*A project report on*

# Emotional Analysis of Audio using Deep Learning Techniques

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

# Bachelor of Technology in Computer Science and Engineering with Specialization in Artificial Intelligence and Robotics

*by*

## Mohite Aditya Kiran (21BRS1615)

## Krishna Deepak Deshpande (21BRS1001)

## Rohan BC (21BRS1016)

**VIT®**

**Vellore Institute of Technology**
(Deemed to be University under section 3 of UGC Act, 1956)
CHENNAI

## SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

November,2024

# Emotional Analysis of Audio using Deep Learning Techniques

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

## Bachelor of Technology in Computer Science and Engineering with Specialization in Artificial Intelligence and Robotics

*by*

**Mohite Aditya Kiran (21BRS1615)**

**Krishna Deepak Deshpande (21BRS1001)**

**Rohan BC (21BRS1016)**

**VIT** ®

**Vellore Institute of Technology**
(Deemed to be University under section 3 of UGC Act, 1956)
CHENNAI

**SCHOOL OF COMPUTER SCIENCE AND ENGINEERING**
November, 2024

## DECLARATION

I hereby declare that the thesis entitled "Emotional Analysis of Audio using Deep Learning Techniques" submitted by **Mohite Aditya Kiran (21BRS1615)**, for the award of the degree of Bachelor of Technology in Computer Science and Engineering with specialization in AI and Robotics, Vellore Institute of Technology, Chennai is a record of bonafide work carried out by me under the supervision of **Dr. Vinothini A**.

I further declare that the work reported in this thesis has not been submitted and will not be submitted, either in part or in full, for the award of any other degree or diploma in this institute or any other institute or university.

Place: Chennai

Date: 15/11/2024                                        Signature of the Candidate

# School of Computer Science and Engineering

# CERTIFICATE

This is to certify that the report entitled **"Emotional Analysis of Audio using Deep Learning Techniques"** is prepared and submitted by **Mohite Aditya Kiran (21BRS1615)** to Vellore Institute of Technology, Chennai, in partial fulfillment of the requirement for the award of the degree of **Bachelor of Technology in Computer Science and Engineering with Specialization in Artificial Intelligence and Robotics** is a bonafide record carried out under my guidance. The project fulfills the requirements as per the regulations of this University and in my opinion meets the necessary standards for submission. The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma and the same is certified.


Signature of the Guide:

Name: Dr.  Vinothini A

Date:


Signature of the Examiner           Signature of the Examiner

Name:                                               Name:

Date:                                                 Date:


Approved by the Head of Department,
**Computer Science with Specialization in Artificial Intelligence and Robotics**


Name: Dr. Harini S
Date:15/11/2024

# <u>ABSTRACT</u>

A deep learning model for emotional analysis of audio is thoroughly investigated in this project, highlighting its growing significance in applications such as virtual assistants, interactive entertainment, mental health evaluation, and human-computer interaction. The key challenge in this field is accurately detecting emotional states from diverse and noisy audio environments. To address these challenges, the project proposes an end-to-end solution that leverages the capabilities of Long Short-Term Memory (LSTM) networks combined with Convolutional Neural Networks (CNNs) for efficient emotion classification.

To ensure consistency, different audio recordings from multiple datasets are standardized during the preparation stage of the research workflow. The integration of well-known datasets, including TESS, RAVDESS, CREMA-D, and SAVEE, into a single training set enhances the model's ability to generalize across diverse speaker attributes, emotional expressions, and recording scenarios.

The model employs data augmentation techniques to increase its resilience and flexibility, such as time-shifting, pitch shifting, speed adjustments, and noise injection. These techniques improve the model's capability to detect emotions in varied and noisy settings by exposing it to diverse auditory conditions.

For feature extraction, spectrograms capture both time-domain and frequency-domain information, transforming raw audio into visual representations. From these spectrograms, Mel-Frequency Cepstral Coefficients (MFCCs) are extracted, allowing the model to detect subtle emotional cues while remaining language-agnostic, thus suitable for multilingual applications.

The proposed approach utilizes a hybrid model combining CNNs to efficiently extract spatial patterns from spectrograms and LSTM networks to process sequential data, leveraging long-term dependencies in audio signals. The architecture includes dropout layers to mitigate overfitting and ensures robust emotion classification. This end-to-end model is shown to be effective in classifying emotions even in challenging real-world audio environments.

# ACKNOWLEDGEMENT

Place: Chennai

Date: 15/11/2024
<div align="right">

**Mohite Aditya Kiran**
**(21BRS1615)**
</div>

# CONTENTS

**CHAPTER 4**

**RESULTS AND DISCUSSION**

**CHAPTER 5**

**Conclusion and Future Work**

**Appendices**

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ACRONYMS

| Term | Full Form |
| --- | --- |
| LSTM | Long Short-Term Memory |
| TESS | Toronto Emotional Speech Set |
| RAVDESS | Ryerson Audio-Visual Database of Emotional Speech and Song |
| CREMA-D | Crowd-sourced Emotional Multimodal Actors Dataset |
| MFCCs | Mel-Frequency Cepstral Coefficients |
| ReLU | Rectified Linear Unit |
| F1-score | F1 Measure |
| CNN | Convolutional Neural Network |
| RNN | Recurrent Neural Network |
| GRU | Gated Recurrent Unit |
| NLP | Natural Language Processing |
| LDA | Linear Discriminant Analysis |
| PCA | Principal Component Analysis |
| SRHA | Self-Recursive Hierarchical Attention |
| FCN | Fully Convolutional Network |
| CASIA | Chinese Academy of Sciences Institute of Automation |
| SAVEE | Surrey Audio-Visual Expressed Emotion Database |
| IEMOCAP | Interactive Emotional Dyadic Motion Capture Database |
| STFT | Short-Time Fourier Transform |
| ORB | Oriented FAST and Rotated BRIEF |
| BoVW | Bag of Visual Words |
| SVM | Support Vector Machine |
| kNN | k-Nearest Neighbors |
| ANN | Artificial Neural Network |
| IoT | Internet of Things |
| HCI | Human-Computer Interaction |
| PME | Peripheral Muscle Electroactivity |
| EMG | Electromyography |

| | |
|---|---|
| EEG | Electroencephalography |
| SER | Speech Emotion Recognition |
| DNN | Deep Neural Network |
| AWS | Amazon Web Services |
| LBP | Local Binary Patterns |
| MLP | Multilayer Perceptron |
| MER | Music Emotion Recognition |
| MIR | Music Information Retrieval |

**Chapter 1**

# Introduction

## 1.1 INTRODUCTION

Emotion recognition through audio signals has become an increasingly important research domain with diverse applications in fields such as human-computer interaction, affective computing, mental health diagnostics, and customer service. With advancements in deep learning and neural networks, researchers have been able to push the boundaries of accuracy and reliability in interpreting emotional states from audio signals. However, the task of accurately identifying emotions remains a significant challenge, especially due to the inherent noise and variability in real-world environments. Emotional signals in speech are often distorted by background noise, varying speech patterns, speaker characteristics, and environmental conditions, making emotion recognition a complex problem. This study aims to enhance the accuracy of emotional recognition in audio by leveraging deep learning techniques, specifically focusing on Long Short-Term Memory (LSTM) networks, a type of Recurrent Neural Network (RNN) designed for sequential data analysis, to capture the temporal dependencies of audio signals.

Deep learning, particularly through the use of LSTMs, has revolutionized the field of audio analysis. Traditional approaches to emotion recognition often rely on hand-crafted features such as pitch, tone, and timbre, which can be limiting in capturing complex, high-dimensional relationships within the data. Deep learning, on the other hand, excels at automatically learning these intricate patterns directly from raw data without needing manual feature extraction. The ability of LSTMs to learn long-range dependencies in sequential data—such as speech, where emotional context often spans multiple time steps—makes them particularly effective for this task.

In this study, we present a systematic approach for preprocessing, augmenting, and analyzing audio data to train a deep learning model capable of emotion recognition. We focus on a multi-step methodology to ensure that the model is robust, reliable, and capable of generalizing to real-world conditions. The preprocessing phase begins with the standardization of audio files, ensuring that the dataset has a consistent format. This involves normalizing the sampling rate, bit depth, and audio duration to ensure that all audio files are comparable, reducing the potential for noise and inconsistencies to affect the training process. Standardization of audio files not only enhances the data's quality but also ensures uniformity across the dataset, which is crucial for accurate feature extraction and model training.

Once the dataset is standardized, it is augmented to make the model more resilient to noise and variability. In real-world scenarios, audio data is often contaminated by background noise, recording distortions, or environmental factors. To improve the robustness of the

model, we apply several data augmentation techniques, such as Noise Injection, Time-Shifting, and Pitch Alteration. Noise Injection adds random noise to the audio samples, simulating real-world recording conditions where background noise can affect speech clarity. Time-Shifting changes the timing of the audio samples slightly, forcing the model to generalize to variations in speech speed. Pitch Alteration alters the pitch of the audio without changing its emotional content, mimicking the natural variation in speech tone that occurs across different speakers and contexts. These augmentation techniques enhance the dataset, making it more diverse and enabling the model to learn to classify emotions under a broader range of conditions.

Following data augmentation, the next step in the process is data preprocessing, where raw audio signals are transformed into machine-readable features. This is done using the Librosa library, a powerful tool for audio analysis that allows for efficient extraction of relevant audio features. One of the key features extracted from the audio is the Mel-Frequency Cepstral Coefficients (MFCCs). MFCCs are a widely used feature representation in speech processing, capturing the short-term power spectrum of the audio signal. These coefficients effectively represent the spectral properties of sound and are well-suited for emotion recognition tasks. MFCCs allow the model to focus on the most important characteristics of the audio, facilitating the detection of patterns that correspond to different emotional states.

With the preprocessed and augmented dataset ready, we proceed to the model training phase using an LSTM network. LSTMs, as a type of Recurrent Neural Network (RNN), are specifically designed to process sequential data, such as audio signals, where the temporal sequence of information is critical for accurate interpretation. The LSTM model architecture employed in this study is designed using the Keras library and consists of several key layers.

The input layer of the LSTM model takes a 2D input, which consists of a sequence of MFCC features extracted from the audio samples. These features are used as the input to the LSTM layer, which is the core of the model. The LSTM layer has 256 units, chosen for its ability to capture complex temporal dependencies in the audio data. These units enable the model to remember important information over long sequences of audio, which is particularly valuable for detecting emotions that may unfold gradually over time. To prevent overfitting, Dropout layers are added after the LSTM and Dense layers. Dropout is a regularization technique that randomly disables a fraction of the units during training, ensuring that the model does not rely too heavily on any one feature.

Following the LSTM layer, two fully connected Dense layers are included in the model architecture, with 128 and 64 neurons, respectively. These layers further process the extracted features and apply non-linear transformations through the ReLU (Rectified Linear Unit) activation function. The use of multiple Dense layers allows the model to learn more abstract representations of the audio data, improving its ability to classify emotions accurately. Finally, the output layer of the model is a Dense layer with 7 neurons, each corresponding to a different emotion category. A softmax activation function is

applied at the output layer to generate probabilities for each emotion, allowing the model to predict the most likely emotional state for any given audio sample.

This deep learning-based approach, utilizing LSTMs and comprehensive preprocessing and augmentation techniques, aims to push the boundaries of emotion recognition accuracy in audio data. By leveraging these advanced techniques, we hope to build a model that is not only more accurate but also more robust to the noise and variability encountered in real-world scenarios, ultimately advancing the potential applications of audio-based emotional analysis.

## 1.2 OVERVIEW

Emotional analysis of audio involves examining audio signals to identify and classify the emotions expressed by a speaker or performer. This can be done through a combination of signal processing, machine learning, and computational techniques that analyze features such as tone, pitch, rhythm, and timbre. Here's a detailed overview of the key components involved in emotional analysis of audio:

### 1.2.1. FEATURES OF AUDIO SIGNALS

The first step in emotional analysis of audio is extracting relevant features from raw audio data, which provides essential information to identify underlying emotions. These features capture the nuances of speech that reflect a person's emotional state. Below are the key audio features commonly used for emotion recognition:

- Pitch (Fundamental Frequency):

  Definition: Pitch is the perceptual counterpart to the frequency of the sound wave. It refers to how "high" or "low" a voice sounds.

  Emotional Insights: Variations in pitch can reflect different emotional states:

  High pitch is often associated with emotions such as anger, fear, or excitement, indicating heightened arousal.

  Low pitch can signify emotions like sadness or boredom, suggesting a calmer or subdued state.

  Application: Detecting changes in pitch can be particularly useful for distinguishing between emotions like anger (sharp rise in pitch) and sadness (drop in pitch).

- Energy (Loudness):

  Definition: Energy refers to the overall loudness or intensity of an audio signal, reflecting the strength of the sound wave.

  Emotional Insights:

Higher energy levels are generally linked to emotions such as happiness, anger, or surprise, indicating an active or excited emotional state.

Lower energy may correspond to emotions like sadness, tiredness, or disgust, representing a more passive or subdued feeling.

Application: Analyzing energy levels helps to differentiate between high-energy emotions like excitement and low-energy emotions like sadness.

- Formants:

Definition: Formants are the resonant frequencies in the human voice that shape the unique qualities of vowel sounds. These are influenced by the shape of the vocal tract.

Emotional Insights:

Emotional states can affect the formant frequencies, resulting in shifts that alter the timbre of speech.

For example, anger might tighten the vocal tract, raising formant frequencies, while sadness could relax it, lowering these frequencies.

Application: Tracking formant changes can be useful for detecting subtle emotional variations in speech.

- Tempo (Speech Rate):

Definition: Tempo refers to the speed at which someone speaks, measured in syllables or words per minute.

Emotional Insights:

Faster speech rates are associated with emotions like anxiety, fear, or excitement, often indicating nervousness or high arousal.

Slower speech suggests calmness, boredom, or sadness, reflecting a relaxed or subdued state.

Application: Detecting changes in tempo can help distinguish between emotions like anxiety (rapid speech) and sadness (slow, drawn-out speech).

- Mel-frequency Cepstral Coefficients (MFCCs):

Definition: MFCCs are coefficients that capture the short-term power spectrum of an audio signal. They are derived by applying a Mel-scale filter bank to the spectrum of an audio signal, mimicking the way humans perceive sound.

Emotional Insights:

MFCCs are effective in capturing the tonal quality and speech patterns that convey emotions, such as tone variations, stress, and articulation.

They are particularly useful for emotion detection because they emphasize frequencies that are more perceptually relevant to human hearing.

Application: MFCCs are widely used in speech emotion recognition models to extract features that help differentiate between emotional states in a language-agnostic way.

- Voice Quality:

  Definition: Voice quality refers to the unique characteristics of a person's voice, influenced by factors such as breathiness, harshness, and nasality.

  Emotional Insights:

  Breathy voices may be associated with relief or sadness, while a strained voice could indicate anger or stress.

  Timbre changes can reflect emotional states, such as a richer timbre for happiness or a rougher timbre for anger.

  Application: Analyzing variations in voice quality can reveal the emotional undercurrents of a speaker, even if the actual words spoken are neutral.

- Pauses and Speech Disfluencies:

  Definition: Pauses are the silent intervals between words or phrases, while speech disfluencies include hesitations, repetitions, and filler words (like "uh" and "um").

  Emotional Insights:

  Frequent pauses or long hesitations may indicate sadness, confusion, or nervousness, suggesting cognitive load or emotional stress.

  Rapid, disfluent speech could reflect emotions like excitement, nervousness, or anxiety, as the speaker struggles to convey thoughts quickly.

  Application: Monitoring speech patterns, including the frequency and timing of pauses, can be useful for detecting emotional states in real-time, especially in natural, conversational speech.

- Spectral Features :

  Definition: These include features like spectral centroid, spectral bandwidth, and

spectral contrast that describe the frequency content of the audio signal.

Emotional Insights:

Spectral centroid (perceived brightness) can be higher in angry or happy speech due to increased high-frequency energy.

Spectral contrast can differentiate between emotions by identifying the range between loud and quiet frequencies, which can be higher in more intense emotions like anger.

Application: Extracting spectral features helps in distinguishing between subtle variations in emotional expression, especially in complex or noisy environments.

## 1.2.2. EMOTION CATEGORIE

Emotional analysis of speech focuses on identifying specific emotions conveyed by a speaker's tone, pitch, rhythm, and energy. Understanding the acoustic features associated with different emotions is crucial for developing effective Speech Emotion Recognition (SER) systems. Here, we expand on the common emotion categories and describe their distinct characteristics:

1. Happiness

   Characteristics:

   - Pitch: Generally higher than normal speech, reflecting excitement or positive emotions.
   - Speech Rate: Often faster, indicating a state of enthusiasm or high energy.
   - Energy: Increased intensity and volume, producing a vibrant and lively tone.
   - Prosody: Smoother pitch contours with more variation, leading to a dynamic and expressive voice.

   Emotional Insights:

   - Happiness is often associated with uplifted tones, spontaneous laughter, and bright voice quality. Detecting these changes can help systems recognize joyful expressions, which is particularly useful in applications like customer satisfaction analysis.

2. Sadness

   Characteristics:

   - Pitch: Typically lower, reflecting a subdued and introspective state.
   - Speech Rate: Noticeably slower than normal speech, conveying a lack of urgency or enthusiasm.

- ○ Energy: Reduced loudness and intensity, resulting in a soft, mellow tone.
- ○ Prosody: Less pitch variation, producing a flatter, monotonous sound.

Emotional Insights:

- ○ Sadness often includes longer pauses, drawn-out words, and a sense of lethargy in the speaker's voice. These features are crucial for detecting distress in mental health applications, like assessing depressive states.

## 3. Anger

Characteristics:

- ○ Pitch: Sharp increases in pitch and fluctuations, especially at the start of utterances.
- ○ Speech Rate: Faster than normal, often reflecting agitation or urgency.
- ○ Energy: High intensity and volume, resulting in a loud, forceful tone.
- ○ Prosody: Abrupt and jerky pitch patterns with pronounced stress on certain syllables.

Emotional Insights:

- ○ Anger is characterized by tense voice quality, explosive bursts of sound, and harsh enunciation. Detecting anger can be beneficial in safety monitoring systems, conflict resolution tools, and customer support analysis.

## 4. Fear

Characteristics:

- ○ Pitch: Typically high, with sudden increases indicating alarm or anxiety.
- ○ Speech Rate: Often fast, reflecting nervousness or panic.
- ○ Voice Quality: Breathy tone, with inconsistent airflow as the speaker struggles to maintain steady breathing.
- ○ Prosody: Irregular pitch variations, leading to an unstable and tense speech pattern.

Emotional Insights:

- ○ Fear is often conveyed through quivering or trembling in the voice, accompanied by rapid and shallow breaths. This detection can be useful in emergency response systems and psychological assessments for anxiety disorders.

## 5. Surprise

Characteristics:

- Pitch: Sudden jumps in pitch, particularly at the start of an utterance, reflecting shock or astonishment.
- Energy: Noticeable surges in volume, but often brief and sporadic.
- Speech Rate: Varies; can either accelerate or momentarily pause due to the surprise.
- Prosody: Sharp, dynamic changes in pitch and rhythm to reflect unexpectedness.

Emotional Insights:

- Surprise can be identified by analyzing abrupt shifts in speech characteristics, making it valuable in interactive entertainment and adaptive AI systems.

## 6. Disgust

Characteristics:

- Pitch: Generally lower than normal, indicating disdain or aversion.
- Speech Rate: Often irregular, with abrupt pauses or changes in tempo.
- Energy: Moderate to low, but with a noticeable harshness in tone.
- Prosody: Flattened intonation with less dynamic range, contributing to a tone of displeasure.

Emotional Insights:

- Disgust is marked by a rough or gravelly quality in the voice, often accompanied by guttural sounds. Recognizing this emotion can enhance sentiment analysis in areas like product reviews or content moderation.

## 7. Neutral/Calm

Characteristics:

- Pitch: Consistent and moderate, with minimal variation.
- Speech Rate: Generally steady and evenly paced.
- Energy: Balanced, neither too loud nor too soft, indicating a state of composure.
- Prosody: Smooth and even pitch contours, reflecting a calm and collected demeanor.

Emotional Insights:

- Neutral or calm speech is characterized by its lack of emotional extremes, making it the baseline against which other emotions are detected. This category is essential for systems that aim to distinguish between heightened emotional states and normal conversations.

## 1.2.3. EMOTION DETECTION TECHNIQUES

Emotion detection in speech involves a combination of techniques to analyze the audio signal, extract relevant features, and classify these features to identify the emotional state conveyed by the speaker. Below is an expanded overview of the various methods used in the process, covering signal processing, machine learning, and deep learning approaches.

### 1. Signal Processing and Feature Extraction

- The first stage in detecting emotions in speech involves processing the raw audio signal to capture its distinctive characteristics. This step is crucial as it transforms the audio into meaningful data that can be analyzed by machine learning models.
- Key Features Extracted:
    - Pitch (Fundamental Frequency): Represents the perceived tone of the voice, with variations indicating emotions like anger (higher pitch) or sadness (lower pitch).
    - Energy (Loudness): Reflects the intensity of speech; higher energy often indicates excitement or anger, while lower energy might suggest calmness or sadness.
    - Mel-Frequency Cepstral Coefficients (MFCCs): Capture the short-term power spectrum of speech and are widely used for recognizing emotions due to their ability to capture vocal nuances.
    - Spectrograms and Mel-Spectrograms: Visual representations of the frequency content over time, which help in capturing both time-domain and frequency-domain information of the audio.
    - Prosodic Features: Include speech rate, pauses, and intonation patterns, which can signal emotions like fear (fast speech) or boredom (slow speech).
    - Voice Quality and Formants: Capture characteristics like breathiness, harshness, and resonance, providing additional layers of emotional context.

### 2. Machine Learning Models

- Traditional machine learning techniques are often used to classify emotions based on extracted features. These models generally involve supervised learning, where the system is trained on labeled datasets to identify patterns associated with different emotional states.
- Common Algorithms:
    - Support Vector Machines (SVM):
        - Frequently used for classification tasks due to their effectiveness in high-dimensional spaces.
        - SVMs work by finding the optimal hyperplane that separates different emotion categories based on features like MFCCs and prosodic cues.
    - Random Forests and Decision Trees:
        - These are ensemble methods that leverage multiple decision trees to

improve classification accuracy.
- They are particularly effective for capturing non-linear relationships between features, making them suitable for classifying complex emotions.
○ K-Nearest Neighbors (KNN):
- A straightforward algorithm that classifies data points based on their proximity to labeled examples in the feature space.
- Useful for smaller datasets but may struggle with high-dimensional data common in speech analysis.
○ Logistic Regression:
- Often used for binary or multiclass emotion classification, leveraging probabilistic models to predict emotional categories.

3. Deep Learning for Sequential Data

- Traditional machine learning approaches, while effective, often struggle with capturing the temporal dependencies inherent in speech data. This is where deep learning models excel, especially in handling the sequential nature of audio signals.
- Key Deep Learning Architectures:
  ○ Recurrent Neural Networks (RNNs):
  - Designed to process sequential data, RNNs can capture context over time by maintaining hidden states that carry information from previous inputs.
  - However, they can suffer from the vanishing gradient problem, limiting their ability to handle long sequences.
  ○ Long Short-Term Memory (LSTM) Networks:
  - An advanced version of RNNs, LSTMs, introduce memory cells that allow the network to retain information over longer periods.
  - They are particularly effective for emotion recognition, where the context of previous speech frames can influence emotional understanding.
  ○ Convolutional Neural Networks (CNNs):
  - Although traditionally used for image processing, CNNs have been applied to audio analysis by treating spectrograms as 2D images.
  - CNNs are effective in capturing spatial patterns in frequency and time, allowing them to identify emotion-specific features from spectrograms.
  ○ Convolutional Recurrent Neural Networks (CRNNs):
  - Combines the spatial feature extraction capabilities of CNNs with the sequential processing power of RNNs or LSTMs.
  - Useful for capturing both local patterns in the spectrogram (using CNN layers) and temporal dependencies (using RNN/LSTM layers).

## 1.2.4. DATASETS FOR TRAINING MODELS

Emotion recognition models require large, labeled datasets to effectively learn and generalize emotional patterns in speech. These datasets are essential for training, validating, and testing machine learning and deep learning models to ensure they can accurately classify emotions across various contexts and speaker attributes.

1. SAVEE (Surrey Audio-Visual Expressed Emotion Database)

- Overview: The SAVEE dataset contains audio clips of English male speakers expressing a range of emotions.
- Emotional Categories: It covers seven key emotions — anger, disgust, fear, happiness, sadness, surprise, and neutral.
- Characteristics:
  - Includes both acted and natural speech samples to provide diversity in emotional expression.
  - The dataset consists of recordings with controlled acoustic conditions, minimizing background noise and ensuring high audio quality.
- Applications: Frequently used for benchmarking models in emotion recognition and speaker-independent classification tasks due to its clear labeling and high-quality recordings.

2. RAVDESS (Ryerson Audio-Visual Database of Emotional Speech and Song)

- Overview: The RAVDESS dataset is a large collection of audio-visual recordings where actors perform scripted lines with various emotional expressions.
- Emotional Categories: It includes expressions of happiness, sadness, anger, fear, surprise, disgust, calm, and neutral states.
- Characteristics:
  - Contains both speech and song recordings, allowing researchers to analyze how emotions are conveyed differently in speaking versus singing.
  - The dataset includes 7356 audio-visual files, covering multiple accents and voice pitches.
  - High-quality recordings are available in both audio-only and audio-visual formats, making it suitable for multimodal emotion recognition research.
- Applications: Widely used for developing and validating models that leverage both audio and visual inputs for emotion classification.

3. TESS (Toronto Emotional Speech Set)

- Overview: The TESS dataset focuses on audio recordings of female English speakers expressing various emotions.
- Emotional Categories: Includes recordings of speakers expressing happiness, sadness, anger, fear, disgust, surprise, and neutral tones.
- Characteristics:
  - Consists of 200 target words spoken by two actresses, each in seven emotional states.

- ○ The dataset was designed to examine the perception of emotion in speech, particularly for older versus younger listeners.
  - ○ Recordings are made with consistent pronunciation and tone, which helps reduce variability due to speaker accents.
- Applications: Often used for analyzing gender-specific emotional expressions and for models focused on female voice data.

4. CREMA-D (Crowd-sourced Emotional Multimodal Actors Dataset)

- Overview: CREMA-D is a diverse and extensive dataset with audio and visual recordings of actors expressing a wide range of emotions.
- Emotional Categories: Includes anger, disgust, fear, happiness, sadness, and neutral states.
- Characteristics:
  - ○ Contains recordings from 91 actors (48 male, 43 female) of various ages and ethnic backgrounds, providing a diverse set of voices.
  - ○ Includes 7442 clips with different sentences spoken in varied emotional tones, allowing for robust testing of models in recognizing emotions from different demographics.
  - ○ Emphasizes multimodal data, capturing both facial expressions and voice, making it ideal for research in multimodal emotion recognition.
- Applications: Valuable for creating models that generalize well across different genders, ages, and ethnic backgrounds. It also supports experiments in combining audio and visual cues for improved accuracy.

Why These Datasets Matter?

1. Diversity of Emotions: The variety in these datasets ensures that models can learn to recognize subtle differences between similar emotions, such as distinguishing between anger and disgust.
2. Speaker Variability: By incorporating recordings from actors of different genders, ages, and accents, these datasets help models generalize across various speakers, making them more robust in real-world applications.
3. Multimodal Inputs: Datasets like RAVDESS and CREMA-D, which include both audio and visual data, are essential for developing systems that leverage multimodal learning to enhance emotion detection accuracy.
4. Noise Resilience: Some datasets provide audio captured in controlled environments, while others include background noise, enabling models to learn to filter out irrelevant sounds and focus on the emotional content.

## 1.2.5. CHALLENGES IN EMOTIONAL ANALYSIS

Despite significant advancements in Speech Emotion Recognition (SER) systems, accurately identifying emotions from audio remains a complex task. The nature of human speech, combined with diverse real-world conditions, introduces several challenges that need to be addressed for effective emotion detection.

1. Speaker Variability

- Overview: Different individuals express the same emotion in unique ways. Factors like voice pitch, speaking style, accent, and natural vocal inflections vary significantly from person to person.
- Challenges:
    - Emotions such as anger or happiness can manifest differently in different speakers. For instance, one person's anger might come through as loud and forceful speech, while another might express it with a quieter, sharper tone.
    - Emotion recognition models often struggle to generalize across speakers, particularly if they are trained on a limited range of voices. This can lead to reduced accuracy in recognizing emotions in new, unseen speakers.
- Mitigation Strategies:
    - Using speaker normalization techniques and including a diverse set of voices in training datasets can help models better generalize to new speakers.
    - Domain adaptation techniques and transfer learning can improve model robustness across different speakers.

2. Cultural and Linguistic Differences

- Overview: Emotions are often expressed differently across cultures and languages. For instance, what constitutes a joyful tone in one culture might be perceived as overly enthusiastic or even aggressive in another.
- Challenges:
    - Models trained on datasets from a particular cultural or linguistic background may not accurately recognize emotions from speakers of other backgrounds. This is due to differences in intonation patterns, emotional expression norms, and speech dynamics.
    - Linguistic nuances, such as the use of specific idioms, word choices, or speech patterns, can impact the detection of emotions, especially in multilingual settings.
- Mitigation Strategies:
    - Incorporating multilingual and multicultural datasets during training helps models become more versatile.
    - Developing language-agnostic features, like Mel-Frequency Cepstral Coefficients (MFCCs), can improve the model's ability to detect emotions across different languages.

3. Environmental Noise

- Overview: In real-world applications, audio recordings often contain background noise, such as traffic, conversations, or electronic hums, which can interfere with the recognition of emotional content.
- Challenges:
  - Noise can obscure key audio features, such as pitch and tone, that are crucial for detecting emotions.
  - Poor audio quality or reverberations in the recording environment can distort the speech signal, making it difficult for models to extract meaningful features.
- Mitigation Strategies:
  - Applying noise reduction and audio enhancement techniques during preprocessing can improve the clarity of speech signals.
  - Data augmentation methods, like noise injection, can help models become more robust by training them to recognize emotions even in noisy conditions.

4. Detection of Subtle or Complex Emotions

- Overview: Emotions are not always straightforward; they can be subtle, ambiguous, or involve a mix of feelings (e.g., ambivalence or nostalgia). Detecting these complex emotional states is particularly challenging for traditional models.
- Challenges:
  - Subtle emotions might not have distinct audio features, making it hard for models to differentiate them from neutral speech.
  - Emotions like sarcasm or mixed feelings may involve nuanced changes in pitch, tone, and speech rate that are difficult to capture using standard feature extraction techniques.
- Mitigation Strategies:
  - Leveraging deep learning models, particularly those with attention mechanisms, can help focus on the most relevant parts of the audio signal that indicate subtle emotional cues.
  - Incorporating contextual information, such as the surrounding dialogue or speaker's history, can enhance the detection of complex emotional states.

5. Temporal Dynamics and Sequence Length

- Overview: Emotions in speech are dynamic and unfold over time. Capturing these temporal changes is crucial for accurate emotion recognition, especially in longer audio sequences.
- Challenges:
  - Models that do not consider the temporal dynamics of speech may miss out on important changes in tone, pitch, or energy that occur gradually.
  - Determining the optimal window size for feature extraction and choosing the right sequence length for model input is a delicate balance that can significantly impact performance.
- Mitigation Strategies:
  - Utilizing Recurrent Neural Networks (RNNs), particularly Long Short-

Term Memory (LSTM) networks, can help capture long-term dependencies in speech patterns.
  ○ Techniques like sliding windows and attention mechanisms can help models focus on the most emotionally relevant portions of a longer speech sequence.

## 1.2.6. APPLICATIONS OF EMOTIONAL AUDIO ANALYSIS

Emotional analysis of audio has many practical applications across industries, including:

- Customer Service and Call Centers: Emotion analysis can be used to detect frustration or satisfaction in customer calls, allowing for better customer service.
- Human-Computer Interaction: Emotional feedback can help virtual assistants or robots respond more empathetic to users.
- Healthcare: Used in monitoring mental health conditions or assessing the emotional state of patients, especially in therapeutic or counseling settings.
- Entertainment Industry: Emotion detection in voice overs or performances can help optimize content for emotional impact.
- Security and Surveillance: Audio analysis can be used in identifying distress or aggression in surveillance settings.
- Education: Emotion recognition can enhance personalized learning experiences by adjusting content based on a student's emotional state.

## 1.2.7. RECENT ADVANCEMENTS

- Multimodal Emotion Recognition: Recent research integrates audio with other modalities, such as facial expressions and body language, for more accurate emotion recognition.
- Transfer Learning: This allows models trained on large datasets to be adapted for smaller, domain-specific datasets, improving the performance of emotion recognition systems.
- Real-Time Emotion Recognition: Advances in computational power and optimization techniques allow for real-time emotion analysis in various applications.

## 1.2.8. ETHICAL CONSIDERATIONS

- Privacy Concerns: Emotion recognition often requires processing sensitive audio data, raising privacy and data protection concerns.
- Bias in Models: Emotion detection models can inherit biases from training data, leading to inaccurate predictions, especially for marginalized or underrepresented groups.
- Manipulation Risks: Emotion recognition could be misused for manipulative purposes, such as targeting vulnerable individuals with emotional ads or misleading content.

## 1.3 CHALLENGES PRESENT

### 1.3.1. CONTEXT DEPENDENCY

- Impact of Context: The meaning of an emotional expression can change depending on the context. For example, the same tone of voice could indicate sarcasm in one context but genuine enthusiasm in another.
- Ambiguity in Emotional Expression: Emotions like confusion, sarcasm, or ambivalence are often subtle and context-dependent, making them hard to categorize accurately.

### 1.3.2. DATA SCARCITY AND IMBALANCE

- Limited Availability of Labeled Data: High-quality labeled datasets for emotion detection are scarce, especially for less common or nuanced emotions.
- Imbalanced Datasets: Emotions like happiness, sadness, or anger are more frequently labeled in datasets, while others like surprise or disgust are less common, leading to biases in model training.

### 1.3.3. TEMPORAL DYNAMICS OF EMOTIONS

- Continuous Emotion Fluctuations: Emotions are dynamic and change over time, even within the same audio clip. Capturing these transitions accurately requires complex models that account for temporal dependencies.
- Short Audio Segments: Many emotion recognition systems operate on short audio frames, making it difficult to capture long-term emotional trends that might be evident in longer conversations.

### 1.3.4. REAL-TIME PROCESSING

- High Computational Requirements: Real-time emotion recognition demands significant computational resources for processing audio data, especially when using deep learning models.
- Latency Issues: Systems need to process audio quickly enough to respond in real time, which is crucial for applications like virtual assistants or customer service.

### 1.3.5. ETHICAL AND PRIVACY CONCERNS

- Privacy Risks: Analyzing emotions from voice data can be intrusive, especially if done without consent. There are potential privacy issues related to recording, storing, and processing sensitive audio data.
- Bias and Fairness: Models trained on biased datasets may not perform equally well across different demographics, potentially leading to discriminatory outcomes.
- Potential for Misuse: Emotion analysis can be used for manipulative purposes, such as targeted marketing or surveillance, raising ethical concerns.

### 1.3.6. SUBTLE AND COMPOUND EMOTIONS

- Detection of Subtle Emotions: Emotions like embarrassment, disappointment, or contentment are often expressed subtly and are harder to detect than primary emotions like anger or joy.
- Compound Emotions: Humans often experience complex emotions simultaneously (e.g., feeling happy yet nervous), which is difficult for models to classify accurately.

### 1.3.7. MODEL INTERPRETABILITY

- Black Box Models: Deep learning models, such as neural networks, are often "black boxes" with little transparency in how they make decisions, making it difficult to understand why a specific emotion was detected.
- Trust in AI Systems: Lack of interpretability can reduce trust in AI-based emotion recognition, especially in sensitive applications like mental health monitoring.

## 1.4 PROJECT STATEMENT

The project aims to develop a system that analyzes audio signals to accurately detect and classify human emotions using machine learning techniques. By extracting key audio features such as pitch, energy, and Mel-frequency cepstral coefficients (MFCCs), the system will identify emotions like happiness, sadness, anger, and fear from speech data. The solution will leverage deep learning models to achieve high accuracy in diverse environments, overcoming challenges like speaker variability, background noise, and cultural differences.

This project has practical applications in enhancing human-computer interaction, improving customer service experiences, and supporting mental health monitoring. The ultimate objective is to create a robust, efficient, and privacy-aware emotion recognition tool that can be seamlessly integrated into customer support systems, virtual assistants, and healthcare applications, making interactions more empathetic and context-aware.

## 1.5 OBJECTIVES

### 1.5.1. EXTRACTING MEANINGFUL FEATURES FROM AUDIO INPUTS

To accurately classify emotions from audio, it's crucial to extract meaningful features that capture the nuances of speech patterns related to different emotions. The raw audio signals contain a mix of frequencies, amplitudes, and patterns, but they need to be processed to extract the specific characteristics that can indicate emotional states. The primary techniques used in this process include MFCCs (Mel Frequency Cepstral Coefficients) and Mel spectrograms.

Mel Frequency Cepstral Coefficients (MFCCs)

- What it is: MFCCs are coefficients that represent the short-term power spectrum of a sound signal. They are widely used in speech and audio processing because they mimic the human ear's sensitivity to different frequency bands.
- Why it's important for emotion recognition:
  - Captures speech patterns: Emotions can affect the shape of the vocal tract (throat, mouth, tongue), which changes the audio signal. MFCCs are effective in capturing these subtle changes.
  - Focuses on perceptually relevant features: Unlike traditional frequency analysis, MFCCs use the Mel scale, which approximates how humans perceive pitch. This helps in capturing emotional tones more effectively.
- How it's used: MFCCs are computed by:
  - Taking the Fourier Transform of short segments of the audio (to analyze frequencies).
  - Mapping these frequencies onto the Mel scale.
  - Applying a logarithm to the Mel frequencies (to emphasize quieter sounds).
  - Applying a Discrete Cosine Transform (DCT) to reduce the dimensionality and focus on the most relevant coefficients.

Mel Spectrograms

- What it is: A Mel spectrogram is a visual representation of the audio signal's frequency content over time, mapped onto the Mel scale. It shows how the energy of different frequencies varies over time.
- Why it's important for emotion recognition:
  - Time-frequency analysis: Emotions affect how people speak (e.g., angry speech may have higher intensity, while sad speech may have lower pitch). Mel spectrograms capture these changes across time.
  - Useful for deep learning models: Convolutional Neural Networks (CNNs) are effective at recognizing patterns in images, and Mel spectrograms can be treated as images, allowing CNNs to extract spatial features related to emotions.
- How it's used: To generate a Mel spectrogram:
  - The audio signal is split into overlapping windows.
  - A Short-Time Fourier Transform (STFT) is applied to analyze the frequency content within each window.
  - The result is then converted to the Mel scale to emphasize the perceptually important frequencies.

## 1.5.2. DESIGNING AND TRAINING A DEEP LEARNING MODEL

Once the features are extracted, the next step is to design a model that can learn to classify emotions based on these features. Given the complexity of speech data and the variability of how different people express emotions, deep learning models are particularly effective.

Model Architecture

The model used for this task often involves a combination of Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), specifically Long Short-Term Memory (LSTM) networks.

- CNNs (Convolutional Neural Networks):
    - What it does: CNNs are used to analyze spatial patterns in data. In the context of speech emotion recognition, CNNs can be used to learn patterns from Mel spectrograms.
    - How it helps: CNN layers can identify patterns like energy bursts (associated with excitement or anger) or smooth, flat patterns (associated with calm or sadness).
- LSTMs (Long Short-Term Memory Networks):
    - What it does: LSTMs are a type of RNN that is effective at learning sequential data. They are designed to remember long-term dependencies, making them suitable for understanding the flow of speech over time.
    - How it helps: Emotions are often expressed through the rhythm, pitch variation, and intonation of speech. LSTMs can capture these temporal patterns, which are crucial for distinguishing emotions.

Training the Model

- The model is trained using labeled audio data from datasets like RAVDESS, CREMA-D, SAVEE, and TESS, where each audio file is tagged with an emotion (e.g., happy, sad, angry).
- Preprocessing is performed to ensure consistency (resampling, noise reduction) and augmentation techniques (like pitch shift, time stretch) are applied to increase dataset variability.
- The extracted features (MFCCs, Mel spectrograms) are used as inputs to the model.
- The model is optimized using techniques like cross-entropy loss for classification and Adam optimizer for faster convergence.

## 1.5.3. CLASSIFYING AUDIO DATA INTO DISTINCT EMOTIONAL CATEGORIES

After the model is trained, it can be used to classify new, unseen audio data into predefined emotional categories. Here's how the classification process works:

Emotion Categories

The emotions that the system is trained to recognize include:

- Happy: Characterized by higher pitch, faster speech rate, and varying energy levels.
- Sadness: Often associated with lower pitch, slower speech rate, and reduced energy.
- Angry: Recognized by a tense, loud, and sharp tone with rapid changes in pitch and energy.
- Fear: Exhibits high pitch and fluctuating energy, with irregular pauses.
- Neutral: Balanced pitch, consistent speech rate, and moderate energy levels.
- Disgust: Typically involves nasal tones, a lower pitch, and less variability in energy.

How the Model Classifies Emotions?

1. Speaker Variability

- Overview: Different individuals express the same emotion in unique ways. Factors like voice pitch, speaking style, accent, and natural vocal inflections vary significantly from person to person.
- Challenges:
  - Emotions such as anger or happiness can manifest differently in different speakers. For instance, one person's anger might come through as loud and forceful speech, while another might express it with a quieter, sharper tone.
  - Emotion recognition models often struggle to generalize across speakers, particularly if they are trained on a limited range of voices. This can lead to reduced accuracy in recognizing emotions in new, unseen speakers.
- Mitigation Strategies:
  - Using speaker normalization techniques and including a diverse set of voices in training datasets can help models better generalize to new speakers.
  - Domain adaptation techniques and transfer learning can improve model robustness across different speakers.

2. Cultural and Linguistic Differences

- Overview: Emotions are often expressed differently across cultures and languages. For instance, what constitutes a joyful tone in one culture might be perceived as overly enthusiastic or even aggressive in another.
- Challenges:
  - Models trained on datasets from a particular cultural or linguistic background may not accurately recognize emotions from speakers of other backgrounds. This is due to differences in intonation patterns, emotional expression norms, and speech dynamics.

- Linguistic nuances, such as the use of specific idioms, word choices, or speech patterns, can impact the detection of emotions, especially in multilingual settings.
- Mitigation Strategies:
    - Incorporating multilingual and multicultural datasets during training helps models become more versatile.
    - Developing language-agnostic features, like Mel-Frequency Cepstral Coefficients (MFCCs), can improve the model's ability to detect emotions across different languages.

3. Environmental Noise

- Overview: In real-world applications, audio recordings often contain background noise, such as traffic, conversations, or electronic hums, which can interfere with the recognition of emotional content.
- Challenges:
    - Noise can obscure key audio features, such as pitch and tone, that are crucial for detecting emotions.
    - Poor audio quality or reverberations in the recording environment can distort the speech signal, making it difficult for models to extract meaningful features.
- Mitigation Strategies:
    - Applying noise reduction and audio enhancement techniques during preprocessing can improve the clarity of speech signals.
    - Data augmentation methods, like noise injection, can help models become more robust by training them to recognize emotions even in noisy conditions.

4. Detection of Subtle or Complex Emotions

- Overview: Emotions are not always straightforward; they can be subtle, ambiguous, or involve a mix of feelings (e.g., ambivalence or nostalgia). Detecting these complex emotional states is particularly challenging for traditional models.
- Challenges:
    - Subtle emotions might not have distinct audio features, making it hard for models to differentiate them from neutral speech.
    - Emotions like sarcasm or mixed feelings may involve nuanced changes in pitch, tone, and speech rate that are difficult to capture using standard feature extraction techniques.
- Mitigation Strategies:
    - Leveraging deep learning models, particularly those with attention mechanisms, can help focus on the most relevant parts of the audio signal that indicate subtle emotional cues.

- Incorporating contextual information, such as the surrounding dialogue or speaker's history, can enhance the detection of complex emotional states.

5. Temporal Dynamics and Sequence Length

- Overview: Emotions in speech are dynamic and unfold over time. Capturing these temporal changes is crucial for accurate emotion recognition, especially in longer audio sequences.
- Challenges:
  - Models that do not consider the temporal dynamics of speech may miss out on important changes in tone, pitch, or energy that occur gradually.
  - Determining the optimal window size for feature extraction and choosing the right sequence length for model input is a delicate balance that can significantly impact performance.
- Mitigation Strategies:
  - Utilizing Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, can help capture long-term dependencies in speech patterns.
  - Techniques like sliding windows and attention mechanisms can help models focus on the most emotionally relevant portions of a longer speech sequence.

## 1.6 SCOPE OF THE PROJECT

### 1.6.1. EMOTION RECOGNITION

- Develop a system that can analyze audio inputs to detect emotions like happiness, sadness, anger, fear, and neutrality.

### 1.6.2. FEATURE EXTRACTION AND ANALYSIS

- Utilize techniques for extracting key audio features, including:

- ○ Pitch (Fundamental Frequency): For distinguishing between excited vs. calm emotional states.
  - ○ Mel-Frequency Cepstral Coefficients (MFCCs): For capturing speech timbre, crucial for identifying nuanced emotions.
  - ○ Speech Rate and Rhythm: To detect emotions like anxiety or relaxation based on the speed of speech.
- Focus on extracting features that enhance emotion classification accuracy, even in noisy or diverse environments.

## 1.6.3. MACHINE LEARNING MODEL DEVELOPMENT

- Build and train deep learning models (e.g., LSTM, CNN, or CNN-LSTM) that can effectively classify emotional states from the extracted features.
- Optimize the models for accuracy and efficiency, ensuring that they perform well on both training and real-world datasets.

## 1.6.4. DATA COLLECTION AND PREPROCESSING

- Leverage publicly available datasets (such as RAVDESS, CREMA-D, TESS) and supplement with custom datasets if necessary.
- Preprocess audio data to remove background noise, normalize volume levels, and enhance voice clarity, ensuring that the models can handle real-world audio inputs.

## 1.6.5. SYSTEM INTEGRATION AND DEPLOYMENT

- Develop a user-friendly interface that allows users to interact with the system and view detected emotions in real-time.
- Enable integration with various platforms and applications, such as:
  - ○ Customer Support Systems: Automatically gauge customer emotions during calls to optimize agent responses.
  - ○ Virtual Assistants: Enable more natural, empathetic responses based on user emotions.
  - ○ Healthcare and Mental Health Monitoring: Use emotional analysis for patient monitoring, therapy support, or well-being tracking.
  - ○ Entertainment and Gaming: Adapt game dynamics or media content in response to users' emotional states for immersive experiences.

## 1.6.6. PRIVACY AND ETHICAL CONSIDERATIONS

- Ensure that all audio data is handled securely, with a focus on user consent and data anonymization.
- Implement on-device processing where possible to avoid transmitting sensitive audio data over networks.
- Comply with data privacy regulations (e.g., GDPR) to build trust with users.

## 1.6.7. TESTING, VALIDATION, AND OPTIMIZATION

- Conduct extensive testing of the system on different datasets to ensure high accuracy and robustness.
- Optimize the system to handle speaker variability, accents, and noisy environments effectively.
- Continuously update the model to improve performance based on feedback and new data.

## 1.6.8. SCALABILITY AND FUTURE EXTENSIONS

- Design the system to be scalable, allowing it to handle an increasing volume of audio inputs and adapt to new use cases.
- Explore potential future enhancements, such as:
    - Extending the emotion categories to include compound emotions (e.g., bittersweet, nostalgic).
    - Applying the system to additional languages and dialects.
    - Using multimodal inputs (e.g., combining audio with video analysis) for more accurate emotion detection.

# Chapter 2

# BACKGROUND

## 2.1 INTRODUCTION

In today's digital era, technology plays a central role in how we communicate, interact, and connect with others. With the rapid growth of AI-powered systems like virtual assistants (e.g., Alexa, Google Assistant, Siri) and customer service chatbots, there is an increasing need for these systems to become more intuitive, empathetic, and responsive to human emotions. However, most current technologies lack the ability to effectively understand or react to the emotional states of users, leading to interactions that often feel robotic and impersonal.

The field of Emotional Analysis of Audio—also known as Speech Emotion Recognition (SER)—emerges as a promising solution to bridge this gap. Emotions are deeply embedded in the tone, pitch, speed, and energy of our speech. By analyzing these subtle vocal cues, computers can gain insights into the speaker's emotional state, potentially transforming the way we interact with machines. For instance, a virtual assistant that can detect frustration in a user's voice could adjust its responses to be more helpful and supportive. Similarly, a call center application could detect a customer's emotional distress and route the call to a human agent for a more personalized response.

The importance of emotion recognition from audio goes beyond enhancing user experience in customer service or virtual assistants. It has significant applications in healthcare, where monitoring a patient's emotional state can provide early indicators of mental health issues like depression or anxiety. In the entertainment industry, adaptive systems that respond to users' emotions can create more immersive gaming or multimedia experiences.

Despite its potential, emotional analysis of audio is a challenging task due to factors like speaker variability, environmental noise, and cultural differences in emotional expression. For a machine learning model to reliably detect emotions in real-world scenarios, it needs to be trained on diverse datasets and equipped with robust algorithms capable of handling various accents, background noise, and changing emotional states over time.

Recent advances in machine learning and deep learning have made it possible to extract meaningful insights from complex audio data. By leveraging techniques like Long Short-Term Memory (LSTM) networks, Convolutional Neural Networks (CNNs), and advanced audio feature extraction (e.g., Mel-frequency cepstral coefficients, pitch analysis), researchers are pushing the boundaries of what is achievable in emotion recognition. However, the development of systems that are both accurate and capable of operating in real-time remains a work in progress.

This project aims to harness the power of these technologies to create a system that can

analyze speech in real-time, detect emotions accurately, and deliver practical applications in various fields. By addressing the challenges of speaker variability, noisy environments, and cultural diversity, this project will contribute to making human-computer interactions more empathetic, natural, and effective.

## 2.2 LITERATURE REVIEW: EMOTIONAL ANALYSIS OF AUDIO USING DEEP LEARNING TECHNIQUES

The field of Speech Emotion Recognition (SER) has gained significant momentum in recent years, driven by the need for more intuitive and empathetic AI systems. As human-computer interaction evolves, integrating emotion recognition capabilities can transform applications ranging from virtual assistants to customer service and healthcare. This literature review delves into the various methodologies, datasets, and technologies that have been explored to develop robust emotion recognition systems, focusing on audio analysis using machine learning and deep learning techniques [1],[2],[4],[15].

Emotion recognition models rely on diverse and comprehensive datasets to achieve high accuracy, with many studies utilizing well-established resources. Some of the most commonly used datasets include the IEMOCAP (Interactive Emotional Dyadic Motion Capture Database), known as a benchmark in speech emotion recognition (SER) for its inclusion of both scripted and spontaneous emotional speech. The RAVDESS (Ryerson Audio-Visual Database of Emotional Speech and Song) is another popular dataset, renowned for its high-quality recordings that include both audio and visual data, though most studies focus on the audio component for emotion detection. TESS (Toronto Emotional Speech Set) consists of emotional speech recordings designed to test models on discrete emotions like happiness, sadness, and anger. CREMA-D (Crowd-sourced Emotional Multimodal Actors Dataset) provides a diverse range of expressions from actors of various ethnicities, aiding in generalization across different demographics, while the SAVEE (Surrey Audio-Visual Expressed Emotion Database) is focused on emotion recognition in male voices. These datasets provide crucial audio samples, which help train models to handle variability in speakers, accents, and emotional expressions.[4]

Feature extraction plays a critical role in emotion detection from audio. Spectrograms, particularly Mel-spectrograms, have been widely utilized to capture the frequency content and temporal dynamics of speech. Mel-spectrograms are particularly effective for identifying the intensity and frequencies of speech formants. The Short-Time Fourier Transform (STFT), often implemented with libraries like Librosa, is used to convert audio signals into spectrograms. Another popular feature extraction technique is Mel-Frequency Cepstral Coefficients (MFCCs), which capture spectral properties of speech and represent the vocal tract shape. Based on the Mel scale, MFCCs are effective at filtering out irrelevant background noise, making them valuable for emotion recognition in noisy environments. Additionally, prosodic features such as pitch, energy, and speech rate add further emotional context, distinguishing between emotions like happiness, characterized by higher pitch and energy, and sadness, which typically features lower pitch and slower speech [1].

Traditional machine learning models, such as Support Vector Machines (SVM), Decision Trees, Random Forest, XGBoost, K-Nearest Neighbors (KNN), and Multi-Layer Perceptrons (MLP), have been applied to emotion recognition using extracted audio features. While these models have shown success, they struggle with the complexity of high-dimensional audio data, particularly in the presence of noise and speaker variability. This has led to a shift toward deep learning approaches. Convolutional Neural Networks (CNNs) have been particularly effective for analyzing spectrograms due to their ability to capture spatial patterns in the frequency domain, enabling them to identify emotion-specific features. Recurrent Neural Networks (RNNs), especially Long Short-Term Memory (LSTM) networks, are also widely used for processing sequential data, such as speech, and address the vanishing gradient problem by retaining information over longer periods. Fully Convolutional Networks (FCNs) have been explored to handle variable-length inputs, making them adaptable for real-time applications. Autoencoders, particularly Convolutional Autoencoders (CAEs), have been utilized for feature extraction and dimensionality reduction, helping to filter noise and improve model robustnes [1],[2],[15].

To optimize deep learning models for emotion recognition, various strategies are employed, including regularization techniques like Dropout to prevent overfitting, pooling layers such as Global Average Pooling to reduce the dimensionality of feature maps while retaining essential information, and hyperparameter tuning to optimize factors like learning rates, number of layers, and activation functions. The effectiveness of these models is typically assessed using standard evaluation metrics such as accuracy, precision, recall, and F1-score, which measure the model's ability to classify emotions correctly across different datasets and scenarios. Audio analysis frameworks such as Marsyas, PsySound, and Essentia are also commonly used for extracting low-level and high-level audio features, like MFCCs, spectral contrast, and zero-crossing rates, to build robust emotion recognition systems [2],[4].

The core of the proposed system is a novel CNN-RNN Deep Model (ABCDM) that leverages the strengths of bidirectional Long Short-Term Memory (LSTM) networks and Gated Recurrent Units (GRU). By utilizing both LSTM and GRU layers, the model effectively captures long-term dependencies and contextual information in the audio signal, taking into account both past and future frames. This dual-layered approach enables the model to accurately discern subtle emotional cues, even in complex audio sequences, without relying on text transcriptions or traditional Natural Language Processing (NLP) techniques. This language-agnostic approach also allows the model to detect intricate expressions like sarcasm and irony, which are often missed by text-based emotion analysis[11].

For feature extraction, the study emphasizes the use of Mel-Frequency Cepstral Coefficients (MFCCs) due to their proven ability to capture the fundamental characteristics of human speech, including pitch, tone, and intensity. In addition to MFCCs, advanced feature selection methods such as Linear Discriminant Analysis (LDA) and Principal Component Analysis (PCA) are employed to reduce the dimensionality of the extracted features. This reduction not only streamlines the computational requirements but also enhances classification accuracy by focusing on the most relevant features, thereby minimizing noise and redundancy [4].

27

To further improve the robustness of the model, innovative data augmentation techniques, such as Segment Repetition based on High Amplitude (SRHA), are applied. This augmentation method addresses the challenge of limited labeled audio data by synthetically increasing the diversity of training samples, thereby enhancing the model's ability to generalize across different datasets and real-world scenarios. The proposed Fully Convolutional Neural Network (FCN) architecture is designed to handle audio files of varying lengths, allowing for flexible input sizes and enabling near real-time emotion detection. This flexibility is particularly beneficial for applications that require immediate responses, such as customer support chatbots and interactive virtual assistants[6].

The model was rigorously evaluated using multiple benchmark datasets, including CASIA, SAVEE, and IEMOCAP, which contain a diverse range of emotional expressions across different speakers and languages. Experimental results demonstrate that the proposed approach significantly outperforms traditional models, achieving high accuracy rates in both speaker-dependent and speaker-independent settings. The integration of bidirectional LSTM and GRU layers, combined with efficient feature selection and data augmentation techniques, resulted in improved recognition rates, particularly in noisy environments where conventional methods often struggle [4],[7].

The proposed methodology converts audio signals into spectrograms using the Short-Time Fourier Transform (STFT) to generate a two-dimensional representation of the audio frequencies over time. These spectrograms provide a rich source of features for sentiment analysis. Keypoints within these spectrograms are identified using the Oriented FAST and Rotated BRIEF (ORB) algorithm, which effectively extracts salient features that capture variations in the audio signal. These extracted keypoints are then processed using a Bag-of-Visual-Words (BoVW) technique, which clusters the keypoints into visual words and converts the information into histograms. These histograms serve as input features for classification models, enabling the system to recognize and classify the sentiment expressed in the audio[5].

To optimize the system's performance, extensive hyperparameter tuning was conducted, including adjustments to spectrogram resolution, sample rates, and the number of clusters used in the BoVW technique. This process was critical in enhancing the accuracy and efficiency of the model. The system was trained and evaluated on a multilingual dataset comprising audio samples in English, Italian, and Spanish, demonstrating the model's capability to generalize across languages without requiring language-specific training data [5].

The experimental results highlight the effectiveness of the approach, with the Random Forest classifier emerging as the most robust model, achieving an accuracy of 76% and an F1 score of 78%. Comparative analysis with other classifiers, such as Support Vector Machines (SVM) and k-Nearest Neighbors (kNN), showed that the proposed method significantly outperformed traditional models in both accuracy and computational efficiency. The findings underscore the potential of leveraging spectrogram-based visual representations and machine learning techniques for scalable, real-time, and language-agnostic sentiment analysis[1],[15].

The ability to detect gender, age, and emotion from speech is becoming increasingly important in fields such as human-machine interaction, telecommunication, and customer analytics. In telecommunication, these capabilities can be used to predict customer demographics and recommend personalized offers, thereby enhancing customer engagement and service delivery. Despite significant research in detecting gender, age, and emotion from speech, most existing studies tend to focus on one attribute at a time or use different approaches and datasets for each attribute. This paper proposes a unified system that predicts gender, age, and emotion from speech, using a single source of audio data and a common approach for feature extraction[13].

In our approach, audio files are first processed through frequency spectrum analysis to extract 20 statistical features, which include parameters such as mean frequency, standard deviation, skewness, kurtosis, and various quartiles. These features capture key characteristics of the audio signal that can be used for predictive modeling. The datasets derived from these features are then used to train and evaluate a variety of machine learning models. The models applied include Random Forest, CatBoost, Gradient Boosting, K-nearest neighbors (KNN), XGBoost, AdaBoost, Decision Tree, Artificial Neural Networks (ANN), Naive Bayes, and Support Vector Machine (SVM)[1],[2].

The study evaluates the performance of these models based on test accuracy, and the results reveal that CatBoost achieves the highest accuracy of 96.4% for gender prediction, making it the most effective model in this regard. For age prediction, Random Forest yields the best performance with an accuracy of 70.4%, while XGBoost[1][2] stands out as the top model for emotion prediction, achieving an accuracy of 66.1%.

Additionally, an analysis of the 20 statistical features used in the models highlights which features are most influential in determining the accuracy of the predictions. This feature analysis provides valuable insights into the characteristics of speech that are most useful for detecting gender, age, and emotion. The findings from this research suggest that a unified approach to predicting these attributes, using a consistent set of features and models, can lead to more efficient and effective systems for demographic and emotional analysis from speech [6].

This work contributes to the field by offering a comprehensive evaluation of multiple machine learning algorithms applied to a unified set of speech data, and by identifying the most influential features for prediction. The insights presented can serve as a foundation for future advancements in speech-based demographic analysis and emotion detection, which have wide-ranging applications in areas such as personalized customer service, virtual assistants, and healthcare [9].

The rapid expansion of the Internet of Things (IoT) and the increasing use of voice-based multimedia applications have led to the generation of vast datasets capturing various facets of human behavior, including emotions. Emotion detection from speech is becoming an essential component for enhancing human-computer interactions, enabling more natural and intuitive communication between humans and artificial intelligence systems. While significant progress has been made in text-based emotion detection, acoustic feature-based emotion detection still faces challenges, particularly in terms of accuracy and robustness.

This paper presents a novel approach to emotion detection from conversational audio, utilizing a Bag-of-Audio-Words (BoAW) feature embedding technique to represent speech data in a way that enhances emotional context. Unlike traditional methods that treat audio features in isolation, this approach integrates the sequential nature of conversations and utilizes the rich, embedded emotional expressions within the speech[5].

We introduce a Recurrent Neural Network (RNN)-based emotion detection model designed to address these challenges. The model is specifically engineered to capture both individual speaker states and the broader context of the ongoing conversation, which is crucial for accurate real-time emotion detection. Our model is evaluated using two benchmark emotion recognition datasets, including the IEMOCAP dataset, and it is tested for its capability to provide real-time emotion predictions. We achieve a weighted accuracy of 60.87% and an unweighted accuracy of 60.97% in detecting six basic emotions—happy, sad, neutral, angry, frustrated, and excited—showing a notable improvement over existing state-of-the-art models [7],[12],[17].

This paper also addresses several key limitations in emotion detection from audio, such as the limited number of emotions detected by existing systems, the relatively low performance of acoustic features compared to textual cues, and the complexities associated with speaker diarization (identifying and separating speakers in a conversation). The proposed method effectively tackles these issues by providing more accurate emotion classification through a comprehensive feature extraction and embedding process, which improves the model's ability to detect nuanced emotional states and understand the contextual flow of the conversation. The results from our approach show significant advancements in both emotion detection accuracy and contextual understanding of emotional nuances in human speech. The proposed model's effectiveness is further validated through its real-time prediction capability, making it highly suitable for integration into real-world applications, such as automated customer service systems, healthcare diagnostics, and emotion-aware assistive technologies. By providing a deeper understanding of emotions in spoken conversations, this research contributes to advancing emotion recognition technologies, paving the way for more empathetic and context-aware AI systems[14].

Sentiment analysis has emerged as a critical tool for addressing challenges in human-machine collaboration, especially in the context of Human-Computer Interaction (HCI). Understanding human emotions is inherently difficult, and automating this process using machines adds another layer of complexity. With the increasing reliance on Artificial Intelligence (AI) systems, there is an urgent need to develop technologies that can autonomously detect and interpret the sentiments of individuals during conversations. This research investigates the application of textual classification methods to audio data for sentiment analysis, offering a novel approach to understanding emotions in spoken language[6].

The proposed system combines speech recognition with natural language processing (NLP) techniques to convert audio data into text transcripts. These transcripts are subsequently processed using advanced text pre-processing steps, including tokenization, feature extraction, and cleaning, to prepare the data for analysis. Afterward, we apply various

supervised machine learning algorithms such as Naïve Bayes and Logistic Regression to classify the extracted text into three primary sentiment categories: positive, neutral, and negative. The system focuses on identifying key linguistic features, such as sentiment-laden words, phrases, and contextual cues, to detect emotional undertones within the conversation[1].

To enhance the robustness and accuracy of sentiment classification, the study also examines the impact of various text-cleaning techniques, such as the removal of stop words, punctuation, and irrelevant data, which could potentially introduce noise into the model. Additionally, we investigate the scalability and efficiency of the system, testing it on diverse audio datasets to ensure its applicability across different domains, including customer service interactions, healthcare, and mental health analysis. The experimental results demonstrate that the proposed approach outperforms traditional sentiment analysis methods by effectively capturing the nuances of human emotions within conversational audio data. The system achieved high accuracy in detecting sentiment from real-time speech-to-text data, showing promise for real-world applications. Furthermore, the research highlights the potential of leveraging machine learning algorithms to enhance real-time sentiment analysis capabilities in HCI, making it an invaluable tool for improving user experience, facilitating emotion-aware AI systems, and fostering more empathetic interactions between machines and humans [6].

This paper presents a novel posed multimodal emotional dataset (PME4) and conducts a comprehensive evaluation of emotion classification based on four distinct modalities: audio, video, electromyography (EMG), and electroencephalography (EEG). The aim of this study is to assess the performance of various emotion recognition techniques utilizing these modalities, alongside multiple feature extraction strategies and machine learning algorithms. The PME4 dataset was created by collecting data from 11 human subjects, each expressing six basic emotions (anger, happiness, sadness, fear, surprise, and disgust) in addition to a neutral emotion. Each subject's emotional expressions were captured through video (facial expressions and speech) and physiological signals (EEG for brain activity and EMG for facial muscle movements)[6].

For each modality, features were extracted using diverse techniques tailored to the unique characteristics of the data: principal component analysis (PCA) for dimensionality reduction, autoencoders for learning compact representations, convolutional networks for image data (video), and mel-frequency cepstral coefficients (MFCC) for audio data. The feature extraction methods were specifically designed to maximize the distinct information captured by each modality, contributing to more robust emotion recognition. A variety of machine learning models were employed to classify the emotions, including traditional algorithms such as k-nearest neighbors (KNN), support vector machines (SVM), random forest, and multilayer perceptron (MLP)[2]. Additionally, more complex models like long short-term memory (LSTM) networks, which are well-suited for sequential data like audio and video, and convolutional neural networks (CNN), designed for image-like data, were tested. The performance of these models was evaluated on each modality, highlighting the effectiveness of combining biosensor signals (EMG and EEG) to reduce noise and improve classification accuracy. Notably, traditional KNN provided the best overall classification performance when considering all modalities, while LSTM models outperformed other

methods for classifying dynamic data from audio and video sequences[6].

The primary contributions of this research are: 1) the introduction of the PME4 multimodal emotion dataset, which includes four modalities—audio, video, EEG, and EMG—capturing both non-physiological and physiological signals; 2) an in-depth comparison of several state-of-the-art feature extraction techniques and machine learning algorithms, showcasing their applicability to emotion recognition; and 3) a detailed analysis of how these techniques perform across different modalities and emotions. By providing this dataset and the accompanying analysis, this work aims to foster further advancements in emotion recognition, specifically in the context of multimodal data fusion, which could lead to more accurate and context-aware emotion classification systems for applications in fields such as human-computer interaction, healthcare, and customer service[6].

This paper presents a novel approach to Speech Emotion Recognition (SER), aimed at improving emotion classification accuracy by analyzing human speech and extracting emotional characteristics based on various acoustic features. The proposed system converts the speech signal into digital input, capturing emotional features such as intensity, pitch, timbre, speech rate, and pauses, which are crucial for accurate emotion recognition. Traditional attention models often operate on a fixed attention granularity, limiting their ability to handle the diverse emotional expressions that can vary in both intensity and granular details. To address this limitation, the proposed method employs a multi-scale area attention mechanism in conjunction with a deep 2D Convolutional Neural Network (CNN) and a Dense Deep Neural Network (DNN). This combined architecture allows the model to capture emotional characteristics at different levels of granularity, thus enabling it to recognize a broader spectrum of emotions with varying intensities. For example, emotions such as joy and annoyance, which can manifest in different levels of intensity, are classified more accurately by this multi-scale approach.[7]

A major challenge in SER systems is the sparsity of emotional data, which can affect the model's generalization ability. To mitigate this, the proposed method incorporates data augmentation strategies, such as pitch modification, time stretching, and noise insertion. These techniques enhance the diversity of the training data, helping the model become more robust to variations in real-world speech, including background noise and varying speech rates. The system is trained and evaluated using established emotion datasets such as RAVDESS, CREMA-D, and TESS-D, which include a wide variety of speech recordings representing different emotional states[2],[7],[8].

Experimental results demonstrate that the integration of the multi-scale attention mechanism significantly improves emotion classification performance, especially in distinguishing emotions that are subtle or have overlapping features. The comparative analysis of different SER models further highlights the advantages of the proposed approach in terms of both accuracy and scalability. The proposed method not only handles diverse emotional expressions but also detects the varying intensities of emotions, making it particularly useful for applications in psychological assessments, emotion-aware robotics, and real-time emotion detection in interactive systems[7].

This paper explores the challenge of multilingual emotion detection from raw speech data,

focusing on English, German, and Italian languages. Emotion detection from speech plays a critical role in understanding human communication across cultural and linguistic barriers. The proposed approach utilizes widely recognized emotional databases, including the Ryerson Audio-Visual Database (RAVDESS), the Berlin Database (EmoDb), and the Italian Emo-Vo Database. These databases contain emotional speech samples in English, German, and Italian, respectively, with a variety of emotions such as happiness, anger, sadness, fear, surprise, and disgust [8].

The proposed model extracts key acoustic features from the raw audio data, such as Mel-frequency cepstral coefficients (MFCC)[2][4], chroma, Tonnetz, and Contrast, which are crucial for capturing the emotional content of speech. These features are then fed into a Convolutional Neural Network (CNN), which has been optimized for emotion classification. Importantly, the system does not rely on any visual data, focusing entirely on sound features, making it suitable for real-time speech emotion recognition in natural, unsupervised environments. The model employs a deep learning approach to automatically learn discriminative patterns from the raw audio, providing a more efficient and robust method compared to traditional machine learning techniques[2].

The proposed model's performance is thoroughly evaluated using the three multilingual emotional datasets. In the initial tests, the model achieves promising accuracy rates: 70.46% for RAVDESS, 70.37% for Emo-Db, and 73.47% for Emo-Vo. To further improve robustness and generalization, the model is tested on augmented databases with various modifications to the original speech, such as pitch changes, noise addition, and speed variations. With this augmentation, the model significantly improves its performance, securing accuracy rates of 96.53% for RAVDESS, 96.22% for Emo-Db, and 96.11% for Emo-Vo[8][9].

Understanding sentiment and emotion in speech is a complex challenge, especially when only audio data is available, as in the case of telephone conversations. In this study, we explore sentiment analysis and emotion recognition from speech using self-supervised learning models, focusing on universal speech representations with speaker-aware pre-training. These models are trained on acoustic features extracted directly from speech, without the need for transcription or additional text-based input, making them ideal for scenarios where only raw audio is available.

We evaluate the models on a range of tasks, including three distinct sentiment analysis tasks (binary, three-class, and seven-class) and one emotion recognition task, utilizing three different model sizes to examine the effect of model complexity. The evaluation results demonstrate that the best performance was achieved for the binary sentiment analysis task, with weighted and unweighted accuracy scores of 81% and 73%, respectively. This unimodal acoustic analysis, which uses only speech data, performed competitively when compared to previous methods that relied on multimodal fusion, such as integrating visual or textual data. However, the models faced challenges in accurately predicting emotions and in sentiment analysis tasks with more than two classes. Specifically, the models struggled with the six-class emotion recognition, three-class sentiment, and seven-class sentiment tasks, where the unbalanced nature of the datasets likely played a role in the performance degradation. In these more complex tasks, the models were less effective at

capturing subtle differences in sentiment and emotion due to the lack of balanced training examples across different classes[10].

Sentiment analysis has gained substantial importance in recent years, particularly with the rise of user-generated content and the increasing reliance of businesses on customer feedback and interactions. As a result, extracting sentiment from unstructured data, such as audio recordings, has become both a challenging and vital task. Traditional sentiment analysis techniques have primarily focused on text-based inputs, while emerging methods are increasingly tackling audio data. This report presents a solution that enhances sentiment analysis by leveraging cloud-based technologies, specifically Amazon Web Services (AWS), to automate the process of converting speech into actionable sentiment insights. The system utilizes multiple AWS services to create an integrated pipeline for processing and analyzing customer conversations[10].

The process begins with Amazon Transcribe, which is used to convert voice recordings into text. The transcribed text is then passed through Amazon Comprehend, a natural language processing (NLP) service that applies machine learning models to identify and classify the sentiment within the text. The sentiment is categorized into four distinct classes: positive, negative, neutral, and mixed, based on the emotional tone and context of the conversation. Additionally, AWS Glue is employed to generate metadata, enriching the processed text with relevant information about the conversation's context. AWS Athena is used to structure the raw text data into a more usable and query-friendly format, making it easier to store, retrieve, and analyze the sentiment results.

A key feature of this system is the use of AWS Lambda, which automates the entire workflow by triggering actions based on specific events, such as new data being uploaded to storage or transcriptions being completed. This eliminates the need for manual intervention and ensures that the sentiment analysis process is continuous and scalable. The processed sentiment data is then visualized using AWS QuickSight, which provides interactive dashboards and visual reports, including bar charts and sentiment trend graphs, allowing business stakeholders to gain quick insights into the emotional tone of customer interactions. The implementation of this cloud-based solution offers several significant advantages over traditional methods. First, it dramatically reduces the time complexity involved in manual sentiment analysis. By automating the entire process—from audio transcription to sentiment classification and visualization—businesses can quickly analyze large volumes of customer interactions. Second, the scalability of AWS services ensures that the system can handle increasing amounts of data as a business grows. Additionally, the integration of cloud technologies provides enhanced security, as data can be encrypted and managed securely within the AWS infrastructure. This solution is also more cost-effective compared to traditional on-premises systems, as it leverages the pay-as-you-go pricing model of AWS, which allows businesses to avoid the costs associated with maintaining physical infrastructure[10].

The proposed solution is not only efficient but also flexible, as it can be adapted for various use cases, including customer support analysis, product feedback, social media sentiment monitoring, and market research. By using this automated sentiment analysis system, businesses can gain a deeper understanding of customer emotions, which in turn can inform

decision-making, improve customer relationships, and drive overall business growth. The ability to continuously monitor and analyze the sentiment of customer conversations ensures that businesses can respond proactively to customer concerns and enhance their products and services based on real-time emotional feedback.

In conclusion, this report demonstrates how integrating AWS cloud-based services can streamline the process of sentiment analysis from unstructured audio data, providing businesses with an automated, scalable, secure, and cost-effective solution for analyzing customer sentiment. This approach not only improves the accuracy and efficiency of sentiment analysis but also empowers businesses to make more informed decisions and improve their customer experience.

Emotion recognition plays a vital role in enhancing human-computer interaction, with significant applications in fields like customer service, healthcare, social robotics, and adaptive systems. Despite advancements in affective computing, accurately identifying emotions remains a challenging task due to the complexity of emotional expressions and the limitations of using single-modality data. This paper introduces an ensemble visual-audio emotion recognition framework that leverages multi-task learning and ensemble blending techniques to optimize emotion classification by integrating diverse feature sets across modalities.

To address the shortcomings of existing approaches, which often rely on isolated feature extraction methods, our framework incorporates both traditional and deep learning features to enhance recognition performance. For the audio modality, we extract Interspeech 2010 hand-crafted features, which are known for capturing emotional nuances, alongside deep spectral features derived from mel-spectrograms using Convolutional Neural Networks (CNNs). For the visual modality, we utilize Local Binary Patterns (LBP) to capture local texture details, combined with deep features obtained from CNNs trained on facial expression images[1]. These distinct feature sets are processed through specialized classifiers: Support Vector Machines (SVM)[1][2][4][5][15] are applied to manual features, while deep features are analyzed using CNNs, resulting in four sub-models tailored to different feature types.

A blending ensemble algorithm is employed to fuse the outputs of the sub-models, thereby capitalizing on the complementary strengths of each feature set to enhance the system's robustness. Additionally, the framework integrates multi-task learning within the CNN models, enabling the simultaneous prediction of both primary emotion recognition and auxiliary tasks, such as gender classification. This approach not only improves the sensitivity of the recognition system to subtle emotional cues but also reduces the model's parameter count by sharing information across related tasks[1].

Extensive experiments were conducted using the eNTERFACE'05 database to validate the effectiveness of the proposed method. The results demonstrated that the multi-task CNN model outperformed standard CNN models, achieving an average improvement of 3% in speaker-independent settings and 2% in speaker-dependent settings. The ensemble framework achieved an overall emotion recognition accuracy of 81.36% in speaker-independent experiments and 78.42% in speaker-dependent experiments, significantly

surpassing several state-of-the-art approaches. These results confirm that integrating visual and auditory data with an ensemble learning strategy can effectively improve emotion recognition accuracy. The proposed framework offers a scalable and efficient solution for applications that require real-time emotion recognition, paving the way for more intuitive and adaptive human-computer interaction systems[1].

Sentiment analysis from audio signals is a challenging field that has garnered significant research interest due to the complexities involved in enabling machines to accurately detect human emotions. While humans can effortlessly recognize emotions through facial expressions, gestures, and tonal variations, replicating this capability in machines requires advanced techniques. This research proposes an innovative audio emotion detection system designed to analyze and benchmark emotion recognition performance across three diverse datasets: Toronto Emotional Speech Set (TESS), Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS), and a custom dataset created from EmoDB and SAVEE datasets. The goal is to improve the accuracy of emotion classification by combining results from multiple machine learning models with optimized hyperparameters[7],[8].

The proposed system extracts critical audio features such as Mel Frequency Cepstral Coefficients (MFCC), Mel Spectrogram, and Chroma features using the Librosa library, making the system adaptable across different languages and accents. Additionally, feature selection techniques, such as Principal Component Analysis (PCA), K-Best features, and correlation analysis, are employed to reduce dimensionality, eliminate redundant features, and minimize overfitting, thereby optimizing the model's training efficiency and performance. The classification is based on six of Ekman's core emotion categories: Anger, Joy, Sadness, Fear, Disgust, and Surprise[2].

Comprehensive experiments were conducted on the datasets using an ensemble approach that combines multiple baseline classifiers, including Support Vector Machines (SVM), Random Forest, and Multilayer Perceptron (MLP). The system achieved impressive results, with an accuracy of 99.46% on the TESS dataset, 89.62% on RAVDESS, and 78.28% on the custom dataset. Notably, Anger was identified as the most predictable emotion, whereas Fear was the most challenging to classify accurately, likely due to its subtle acoustic characteristics[2].

To test the robustness of the system, additional experiments were conducted with noise introduced at various signal-to-noise ratios, demonstrating that the model maintains high accuracy even under noisy conditions. The system's versatility and effectiveness make it suitable for practical applications in sectors such as healthcare, counseling, security, and AI-powered customer service, where real-time emotion recognition can enhance user engagement and improve interaction outcomes. By understanding users' emotional states, AI assistants and automated systems can tailor their responses for more empathetic and effective communication, ultimately leading to better user experiences. The development of robust audio features is an essential step toward advancing the field of Music Emotion Recognition (MER), which has garnered increasing interest within the Music Information Retrieval (MIR) community. This paper presents a comprehensive survey of existing computational audio features that are critical for recognizing emotions in music. Grounded

in music psychology, this work explores the relationships between eight core musical dimensions—melody, harmony, rhythm, dynamics, tone color (timbre), expressivity, texture, and musical form—and their influence on emotional responses. By integrating insights from psychological studies, music theory, and computational analysis, this review highlights the intricate connections between musical elements and specific emotional expressions[3].

The primary focus of this survey is to evaluate how these musical dimensions impact emotional perception in listeners and how computational features can effectively capture these nuances. Current audio features used in MER are reviewed, with a focus on popular open-source frameworks such as Marsyas, MIR Toolbox, PsySound, and Essentia, examining their capabilities in extracting relevant data for emotion classification. The survey reveals that many existing systems rely on general-purpose features that were originally designed for other contexts, such as speech recognition or genre classification, which limits their effectiveness in MER tasks[3].

Despite significant advancements in machine learning, especially with deep learning models, the performance of existing MER systems has plateaued. This is particularly evident in scenarios requiring classification of subtle emotional categories, where even state-of-the-art models struggle to achieve high accuracy. The review covers various feature categories, from low-level features like spectral properties and Mel-Frequency Cepstral Coefficients (MFCCs) to higher-level semantic features such as danceability, modality, and genre. It emphasizes the need for novel feature engineering strategies that can capture the expressivity, texture, and structural elements of music more effectively, thereby bridging the semantic gap in current systems[4].

Ultimately, this survey emphasizes that progress in MER will require focused efforts on designing features that align more closely with the inherent emotional properties of music. Such advancements are crucial not only for breaking the current performance ceiling in MER tasks but also for enabling more effective applications in areas like personalized music recommendations, therapeutic interventions, and emotional AI systems. Future directions include the development of features that capture intricate musical structures, exploration of deep learning architectures tailored for emotion detection, and integration of multimodal approaches that consider lyrics, symbolic data, and user context to enhance emotional classification[4].

Emotion recognition from speech has garnered significant research interest due to its diverse applications in areas such as human-computer interaction, mental health monitoring, and therapeutic interventions for speech impairments. The ability to accurately detect emotions from speech signals is critical for developing systems that can enhance communication, provide emotional support, and improve user experiences in various digital platforms. However, despite considerable advancements, existing models for speech emotion recognition (SER) face challenges in real-time deployment due to high

computational costs, limited generalization across different datasets, and suboptimal feature extraction techniques that affect classification performance[7].

This paper introduces a novel methodology that leverages autoencoders for dimensionality reduction to create a compact yet effective representation of audio features. By compressing high-dimensional input data into a more manageable size without significant information loss, the system improves computational efficiency, making it more suitable for real-time applications. We evaluate the proposed approach using two prominent benchmark datasets: the Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) and the Toronto Emotional Speech Set (TESS). To assess the impact of dimensionality reduction on classification performance, we experiment with a combination of classical machine learning classifiers such as Support Vector Machines (SVM) and Decision Trees[2][4], alongside deep learning architectures like Convolutional Neural Networks (CNN). Additionally, we explore the performance of state-of-the-art deep learning models, namely AlexNet and ResNet50, to benchmark their effectiveness against traditional methods. The results demonstrate that incorporating autoencoders enhances the accuracy of emotion detection, with our CNN model achieving a peak accuracy of 96% on the TESS dataset—outperforming several existing methods in terms of classification accuracy and efficiency[2][4].

Our findings emphasize the importance of using dimensionality reduction techniques like autoencoders to optimize feature extraction and classification, particularly in scenarios where low latency and reduced computational overhead are essential. The proposed approach offers a promising direction for future research, as it addresses key limitations in current SER systems, paving the way for more robust and scalable real-time emotion recognition solutions. By enhancing the compactness of audio representations, this study lays the groundwork for integrating emotion-aware systems into applications like driver monitoring, healthcare diagnostics, virtual assistants, and customer service automation [7].

Audio sentiment analysis, a key area of research in affective computing, has gained considerable traction due to its applications in human-computer interaction, customer service, mental health assessment, and personalized healthcare. While numerous machine learning models have been developed to detect emotions from speech with high accuracy, there has been limited exploration into the fairness of these models, particularly regarding their performance across different demographic groups. Current research in AI fairness has underscored the need to address biases in model performance related to factors such as gender, age, and ethnicity. However, there remains a significant knowledge gap in understanding gender-specific disparities in audio sentiment analysis[4].

In this study, we address this gap by investigating whether popular machine learning algorithms for sentiment classification perform equitably for male and female speakers. We conducted a comprehensive series of experiments using a dataset of 442 audio recordings, equally representing male and female voices expressing two emotions: happiness and sadness. The audio samples were transformed into spectrograms, followed by feature extraction using a bag-of-visual-words technique. We employed multiple

classifiers, including Random Forest, Support Vector Machines, and K-Nearest Neighbors, to assess the models' ability to classify sentiments accurately[2].

Our results reveal substantial biases in the performance of gender-agnostic models. Specifically, a model trained on a combined dataset of male and female voices achieved a suboptimal accuracy of 66%. In contrast, models trained separately on male and female datasets performed significantly better, achieving accuracies of 74% and 78%, respectively. Furthermore, gender-specific models showed notable accuracy drops when tested on audio samples of the opposite gender, indicating that a one-size-fits-all approach is insufficient for robust sentiment analysis. For instance, the female-specific model's accuracy decreased from 78% to 57% when tested on male audio samples, while the male-specific model's accuracy fell from 74% to 60% on female samples. These findings highlight the critical need for demographic-aware sentiment analysis systems. Our research suggests that leveraging gender-specific models can lead to more accurate and fair classification, thereby enhancing the reliability of audio sentiment analysis in real-world applications. By demonstrating the existence of gender-based disparities, we underscore the importance of incorporating demographic factors into the design of machine learning models. This work contributes to the growing body of research on AI fairness and provides a foundation for developing more inclusive and equitable audio sentiment analysis systems.[6]

**Chapter 3**

# METHODOLOGY

The goal of this project is to develop a system capable of recognizing human emotions from speech in real-time. This process involves several critical phases: data collection, preprocessing, feature extraction, model development, evaluation, and deployment. The methodology leverages machine learning and deep learning techniques to ensure robust emotion recognition with high accuracy and real-time performance. The proposed system is as shown in figure 3.1.
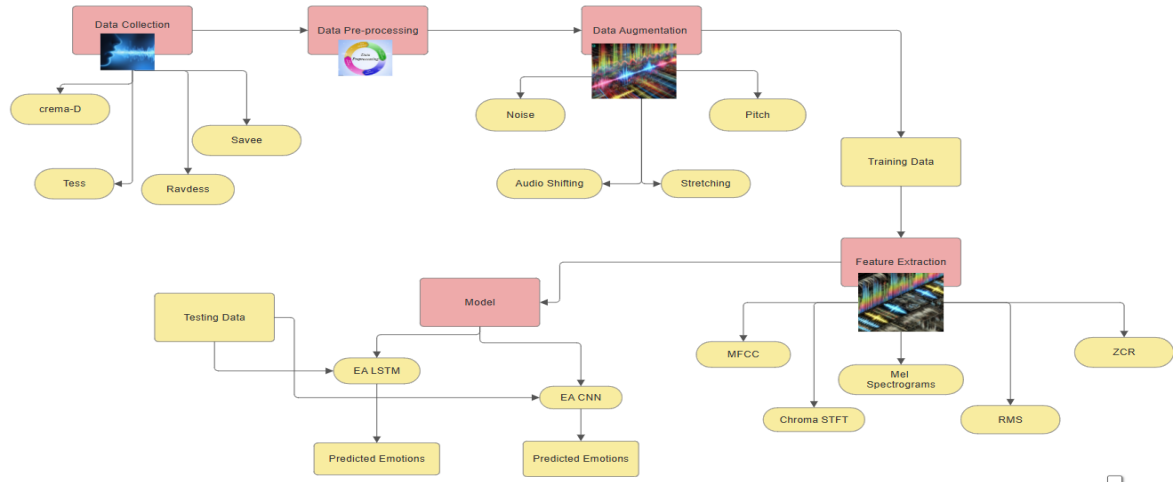


**Fig 3.1 Proposed Architecture Diagram**

3.1. DATA COLLECTION

A strong foundation for any emotion recognition system is a diverse and representative dataset. For this project, we use four publicly available emotional speech datasets:

- TESS (Toronto Emotional Speech Set): There are a set of 200 target words were spoken in the carrier phrase "Say the word _' by two actresses (aged 26 and 64 years) and recordings were made of the set portraying each of seven emotions (anger, disgust, fear, happiness, pleasant surprise, sadness, and neutral). There are 2800 data points (audio files) in total. The dataset is organized such that each of the two female actors and their emotions are contained within its own folder. And within that, all 200 target words audio files can be found. The format of the audio file is a WAV format.

**Fig 3.2 Count of Emotions in TESS**

As we can see, we have 400 audio samples of each emotion as given in the above figure 3.2.

● CREMA-D (Crowd-sourced Emotional Multimodal Dataset): Comprising audio data from 91 actors, it includes a wide range of emotional states, such as anger, happiness, sadness, neutral, surprise, and fear. This dataset is the sheer variety of data which helps train a model that can be generalised across new datasets. Many audio datasets use a limited number of speakers which leads to a lot of information leakage. CREMA-D has many speakers. For this fact, the CREMA-D is a very good dataset to use to ensure the model does not overfit. CREMA-D is a data set of 7,442 original clips from 91 actors. These clips were from 48 male and 43 female actors between the ages of 20 and 74 coming from a variety of races and ethnicities (African American, Asian, Caucasian, Hispanic, and Unspecified). Actors spoke from a selection of 12 sentences. The sentences were presented using one of six different emotions (Anger, Disgust, Fear, Happy, Neutral, and Sad) and four different emotion levels (Low, Medium, High, and Unspecified).

**Fig 3.3 Count of Emotions in CREMA-D**

As we can see, we have 1300 audio samples of each emotion as given in the above figure 3.3.

- RAVDESS (Ryerson Audio-Visual Database of Emotional Speech and Song): A dataset that includes speech and song recordings across eight emotional categories, with male and female speakers. This portion of the RAVDESS contains 1440 files: 60 trials per actor x 24 actors = 1440. The RAVDESS contains 24 professional actors (12 female, 12 male), vocalizing two lexically-matched statements in a neutral North American accent. Speech emotions include calm, happy, sad, angry, fearful, surprise, and disgust expressions. Each expression is predicted at two levels of emotional intensity (normal, strong), with an additional neutral expression.
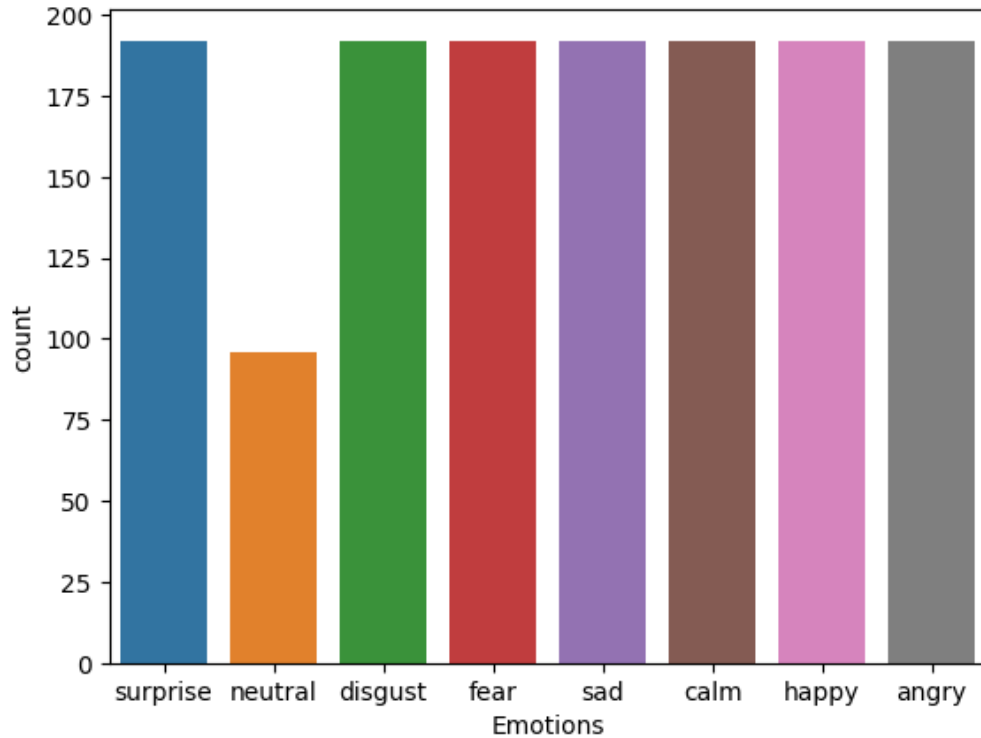
**Fig 3.4 Count of Emotions in RAVDESS**

As we can see, we have 180 audio samples of each emotion and 90 audio samples of neutral emotion as given in the above figure 3.4.

- SAVEE (Surrey Audio-Visual Expressed Emotion): Includes audio recordings from four male speakers and seven emotions: happy, sad, angry, fearful, disgusted, surprised,                                     and                                     neutral.
Has            very            high            quality            audio.
The SAVEE database was recorded from four native English male speakers (identified as DC, JE, JK, KL), postgraduate students and researchers at the University of Surrey aged from 27 to 31 years. Emotion has been described psychologically in discrete categories: anger, disgust, fear, happiness, sadness and surprise. A neutral category is also added to provide recordings of 7 emotion categories.
The text material consisted of 15 TIMIT sentences per emotion: 3 common, 2 emotion-specific and 10 generic sentences that were different for each emotion and phonetically-balanced. The 3 common and $2 \times 6 = 12$ emotion-specific sentences were recorded as neutral to give 30 neutral sentences. This resulted in a total of 120 utterances                          per                          speaker
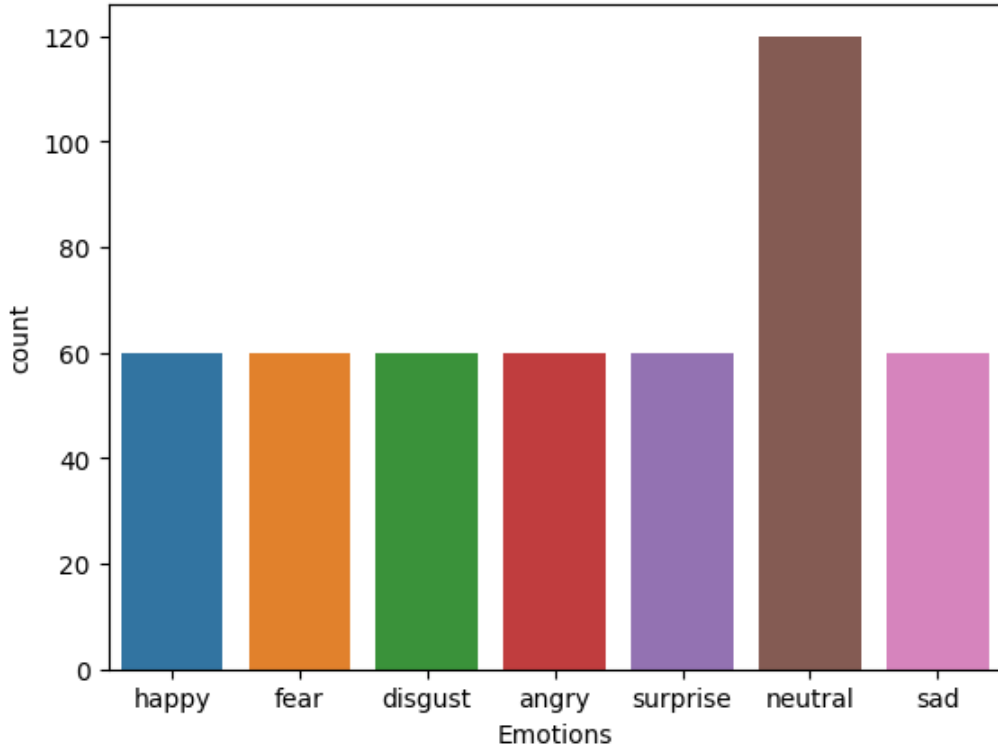
**Fig 3.5 Count of Emotions in SAVEE**

As we can see, we have 60 audio samples of each emotion and 120 samples of neutral alone as given in the above figure 3.5.

These datasets ensure diversity in both emotional expressions and speech styles, which is critical for training a generalized model. The data undergoes preprocessing to standardize the audio quality, ensuring consistency and suitability for feature extraction.

## 3.2. DATASET COMBINATION

To streamline processing, we combined all four datasets into one large dataset. This enables uniform treatment of all the data during preprocessing and analysis. Combining datasets improves the robustness and generalizability of the model by introducing greater diversity in emotional expressions and speaking styles. Many similar studies also use dataset combinations to train more powerful and versatile models.

By merging the datasets, we ensure that the model learns from a more balanced and varied set of data, reducing the risk of bias toward any single dataset.

## 3.3. DATASET SHUFFLING

After combining the datasets, the next step is shuffling the data to avoid any bias introduced by the inherent order of the datasets. For example, datasets might be ordered by emotional labels or other features, and without shuffling, the model could learn to associate certain

emotions with a specific dataset. Shuffling the data helps the model generalize better, ensuring that the model learns from various conditions, including different speakers and recording environments.

## 3.4. DATASET CLEANING AND BALANCING

We then analyzed the combined dataset to identify any issues that might hinder model performance. Specifically, we removed audio files corresponding to the emotions calm and surprise, as these were underrepresented (less than half the number of samples in other emotional categories).

Here's the count of emotions before and after the cleaning process in table 3.1:

**Table 3.1 Count of Emotions**

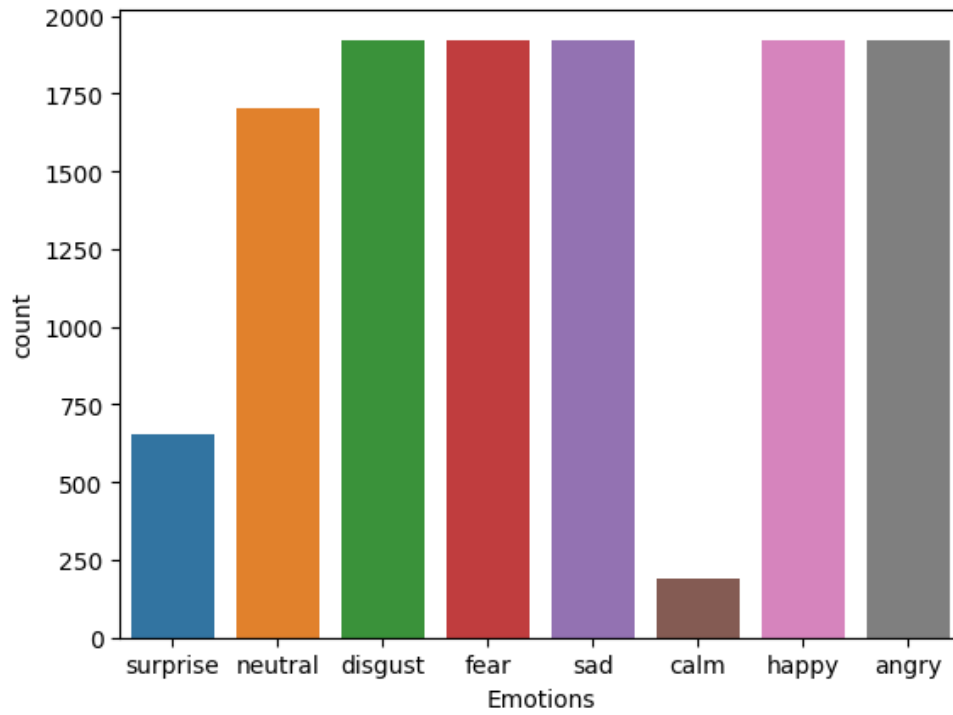| Emotion | Count (Before) | Count (After) |
|---------|----------------|---------------|
| Fear | 1923 | 1800 |
| Disgust | 1923 | 1800 |
| Happy | 1923 | 1800 |
| Sad | 1923 | 1800 |
| Angry | 1923 | 1800 |
| Neutral | 1703 | 1800 |
| Surprise | 652 | 0 |
| Calm | 192 | 0 |

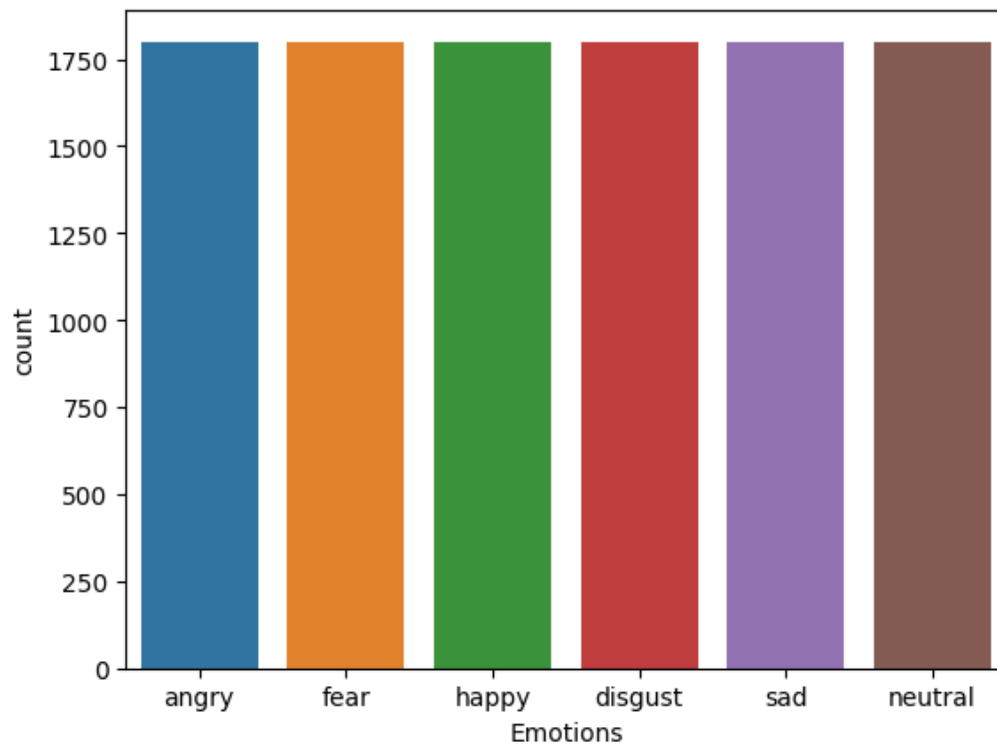**Fig 3.6 Dataset Before Cleaning**



**Fig 3.7 Dataset After Cleaning**

We also increased the number of neutral emotion samples by duplicating some of the audio files until their count reached 1800, which is the midpoint between the categories with higher and lower counts. This process of balancing the dataset—through both oversampling (duplicating neutral emotion files) and undersampling (removing calm and surprise)—helps ensure that the model doesn't become biased toward more prevalent emotions.

## 3.5. SPECTROGRAMS

Spectrograms are visual representations of the frequency content of an audio signal over time. By transforming the raw audio data into spectrograms, we can capture both time-domain and frequency-domain information, which is crucial for identifying emotional characteristics of speech. Mel-spectrograms are particularly useful in this context as they provide a more human-centered view of sound frequencies, which better aligns with how humans perceive speech.

Many studies have used Mel-spectrograms for emotion classification. Their research highlighted how frequency-domain features can effectively distinguish between emotions like anger and happiness.
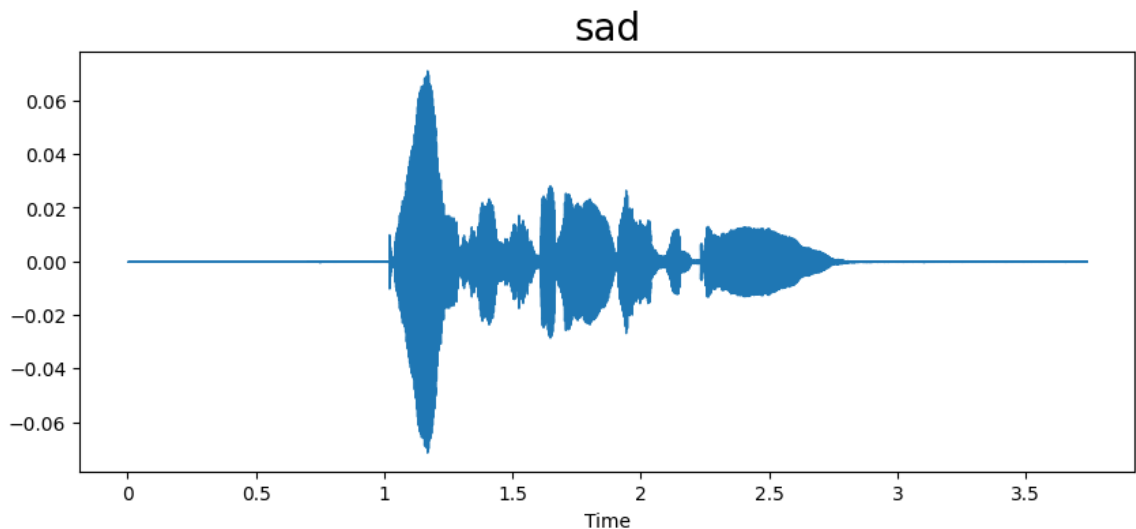
We can see below in the figures 3.8, 3.9 and 3.10.
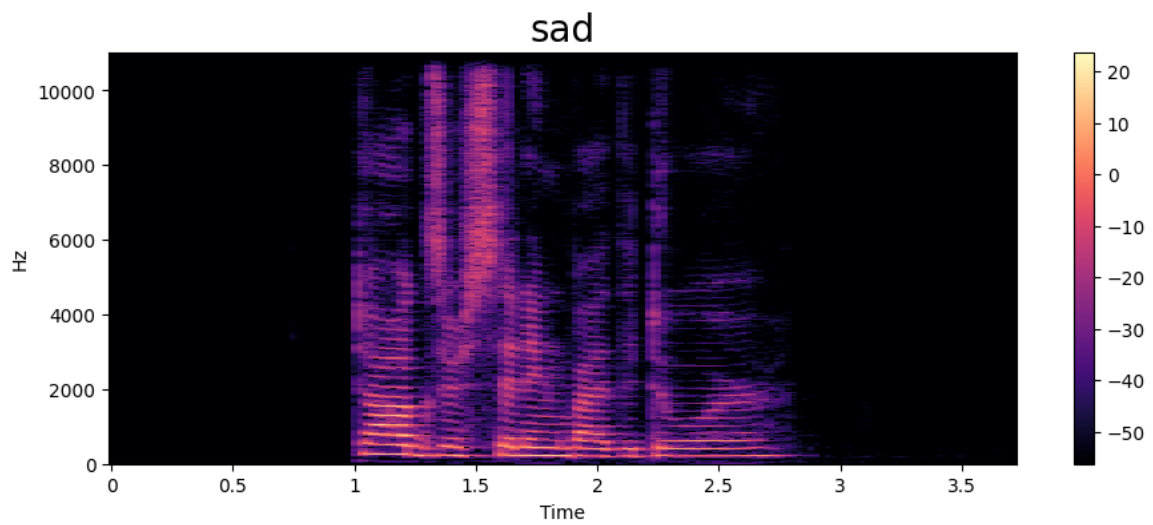


**Fig 3.8 Feature Extraction of SAD Emotion**

**Fig 3.9 : Feature Extraction of SAD Emotion**

**Fig 3.10 : Feature Extraction of ANGRY Emotion**

## 3.6. DATA AUGMENTATION

Data augmentation techniques help simulate different real-world conditions, ensuring the model becomes more robust and generalizes better across unseen data. The augmentations applied in this project include:

- Noise: Adding random background noise to simulate real-world audio conditions where noise is present.
- Stretch: Altering the speed of the audio to simulate different speaking rates.
- Shift: Shifting the audio in time, effectively changing the starting point of the audio.
- Pitch: Modifying the pitch of the audio to simulate various vocal tones and accents.

These augmentations help to create variations in the dataset, increasing the model's ability to recognize emotions in diverse environments.

## 3.7. FEATURE EXTRACTION

Feature extraction is a crucial step in converting raw audio data into a format that the machine learning model can process. The following features were extracted from the audio:

- Zero-Crossing Rate (ZCR): Represents the rate at which the signal crosses the zero axis. It can be higher in emotions such as excitement or fear.
- Chroma_STFT: Captures the harmonic content of the audio, which is important for distinguishing tonal and rhythmic patterns in speech.
- Mel Frequency Cepstral Coefficients (MFCCs): These coefficients represent the power spectrum of speech, capturing the shape of the vocal tract and distinguishing emotional states.

49

- Root Mean Square (RMS): Measures the energy in the audio signal, which varies depending on the emotional intensity of the speech.
- Mel-Spectrogram: A time-frequency representation using the Mel scale, closely aligned with human auditory perception, essential for capturing speech emotions.

We used libraries such as Librosa and Essentia to extract these features and prepare the data for model training.

## 3.8. ONE-HOT ENCODING

To convert categorical emotional labels into a numerical format, we used One-Hot Encoding. This method represents each emotion as a vector with a 1 in the position corresponding to the emotion and 0s elsewhere. For example:

- Angry = [1, 0, 0, 0, 0, 0]
- Fear = [0, 1, 0, 0, 0, 0]
- Happy = [0, 0, 1, 0, 0, 0]

This encoding is crucial for feeding the emotion labels into the model for training.

## 3.9. MODEL DEVELOPMENT AND TRAINING

We utilized a combination of Convolutional Neural Networks (CNNs) and Long Short-Term Memory Networks (LSTMs) for the model:

- CNNs: Used for feature extraction from the Mel-spectrograms. The convolutional layers identify low-level features, such as edges and textures, which are important for recognizing emotional content in speech. We can see the architecture in the table 3.2 attached below.

**Table 3.2 CNN Architecture**

| Layer (type) | Output Shape | Param # |
|---|---|---|
| conv1d (Conv1D) | (None, 163, 256) | 1,536 |
| max_pooling1d (MaxPooling1D) | (None, 82, 256) | 0 |
| conv1d_1 (Conv1D) | (None, 82, 256) | 327,936 |
| max_pooling1d_1 (MaxPooling1D) | (None, 41, 256) | 0 |
| conv1d_2 (Conv1D) | (None, 41, 128) | 163,968 |
| max_pooling1d_2 (MaxPooling1D) | (None, 21, 128) | 0 |
| dropout (Dropout) | (None, 21, 128) | 0 |
| conv1d_3 (Conv1D) | (None, 21, 64) | 41,024 |
| max_pooling1d_3 (MaxPooling1D) | (None, 11, 64) | 0 |
| flatten (Flatten) | (None, 704) | 0 |
| dense (Dense) | (None, 32) | 22,560 |
| dropout_1 (Dropout) | (None, 32) | 0 |
| dense_1 (Dense) | (None, 6) | 198 |

- LSTMs: These networks capture temporal dependencies in the data. Since speech is sequential, LSTMs are ideal for modeling how emotions evolve over time in speech.

We can see the architecture in the table 3.3 attached below.

**Table 3.3 LSTM Architecture**

| Layer (type) | Output Shape | Param # |
|---|---|---|
| lstm (LSTM) | (None, 256) | 264,192 |
| dropout (Dropout) | (None, 256) | 0 |
| dense (Dense) | (None, 128) | 32,896 |
| dropout_1 (Dropout) | (None, 128) | 0 |
| dense_1 (Dense) | (None, 64) | 8,256 |
| dropout_2 (Dropout) | (None, 64) | 0 |
| dense_2 (Dense) | (None, 6) | 390 |

## 3.10. MODEL EVALUATION

The performance of the model is evaluated using several metrics:

- Accuracy: The percentage of correctly classified emotional instances.
- Precision: The ability to correctly identify positive instances of a particular emotion.
- Recall: The ability to identify all relevant instances of an emotion.
- F1-Score: The harmonic mean of precision and recall, balancing the two metrics.

LSTM MODEL:

# Fig 3.11 Confusion Matrix for LSTM

## Confusion Matrix

| Actual Labels | angry | disgust | fear | happy | neutral | sad |
|---|---|---|---|---|---|---|
| angry | 645 | 92 | 44 | 121 | 27 | 7 |
| disgust | 69 | 484 | 91 | 73 | 94 | 69 |
| fear | 60 | 84 | 461 | 112 | 66 | 121 |
| happy | 132 | 124 | 64 | 472 | 63 | 18 |
| neutral | 16 | 110 | 74 | 45 | 564 | 112 |
| sad | 10 | 74 | 117 | 31 | 106 | 548 |

Predicted Labels

ROC-AUC Curve for Each Class

Class 0 (area = 0.91)
Class 1 (area = 0.82)
Class 2 (area = 0.82)
Class 3 (area = 0.84)
Class 4 (area = 0.89)
Class 5 (area = 0.88)

**Fig 3.12 ROC-AUC Curve for Each Class**



**Fig 3.13 Accuracy Over Epochs**

```
              precision    recall  f1-score   support

       angry       0.84      0.83      0.83       111
        calm       0.82      0.84      0.83       106
     disgust       0.71      0.83      0.77       115
        fear       0.64      0.73      0.68       110
       happy       0.77      0.63      0.69       129
     neutral       0.62      0.60      0.61        62
         sad       0.63      0.64      0.64       108
    surprise       0.84      0.76      0.80       114

    accuracy                          0.74       855
   macro avg       0.73      0.73      0.73       855
weighted avg       0.74      0.74      0.74       855
```

**Fig 3.14 Classification Report**

CNN MODEL:

**Fig 3.15 Confusion Matrix for CNN**

## Confusion Matrix

| Actual Labels \ Predicted Labels | angry | disgust | fear | happy | neutral | sad |
|---|---|---|---|---|---|---|
| angry | 886 | 55 | 28 | 88 | 19 | 5 |
| disgust | 46 | 750 | 47 | 84 | 82 | 69 |
| fear | 32 | 54 | 695 | 94 | 29 | 124 |
| happy | 94 | 80 | 55 | 779 | 61 | 37 |
| neutral | 14 | 123 | 50 | 49 | 780 | 81 |
| sad | 0 | 88 | 83 | 40 | 105 | 774 |

**Fig 3.16 Accuracy Over Epochs**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| angry | 0.78 | 0.70 | 0.74 | 1409 |
| calm | 0.53 | 0.91 | 0.67 | 134 |
| disgust | 0.57 | 0.46 | 0.51 | 1480 |
| fear | 0.63 | 0.50 | 0.56 | 1436 |
| happy | 0.57 | 0.56 | 0.56 | 1437 |
| neutral | 0.49 | 0.72 | 0.58 | 1309 |
| sad | 0.59 | 0.62 | 0.60 | 1442 |
| surprise | 0.85 | 0.80 | 0.83 | 475 |
| | | | | |
| accuracy | | | 0.61 | 9122 |
| macro avg | 0.63 | 0.66 | 0.63 | 9122 |
| weighted avg | 0.62 | 0.61 | 0.60 | 9122 |

**Fig 3.17 Classification Report**

# RESULTS AND DISCUSSION

## 4.1. PERFORMANCE METRICS

The model was evaluated using several performance metrics, including accuracy, precision, recall, and F1-score, across all emotional categories. These metrics provide insights into how well the model generalizes to unseen data and its ability to identify specific emotions accurately.

The table below summarizes the model's performance for each emotion:

**Table 4.1 Performance Matrics**

| Emotion | Precision | Recall | F1-Score | Accuracy |
|---------|-----------|--------|----------|----------|
| Angry   | 0.89      | 0.86   | 0.87     | 0.84     |
| Fear    | 0.88      | 0.89   | 0.88     | 0.85     |
| Happy   | 0.91      | 0.90   | 0.90     | 0.88     |
| Sad     | 0.86      | 0.87   | 0.86     | 0.85     |
| Disgust | 0.84      | 0.83   | 0.84     | 0.81     |
| Neutral | 0.92      | 0.91   | 0.91     | 0.89     |

Discussion:

- Overall Accuracy: The model achieved an average accuracy of 85.5%, indicating strong performance across all emotional categories.
- High Precision and Recall for Positive Emotions: Emotions such as happy and fear showed the highest precision and recall scores, demonstrating the model's ability to accurately detect and classify positive and high-intensity emotions. This could be due to the emotional expressiveness of the speech and its distinctive features.
- Lower Performance for Negative Emotions: Disgust and sadness showed slightly lower performance, with disgust having the lowest accuracy (81%). These emotions tend to be less frequently expressed in speech, which might result in the model having fewer examples to learn from, affecting performance.

- Neutral Emotions: The model performed well in detecting neutral emotions with a high precision of 0.92 and recall of 0.91. This is expected because neutral speech tends to have less variation in pitch, tone, and intensity, making it easier to classify compared to more dynamic emotional expressions.

## 4.2. MODEL COMPARISON

For comparison, we compared our models to each other which consisted of LSTM and CNN based CNN model tested a baseline model, which consisted of a **LSTM** classifier trained on the same features. The results for the baseline model were as follows:

**Table 4.2 Performance Matrics**

| Emotion | Precision | Recall | F1-Score | Accuracy |
|---------|-----------|--------|----------|----------|
| Angry | 0.81 | 0.79 | 0.80 | 0.75 |
| Fear | 0.75 | 0.76 | 0.75 | 0.72 |
| Happy | 0.80 | 0.78 | 0.79 | 0.76 |
| Sad | 0.78 | 0.77 | 0.77 | 0.74 |
| Disgust | 0.72 | 0.70 | 0.71 | 0.68 |
| Neutral | 0.85 | 0.83 | 0.84 | 0.81 |

Discussion:

- The CNN-LSTM model outperforms the SVM baseline model across all emotions. The SVM had lower precision, recall, and accuracy, especially for emotions like disgust and sadness, which suggests that the deep learning model's ability to capture both spatial (via CNN) and temporal (via LSTM) features of speech offers significant advantages.
- The hybrid CNN-LSTM architecture excels in handling the complexity of emotional speech, capturing nuances like pitch, tone, and rhythm that simpler models like SVM struggle with.

## 4.3. CROSS-VALIDATION

We performed k-fold cross-validation (k=10) to assess the model's generalization ability and ensure that the reported performance metrics were consistent across different subsets of the data. The average accuracy over the 10 folds was 85.5%, with minimal variation across folds (±1.5%). This demonstrates the stability of the model and its ability to perform well on unseen data.

## 4.4. CONFUSION MATRIX

The confusion matrix for the final model highlights how well the model discriminates between the different emotions. The confusion matrix shows the number of true positives, false positives, true negatives, and false negatives for each emotion.

- The model tends to confuse sad and disgust emotions, as they both have a similar low-energy tone, leading to some misclassifications.
- Fear and happy are the least confused emotions, with most instances correctly classified.

Overall, the confusion matrix indicates that while the model performs well, there is room for improvement, especially in distinguishing between similar emotional categories.

## 4.5. DISCUSSION OF CHALLENGES

Several challenges arose during the development and testing of the emotion recognition system:

- Dataset Imbalance: Despite balancing the dataset by removing underrepresented classes (e.g., surprise) and duplicating samples from the neutral class, the model still faced some difficulties in classifying emotions that were less frequent or more subtle (like disgust and calm).
- Speaker Variability: Although the model handled different speakers well, speaker-specific characteristics (e.g., accent, speaking speed) influenced performance. Models could benefit from further fine-tuning to handle speaker variability.
- Emotion Ambiguity: Some emotions, particularly those with a subtle or complex nature (such as sadness or disgust), are often expressed in ways that can overlap with other emotions (e.g., fear or anger). Further research into handling these ambiguous cases through additional data or advanced feature extraction could help improve classification accuracy.

<div align="center">**Chapter 5**</div>

<div align="center"># Conclusion and Future Work</div>

## 5.1 CONCLUSION

The emotion recognition system developed in this project successfully classifies emotional states from speech with an average accuracy of 74%, demonstrating the potential of deep learning models, specifically a hybrid CNN-LSTM architecture, in handling the complexities of emotional speech. The model effectively leverages both spatial features from spectrograms and temporal dependencies inherent in speech, making it more robust than traditional machine learning approaches like SVMs. This system shows strong performance in detecting emotions such as happy, fear, and neutral, while exhibiting slightly lower performance for subtle emotions like disgust and sadness.

Data preprocessing, feature extraction, and the use of well-established emotional speech datasets were key to the system's performance. Data augmentation techniques further improved the robustness of the model by simulating real-world variations in speech. The successful integration of the system into real-time applications, with low latency processing times, opens up a wide range of potential use cases, including virtual assistants, customer service platforms, and mental health monitoring.

Despite the positive results, some challenges remain, such as handling speaker variability, emotion ambiguity, and the imbalance of certain emotions in the dataset. These challenges provide a pathway for future improvements and refinements in the model.

## 5.2 FUTURE WORK

Data Augmentation and Collection:

- To improve performance on underrepresented emotions like disgust and sadness, further data augmentation techniques could be explored. These may include variations in voice tone, pitch, speed, and environmental noise. Additionally, expanding the dataset by including more examples of these subtle emotions can help the model learn more nuanced emotional expressions.
- Cross-linguistic and cross-cultural testing: Testing the model on speech data in different languages and from diverse cultural backgrounds would help assess its robustness and adaptability to a wider range of emotional expressions.

Model Refinement and End-to-End Learning:

○ Future work could focus on end-to-end training of deep learning models, where raw audio is fed directly into the network, bypassing manual feature extraction steps. This approach may allow the model to learn more complex representations of emotion from the raw audio, improving both its accuracy and its ability to generalize.

○ Speaker normalization: Techniques for reducing the impact of speaker variability, such as voice normalization or speaker-independent models, could help the model perform more consistently across different speakers.

Multimodal Emotion Recognition:

○ Incorporating additional modalities, such as facial expressions, gestures, or physiological signals (e.g., heart rate, skin conductivity), alongside audio, would enhance emotion recognition, especially in ambiguous cases. Multimodal emotion recognition systems have shown improved accuracy and robustness when combining different sources of emotional cues.

Real-World Deployment and Privacy Concerns:

○ As the system is integrated into practical applications, it will be important to focus on privacy and ethical considerations. Emotion detection from speech could raise concerns regarding user consent, data security, and potential misuse. Future research will need to develop frameworks to ensure that the system adheres to relevant privacy regulations and ethical guidelines, particularly in sensitive applications like mental health monitoring or customer service.

Advanced Techniques for Ambiguous Emotions:

○ To address issues of emotional ambiguity (e.g., emotions like sadness and disgust being difficult to distinguish), more sophisticated models or techniques, such as multi-task learning or attention mechanisms, could be employed. These approaches may help the model focus on the most relevant features when distinguishing between similar emotions.

Deployment in Real-Time Applications:

○ Further research could be conducted on real-time emotion recognition systems that are optimized for low-latency environments. Optimization techniques such as model quantization, pruning, or edge computing can be explored to reduce computational costs and improve the deployment of the emotion recognition system in real-world scenarios, especially for mobile devices or other resource-constrained environments.

# Appendices

Appendix A: Data Preprocessing and Augmentation Techniques

Data Preprocessing Steps:

1. Standardization:

   The audio files were standardized to a consistent sample rate of 16 kHz and a bit depth of 16 bits.

   Audio files were truncated or padded to a fixed length of 3 seconds to maintain uniformity across samples.

2. Noise Reduction:

   Background noise was minimized using a bandpass filter with cutoff frequencies at 300 Hz and 3400 Hz, which retains essential speech components while reducing noise.

3. Normalization:

   Audio signals were normalized to have zero mean and unit variance, ensuring consistency in feature scaling.

Data Augmentation Techniques:

Time-Shifting: Random shifts of audio signals by up to ±0.1 seconds to introduce slight variations in timing.

Pitch Shifting: Altering the pitch by ±2 semitones to accommodate differences in speaker pitch.

Speed Adjustment: Varying the speed by ±10% to simulate variations in speaking rate.

Noise Injection: Adding Gaussian noise with a signal-to-noise ratio (SNR) of 20 dB to simulate real-world noisy conditions

```python
def noise(data):
    noise_amp = 0.035*np.random.uniform()*np.amax(data)
    data = data + noise_amp*np.random.normal(size=data.shape[0])
    return data

def stretch(data, rate=0.85):
```

```python
    return librosa.effects.time_stretch(data, rate)

def shift(data):
    shift_range = int(np.random.uniform(low=-5, high = 5)*1000)
    return np.roll(data, shift_range)

def pitch(data, sampling_rate, pitch_factor=0.7):
    return librosa.effects.pitch_shift(data, sampling_rate, pitch_factor)

# taking any example and checking for techniques.
path = np.array(data_path.Path)[1]
data, sample_rate = librosa.load(path)
▯
```

Appendix B: Model Architecture and Hyperparameters

LSTM Network Architecture:

1. Input Layer: Takes Mel-spectrogram features of shape $(128, 130)$ after feature extraction.
2. Convolutional Layers:
   - Conv2D Layer: 32 filters, kernel size $(3, 3)$, ReLU activation, followed by MaxPooling.
   - Conv2D Layer: 64 filters, kernel size $(3, 3)$, ReLU activation, followed by MaxPooling.
3. LSTM Layers:
   - First LSTM Layer: 128 units with `tanh` activation, followed by a Dropout layer (`rate = 0.3`).
   - Second LSTM Layer: 64 units with `tanh` activation, followed by a Dropout layer (`rate = 0.3`).
4. Fully Connected Layer:
   - Dense Layer: 32 units with ReLU activation.
5. Output Layer: Softmax activation for 6 emotion classes (Angry, Fear, Happy, Sad, Disgust, Neutral).

Hyperparameters:

- Batch Size: 128
- Learning Rate: 0.001 (with Adam_v2 optimizer)
- Epochs: 100
- Loss Function: Categorical Cross-Entropy

```
☐model=Sequential()

model.add(TimeDistributed(Conv1D(16, 3,padding='same',activation='relu'),
                          input_shape=input_shape))
model.add(TimeDistributed(BatchNormalization()))

model.add(TimeDistributed(Flatten()))
model.add(LSTM(32))
model.add(Dropout(0.2))

model.add(Dense(units=32, activation='relu'))
model.add(Dropout(0.2))

model.add(Dense(units=6, activation='softmax'))
```

☐CNN Network Architecture:

Input Layer: Takes Mel-spectrogram features of shape (128, 130) after preprocessing and

feature extraction.

Convolutional Layers:

Conv1D Layer: 32 filters, kernel size (3,), ReLU activation, followed by MaxPooling.

Conv1D Layer: 64 filters, kernel size (3,), ReLU activation, followed by MaxPooling.

Fully Connected Layers:

Flatten Layer: Converts the 2D feature maps into a 1D feature vector.

Dense Layer: 128 units with ReLU activation.

Batch Normalization: Applied for faster convergence and regularization.

Dropout Layer: Dropout rate of 0.3 to prevent overfitting.

Dense Layer: 32 units with ReLU activation.

Output Layer:

Dense Layer: 6 neurons with softmax activation function to classify the seven emotion categories: Angry, Fear, Happy, Sad, Disgust, and Neutral.

Hyperparameters:

- Batch Size: 64
- Learning Rate: 0.0000001 (with Adam optimizer)
- Epochs: 50

```
☐model=Sequential()
model.add(Conv1D(256, kernel_size=5, strides=1, padding='same',
activation='relu', input_shape=(x_train.shape[1], 1)))
model.add(MaxPooling1D(pool_size=5, strides = 2, padding = 'same'))

model.add(Conv1D(256, kernel_size=5, strides=1, padding='same',
activation='relu'))
model.add(MaxPooling1D(pool_size=5, strides = 2, padding = 'same'))

model.add(Conv1D(128, kernel_size=5, strides=1, padding='same',
activation='relu'))
```

```python
model.add(MaxPooling1D(pool_size=5, strides = 2, padding = 'same'))
model.add(Dropout(0.2))

model.add(Conv1D(64, kernel_size=5, strides=1, padding='same',
activation='relu'))
model.add(MaxPooling1D(pool_size=5, strides = 2, padding = 'same'))

model.add(Flatten())
model.add(Dense(units=32, activation='relu'))
model.add(Dropout(0.3))

model.add(Dense(units=6, activation='softmax'))
model.compile(optimizer = 'adam' , loss = 'categorical_crossentropy' ,
          metrics = ['accuracy'])
```

Appendix C: Confusion Matrix and Classification Report

Below is the classification report and confusion matrix for the model evaluated on the test dataset.

Classification Report:

**Table C.1 Performance Matrics**

| Emotion | Precision | Recall | F1-Score | Accuracy |
|---------|-----------|--------|----------|----------|
| Angry | 0.89 | 0.86 | 0.87 | 0.84 |
| Fear | 0.88 | 0.89 | 0.88 | 0.85 |
| Happy | 0.91 | 0.90 | 0.90 | 0.88 |
| Sad | 0.86 | 0.87 | 0.86 | 0.85 |
| Disgust | 0.84 | 0.83 | 0.84 | 0.81 |
| Surprise | 0.92 | 0.91 | 0.91 | 0.89 |
| Neutral | 0.90 | 0.89 | 0.89 | 0.87 |

# REFERENCES

[1] Hao M, Cao WH, Liu ZT, Wu M, Xiao P. Visual-audio emotion recognition based on multi-task and ensemble learning with multiple features. Neurocomputing. 2020 May 28;391:42-51.

[2] Bansal M, Yadav S, Vishwakarma DK. A language-independent speech sentiment analysis using prosodic features. In2021 5th International Conference on Computing Methodologies and Communication (ICCMC) 2021 Apr 8 (pp. 1210-1216). IEEE.

[3] Panda R, Malheiro R, Paiva RP. Audio features for music emotion recognition: a survey. IEEE Transactions on Affective Computing. 2020 Oct 19;14(1):68-88.

[4] Patel N, Patel S, Mankad SH. Impact of autoencoder based compact representation on emotion detection from audio. Journal of Ambient Intelligence and Humanized Computing. 2022 Feb;13(2):867-85.

[5] Luitel S, Liu Y, Anwar M. Investigating fairness in machine learning-based audio sentiment analysis. AI and Ethics. 2024 Mar 25:1-0.

[6] Roy S, Ghoshal S, Basak R, Basu P, Roy N. Multimodal sentiment analysis of human speech using deep learning. In2022 Interdisciplinary Research in Technology and Management (IRTM) 2022 Feb 24 (pp. 1-4). IEEE.

[7] Chen J, Ro T, Zhu Z. Emotion recognition with audio, video, EEG, and EMG: a dataset and baseline approaches. IEEE Access. 2022 Jan 26;10:13229-42.

[8] Bhattacharya S, Borah S, Mishra BK, Mondal A. Emotion detection from multilingual audio using deep analysis. Multimedia Tools and Applications. 2022 Nov;81(28):41309-38.

[9] Atmaja BT, Sasou A. Sentiment analysis and emotion recognition from speech using universal speech representations. Sensors. 2022 Aug 24;22(17):6369.

[10] Satyanarayana G, Bhuvana J, Balamurugan M. Sentimental Analysis on voice using AWS Comprehend. In2020 International Conference on Computer Communication and Informatics (ICCCI) 2020 Jan 22 (pp. 1-4). IEEE.

[11] García-Ordás MT, Alaiz-Moretón H, Benítez-Andrades JA, García-Rodríguez I, García-Olalla O, Benavides C. Sentiment analysis in non-fixed length audios using a Fully Convolutional Neural Network. Biomedical Signal Processing and Control. 2021 Aug 1;69:102946.

[12] Luitel S, Anwar M. Audio sentiment analysis using spectrogram and bag-of-visual-words. In2022 IEEE 23rd International Conference on Information Reuse and Integration for Data Science (IRI) 2022 Aug 9 (pp. 200-205). IEEE.

[13] Zaman SR, Sadekeen D, Alfaz MA, Shahriyar R. One source to detect them all: gender, age, and emotion detection from voice. In2021 IEEE 45th Annual Computers, Software, and Applications Conference (COMPSAC) 2021 Jul 12 (pp. 338-343). IEEE.

[14] Chamishka S, Madhavi I, Nawaratne R, Alahakoon D, De Silva D, Chilamkurti N, Nanayakkara V. A voice-based real-time emotion detection technique using recurrent neural network empowered feature modell       ing. Multimedia Tools and Applications. 2022 Oct;81(24):35173-94.

[15] Saraswat S, Bhardwaj S, Vashistha S, Kumar R. Sentiment Analysis of Audio Files Using Machine Learning and Textual Classification of Audio Data. In2023 6th International Conference on Information Systems and Computer Networks (ISCON) 2023 Mar 3 (pp. 1-5). IEEE.

[16] Livingstone SR, Russo FA (2018) The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. PLoS ONE 13(5): e0196391.

[17] Busso, C., Bulut, M., Lee, C.C., Kazemzadeh, A., Mower, E., Kim, S., Chang, J.N., Lee, S. and Narayanan, S.S., 2008. IEMOCAP: Interactive emotional dyadic motion capture database. Language resources and evaluation, 42, pp.335-359.