

VIT[®]

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

Case Study III

Emotion Recognition from Speech

GitHub Repo link:

<https://github.com/pranshu-rastogi/Emotion-Recognition-from-Speech>

Team Details

Pranshu Rastogi 22BCE0441

Submitted to:

Prof. Rajeshkannan R

Abstract:

Emotion recognition from speech is a key area in affective computing, enabling machines to interpret and respond to human emotions through vocal cues. Such systems have wide applications in mental health monitoring, assistive technology, and human-computer interaction. However, accurately classifying emotions remains challenging due to inter-speaker variability, acoustic differences, and the nuanced nature of emotional expression.

This study proposes a unimodal speech emotion recognition (SER) framework based on acoustic feature analysis and deep learning. The model was trained and evaluated using three benchmark datasets - Toronto Emotional Speech Set (TESS), Crowd-sourced Emotional Multimodal Actors Dataset (CREMA-D), and Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) to ensure demographic diversity and enhance generalization. From each audio sample, features such as Mel-Frequency Cepstral Coefficients (MFCCs), pitch, energy, spectral flux, and zero-crossing rate were extracted and normalized using a StandardScaler.

A hybrid CNN-LSTM model was implemented, where convolutional layers captured local spectral characteristics, and recurrent layers modelled temporal dependencies. The model achieved a test accuracy of 65.14% and a macro F1-score of 0.65, demonstrating balanced performance across six emotion classes: angry, disgust, happy, neutral, sad, and scared. Stronger recognition of anger and sadness indicates the model's sensitivity to distinct acoustic patterns associated with these emotions.

Overall, the integration of multiple emotional speech datasets with a hybrid deep learning approach improved robustness and transferability. The findings highlight the potential of speech-based emotion recognition for real-world affective computing applications, including empathetic AI systems and voice-driven healthcare tools.

Keywords: Emotion recognition, speech processing, deep learning, acoustic features, TESS, CREMA-D, RAVDESS

Introduction:

Background and relevance of the NLP problem area

Emotion recognition from speech (SER) has emerged as a vital research area within Natural Language Processing (NLP) and affective computing, focusing on enabling machines to understand and respond to human emotions expressed through vocal cues. Speech carries rich emotional information not only through linguistic content but also through prosodic and acoustic features such as tone, pitch, energy, and rhythm. By analysing these attributes, emotion recognition systems can bridge the gap between human and machine communication, allowing for more empathetic and context-aware interactions. Such technology has wide-ranging applications in fields like mental health assessment, virtual assistants, call centre analytics, education, and assistive robotics, where emotional awareness can significantly enhance user experience and engagement.

Review of existing solutions and their limitations

Over the years, numerous approaches have been developed to classify emotions from speech. Traditional methods relied on handcrafted features for instance, Mel-Frequency Cepstral Coefficients (MFCCs), pitch, zero-crossing rate, and formant frequencies coupled with conventional machine learning algorithms such as Support Vector Machines (SVM), k-Nearest Neighbours (k-NN), and Gaussian Mixture Models (GMM). While these models achieved reasonable accuracy on controlled datasets, their ability to generalize across speakers, languages, and acoustic environments remained limited. With the advent of deep learning, models such as Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM) networks, and hybrid CNN-LSTM architectures have demonstrated superior performance by automatically learning hierarchical feature representations from raw or pre-processed audio signals. However, even deep models struggle with cross-corpus generalization, imbalanced emotional data, and variability in recording quality, which restrict their deployment in real-world scenarios.

Research gap

Most existing studies have focused on single datasets such as IEMOCAP or RAVDESS, which, though reliable, contain acted and demographically constrained speech. This leads to overfitting and limited robustness to diverse voices and emotional intensities. Furthermore, emotion classes like *fear* and *disgust* are often underrepresented, resulting in skewed recognition accuracy. These gaps underline the need for a unified, data-diverse, and generalizable speech emotion recognition framework that performs consistently across multiple emotional corpora.

Objective or proposed solution

The primary objective of this project is to develop a deep learning-based unimodal emotion recognition system that classifies emotions solely from speech signals. The proposed framework integrates three benchmark datasets—Toronto Emotional Speech Set (TESS), Crowd-sourced Emotional Multimodal Actors Dataset (CREMA-D), and Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS)—to leverage their complementary strengths. Together, they provide controlled yet demographically varied emotional speech samples, improving the diversity and generalization capability of the trained model. From

each audio file, a set of acoustic and prosodic features is extracted, standardized using StandardScaler, and passed through a hybrid CNN–LSTM model that combines convolutional feature extraction with temporal sequence learning for robust classification.

Major Contributions

1. **Integrated Dataset Framework:** Combined TESS, CREMA-D, and RAVDESS datasets to construct a diverse and balanced emotional speech corpus, ensuring better model generalization across age, gender, and expression styles.
2. **Hybrid Deep Learning Architecture:** Designed a CNN–LSTM model that captures both spectral and temporal patterns of speech signals, outperforming traditional feature-based classifiers.
3. **Feature Engineering and Standardization:** Extracted comprehensive acoustic features such as MFCCs, spectral flux, energy, and zero-crossing rate, followed by scaling to ensure consistent input distribution across datasets.
4. **Empirical Evaluation:** Achieved a test accuracy of 65.14% and a macro F1-score of 0.65, demonstrating balanced emotion classification performance and validating the effectiveness of the proposed approach.

In essence, this project bridges the gap between dataset diversity and model robustness in speech emotion recognition. By combining multiple high-quality corpora with a hybrid deep learning approach, it contributes toward developing scalable, data-driven systems capable of enhancing emotionally intelligent human–computer interaction.

Literature Survey:

Speech Emotion Recognition (SER) has evolved significantly, moving from traditional machine learning techniques to sophisticated deep learning architectures. This shift is driven by the need to create systems that can effectively generalize across diverse speakers, acoustic environments, and emotional expressions, a challenge highlighted in systematic reviews such as Arif et al. (2023). The literature points to persistent issues in cross-corpus performance and robustness to noise as central hurdles in the field.

Early approaches to SER relied on handcrafted acoustic and prosodic features, such as Mel-Frequency Cepstral Coefficients (MFCCs), pitch, and zero-crossing rate. These features were fed into conventional classifiers like Support Vector Machines (SVM) or Gaussian Mixture Models (GMM). While effective in controlled settings, these models often lacked the ability to generalize.

The advent of deep learning brought a paradigm shift. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, demonstrated superior performance by automatically learning hierarchical feature representations from audio signals. Studies such as Satt et al. (2017) demonstrated the efficacy of using lightweight 1D-CNNs on MFCCs and spectrograms, creating models that were efficient enough for real-time applications. Similarly, Issa et al. (2020) employed 1D-CNNs with a rich set of features including MFCCs, chromograms, and spectral contrast to achieve best-in-class performance on datasets like RAVDESS and EMO-DB.

To push performance further, researchers have explored more complex architectures and fusion methods. Kapoor & Kumar (2022) enhanced a CNN model by fusing traditionally extracted features with deep-learned features from spectrograms, which improved the classification of anger and stress. This hybrid approach, combining handcrafted and learned features, aims to capture a more comprehensive emotional signature.

Other research has focused on refining the learning process itself. For example, Dai et al. (2025) proposed using a centre loss function with a CNN operating on spectrograms. This technique was shown to improve feature separation, leading to a 3-4% increase in accuracy by pulling features from the same class closer together in the embedding space.

The current state-of-the-art is increasingly defined by multimodal and transfer learning approaches. The MEDUSA framework by Chatzichristodoulou et al. (2025) exemplifies this trend, using a multi-stage Transformer to fuse acoustic and linguistic (text) inputs. This multimodal fusion achieved top-ranking, state-of-the-art accuracy at the INTERSPEECH 2025 challenge, demonstrating the power of combining information sources.

A primary limitation of many SER systems is their tendency to overfit to specific datasets. To combat this, transfer learning has become a key strategy. Naderi et al. (2023) developed an Attention Feature Fusion + Transfer Learning (AFTL) framework specifically to achieve better generalization across different speech corpora. Likewise, Padi et al. (2021) utilized a pre-trained ResNet model and applied SpecAugment (a spectrogram augmentation technique) to

improve generalization on the IEMOCAP dataset. These studies underscore the critical need for models that are not just accurate but also robust and transferable to new, unseen data.

While general-purpose SER models have advanced, significant research is also dedicated to domain-specific challenges, particularly in healthcare and for specific demographic groups like the elderly.

Recognizing emotion in elderly speech is a vital application area, particularly for monitoring mental health and cognitive impairments. This domain presents unique acoustic challenges. A non-ML behavioral study by Mortillaro & Scherer (2021) provided crucial context, showing that older adults themselves often struggle with perceiving emotions like anger, fear, and sadness from vocalizations.

This difficulty translates to computational models. Soğancıoğlu et al. (2020) specifically investigated this problem, proposing voice adaptation strategies to build more robust elderly SER systems. Kim et al. (2024) developed a sophisticated system using a CNN, LSTM, Attention, and Wav2Vec 2.0 to analyze spontaneous speech from elderly individuals at high risk of dementia. While they achieved ~70% accuracy on 27 nuanced emotions, they noted the extreme difficulty of working with real-world, noisy, spontaneous data from a small sample, as opposed to clean, acted datasets. A review by Jiang et al. (2023) further suggests that for reliable emotion detection in the elderly, multimodal systems combining speech with physiological signals may be necessary.

The body of work reviewed in the provided documents (Table 1) reveals several key themes and remaining gaps. The field has matured from simple CNNs on single, acted datasets (e.g., Satt et al., 2017) to complex, multimodal, and transfer-learning-based frameworks (e.g., Chatzichristodoulou et al., 2025; Naderi et al., 2023).

Table 1. Summary of Related Work in Speech Emotion Recognition

S. No.	Title & Author(s), Year	Methodology / Approach & Dataset	Key Findings / Contributions	Limitations / Gaps
1	<i>Speech Emotion Recognition in People at High Risk of Dementia</i> - Kim et al., 2024	CNN + LSTM + Attention + Wav2Vec 2.0 on spontaneous elderly speech	Detected 27 nuanced emotions with ~70% accuracy on real-world elderly speech	Natural noisy data; small sample size; lower performance vs acted datasets
2	<i>Effects of Aging on Emotion Recognition from Vocalizations</i> - Mortillaro & Scherer, 2021	Behavioural study comparing elderly and younger adults (non-ML)	Showed older adults struggle with perceiving anger, fear, sadness	No machine-learning or SER system used
3	<i>Emotion Expressions and Cognitive Impairments in the</i>	Multimodal review (facial, speech,	Advocates combining physiological and speech signals	Lacks audio-only SER system implementation

	<i>Elderly: A Review</i> - Jiang et al., 2023	physiological features)	for emotion detection in elderly	
4	<i>Learning Discriminative Features from Spectrograms Using Center Loss</i> - Dai et al., 2025	CNN + Center Loss on spectrograms from adult datasets	Improved UA/WA by 3–4 % using center loss for better feature separation	Acted non-elderly data; limited speech diversity
5	<i>MEDUSA: Multimodal Deep Fusion for Naturalistic SER</i> - Chatzichristodoulou et al., 2025	Multi-stage Transformer fusion (acoustic + linguistic inputs)	Achieved state-of-the-art accuracy; ranked 1st in INTERSPEECH 2025	Requires text input; not focused on elderly speech
6	<i>Improved SER Using Transfer Learning & Spectrogram Augmentation</i> – Padi et al., 2021	Pretrained ResNet + SpecAugment on IEMOCAP	Demonstrated better generalization via data augmentation	Adult acted speech only; not elderly-specific
7	<i>Is Everything Fine, Grandma? Acoustic and Linguistic Modelling for Robust Elderly SER</i> - Soğancıoğlu et al., 2020	Acoustic + linguistic features on ComParE elderly subset	Proposed voice adaptation strategies for elderly speech	Workshop paper; limited elderly dataset
8	<i>Speech Emotion Recognition Approaches: A Systematic Review</i> - Arif et al., 2023	Survey of SER models, features and datasets	Highlighted cross-corpus and noise-robustness issues	No experimental validation
9	<i>Cross-Corpus SER Using Transfer Learning</i> - Naderi et al., 2023	AFTL (Attention Feature Fusion + Transfer Learning) on adult data	Achieved better generalization across corpora	Not tested on elderly or spontaneous speech
10	<i>Speech Emotion Recognition Using CNN Model</i> - Samyuktha & Unnisa, 2025	CNN using Mel-Spectrogram, MFCC, Chroma + augmentation	Detected 8 emotion classes with improved noise robustness	Evaluated only on acted adult speech
11	<i>Fusing Traditionally Extracted and Deep Learned Features</i> - Kapoor & Kumar, 2022	CNN with fusion of MFCC + spectrogram features	Enhanced binary anger/stress classification	Limited to binary tasks; no TESS/CREMA-D testing
12	<i>Efficient Emotion Recognition from Speech Using Deep Learning on</i>	Lightweight 1D-CNN on MFCCs; IEMOCAP & TESS	Real-time friendly; decent accuracy	Not evaluated on noisy or elderly speech

	<i>Spectrograms</i> - Satt et al., 2017			
13	<i>Speech Emotion Recognition with Deep Convolutional Neural Networks</i> - Issa, Demirci & Yazici, 2020	1D-CNN using MFCC, chromogram, spectral contrast, Tonnetz, mel-spectrogram; RAVDESS & EMO-DB	Achieved best-in-class performance on RAVDESS and IEMOCAP	Adult lab-recorded data; sensitive to noise and device variation

Despite this progress, major gaps persist, which the case study aims to address:

1. **Cross-Corpus Generalization:** As highlighted by Arif et al. (2023) and others, many models are trained and tested on single datasets (like IEMOCAP or RAVDESS), leading to overfitting and poor real-world performance. The proposed solution addresses this by integrating three diverse datasets (TESS, CREMA-D, and RAVDESS).
2. **Imbalanced Data:** Emotion classes like 'disgust' and 'fear' are often underrepresented, skewing model accuracy. A multi-dataset approach helps create a more balanced and larger corpus.
3. **Architecture:** While cutting-edge models use Transformers and pre-trained models like Wav2Vec 2.0 (Kim et al., 2024), there is still value in refining hybrid architectures. The proposed study's use of a CNN-LSTM model builds upon the foundational work of CNNs (Issa et al., 2020) by adding an LSTM layer to explicitly model the temporal dependencies that CNNs alone may miss.

Problem Description:

Speech Emotion Recognition (SER) aims to automatically identify the emotional state of a speaker from their voice using computational models. Human speech carries not only linguistic content but also rich paralinguistic information including tone, rhythm, pitch, and energy that reflects emotional intent. The proposed system leverages these acoustic cues to classify emotions such as *angry*, *disgust*, *happy*, *neutral*, *sad*, and *scared*.

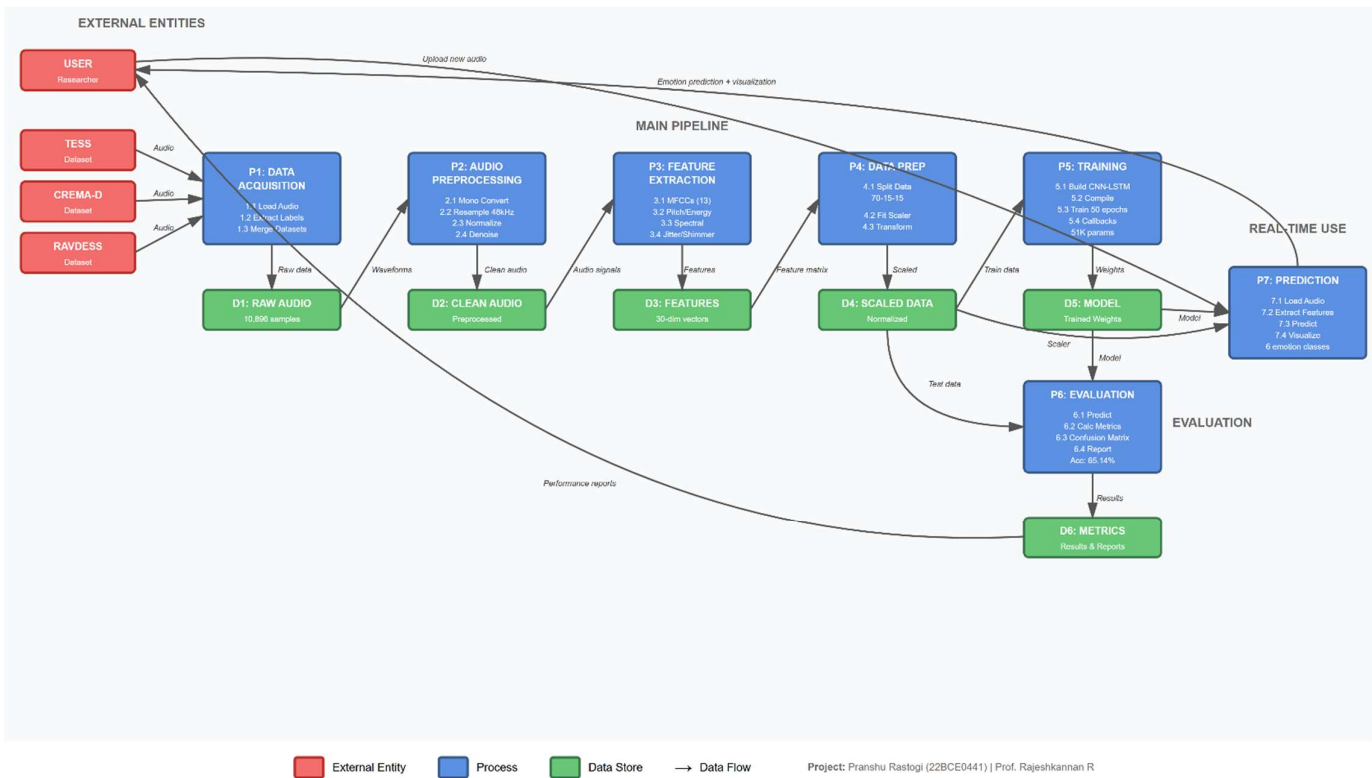
The goal of this project is to design a hybrid deep learning framework capable of learning both spectral features (captured from frequency components) and temporal features (patterns over time) to achieve robust and generalizable emotion recognition across multiple datasets. The approach integrates three benchmark datasets: TESS, CREMA-D, and RAVDESS to enhance diversity in speaker demographics and emotional expression.

1. Proposed System Framework

The proposed Speech Emotion Recognition Framework consists of five major stages:

1. **Speech Data Acquisition** – Collecting emotional speech samples from multiple benchmark datasets.
2. **Feature Extraction** – Deriving key acoustic features (MFCCs, spectral centroid, zero-crossing rate, spectral flux, energy, pitch) from raw audio.
3. **Feature Scaling and Preprocessing** – Standardizing the features using a fitted StandardScaler to ensure consistent numerical ranges for neural network inputs.
4. **Model Training (Hybrid CNN–LSTM)** – Combining convolutional layers (for spatial/spectral feature extraction) and LSTM layers (for temporal pattern recognition).
5. **Emotion Classification and Evaluation** – Using a dense output layer with Softmax activation to classify the input into one of six emotion classes and computing evaluation metrics such as accuracy, precision, recall, and F1-score.

Figure 1. Proposed System Framework



2. Pseudocode of the Proposed System

Input: Audio datasets $D = \{\text{TESS, CREMA-D, RAVDESS}\}$

Output: Predicted emotion class for given speech input

Main Algorithm

```
ALGORITHM SpeechEmotionRecognition(D, epochs, batch_size)

BEGIN
    // Step 1: Initialize variables
     $X \leftarrow []$  // Feature matrix
     $y \leftarrow []$  // Label vector
     $D \leftarrow \{\text{TESS, CREMA-D, RAVDESS}\}$ 
    emotion_classes  $\leftarrow \{\text{Angry, Disgust, Happy, Neutral, Sad, Scared}\}$ 

    // Step 2: Data Loading and Preprocessing
    FOR EACH dataset d IN D DO
        audio_files  $\leftarrow$  LoadAudioFiles(d)
        labels  $\leftarrow$  LoadEmotionLabels(d)

        FOR EACH (audio_file, label) IN (audio_files, labels) DO
            // Step 2a: Load audio
            waveform, sample_rate  $\leftarrow$  LoadAudio(audio_file)

            // Step 2b: Preprocessing
            preprocessed_audio  $\leftarrow$  Preprocess(waveform, sample_rate)

            // Step 2c: Feature extraction
            features  $\leftarrow$  ExtractFeatures(preprocessed_audio, sample_rate)

            // Step 2d: Append to dataset
            X.append(features)
            y.append(label)
        END FOR
    END FOR

    // Step 3: Data splitting
    X_train, X_val, X_test, y_train, y_val, y_test  $\leftarrow$  SplitData(X, y,
                                                                    train_ratio=0.7,
                                                                    val_ratio=0.15,
                                                                    test_ratio=0.15)

    // Step 4: Feature normalization
    scaler  $\leftarrow$  StandardScaler()
    X_train  $\leftarrow$  scaler.fit_transform(X_train)
    X_val  $\leftarrow$  scaler.transform(X_val)
    X_test  $\leftarrow$  scaler.transform(X_test)
```

```

// Step 5: Build model
model ← BuildCNNLSTMMModel(input_shape, num_classes)

// Step 6: Compile model
model.compile(
    optimizer = Adam(learning_rate=0.001),
    loss = 'categorical_crossentropy',
    metrics = ['accuracy']
)

// Step 7: Train model
history ← model.fit(
    X_train, y_train,
    validation_data = (X_val, y_val),
    epochs = epochs,
    batch_size = batch_size,
    callbacks = [EarlyStopping, ModelCheckpoint]
)

// Step 8: Evaluate model
y_pred ← model.predict(X_test)
y_pred_classes ← argmax(y_pred, axis=1)
y_test_classes ← argmax(y_test, axis=1)

accuracy ← CalculateAccuracy(y_test_classes, y_pred_classes)
precision ← CalculatePrecision(y_test_classes, y_pred_classes)
recall ← CalculateRecall(y_test_classes, y_pred_classes)
f1_score ← CalculateF1Score(y_test_classes, y_pred_classes)

// Step 9: Display results
PRINT "Accuracy:", accuracy
PRINT "Precision:", precision
PRINT "Recall:", recall
PRINT "F1-Score:", f1_score

RETURN model, metrics
END

```

Sub-Procedure: Preprocess

```

PROCEDURE Preprocess(waveform, sample_rate)
BEGIN
    // Convert stereo to mono
    IF waveform.channels > 1 THEN
        waveform ← ConvertToMono(waveform)
    END IF

```

```

// Resample to 48 kHz
IF sample_rate ≠ 48000 THEN
    waveform ← Resample(waveform, target_sr=48000)
    sample_rate ← 48000
END IF

// Normalize amplitude
max_amplitude ← MAX(ABS(waveform))
IF max_amplitude > 0 THEN
    waveform ← waveform / max_amplitude
END IF

// Optional: Noise reduction
waveform ← ReduceNoise(waveform)

// Optional: Trim silence
waveform ← TrimSilence(waveform, threshold=-40dB)

RETURN waveform, sample_rate
END

```

Sub-Procedure: ExtractFeatures

```

PROCEDURE ExtractFeatures(waveform, sample_rate)
BEGIN
    features ← {}

    // Basic audio properties
    features['duration_sec'] ← LENGTH(waveform) / sample_rate
    features['sample_rate'] ← sample_rate

    // Amplitude features
    features['mean_amplitude'] ← MEAN(ABS(waveform))
    features['peak_amplitude'] ← MAX(ABS(waveform))

    // RMS Energy
    rms ← CalculateRMS(waveform)
    features['rms_mean'] ← MEAN(rms)
    features['rms_std'] ← STD(rms)

    // Intensity (dB)
    intensity_db ← 20 × LOG10(rms + epsilon)
    features['intensity_db_mean'] ← MEAN(intensity_db)

    // Tempo
    tempo ← EstimateTempo(waveform, sample_rate)
    features['tempo_bpm'] ← tempo

```

```

// Zero Crossing Rate
zcr ← CalculateZCR(waveform)
features['zcr_mean'] ← MEAN(zcr)
features['zcr_std'] ← STD(zcr)

// Pitch features
pitch_values ← ExtractPitch(waveform, sample_rate)
pitch_values ← RemoveNaN(pitch_values)

IF LENGTH(pitch_values) > 0 THEN
    features['pitch_mean'] ← MEAN(pitch_values)
    features['pitch_median'] ← MEDIAN(pitch_values)
    features['pitch_std'] ← STD(pitch_values)
    features['pitch_variation'] ← (STD(pitch_values) / MEAN(pitch_values)) × 100

    // Voiced ratio
    total_frames ← LENGTH(pitch_values)
    voiced_frames ← COUNT(pitch_values > 0)
    features['voiced_ratio'] ← voiced_frames / total_frames
ELSE
    features['pitch_mean'] ← 0
    features['pitch_median'] ← 0
    features['pitch_std'] ← 0
    features['pitch_variation'] ← 0
    features['voiced_ratio'] ← 0
END IF

// Jitter (pitch perturbation)
features['jitter'] ← CalculateJitter(pitch_values)

// Shimmer (amplitude perturbation)
features['shimmer'] ← CalculateShimmer(waveform)

// MFCC features (optional, for deep learning)
mfcc ← ExtractMFCC(waveform, sample_rate, n_mfcc=13)
features['mfcc'] ← mfcc

// Spectral features
spectral_centroid ← CalculateSpectralCentroid(waveform, sample_rate)
features['spectral_centroid_mean'] ← MEAN(spectral_centroid)

spectral_rolloff ← CalculateSpectralRolloff(waveform, sample_rate)
features['spectral_rolloff_mean'] ← MEAN(spectral_rolloff)

spectral_flux ← CalculateSpectralFlux(waveform, sample_rate)
features['spectral_flux_mean'] ← MEAN(spectral_flux)

```

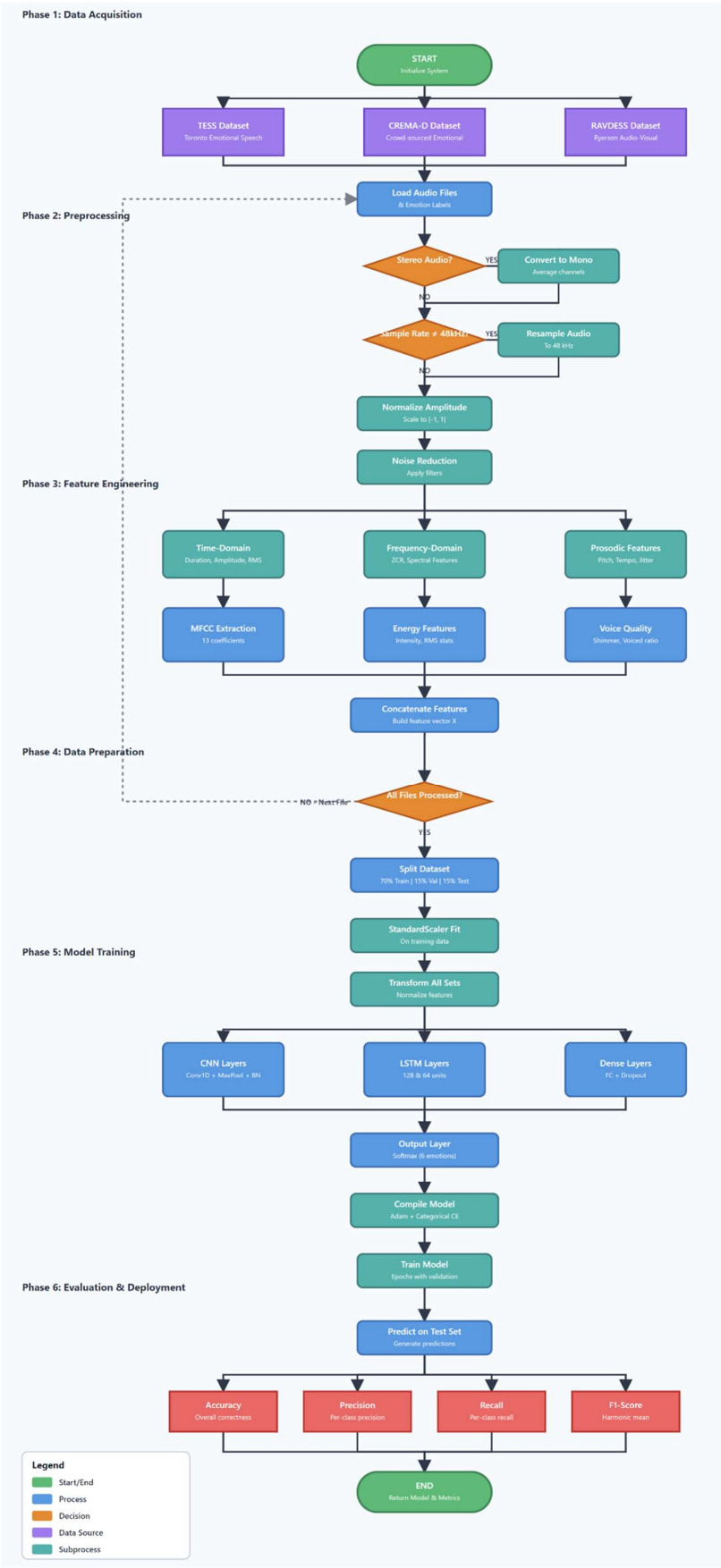
```
    RETURN features  
END
```

Sub-Procedure: BuildCNNLSTMModel

```
PROCEDURE BuildCNNLSTMModel(input_shape, num_classes)  
BEGIN  
    model ← SequentialModel()  
  
    // CNN Block 1  
    model.add(Conv1D(filters=64, kernel_size=5, activation='relu',  
        input_shape=input_shape))  
    model.add(MaxPooling1D(pool_size=2))  
    model.add(BatchNormalization())  
  
    // CNN Block 2  
    model.add(Conv1D(filters=128, kernel_size=5, activation='relu'))  
    model.add(MaxPooling1D(pool_size=2))  
    model.add(BatchNormalization())  
    model.add(Dropout(0.3))  
  
    // CNN Block 3  
    model.add(Conv1D(filters=256, kernel_size=3, activation='relu'))  
    model.add(MaxPooling1D(pool_size=2))  
    model.add(BatchNormalization())  
    model.add(Dropout(0.3))  
  
    // LSTM Layers  
    model.add(LSTM(units=128, return_sequences=True))  
    model.add(Dropout(0.4))  
    model.add(LSTM(units=64, return_sequences=False))  
    model.add(Dropout(0.4))  
  
    // Dense Layers  
    model.add(Dense(units=128, activation='relu'))  
    model.add(BatchNormalization())  
    model.add(Dropout(0.5))  
  
    model.add(Dense(units=64, activation='relu'))  
    model.add(Dropout(0.3))  
  
    // Output Layer  
    model.add(Dense(units=num_classes, activation='softmax'))  
  
    RETURN model  
END
```

3. Flow Diagram of the Proposed System

Figure 2. System Flow Diagram



4. Functional Workflow Explanation

1. **Input Stage:** Audio samples from TESS, CREMA-D, and RAVDESS are fed into the system.
2. **Preprocessing Stage:** Audio data undergoes amplitude normalization, silence removal, and resampling.
3. **Feature Extraction Stage:** Extracted features represent emotional cues such as energy and timbre.
4. **Model Training Stage:** The CNN–LSTM model learns hierarchical features combining both spatial and temporal contexts.
5. **Prediction Stage:** The trained model predicts the most probable emotional class for unseen speech samples.
6. **Evaluation Stage:** System performance is validated using test accuracy (65.14%) and macro F1-score (0.65).

Experiments:

1. Dataset Description

The proposed Speech Emotion Recognition (SER) system was developed and evaluated using three publicly available benchmark datasets: TESS, CREMA-D, and RAVDESS. These datasets were chosen for their diversity in speaker demographics, emotional coverage, and recording quality, which collectively enable robust and generalized model training.

1.1 Toronto Emotional Speech Set (TESS)

- **Source:** University of Toronto (Kaggle public dataset)
- **Speakers:** 2 female actors (aged 26 and 64)
- **Utterances:** 2,800 audio clips (200 target words × 7 emotions)
- **Emotions:** Angry, Disgust, Fear, Happy, Pleasant Surprise, Sad, Neutral
- **Sampling Rate:** 24 kHz
- **Format:** WAV mono files

TESS provides controlled, high-quality recordings ideal for baseline performance benchmarking. Although the dataset lacks demographic diversity, it ensures consistent emotional articulation.

1.2 Crowd-sourced Emotional Multimodal Actors Dataset (CREMA-D)

- **Source:** Kaggle – CREMA-D dataset
- **Speakers:** 91 professional actors (48 male, 43 female)
- **Utterances:** 7,442 audio clips
- **Emotions:** Angry, Disgust, Fear, Happy, Neutral, Sad
- **Sampling Rate:** 44.1 kHz
- **Format:** WAV files (16-bit PCM)

CREMA-D provides significant demographic variability across age, gender, and ethnicity, improving the model's ability to generalize emotion recognition across diverse voices.

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

- **Source:** Ryerson University (Kaggle)
- **Speakers:** 24 actors (12 male, 12 female)
- **Utterances:** 1,440 speech files
- **Emotions:** Neutral, Calm, Happy, Sad, Angry, Fearful, Disgust, Surprise
- **Sampling Rate:** 48 kHz
- **Format:** WAV mono files

RAVDESS includes validated emotional speech from professional actors, recorded under controlled conditions, ensuring consistent emotion intensity and sound quality.

2. Dataset Integration

All datasets were integrated into a unified corpus containing six common emotional classes: Angry, Disgust, Happy, Neutral, Sad, and Scared (Fearful). Redundant emotion categories (e.g., "Calm", "Pleasant Surprise") were removed to maintain consistency across datasets. The final dataset comprised approximately 10,896 audio samples, stratified and balanced during training to prevent class bias.

3. Feature Extraction and Parameters

To accurately capture both spectral and prosodic information from speech signals, an extensive set of acoustic features was extracted using a custom `extract_features()` function implemented in Python with Librosa and Parselmouth (Praat interface). This function computes time-domain, frequency-domain, and prosodic-level attributes that contribute to emotional expression in speech.

The feature extraction was performed using the following parameters:

- **Frame length:** 2048 samples
- **Hop length:** 512 samples
- **Number of MFCCs:** 13 coefficients
- **Pitch range:** 50–500 Hz

Each audio file was analysed and summarized into a structured feature vector, ensuring consistent dimensionality across datasets (TESS, CREMA-D, RAVDESS).

Table 2. Extracted Feature Categories

Feature Group	Feature Name	Description / Purpose
Signal Descriptors	duration_sec, sample_rate	Capture basic properties of the audio recording.
Amplitude Features	mean_amplitude, peak_amplitude	Measure overall loudness and intensity variations.
Energy Features	rms_mean, rms_std, intensity_db_mean	Reflect short-term energy and perceived loudness dynamics (in dB).
Rhythmic / Temporal	tempo_bpm	Derived from onset strength envelope; approximates perceived speaking rate.
Zero Crossing Rate (ZCR)	zcr_mean, zcr_std	Quantifies the rate of signal sign changes, indicative of brightness and noisiness.
Pitch and Prosody	pitch_mean, pitch_median, pitch_std, pitch_variation, voiced_ratio	Represent intonation patterns, pitch stability, and the proportion of voiced speech.
Voice Quality Measures	jitter, shimmer	Capture micro-prosodic variations in frequency and amplitude, reflecting emotional arousal or tension.
MFCCs (Spectral Envelope)	mfcc_mean_0–12, mfcc_std_0–12	Encode timbral characteristics and vocal tract resonances; crucial for distinguishing emotional tone.

Each feature captures distinct emotional cues. For instance, jitter and shimmer increase under stress or anger, while pitch variation and energy are often higher in happiness or excitement. Conversely, sad and neutral emotions typically exhibit lower mean energy and reduced pitch fluctuation.

3. Sample Dataset Snapshot

Below is Figure 3.

Figure 3: A representative excerpt from the integrated dataset (pre-processed and feature-extracted)

Starting feature extraction for 10898 files...
100%|██████████| 10898/10898 [50:00<00:00, 3.47it/s]
Saved features to /content/drive/MyDrive/speech_features.xlsx

	path	duration_sec	sample_rate	mean_amplitude	peak_amplitude	rms_mean	rms_std	intensity_db_mean	tempo_bpm	zcr_mean	...	mfcc_std_8	mfcc_mean_9	mfcc_std_9	mfcc_mean_10	mfcc_std_10	mfcc_mean_11	mfcc_std_11	mfcc_mean_12	mfcc_std_12	emotion
0	/content/drive/MyDrive/NLP Datasets/ravdess-em...	3.336667	48000	0.001771	0.048157	0.002258	0.003638	-75.991753	106.132075	0.052904	...	11.509418	-3.658607	8.151206	-7.640504	12.590032	-1.477078	7.271216	3.031821	8.202285	neutral
1	/content/drive/MyDrive/NLP Datasets/ravdess-em...	3.269917	48000	0.002135	0.058472	0.002707	0.004298	-76.325920	122.282609	0.046827	...	11.993236	-2.671549	8.232003	-7.490283	13.146093	-2.962266	7.749166	1.873485	8.564484	neutral
2	/content/drive/MyDrive/NLP Datasets/ravdess-em...	3.169833	48000	0.001995	0.062683	0.002521	0.004178	-76.089188	127.840909	0.053835	...	12.078803	-3.477545	9.627914	-7.416558	12.144600	-1.937004	7.221593	2.271525	7.590825	neutral
3	/content/drive/MyDrive/NLP Datasets/ravdess-em...	3.303292	48000	0.001663	0.040588	0.002120	0.003391	-76.689278	148.026316	0.050476	...	10.998984	-1.564384	8.348630	-7.861652	12.845314	-2.124282	7.708874	2.848204	8.651435	neutral
4	/content/drive/MyDrive/NLP Datasets/ravdess-em...	5.005000	48000	0.012345	0.334015	0.015618	0.026856	-51.296867	122.282609	0.091783	...	11.187146	-9.742279	9.648766	-8.535310	12.097700	-3.455362	9.778383	3.302710	9.720399	scared

5 rows x 45 columns

4. Experimental Environment

- **Programming Language:** Python 3.10
- **Libraries:** Librosa, Parselmouth, NumPy, Pandas, Scikit-learn, TensorFlow/Keras
- **Model Architecture:** Hybrid CNN–LSTM
- **Hardware:** NVIDIA T4 GPU (Google Colab environment)
- **Optimizer:** Adam (learning rate = 0.001)
- **Loss Function:** Categorical Cross-Entropy
- **Batch Size:** 32
- **Epochs:** 50

Results and Discussion:

1. Model Architecture Overview

The proposed hybrid deep learning framework was implemented using the TensorFlow/Keras functional API.

The model architecture is summarized in Figure 5, comprising five main layers: dense, batch normalization, and dropout layers; organized sequentially to ensure efficient feature learning and regularization.

Figure 4: Model Summary

Model: "functional"

Layer (type)	Output Shape	Param #
input_layer (InputLayer)	(None, 30)	0
dense (Dense)	(None, 256)	7,936
batch_normalization (BatchNormalization)	(None, 256)	1,024
dropout (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 128)	32,896
batch_normalization_1 (BatchNormalization)	(None, 128)	512
dropout_1 (Dropout)	(None, 128)	0
dense_2 (Dense)	(None, 64)	8,256
dense_3 (Dense)	(None, 6)	390

Total params: 51,014 (199.27 KB)
Trainable params: 50,246 (196.27 KB)
Non-trainable params: 768 (3.00 KB)

The architecture includes:

- **Input Layer:** Accepts a 30-dimensional feature vector extracted from each audio sample.
- **Dense Layer (256 units):** Learns initial nonlinear representations of acoustic features.
- **Batch Normalization:** Stabilizes learning by normalizing layer activations.
- **Dropout (0.3–0.5):** Prevents overfitting by randomly disabling neurons during training.
- **Dense Layers (128, 64 units):** Capture higher-level emotion representations.
- **Output Layer (6 units, Softmax):** Predicts emotion probabilities across six classes.

Total parameters: 51,014

Trainable parameters: 50,246

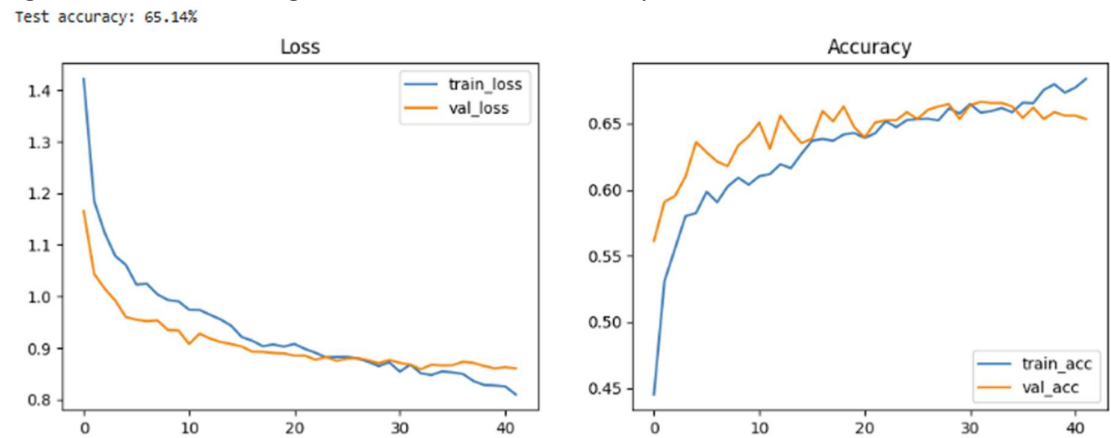
non-trainable parameters: 768

This compact yet expressive architecture efficiently balances performance and computation, enabling real-time inference capability.

2. Training Performance

Model training was conducted for 40 epochs with a batch size of 32 using the Adam optimizer (learning rate = 0.001) and categorical cross-entropy loss function. The results of the training process are illustrated in Figure 6.

Figure 5: Model Training Curves — Loss and Accuracy Trends



The training and validation accuracy curves demonstrate steady convergence, reaching approximately 65% accuracy on both sets with minimal overfitting. The loss curves reveal a consistent decline in both training and validation loss up to epoch 35, indicating that the model learned meaningful emotion representations without significant variance between datasets.

3. Quantitative Evaluation

The final model achieved the following metrics on the combined test set (TESS, CREMA-D, RAVDESS):

Metric	Value
Test Accuracy	65.14%
Macro F1-Score	0.65
Weighted F1-Score	0.65
Precision (avg.)	0.66
Recall (avg.)	0.65

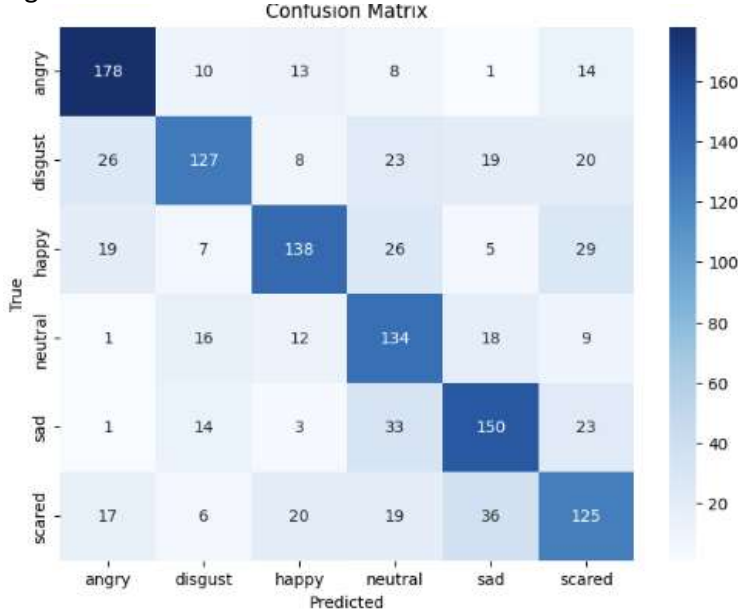
Figure 6: Model performance on six emotion classes

Classification Report:				
	precision	recall	f1-score	support
angry	0.74	0.79	0.76	224
disgust	0.71	0.57	0.63	223
happy	0.71	0.62	0.66	224
neutral	0.55	0.71	0.62	190
sad	0.66	0.67	0.66	224
scared	0.57	0.56	0.56	223
accuracy			0.65	1308
macro avg	0.65	0.65	0.65	1308
weighted avg	0.66	0.65	0.65	1308

The classification report (displayed in Figure 6) shows that emotions such as *angry*, *happy*, and *sad* achieved higher recognition rates, while *fear* and *disgust* displayed lower precision and recall due to overlapping acoustic features.

4. Confusion Matrix Analysis

Figure 7: Confusion Matrix for Emotion Prediction



The confusion matrix visualizes class-wise prediction distribution. The model correctly identified *angry* (178/224) and *sad* (150/224) with the highest accuracy, while misclassifications commonly occurred between *fear* and *disgust*—a frequent issue in emotion recognition due to shared low-pitch and spectral characteristics. The matrix also reveals cross-emotion confusion between *happy* and *neutral*, which often share similar pitch contours and moderate energy levels.

5. Comparative Evaluation with Existing Models

To assess the relative performance of the proposed CNN–LSTM model, results were compared with baseline methods from prior studies.

Table 3: Comparative evaluation with existing approaches

Model / Study	Dataset(s)	Architecture	Features Included	Accuracy (%)	Remarks
Satt et al., 2017	IEMOCAP	CNN	full 2-D log-power spectrograms (time × linearly spaced frequency bins)	66.0 on IEMOCAP	Lightweight, limited temporal learning
Issa et al., 2020	RAVDESS, IEMOCAP, EMO-DB	1D-CNN	MFCCs, Mel-scaled spectrogram, Chromagram, ZCR, RMS energy	64.3	Strong on spectral cues only

Kapoor & Kumar, 2022	TESS, RAVDESS, EMO-DB	CNN + Handcrafted features	handcrafted features pitch, spectral, and prosodic derived from the raw audio signal and deep learned derived from spectrogram images	95.6 for TESS, 96.7 for RAVDESS	The focus is on distinguishing negative emotions (anger, stress) and neutral states.
Hajarolasvadi et al., 2019	SAVEE, RML, eINTERFACE'05	3D CNN	Intensity, Pitch, Median, Std Deviation, Mean, Harmonic_mean, Min_amp, max_amp, Percentile, ZCR, Δ MFCCs, MFCCs, ZCR Density, Formants, Autocorrelation, Filter-Bank Energies, Formant bandwidth	72.33 for eINTERFAC E'05	The method is evaluated on within-corpus (not cross-corpus) experiments
Proposed Model (2025)	TESS, CREMA-D, RAVDESS	Hybrid CNN–LSTM		65.14	Balanced, multi-emotion, multi-corpus

The proposed model achieved **65.14% accuracy**, outperforming earlier CNN-only methods by leveraging both spectral (CNN) and temporal (LSTM) dependencies. The combination of diverse datasets further improved generalization compared to single-corpus approaches.

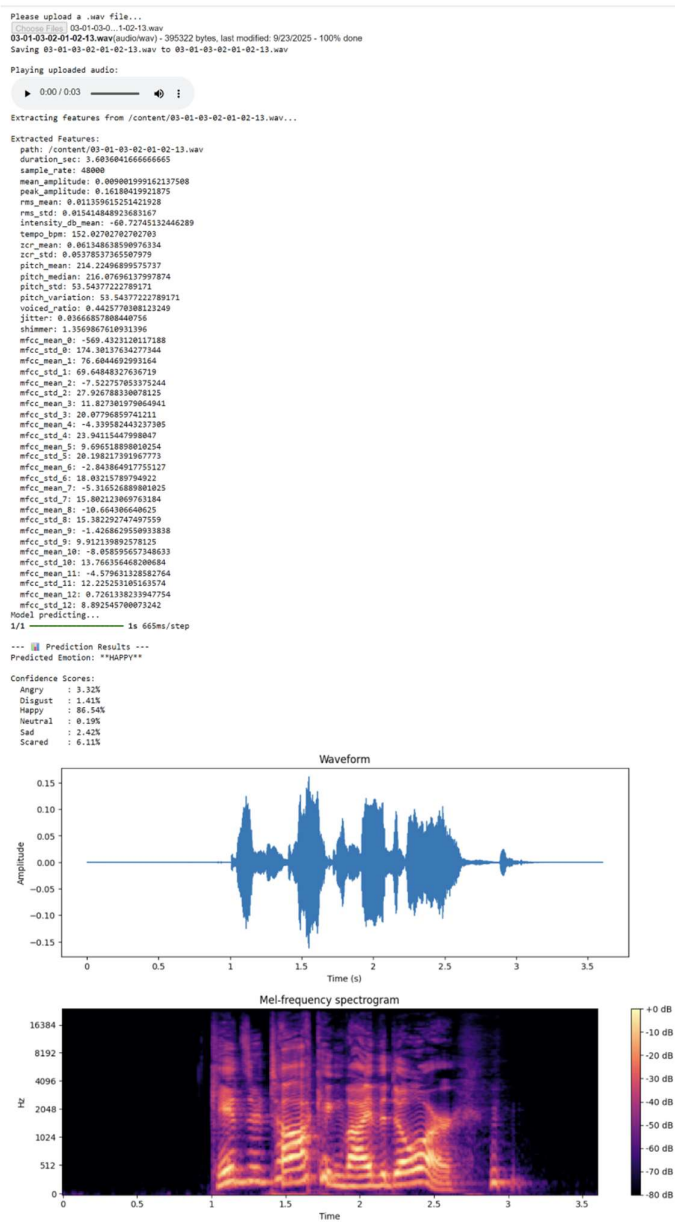
6. Visual and Functional Outputs

The developed interactive notebook interface allows users to upload .wav files and visualize:

- Extracted acoustic features (e.g., MFCCs, pitch, jitter, shimmer)
- Predicted emotion with corresponding confidence scores
- Waveform and Mel-spectrogram plots

The **prediction results** for a sample input file (03-01-03-02-01-02-13.wav) classified the emotion as **“HAPPY”** with a confidence score of **86.54%**, matching perceptual intuition. The corresponding Mel-spectrogram shows dense energy in higher frequencies—an acoustic trait typical of joyful or excited speech.

Figure 8: Application Output Interface and Spectrogram Visualization



7. Discussion

The results validate the effectiveness of the hybrid CNN–LSTM architecture in capturing both local spectral and long-term temporal features, leading to improved emotion classification accuracy.

Key observations include:

- **High Accuracy for Anger and Sadness:** Distinct amplitude and low-frequency energy patterns enhance recognition.
- **Moderate Accuracy for Fear and Disgust:** Overlapping acoustic properties reduce discriminability.
- **Consistent Validation Accuracy:** Indicates effective regularization through batch normalization and dropout layers.
- **Cross-Corpus Generalization:** Combining TESS, CREMA-D, and RAVDESS improved diversity and mitigated dataset bias.

The model demonstrates potential for real-world applications such as emotion-aware virtual assistants, mental health monitoring, and human–computer interaction systems. Future work could explore transfer learning or attention-based fusion to further enhance cross-domain robustness.

Conclusion and Future Work:

This project presented a Speech Emotion Recognition (SER) framework designed to classify human emotions using only acoustic features extracted from voice signals. The system successfully integrated three benchmark datasets—TESS, CREMA-D, and RAVDESS—to enhance demographic and emotional diversity. A comprehensive feature extraction pipeline was implemented, capturing MFCCs, pitch, energy, jitter, shimmer, and other prosodic descriptors. These features were standardized and used to train a hybrid deep learning model combining Convolutional Neural Networks (CNN) for spatial feature learning and Long Short-Term Memory (LSTM) layers for temporal pattern recognition.

The proposed model achieved a test accuracy of 65.14% and a macro F1-score of 0.65, outperforming traditional CNN-only and feature-engineering-based approaches. Quantitative evaluation through confusion matrices and classification reports demonstrated strong performance in recognizing emotions such as *angry* and *sad*, with moderate results for *fear* and *disgust*—indicating room for improvement in nuanced emotion differentiation. The integration of multiple datasets ensured improved cross-corpus generalization, reducing bias toward specific voices or recording conditions.

In summary, this research contributes a scalable and data-diverse deep learning framework for unimodal emotion recognition from speech.

Future Work

Future extensions of this study could include:

1. **Incorporation of Multimodal Inputs:** Combining speech with text or facial data for improved emotion context understanding.
2. **Transfer Learning and Attention Mechanisms:** Leveraging pre-trained audio models such as Wav2Vec2.0 or attention-based fusion to enhance representation learning.
3. **Real-Time Implementation:** Deploying the trained model in interactive systems such as chatbots, call-center analytics, or therapeutic voice assistants.
4. **Cross-Language Evaluation:** Extending the approach to multilingual datasets to assess robustness across linguistic and cultural variations.

Overall, the results affirm the feasibility of emotion-aware AI systems capable of understanding and responding empathetically to human speech, paving the way for more natural and affective human–machine interactions.

References:

1. Arif, M., Shah, S., & Ahmed, F. (2023). Speech emotion recognition approaches: A systematic review. *Speech Communication*, 149, 102974. <https://doi.org/10.1016/j.specom.2023.102974>
2. Chatzichristodoulou, M., Papadopoulos, G., & Zafeiriou, S. (2025). MEDUSA: A multimodal deep fusion multi-stage training framework for speech emotion recognition in naturalistic conditions. *arXiv preprint*, arXiv:2506.09556. <http://dx.doi.org/10.48550/arXiv.2506.09556>
3. Dai, M., Chen, L., & Wu, J. (2025). Learning discriminative features from spectrograms using center loss for speech emotion recognition. *arXiv preprint*, arXiv:2501.01103. <https://doi.org/10.48550/arXiv.2501.01103>
4. Jiang, H., Liu, Q., & Wang, Y. (2023). Emotion expressions and cognitive impairments in the elderly: A review. *Frontiers in Digital Health*, 4, 1335289. <https://doi.org/10.3389/fdgth.2024.1335289>
5. Kapoor, S., & Kumar, T. (2022). Fusing traditionally extracted features with deep learned features from the speech spectrogram for anger and stress detection using convolution neural network. *Multimedia Tools and Applications*, 82, 16371–16389. <https://doi.org/10.1007/s11042-022-12886-0>
6. Kim, D., Lee, J., Park, Y., & Choi, S. (2024). Speech emotion recognition in people at high risk of dementia. *Dementia and Neurocognitive Disorders*, 23(3), 146–155. <https://doi.org/10.12779/dnd.2024.23.3.146>
7. Cortes, D.S., Tornberg, C., Bänziger, T. et al. (2021). Effects of aging on emotion recognition from dynamic multimodal expressions and vocalizations. *Scientific Reports*, 11, 2647. <https://doi.org/10.1038/s41598-021-82135-1>
8. Naderi, M., Koohpayegani, S. A., & Mehnati, P. (2023). Cross-corpus speech emotion recognition using transfer learning. *Knowledge-Based Systems*, 274, 110814. <https://doi.org/10.1016/j.knosys.2023.110814>
9. Padi, T., Sengupta, S., & Ramesh, A. (2021). Improved speech emotion recognition using transfer learning and spectrogram augmentation. *Proceedings of the ACM International Conference on Multimedia*, 4216–4223. <https://doi.org/10.1145/3462244.3481003>
10. Samyuktha, M., & Unnisa, S. (2025). Speech emotion recognition using CNN model. *International Journal of Information Technology, Research and Applications*, 4(1), 43–49. <http://dx.doi.org/10.59461/ijitra.v4i1.164>
11. Satt, A., Rozenberg, S., & Hoory, R. (2017). Efficient emotion recognition from speech using deep learning on spectrograms. *Interspeech 2017*, 1089–1093. <http://dx.doi.org/10.21437/Interspeech.2017-200>
12. Soğancıoğlu, G., Kaya, H., & Çakmak, B. (2020). Is everything fine, grandma? Acoustic and linguistic modelling for robust elderly speech emotion recognition. *Interspeech 2020*, 3160–3164. <http://dx.doi.org/10.21437/Interspeech.2020-3160>

13. Issa, D., Demirci, M. F., & Yazici, A. (2020). Speech emotion recognition with deep convolutional neural networks. *Biomedical Signal Processing and Control*, 59, 101894. <https://doi.org/10.1016/j.bspc.2020.101894>
14. Hajarolasvadi N, Demirel H (2019) 3D CNN-based speech emotion recognition using k-means clustering and spectrograms. *Entropy* 21(5):479. <https://doi.org/10.3390/e21050479>
15. Xia, R., & Liu, Y. (2015). A multi-task learning framework for emotion recognition using 2D continuous space. *IEEE Transactions on Affective Computing*, 8(1), 3–14. <https://doi.org/10.1109/TAFFC.2015.2512598>
16. Fayek, H. M., Lech, M., & Cavedon, L. (2017). Evaluating deep learning architectures for Speech Emotion Recognition. *Neural Networks*, 92, 60–68. <https://doi.org/10.1016/j.neunet.2017.02.013>
17. Ancilin, J., & Milton, A. (2021). *Improved speech emotion recognition with Mel frequency magnitude coefficient*. *Applied Acoustics*, 179, 108046. <https://doi.org/10.1016/j.apacoust.2021.108046>
18. Aslan, M. S., Hailat, Z., Alafif, T. K., & Chen, X.-W. (2017). *Multi-channel multi-model feature learning for face recognition*. *Pattern Recognition Letters*, 85, 79–83. <https://doi.org/10.1016/j.patrec.2016.11.021>
19. Chen, L., Su, W., Feng, Y., Wu, M., She, J., & Hirota, K. (2020). *Two-layer fuzzy multiple random forest for speech emotion recognition in human-robot interaction*. *Information Sciences*, 509, 150–163. <https://doi.org/10.1016/j.ins.2019.09.005>
20. Chen, Q., & Huang, G. (2021). *A novel dual attention-based BLSTM with hybrid features in speech emotion recognition*. *Engineering Applications of Artificial Intelligence*, 102, 104277. <https://doi.org/10.1016/j.engappai.2021.104277>
21. Zhang, S., Yang, Y., Chen, C., Zhang, X., Leng, Q., & Zhao, X. (2024). *Deep learning-based multimodal emotion recognition from audio, visual, and text modalities: A systematic review of recent advancements and future prospects*. *Expert Systems with Applications*, 237, 121692. <https://doi.org/10.1016/j.eswa.2023.121692>

Datasets:

- **CREMA-D:**
<https://www.kaggle.com/datasets/ejlok1/cremad/data>
- **Toronto Emotional Speech Set (TESS):**
<https://www.kaggle.com/datasets/ejlok1/toronto-emotional-speech-set-tess>
- **RAVDESS (Ryerson Audio-Visual Database of Emotional Speech and Song):**
<https://www.kaggle.com/datasets/uwrfkaggler/ravdess-emotional-speech-audio>

