.

- **Data Augmentation (Optional):** To enhance the robustness of the model, data augmentation techniques such as pitch shifting, speed variation, and background noise injection will be applied to the audio data to simulate different acoustic conditions.

- **Train-Test Split:** The dataset will be split into training and testing sets, ensuring that no overlap occurs between the sets. This helps in evaluating the model's generalization ability.

## 4.3 Data Quality Assessment

Before proceeding to model development, it is essential to ensure the quality and integrity of the collected datasets. The following steps are taken for data quality assessment:

- **Class Balance Check**:

    The datasets are inspected to verify that there is a sufficient and balanced number of real and fake audio samples. Any significant class imbalance could lead to biased model performance. Techniques such as oversampling (SMOTE) or undersampling are considered if necessary.

- **Audio Integrity Verification**:

    Each audio file is checked for corruption or loading errors. Files that cannot be properly loaded, decoded, or that are found to be empty are excluded from the dataset.

- **Duration Consistency**:

    Audio recordings are checked to ensure reasonable duration limits. Extremely short or excessively long audio files are removed or trimmed to ensure that the model receives consistent-length inputs for training.

- **Signal-to-Noise Ratio (SNR) Analysis**:

    SNR is measured for each audio file to assess the quality of recordings. Files with very poor SNR (high noise, low speech content) are flagged and either enhanced through preprocessing or removed.

- **Speaker Diversity Analysis**:

    The datasets are analyzed to ensure a diverse range of speakers (in terms of age, gender, accent) are included. This prevents the model from overfitting to voice characteristics and enhances its ability to generalize to unseen data.

## Data Preprocessing

Raw audio data often contains inconsistencies such as varying sampling rates, background noise, and silent periods. To handle these issues:

- **Audio Resampling**:

  All audio files were resampled to a uniform sampling rate of **16 kHz** to maintain consistency across the dataset.

- **Silence Removal**:

  Using **Voice Activity Detection (VAD)** algorithms, silent segments in the audio files were identified and removed to focus only on the parts containing speech.

- **Noise Reduction**:

  **Spectral gating** and **noise filtering** techniques were applied to minimize background noise and enhance the clarity of the audio signals.

This preprocessing ensures that the machine learning models receive high-quality and standardized inputs, improving their ability to learn meaningful patterns.

## Feature Extraction

Instead of feeding raw audio directly to models, **audio features** that represent meaningful characteristics were extracted:

- **Mel-Frequency Cepstral Coefficients (MFCCs)**:

  Capture the power spectrum of audio signals and are crucial for speech and speaker recognition tasks.

- **Chroma Features**:

  Represent the distribution of energy across the 12 distinct pitch classes, helping differentiate tonal properties in real and fake speech.

- **Spectral Centroid**:

  Represents the "center of mass" of the spectrum, often correlating with the brightness of a sound.

**Model Selection and Training**

After preparing the feature set, four machine learning algorithms were selected:

- **Random Forest Classifier**

- **Support Vector Machine (SVM)**

- **XGBoost Classifier**

- **Gradient Boosting Classifier**

Each model was trained to learn the distinction between real and deepfake audio samples based on extracted features

**Training and Validation**

- The dataset was split into **80% for training** and **20% for testing**.

- To ensure that the models generalize well and avoid overfitting, **5-Fold Cross-Validation** was employed:

    i. The training set was divided into 5 subsets.
    ii. In each iteration, 4 subsets were used for training and 1 subset for validation.
    iii. The final model performance was averaged over the 5 folds.

**Model Evaluation**

The models were evaluated using important classification metrics:

- **Accuracy**:
    Measures the proportion of correctly predicted instances over the total predictions.

- **Precision**:
    Indicates how many of the instances predicted as "deepfake" are deepfake.

- **Recall**:
    Measures the ability of the model to find all the deepfake instances in the dataset.

- **F1-Score**:
    The harmonic mean of precision and recall, providing a balanced evaluation metric.

**Deployment Strategy**

Once the best-performing model was selected:

- **Backend**:

  - Developed using **FastAPI**, a modern and fast (high-performance) Python web framework.
  - It handled model inference, file uploads, and prediction routing.

- **Frontend**:

  - Created using **Streamlit**, providing a user-friendly interface where users could upload audio files and receive real-time predictions about whether the audio is real or deepfake.
  - This integrated approach ensured a smooth transition from model development to real-world deployment.

## 5.2 Algorithms and Techniques

The following algorithms and techniques were utilized in depth:

**Random Forest Classifier**

- **Description**:
  An ensemble learning method that builds multiple decision trees during training and outputs the class that is many of the classes predicted by individual trees.

- **Key Features**:

  - Handles large datasets efficiently.
  - Reduces overfitting by averaging multiple trees.
  - Robust against noisy data and outliers.

- **Why Used**:

  Effective in handling complex patterns in audio features and providing strong baseline performance.

**Support Vector Machine (SVM)**

- **Description**:

  A supervised learning model that finds the optimal hyperplane that best separates the two classes (real and fake).

- **Key Features**:

  - Effective in high-dimensional spaces.
  - Uses different kernels (linear, RBF) to adapt to feature patterns.

- **Why Used**:

  SVMs are particularly good at binary classification problems with clear margins of separation.

**XGBoost (Extreme Gradient Boosting)**

- **Description**:

  A scalable, accurate, and fast implementation of gradient boosting algorithms designed for speed and performance.

- **Key Features**:

  - Efficiently handles imbalanced data.
  - Regularization (L1 and L2) improves generalization.
  - Parallelized tree building improves computational speed.

- **Why Used**:

  XGBoost often delivers state-of-the-art performance in many machine learning competitions and real-world tasks.

**Gradient Boosting Classifier**

- **Description**:

  - Another powerful boosting method that builds models iteratively.

- **Key Features**:

  - Minimizes loss by combining weak learners to create a strong overall model

- **Why Used**:

- Particularly effective when datasets have complex patterns.

**Feature Scaling (Standardization)**

- **Description**:
  StandardScaler was used to normalize feature distributions.

- **Key Features**:

  - Centers features at mean zero and unit variance.
  - Essential for SVM and Gradient Descent models for better convergence.

## 5.3 Justification for Chosen Approach

The chosen approach combines extensive audio feature extraction with multiple machine learning classifiers to ensure robust and reliable detection of deepfake audio. Feature extraction techniques like MFCCs, Chroma features, and Spectral Centroid were selected because they capture essential acoustic and harmonic properties that differentiate real human speech from synthetically generated audio. Instead of relying solely on raw audio signals, these features offer a structured numerical representation that enhances model learning and performance.

Multiple machine learning models—Random Forest, Support Vector Machine (SVM), XGBoost, and Gradient Boosting—were trained and evaluated to identify the most effective classifier for the task. This multi-model strategy was adopted to avoid over-reliance on a single algorithm and to increase the system's resilience across different types of deepfake audio. Gradient Boosting was ultimately selected based on its superior accuracy and ability to capture subtle differences in audio features.

Additionally, the deployment architecture using FastAPI and Streamlit was justified to ensure the solution is accessible, scalable, and user-friendly. FastAPI offers high-performance backend support for real-time audio analysis, while Streamlit provides an intuitive frontend interface for non-technical users.

# CHAPTER – 6

# FEATURE ENGINEERING AND ANALYSIS

## 6.1 Feature Extraction Process

In deepfake audio detection, feature engineering is crucial to transforming raw audio data into structured numerical features that machine learning models can process. The features extracted in this project help characterize audio samples and distinguish between real and deepfake content. This process is essential for the success of model training and evaluation.

The following features were extracted from the audio files:

- **Mel Frequency Cepstral Coefficients (MFCCs):**

  - MFCCs are widely used in speech and audio signal processing to represent the power spectrum of sound. They capture the short-term spectral properties of an audio signal.
  - 13 MFCC coefficients were extracted from each audio frame, representing the frequency content of the signal in a compact form. The MFCCs are often used because they capture the most relevant auditory characteristics of speech and sound.
  - Use: These coefficients help identify phonetic content, speaker characteristics, and tonal qualities, making them highly useful for distinguishing real and manipulated audio.

- **Chroma Features:**

  - Chroma features represent the distribution of energy across the 12 pitch classes of the chromatic scale (semitones in an octave). These features help capture harmonic and tonal information in the audio signal.
  - Use: Chroma features are important for music and speech analysis where pitch and tonality play a significant role in the audio's content. Deepfake audio typically lacks the natural harmonic richness of real human speech, making this feature useful for classification.

- **Spectral Centroid:**
  - The spectral centroid is a measure of the "brightness" of the sound. It represents the center of mass of the frequency spectrum.

- A higher spectral centroid typically indicates a brighter sound, often associated with higher-pitched frequencies.
- Use: This feature helps in capturing the perceptual characteristics of sound that distinguish real speech from synthetic or manipulated speech. The spectral centroid reflects how "sharp" or "dull" the audio feels, which can be a distinguishing factor between real and deepfake audio.

Each audio sample was converted into a feature vector consisting of the values derived from these features. These vectors serve as the input for machine learning models, allowing them to learn patterns in real versus fake audio.

## 6.2 Data Standardization

Data standardization is a critical step before applying machine learning algorithms. Raw features, such as MFCCs, Chroma features, and Spectral Centroid, often have different ranges and distributions, which can lead to poor model performance.

- **Standardization Method:**

  Standard Scaler was used to standardize the extracted features. This process ensures that:

  - All features have a mean of 0 and a standard deviation of 1.
  - This ensures that no feature dominates over others, allowing the models to focus on learning the relationships between features rather than their scale.

**Why Standardization?**

Standardization is particularly important for models like Support Vector Machines (SVMs) and Gradient Descent-based models (e.g., Logistic Regression), which are sensitive to the scale of features. It also helps speed up convergence during training.

## 6.3 Analysis of Extracted Features

Once the features were extracted and standardized, statistical analyses were conducted to identify how real and deepfake audio samples differ based on these features. The key insights from these analyses are summarized below:

| Feature Name | Real Audio (Mean ± Std) | Deepfake Audio (Mean ± Std) | Key Observations |
|---|---|---|---|
| MFCCs | ~ [-400, 400] | ~[-500, 500] | Deepfake audio showed slightly higher variations in MFCCs. |
| Chroma Features | ~ [0.1, 0.9] | ~ [0.05, 1.0] | Real audio had richer harmonic content compared to deepfake audio |
| Spectral Centroid | ~3000–4000Hz | ~2500–4500 Hz | Spectral centroid in deepfakes was slightly more variable. |

**Table-6.3.1 Audio Feature Statistics: Real vs Deepfake**

## Real vs Fake Separation:

The differences in these features between real and deepfake audio samples provided useful indicators for classification. Deepfake audio generally showed:

- Higher variations in MFCCs, likely due to synthetic artifacts.

- Lower harmonic richness in Chroma features.

- Wider range of Spectral Centroid values, indicating a less consistent frequency profile compared to real audio.

## Feature Correlations:

There were mild correlations observed between spectral features (MFCC and Spectral Centroid), which highlighted that certain frequency-based characteristics are closely related. For example, a high spectral centroid often aligned with certain MFCC patterns, showing how both features could work together to help distinguish real from deepfake audio.

## 6.4 Challenges in Feature Engineering

While feature extraction helped convert raw audio data into a usable form, several challenges arose during this process:

- **Background Noise and Recording Quality:**

  - Real-world audio recordings often contain background noise or distortions from different recording environments. These factors can affect the accuracy of feature extraction, especially for MFCC and Spectral Centroid.
  - Mitigation: Noise reduction techniques and Voice Activity Detection (VAD) were used to remove silent or non-speech segments.

- **Variability in Speaking Style:**

  - Variations in how individuals speak — such as accent, pitch, tone, and speed — introduced noise in the features.
  - This made it harder to distinguish between real and deepfake audio samples, as both types could share similar speaking styles.
  - Solution: Data preprocessing and feature extraction techniques were designed to capture these nuances while minimizing their impact on classification.

- **Balancing Real and Fake Samples:**

  - To avoid introducing bias, the dataset was checked for a balanced representation of real and deepfake audio samples.
  - Solution: Oversampling or under-sampling techniques could be employed if needed to ensure equal distribution between real and fake classes.

# CHAPTER – 7

## 7.Model Development

## 7.1 Model Selection and Rationale

In this chapter, focus on the selection and development of machine learning models for deepfake audio detection. Several well-known algorithms were evaluated to determine which would provide the best performance for the task.

The following models were chosen:

- **Random Forest Classifier**:

    A robust ensemble model that aggregates multiple decision trees to make predictions, reducing overfitting and improving classification accuracy.

- **Support Vector Machine (SVM)**:

    A powerful method for binary classification that aims to find the optimal hyperplane to separate data points of different classes.

- **XGBoost (Extreme Gradient Boosting)**:

    An efficient and highly accurate boosting algorithm that sequentially builds trees and focuses on reducing errors from previous models. It's known for handling imbalanced datasets well.

- **Gradient Boosting Classifier**:

    A gradient-based ensemble method that improves classification accuracy by fitting weak learners to residual errors of previous iterations.

## 7.2 Model Training and Validation

- **Data Split**:

    The dataset was split into training and testing sets using an **80:20 ratio**, meaning 80% of the data was used for training the models, and 20% was kept aside for testing. This helped in validating the models' generalization ability on unseen data.

- **Cross-Validation**:

    To further validate model performance, **5-fold cross-validation** was applied during training. This method splits the training data into 5 parts, where each part gets a chance to serve as the validation set while the other four are used for training. The results are then averaged to provide a more reliable estimate of model performance.

- **Hyperparameter Tuning**:

    Each model underwent hyperparameter tuning to optimize performance. Grid search and random search techniques were used to explore different values for hyperparameters such as:

    - **Random Forest**: Number of trees, max depth, min samples per leaf
    - **SVM**: Kernel type, C (regularization parameter), gamma
    - **XGBoost**: Learning rate, number of estimators, max depth
    - **Gradient Boosting**: Learning rate, n_estimators, max depth

## 7.3 Model Evaluation Metrics

- **Accuracy**:

    The percentage of correctly classified samples. While simple, accuracy provides a general sense of how well a model is performing.

- **Precision**:

    Measures the proportion of positive predictions that are correct. This is critical when false positives (classifying real audio as fake) can be problematic.

- **Recall**:

    Measures the proportion of actual positives correctly identified by the model. This metric is especially important when false negatives (failing to detect a deepfake) need to be minimized.

- **F1-Score**:

    The harmonic mean of precision and recall. The F1-score balances precision and recall, making it an important metric for tasks where both false positives and false negatives are costly.

# CHAPTER – 8

## 8. DISCUSSION

## 8.1 Presentation of Results

The evaluation of different machine learning models—XGBoost, Random Forest, SVM, and Gradient Boosting—yielded insights into their capabilities for deepfake audio detection. These models were evaluated based on several performance metrics, including accuracy, precision, recall, and F1-score. Among the models tested, Gradient Boosting exhibited superior performance, achieving an accuracy of 92.31%, the highest among all models.

A classification report based on precision, recall, and F1-score showed that the Gradient Boosting model excelled at identifying FAKE audio samples, with high precision and recall. Precision reflects how many of the predicted fake audio samples were correctly classified as fake, while recall measures how many of the actual fake audio samples were correctly identified. High precision and recall for the fake class suggest that the model effectively captured the characteristics of artificial audio signals. However, despite its excellent performance on fake audio detection, the REAL class posed a challenge for the model. The recall for REAL audio was relatively low, indicating that the model struggled to detect real audio samples accurately.

This difficulty in detecting real audio is not unique to this project but is a known issue in deepfake detection tasks. The challenge is attributed to the imbalanced nature of the dataset used for training. The real audio class is often underrepresented in deepfake detection datasets, leading to a bias where the model favors identifying fake audio over real audio. As a result, while the Gradient Boosting model performed excellently in distinguishing fake from real audio in terms of FAKE detection, it faltered in correctly classifying real audio.

## Summary of Results:

- Gradient Boosting: Best-performing model with 92.31% accuracy, excelling in fake audio detection.

- XGBoost, Random Forest, and SVM: All achieved similar performance with 84.6% accuracy, but struggled with low recall for the REAL audio class.

## 8.2 Comparison with Expectations and Literature

The results observed in this study, particularly the performance of the Gradient Boosting model, align closely with findings from existing literature on deepfake audio detection. Boosting algorithms, such as Gradient Boosting and XGBoost, have been widely praised for their ability to handle complex patterns in datasets, especially when distinguishing between subtle differences in classes like real and fake audio. Previous studies have shown that these algorithms, due to their ensemble nature, are highly effective at learning non-linear patterns and fine-grained features, which makes them well-suited for deepfake detection tasks where the difference between real and fake samples can be minimal.

For instance, XGBoost and Random Forest have frequently been used in classification problems due to their robustness and ability to handle imbalanced datasets. This is consistent with the findings of this study, where these models exhibited similar performance, both achieving 84.6% accuracy. However, the challenge of low recall for the REAL audio class in these models reflects a broader issue in deepfake detection. Imbalanced datasets have been a recurring problem in many deepfake audio studies, as the fake class tends to dominate the training process, thereby limiting the model's ability to accurately classify real audio samples.

One key insight from the literature is the widespread use of dataset balancing techniques or data augmentation to improve model performance, particularly in real audio detection. Techniques like SMOTE (Synthetic Minority Over-sampling Technique) or random sampling have been reported to enhance the ability of models to classify underrepresented classes, such as REAL audio in this case. Additionally, research has emphasized the need for more advanced feature extraction techniques, such as voice biometrics or acoustic analysis, to capture more detailed characteristics of real audio that might be missed by standard features like MFCCs or Spectral Centroid.

## Key Literature Insights:

- Ensemble models like Random Forest and Gradient Boosting are commonly used in classification tasks due to their ability to process complex data patterns.

- Deepfake audio detection remains challenging due to the subtle variations in real human speech, which are difficult to capture with traditional audio features.

- Imbalanced datasets are a recurring challenge, with real audio often being underrepresented, leading to difficulties in training models to detect it accurately.

## 8.3 Interpretation and Implications

The superior performance of the Gradient Boosting model highlights its potential in real-time deepfake audio detection applications, especially for identifying FAKE audio with high precision. This high performance underscores the power of ensemble methods, which combine the output of multiple weak learners to produce a stronger, more accurate model. The model's ability to detect fake audio reliably makes it an excellent candidate for real-world deepfake detection tasks, such as content moderation, forensic audio analysis, and security systems.

However, despite its success in detecting fake audio, the REAL class continues to present challenges. The low recall for real audio detection indicates that the model struggles to distinguish real audio from other types of audio. This limitation highlights an inherent issue in the detection of deepfake audio: real audio samples tend to have subtle, complex variations in tone, intonation, and speech patterns, which are difficult to model with standard feature extraction techniques. Furthermore, the imbalanced nature of the dataset exacerbates this challenge, as the model has been primarily trained to identify fake audio samples, which are often more pronounced and easier to classify.

## Implications for Future Work:

- Dataset Imbalance: The low recall for real audio suggests the need for more balanced datasets or techniques like SMOTE to generate synthetic real audio samples. Addressing this imbalance will likely improve the model's recall and F1-score for real audio.

- Advanced Feature Engineering: Incorporating more sophisticated features, such as voice biometrics, intonation analysis, and prosody features, could provide a richer understanding of real audio. These additional features might help the model better distinguish between real and fake samples.

- Real-Time Detection Optimization: Although the model performed well in offline evaluation, deploying the system for real-time detection poses additional challenges. Optimizing the model to handle large datasets and ensuring low-latency responses are

essential for real-world applications. This may involve optimizing the model for speed, reducing its computational requirements, and improving its ability to scale.

## 8.4 Limitations of the Study

While this study has contributed valuable insights into the use of machine learning for deepfake audio detection, several limitations must be addressed in future research. One of the most significant challenges encountered was the dataset imbalance, which negatively affected the recall and F1-score for the REAL audio class. This limitation highlights a fundamental issue in deepfake detection, where real audio samples are often underrepresented compared to fake samples. The imbalanced dataset results in the model being biased toward detecting fake audio, leading to a loss in real audio detection performance.

In addition to dataset imbalance, the limited set of audio features used in this study may have constrained the model's ability to capture the full range of characteristics present in both real and fake audio samples. While MFCCs, Chroma, and Spectral Centroid are commonly used features in audio classification, more advanced techniques like prosody analysis or intonation-based features could have provided deeper insights into the nuances of real human speech, further improving the model's performance.

Another important consideration is the risk of overfitting. The Gradient Boosting model, while highly effective, may be at risk of overfitting due to the imbalance in the dataset. Overfitting occurs when a model learns too much from the training data, capturing noise and irrelevant patterns that do not generalize well to new, unseen data. To mitigate this risk, more cross-validation and data augmentation techniques should be employed to ensure the model is robust and can generalize to a variety of real-world scenarios.

Finally, while the models were evaluated offline, real-time detection presents new challenges, particularly in terms of processing power and latency. Real-time deepfake audio detection requires high computational resources and quick response times, especially when deployed in applications that require instant feedback. Optimizing the model for real-time deployment without sacrificing accuracy is a critical area for future development.

# CHAPTER – 9

## 9. DEPLOYMENT AND USER INTERFACE

## 9.1 Backend Using Fast API & Frontend Using Streamlit

The integration of the backend and frontend is a crucial part of deploying any machine learning model in a user-friendly manner. For this project, **FastAPI** was selected for backend integration, while **Streamlit** was used for creating the frontend interface.

**Backend Integration with FastAPI:**

FastAPI is an asynchronous web framework for building APIs with Python, ideal for our machine learning model deployment. The primary role of FastAPI in this project was to create the API that interacts with the trained **Gradient Boosting** model. This API handles several key tasks:

- **Audio File Upload:** FastAPI provides an endpoint that allows users to upload their audio files for deepfake detection.

- **Feature Extraction and Prediction:** Once an audio file is uploaded, FastAPI extracts the required features (such as MFCCs, Chroma Features, and Spectral Centroid) and passes them to the trained **Gradient Boosting** model.

- **Prediction Response:** After processing the audio, the backend returns a prediction (whether the audio is real or fake) to the frontend.

FastAPI's high performance and ease of use made it the perfect choice for our backend, allowing fast asynchronous responses, which is critical for real-time applications.

**Frontend Development with Streamlit:**

Streamlit was chosen for its simplicity and ability to quickly create interactive web applications. It allows users to interact with the system without needing to write any code. The frontend performs several tasks:

- **Audio Upload:** The Streamlit interface has a file uploader widget, allowing users to drag and drop audio files directly into the app.

- **Prediction Display:** After receiving the prediction from the backend, Streamlit displays the result (whether the audio is real or fake) to the user.

- **User Experience:** Streamlit also allows the display of visualizations and feedback in a clean, easy-to-understand format, enhancing the user experience.

## 9.2 Real-Time Workflow of Detection System

The real-time workflow of the deepfake detection system is designed to be fast, efficient, and user-friendly. The following steps describe the process from the moment the user interacts with the system to receiving the result:

**Step 1: User Uploads Audio**

The first step in the workflow is the user interaction. The user uploads an audio file through the **Streamlit** interface. Streamlit provides a simple file uploader widget where the user can either drag and drop the audio file or select it manually from their system. The system supports common audio formats such as WAV, MP3, etc.

**Step 2: Audio File Sent to FastAPI Backend**

Once the user uploads the audio file, it is sent to the backend, which is powered by **FastAPI**. The file is sent via an HTTP POST request, where FastAPI handles the audio data, processes it, and prepares it for classification. The backend is designed to handle multiple types of audio files, ensuring flexibility in real-world applications.

**Step 3: Feature Extraction and Prediction**

After the audio file is received by the FastAPI backend, it undergoes feature extraction. The relevant features (MFCCs, Chroma Features, and Spectral Centroid) are extracted from the audio data using libraries such as **Librosa**. These features are then passed through the **Gradient Boosting** model, which has been trained to differentiate between real and fake audio. The model performs the classification task by analyzing the audio's characteristics and determining whether it belongs to the real or fake category.

The **Gradient Boosting** model was selected for its ability to handle complex relationships between features, and its superior accuracy in classifying audio samples. The prediction process is fast enough to provide near-instant feedback to the user, allowing for a real-time experience.

**Step 4: Display Results to User**

Once the backend completes the prediction, the result is sent back to the **Streamlit** frontend, where it is displayed to the user. The result typically shows a simple message indicating whether the audio is "Real" or "Fake." Additionally, more detailed information can be provided, such as the confidence level of the prediction, which further enhances the transparency and trust in the system.

**Step 5: Continuous Real-Time Detection**

The workflow ensures that each audio sample is processed independently and quickly. The system can handle many users simultaneously without significant delays, thanks to the asynchronous nature of FastAPI and the efficient design of Streamlit. The real-time interaction model allows users to test multiple audio files one after the other, with each file being processed and analyzed instantly.

By streamlining the workflow in this way, the system can be deployed in environments where quick audio classification is necessary, such as content moderation platforms, security systems, or any real-time application that requires deepfake detection.

# CHAPTER 10

## 10. CONCLUSION AND FUTURE WORK

## 10.1 Summary of Key Findings

This study focused on building and evaluating machine learning models for the detection of deepfake audio. Several models were explored, including **XGBoost**, **Random Forest**, **SVM**, and **Gradient Boosting**. Among these, **Gradient Boosting** emerged as the most effective, with an impressive accuracy of **92.31%**. It demonstrated superior performance in detecting **FAKE** audio samples, especially when compared to the other models, which all achieved around **84.6% accuracy**. However, despite its strength in identifying fake audio, the model faced significant challenges when attempting to detect **REAL** audio, exhibiting a low recall for this class due to the imbalanced nature of the training dataset.

Key findings include:

- **Gradient Boosting** showed the highest accuracy and robust performance in detecting fake audio, making it the best model for deepfake audio detection.

- **XGBoost**, **Random Forest**, and **SVM** performed similarly but struggled with detecting real audio due to their inability to generalize well on underrepresented real audio samples.

- The dataset imbalance, with a disproportionate number of fake samples, was a key factor in the low recall and F1-score for real audio classification.

- Despite its impressive offline performance, challenges in real-time detection remain, particularly in terms of latency, computational power, and the scalability of the models.

## 10.2 Conclusion

In conclusion, this study underscores the potential of machine learning models, particularly **Gradient Boosting**, in the detection of deepfake audio. The **Gradient Boosting** model's ability to identify fake audio with high precision demonstrates the viability of ensemble methods in tackling complex audio classification tasks. However, it also highlights critical challenges in deepfake detection, particularly in balancing the dataset and improving the detection of real audio samples. The low recall for the **REAL** class emphasizes the

difficulty of distinguishing subtle differences in human speech patterns, which are often hard to capture in training data.

The study's findings suggest that while deepfake detection in audio can be highly effective with the right models, further improvements are needed to achieve a balanced system that can accurately classify both fake and real audio samples. As the research progresses, addressing the imbalance in datasets, employing more advanced feature extraction techniques, and optimizing models for real-time deployment will be key to improving overall performance.

Ultimately, this research lays the groundwork for future deepfake audio detection systems that could be deployed in various sectors, such as security, media, and content moderation.

## 10.3 Practical Applications

The practical applications of deepfake audio detection are vast, given the rapid rise of artificial intelligence and the increasing sophistication of deepfake technologies. This research has important implications for a range of industries where the authenticity of audio content is crucial. Some of the key practical applications include:

- **Security and Forensics**: Detecting deepfake audio could be critical in security and forensic applications, such as identifying fraudulent audio used in scams, voice-based authentication systems, and legal investigations. For example, audio deepfakes might be used to impersonate voices in a corporate setting or to manipulate evidence in criminal cases. A robust deepfake detection system could prevent such misuse of audio content.

- **Content Moderation**: With the rise of platforms hosting user-generated content, deepfake detection systems could be used to flag fake or manipulated audio in videos or podcasts. Content platforms like social media, news outlets, and streaming services could use these systems to maintain the integrity of the content they distribute.

- **Voice Biometric Authentication**: As voice-based authentication systems become more popular, ensuring the legitimacy of the audio is paramount. The ability to detect deepfake voices would bolster security measures by ensuring that impersonators cannot bypass voice recognition systems.

- **Media and Entertainment**: The entertainment industry could use deepfake detection tools to maintain the credibility of audio content, ensuring that fake audio clips or voiceovers do not mislead audiences. Additionally, content creators can use such detection tools to verify the authenticity of user-generated content in a more transparent way.

- **Education and Awareness**: In the context of educational content, it is important to ensure that the voice recordings used for teaching materials or online courses are authentic. Deepfake detection models could be applied to verify such content, especially in sensitive contexts where misinformation could have significant repercussions.

## 10.4 Suggestions for Future Research and Improvement

While this study provides valuable insights, several avenues for future research and improvement could help enhance the accuracy and applicability of deepfake audio detection systems. Key suggestions for future work include:

- **Dataset Expansion and Balancing**: A major limitation of this study was the imbalanced dataset, which negatively impacted the recall and F1-score for real audio detection. Future research should focus on developing more balanced datasets, either by collecting more real audio samples or using techniques like **SMOTE** to generate synthetic real samples. Expanding the dataset to include various accents, languages, and background noises would also make the models more robust and generalizable.

- **Advanced Feature Engineering**: While **MFCC**, **Chroma**, and **Spectral Centroid** provided a good starting point for feature extraction, future work should consider incorporating more sophisticated features. **Voice biometrics**, **prosody**, **intonation**, and **acoustic analysis** could offer more detailed insights into the characteristics of real and fake audio. Such features might help the model better differentiate between human speech and artificially generated voices.

- **Model Optimization for Real-Time Detection**: The models used in this study were evaluated in an offline setting. However, deploying these models for real-time detection, such as in live video or phone calls, requires optimizing the model's performance to handle large datasets quickly. Real-time detection will require techniques to reduce latency, improve inference speed, and manage computational

requirements, possibly by leveraging more efficient architectures or edge computing solutions.

- **Cross-Platform Testing**: Deepfake detection systems should be tested across different audio sources and platforms to ensure they generalize well to real-world conditions. Future research should explore the generalization of deepfake detection models to various accents, voice tones, recording devices, and environments. This would help ensure the models' effectiveness in diverse, real-world scenarios.

- **Hybrid Models and Transfer Learning**: Combining multiple machine learning approaches, such as integrating **deep learning** with traditional models like **Gradient Boosting**, could offer better performance. Transfer learning, where models pre-trained on large datasets are fine-tuned on specific tasks, could also be explored. This approach may provide an efficient way to overcome the limitations of small or imbalanced datasets.

- **Explainability and Interpretability**: One important area of improvement lies in the interpretability of machine learning models. While models like **Gradient Boosting** can achieve high accuracy, understanding why the model makes a particular prediction is crucial, especially in critical applications like forensics or security. Future work should explore **explainable AI** techniques that provide insights into how models classify audio samples as real or fake.

# Chapter -11

## References

1. Chorowski, J., et al. (2023). *Attention-based Models for Speech Recognition*. In Proceedings of the Advances in Neural Information Processing Systems (NeurIPS).

   o This paper discusses attention mechanisms in speech recognition and provides foundational concepts in audio processing that may be applied to deepfake detection.

2. Goodfellow, I., et al. (2023). *Generative Adversarial Nets*. In Advances in Neural Information Processing Systems (NeurIPS).

   o The seminal paper on Generative Adversarial Networks (GANs), which are a primary method for generating deepfakes.

3. Zhou, L., & Zhang, Z. (2023). *Deepfake Detection in Audio: An Overview*. Journal of Digital Media.

   o This paper provides an overview of the challenges and methodologies used in deepfake detection, with a particular focus on audio.

4. Kumar, A., & Gupta, R. (2022). *XGBoost for Deepfake Detection: A Comparative Study*. International Journal of Machine Learning.

   o This study compares various machine learning models, including XGBoost, for deepfake detection tasks, and discusses performance metrics.

5. Afchar, D., et al. (2023). *MFC-HMM: Audio Deepfake Detection using Modified MFCC Features*. International Conference on Audio, Speech, and Signal Processing (ICASSP).

   o This paper explores the use of Modified Mel Frequency Cepstral Coefficients (MFCC) features in deepfake audio detection.

6. Nikolov, R., & Vassilev, T. (2022). *Real-time Deepfake Audio Detection using Gradient Boosting Machines*. Journal of AI Research.

o This study focuses on real-time deepfake audio detection using ensemble methods such as Gradient Boosting, which is relevant to the findings of this study.

7. Barton, D., & Li, Y. (2022). *Machine Learning for Audio Classification: Techniques and Applications*. Springer.

    o Provides comprehensive coverage on the various machine learning techniques applied to audio classification, including Random Forest, Gradient Boosting, and SVM, all of which were used in the present study.

8. Shen, D., et al. (2022). *Audio Deepfake Detection using Convolutional Neural Networks*. Journal of Computer Vision and Pattern Recognition.

    o This paper investigates the use of Convolutional Neural Networks (CNNs) for detecting deepfake audio, offering an alternative to traditional machine learning approaches.

9. Nguyen, P., & Zisserman, A. (2023). *Deepfake Detection in Audio and Video Streams*. IEEE Transactions on Information Forensics and Security.

    o The paper addresses the broader issue of deepfake detection in both audio and video streams, discussing challenges and solutions applicable across modalities.

10. Zhang, S., & Wang, J. (2023). *Using SMOTE for Handling Imbalanced Datasets in Audio Classification*. Journal of Machine Learning Research.

- This article discusses techniques like SMOTE for balancing datasets in audio classification tasks, which is relevant for addressing the dataset imbalance in the present study.

11. Li, H., & Liu, X. (2023). *Real-time Audio Deepfake Detection: Challenges and Solutions*. In Proceedings of the ACM Conference on Security and Privacy.

- This paper explores the challenges of deploying deepfake detection systems in real-time environments and suggests methods for improving computational efficiency.

12. Wang, Y., et al. (2023). *Advancements in Audio Deepfake Detection: A Comprehensive Review*. Computer Vision and Artificial Intelligence Journal.

# CHAPTER -12

## Appendix A: Code Snippets

## Backend

## main.py File

```python
# Import necessary libraries
import os
import numpy as np
import librosa
import joblib
from typing import List, Dict, Any
from fastapi import FastAPI, UploadFile, File, Form, HTTPException
from fastapi.middleware.cors import CORSMiddleware
from pydantic import BaseModel
import uvicorn
from xgboost import XGBClassifier
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report


app = FastAPI(title="Deep Fake Voice Detection API")

# Add CORS middleware to allow cross-origin requests
app.add_middleware(
    CORSMiddleware,
    allow_origins=["*"],
    allow_credentials=True,
    allow_methods=["*"],
    allow_headers=["*"],
)
```

```python
# ================================================
# SET YOUR FOLDER PATHS HERE
# ================================================

REAL_FOLDER_PATH = "D:\Deepfake_project\KAGGLE\AUDIO\FAKE"
FAKE_FOLDER_PATH = "D:\Deepfake_project\KAGGLE\AUDIO\REAL"
# ================================================

# Paths to save the trained models
MODEL_DIR = "models"
os.makedirs(MODEL_DIR, exist_ok=True)

# Function to extract audio features
def extract_features(file_path, mfcc=True, chroma=True, mel=True):
    """Extract features from audio file for classification"""
    try:
        # Load audio file
        y, sr = librosa.load(file_path, sr=None)

        # Set duration to 3 seconds (trim or pad)
        if len(y) > sr * 3:
            y = y[:sr * 3]
        else:
            y = np.pad(y, (0, max(0, sr * 3 - len(y))), 'constant')

        features = []

        # Extract features
        if mfcc:
            mfccs = np.mean(librosa.feature.mfcc(y=y, sr=sr, n_mfcc=13).T, axis=0)
            features.extend(mfccs)

        if chroma:
            chroma = np.mean(librosa.feature.chroma_stft(y=y, sr=sr).T, axis=0)
```

```python
            features.extend(chroma)

        if mel:
            mel = np.mean(librosa.feature.melspectrogram(y=y, sr=sr).T, axis=0)
            features.extend(mel)

        return features

    except Exception as e:
        print(f"Error extracting features from {file_path}: {str(e)}")
        return None

# Function to process all audio files in a directory
def process_audio_directory(directory):
    features = []
    labels = []

    # Determine if this is a "REAL" or "FAKE" directory
    is_real = "REAL" in directory
    label = 1 if is_real else 0

    for root, _, files in os.walk(directory):
        for file in files:
            if file.endswith(('.wav', '.mp3', '.flac', '.ogg')):
                file_path = os.path.join(root, file)

                # Extract features from the audio file
                audio_features = extract_features(file_path)

                if audio_features is not None:
                    features.append(audio_features)
                    labels.append(label)

    return features, labels
```

```python
def convert_numpy_types(obj):
    """Convert numpy types to Python native types for JSON serialization"""
    if isinstance(obj, (np.int8, np.int16, np.int32, np.int64,
                        np.uint8, np.uint16, np.uint32, np.uint64)):
        return int(obj)
    elif isinstance(obj, (np.float16, np.float32, np.float64)):
        return float(obj)
    elif isinstance(obj, np.ndarray):
        return obj.tolist()
    elif isinstance(obj, dict):
        return {k: convert_numpy_types(v) for k, v in obj.items()}
    elif isinstance(obj, list):
        return [convert_numpy_types(i) for i in obj]
    else:
        return obj


@app.post("/train")
async def train_model(real_folder: str = Form(None), fake_folder: str = Form(None)):
    """
    Train multiple models using Real and Fake audio folders and compare their performance.
    If no folders are provided, use the default paths set at the top of the file.
    """
    try:
        # Use provided folders or fall back to defaults
        real_dir = real_folder if real_folder else REAL_FOLDER_PATH
        fake_dir = fake_folder if fake_folder else FAKE_FOLDER_PATH

        # Validate directories
        if not os.path.isdir(real_dir):
            raise HTTPException(status_code=400, detail=f"Real folder path '{real_dir}' is not a
valid directory")
        if not os.path.isdir(fake_dir):
```

```python
        raise HTTPException(status_code=400, detail=f"Fake folder path '{fake_dir}' is not a
valid directory")

    # Process real audio files
    print(f"Processing real audio files from: {real_dir}")
    real_features, real_labels = process_audio_directory(real_dir)

    # Process fake audio files
    print(f"Processing fake audio files from: {fake_dir}")
    fake_features, fake_labels = process_audio_directory(fake_dir)

    # Combine datasets
    X = np.array(real_features + fake_features)
    y = np.array(real_labels + fake_labels)

    if len(X) == 0:
        raise HTTPException(status_code=400, detail="No valid audio files found in the
provided directories")

    print(f"Total samples: {len(X)}, Real: {sum(y)}, Fake: {len(y) - sum(y)}")

    # Split data for training and testing
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

    # Define models to train
    models = {
        "xgboost": XGBClassifier(
            n_estimators=100,
            learning_rate=0.1,
            max_depth=5,
            random_state=42,
            use_label_encoder=False,
            eval_metric='logloss'
        ),
```

```python
        "random_forest": RandomForestClassifier(
            n_estimators=100,
            max_depth=10,
            random_state=42
        ),
        "svm": SVC(
            kernel='rbf',
            probability=True,
            random_state=42
        ),
        "gradient_boosting": GradientBoostingClassifier(
            n_estimators=100,
            learning_rate=0.1,
            max_depth=5,
            random_state=42
        )
    }

    # Results dictionary
    results = {}

    # Train and evaluate each model
    best_model = None
    best_accuracy = 0.0

    for model_name, model in models.items():
        print(f"Training {model_name}...")
        model.fit(X_train, y_train)

        # Evaluate the model
        y_pred = model.predict(X_test)
        accuracy = float(accuracy_score(y_test, y_pred))
        report = classification_report(y_test, y_pred, target_names=["FAKE", "REAL"],
output_dict=True)
```

```python
    # Save the model
    model_path = os.path.join(MODEL_DIR, f"{model_name}_model.pkl")
    joblib.dump(model, model_path)

    # Store results
    results[model_name] = {
        "accuracy": accuracy,
        "classification_report": convert_numpy_types(report)
    }

    # Check if this is the best model so far
    if accuracy > best_accuracy:
        best_accuracy = accuracy
        best_model = model_name

# Save the best model as the default model
best_model_path = os.path.join(MODEL_DIR, f"{best_model}_model.pkl")
default_model_path = os.path.join(MODEL_DIR, "default_model.pkl")
# Copy the best model to the default model path
best_model_obj = joblib.load(best_model_path)
joblib.dump(best_model_obj, default_model_path)

# Return comprehensive results
return {
    "message": "Multiple models trained successfully",
    "total_samples": int(len(X)),
    "real_samples": int(sum(y)),
    "fake_samples": int(len(y) - sum(y)),
    "real_folder_used": real_dir,
    "fake_folder_used": fake_dir,
    "results": results,
    "best_model": {
        "name": best_model,
```

```python
                "accuracy": float(best_accuracy)
            },
            "models_saved": list(models.keys())
        }


    except Exception as e:
        raise HTTPException(status_code=500, detail=f"Training error: {str(e)}")


@app.post("/predict")
async def predict_audio(file: UploadFile = File(...), model_name: str = Form("default")):
    """

    Predict if the uploaded audio file is real or fake using the specified model
    """

    try:
        # Determine which model to use
        if model_name == "default":
            model_path = os.path.join(MODEL_DIR, "default_model.pkl")
        else:
            model_path = os.path.join(MODEL_DIR, f"{model_name}_model.pkl")


        # Check if model exists
        if not os.path.exists(model_path):
            available_models = [f.replace("_model.pkl", "") for f in os.listdir(MODEL_DIR) if
f.endswith("_model.pkl")]
            if not available_models:
                raise HTTPException(status_code=400, detail="No models trained yet. Please train
models first.")
            else:
                raise HTTPException(status_code=400, detail=f"Model '{model_name}' not found.
Available models: {', '.join(available_models)}")


        # Save the uploaded file temporarily
        temp_file = f"temp_{file.filename}"
        with open(temp_file, "wb") as f:
```

```python
        f.write(await file.read())

    # Extract features
    features = extract_features(temp_file)

    if features is None:
        os.remove(temp_file)
        raise HTTPException(status_code=400, detail="Could not extract features from the
uploaded file")

    # Load model and make prediction
    clf = joblib.load(model_path)
    features_array = np.array(features).reshape(1, -1)
    prediction = int(clf.predict(features_array)[0])
    probabilities = clf.predict_proba(features_array)[0]

    # Clean up temp file
    os.remove(temp_file)

    # Return prediction results
    result_label = "REAL" if prediction == 1 else "FAKE"
    confidence = float(probabilities[prediction])

    return {
        "filename": file.filename,
        "model_used": model_name,
        "prediction": result_label,
        "confidence": confidence,
        "probability_fake": float(probabilities[0]),
        "probability_real": float(probabilities[1])
    }

except Exception as e:
    # Clean up in case of error
```

```python
        if os.path.exists(f"temp_{file.filename}"):
            os.remove(f"temp_{file.filename}")
        raise HTTPException(status_code=500, detail=f"Prediction error: {str(e)}")


if __name__ == "__main__":
    print(f"Starting Deep Fake Voice Detection API with multiple ML classifiers")
    print(f"Default REAL folder path: {REAL_FOLDER_PATH}")
    print(f"Default FAKE folder path: {FAKE_FOLDER_PATH}")
    print(f"Models will be saved to: {MODEL_DIR}")
    uvicorn.run(app, host="127.0.0.1", port=8000)
```

## Frontend

```python
import streamlit as st

import requests

import numpy as np

import librosa

import librosa.display

import matplotlib.pyplot as plt

# Set page title and layout

st.set_page_config(page_title="Deep Fake Voice Detection", page_icon="🔊",
layout="wide")

# Custom CSS for styling

st.markdown(

    """

    <style>

        .stAudio {

            background-color: white;

            padding: 5px;
```

```
        border-radius: 5px;

    }

    .stButton button {

        background-color: #007BFF !important;

        color: white !important;

        font-size: 16px !important;

        padding: 8px 20px !important;

        border-radius: 10px !important;

        display: flex;

        margin: 0 auto;

    }

</style>

""",

unsafe_allow_html=True

)

# Title

st.markdown("<h2 style='text-align: center; color: #007BFF;'>🔊 Deep Fake Voice
Detection</h2>", unsafe_allow_html=True)

# File uploader

uploaded_file = st.file_uploader("📂 Upload an audio file", type=["wav", "mp3", "flac",
"ogg"])

API_URL = "http://127.0.0.1:8000/predict"  # FastAPI backend URL

if uploaded_file:

    st.markdown("### 🎵 Audio Preview")

    st.audio(uploaded_file, format='audio/wav')
```

```python
# Convert uploaded file to numpy array for analysis

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


# Display waveform & spectrogram side by side

col1, col2 = st.columns(2)

with col1:

    st.markdown("### 📊 Waveform")

    fig, ax = plt.subplots(figsize=(4, 2))  # Smaller size

    librosa.display.waveshow(y, sr=sr, ax=ax, color='blue')

    ax.set_xlabel("Time (s)", fontsize=8)

    ax.set_ylabel("Amplitude", fontsize=8)

    ax.set_title("Waveform", fontsize=10)

    ax.tick_params(axis='both', labelsize=6)

    st.pyplot(fig)


with col2:

    st.markdown("### 🎭 Spectrogram")

    fig, ax = plt.subplots(figsize=(4, 2))  # Smaller size

    S = librosa.amplitude_to_db(np.abs(librosa.stft(y)), ref=np.max)

    librosa.display.specshow(S, sr=sr, x_axis='time', y_axis='log', ax=ax, cmap='coolwarm')

    ax.set_title("Spectrogram", fontsize=10)

    ax.tick_params(axis='both', labelsize=6)

    st.pyplot(fig)
```

```python
# Centered Analyze button

st.markdown("<div style='text-align: center;'>", unsafe_allow_html=True)

analyze = st.button("🚀 Analyze Audio")

st.markdown("</div>", unsafe_allow_html=True)

if analyze:

    with st.spinner("Processing... 🔄"):

        files = {"file": uploaded_file.getvalue()}

        response = requests.post(API_URL, files=files)

    if response.status_code == 200:

        result = response.json()

        prediction = result['prediction']

        # Display prediction

        st.success(f"✅ Prediction: {prediction}")

    else:

        st.error("❌ Error processing the file.")
```