# Telugu Accent Classifier

A TERM PROJECT REPORT SUBMITTED IN PARTIAL FULFILMENT OF THE
REQUIREMENTS FOR THE AWARD OF THE DEGREE OF
**BACHELOR OF TECHNOLOGY IN COMPUTER SCIENCE AND
ENGINEERING**



By
**Batch-A12**

| | |
|---|---|
| **T. Lavanya** | **K. Niharika** |
| (21JG1A0560) | (21JG1A0543) |
| **G. Aakruthi** | **D. Harshini** |
| (21JG1A0527) | (22JG5A0503) |

**Under the esteemed guidance of**

**Mrs. D. Indu**

(Asst. Professor)

CSE Department

**Department of Computer Science and Engineering**
**GAYATRI VIDYA PARISHAD COLLEGE OF ENGINEERING FOR WOMEN**
**(Affiliated to Jawaharlal Nehru Technological University, Kakinada)**
**2021-25**

# GAYATRI VIDYA PARISHAD COLLEGE OF ENGINEERING FOR WOMEN
## Department of Computer Science and Engineering



## CERTIFICATE

This is to certify that the project report titled **"Telugu Accent Classifier"** is a bonafide work of following III B.Tech students in the Department of Computer Science and Engineering, Gayatri Vidya Parishad College of Engineering for Women affiliated to JNT University, Kakinada during the academic year 2023-24, in partial fulfilment of the requirement for the award of the degree of Bachelor of Technology of this university.

**Ms T. Lavanya** (21JG1A0560)          **Ms K. Niharika** (21JG1A0543)

**Ms G. Aakruthi** (21JG1A0527)          **Ms D. Harshini** (22JG5A0503)

Internal Guide                                              Head of the Department

**Mrs. D. Indu**                                           **Dr. P. V. S. L. Jagadamba**

Asst. Professor                                            Professor

Dept. of CSE                                               Dept. of CSE

# ACKNOWLEDGEMENT

The satisfaction that accompanies the successful completion of any task would be incomplete without the mention of people who made it possible and whose constant guidance and encouragement crown all the efforts with success.

We feel elated to extend our sincere gratitude to our guide **Mrs. D. Indu**, Assistant Professor for encouragement all the way during analysis of the project. Her annotations, insinuations and criticisms are the key behind the successful completion of the thesis and for providing us all the required facilities.

We would like to take this opportunity to extend our gratitude to Project Coordinator, **Dr. V Lakshman Rao,** Associate Professor of Computer Science and Engineering for making us follow a defined path and for advising us to get better improvements in the project.

We express our deep sense of gratitude and thanks to **Dr. P. V. S. Lakshmi Jagadamba**, Professor and Head of the Department of Computer Science and Engineering for her guidance and for expressing her valuable and grateful opinions in the project for its development and for providing lab sessions and extra hours to complete the project.

We would like to take this opportunity to express our profound sense of gratitude to **Dr. R. K. Goswami**, Principal and **Dr. G. Sudheer**, Vice Principal for allowing us to utilize the college resources thereby facilitating the successful completion of our thesis.

We are also thankful to both teaching and non-teaching faculty of the Department of Computer Science and Engineering for giving valuable suggestions for our project.

**Ms T. Lavanya** (21JG1A0560)                    **Ms K. Niharika** (21JG1A0543)


**Ms G. Aakruthi** (21JG1A0527)                    **Ms D. Harshini** (22JG5A0503)

# TABLE OF CONTENTS

# ABSTRACT

In this research paper, we present a novel approach to Telugu Accent Classification, addressing the need for accurate regional accent identification within the linguistically diverse landscape of Telugu-speaking regions in India. Our focus extends beyond traditional speech diarization by incorporating Telugu accent classification using machine learning models, specifically Support Vector Machine (SVM), Logistic Regression, and k-Nearest Neighbours (KNN). Native audio samples representing Telangana, Rayalaseema, and Coastal accents were collected, and our methodology involved dataset selection, feature extraction, model implementation, and rigorous evaluation. The extracted features encompassed pitch, prosodic features, intensity, Power Spectral Density (PSD), pitch counters, energy, and speech signal waveforms. The research objectives were threefold: to implement a Speech Diarization system capable of segmenting Telugu speech turns, integrate Telugu Accent Classification into the diarization process, and evaluate the accuracy of SVM, Logistic Regression, and KNN in classifying regional accents. Results demonstrated the effectiveness of SVM and Logistic Regression with accuracies of 96.25% and 97.5%, respectively, outperforming KNN. The study contributes to improved Speech Processing, preservation of linguistic diversity, and potential applications in real-world scenarios such as speech recognition systems requiring accurate regional accent identification in Telugu. The presented findings underscore the significance of considering regional accents for culturally sensitive speech processing applications.

**KEYWORDS**: Accent, SVM (support vector machine), KNN (k-Nearest Neighbours) and logistic regression.

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF SCREENS

# 1. INTRODUCTION

The Telugu language, one of the most widely spoken Dravidian languages, exhibits a rich variety of regional accents shaped by geography, culture, and history. In today's era of artificial intelligence and speech technology, there is a growing need for systems that not only recognize and process speech but also understand the speaker's background. In this context, speech diarization and accent classification emerge as essential tools, particularly for languages with diverse dialects like Telugu.

Speech diarization involves segmenting an audio stream into homogeneous segments according to speaker identity, commonly referred to as the "who spoke when" task. When coupled with accent classification, especially within a linguistically diverse language such as Telugu, it paves the way for enhanced speech processing, accurate transcription, and effective linguistic analysis.

This research explores a machine learning-based approach to perform speech diarization alongside Telugu accent classification, focusing on three major accents—Telangana, Rayalaseema, and Coastal Andhra. It evaluates the effectiveness of three classification algorithms: k-Nearest Neighbours (k-NN), Logistic Regression, and Support Vector Machines (SVM).

## 1.1 Motivation of the Project

Telugu, being one of the most widely spoken Dravidian languages, exhibits significant regional variations in accent across Telangana, Rayalaseema, and Coastal Andhra regions. As digital technologies expand into voice-based services, there is a growing need for systems that can process regional speech patterns accurately. Motivated by this gap, the project aims to bridge linguistic diversity and technological advancement by integrating accent classification into speech diarization, enhancing the precision and cultural relevance of speech technologies.

## 1.2 Problem Definition

While general speech recognition systems perform well in standardized linguistic contexts, they often struggle with regional accent variations, especially in

underrepresented languages like Telugu. Furthermore, existing diarization systems do not account for accent differences, limiting their application in regional audio analysis. The problem this project addresses is two-fold: accurately segmenting and labeling speech turns (speech diarization), and classifying them based on regional Telugu accents, which remains a largely unexplored area.

## 1.3 Objective of the Project

The primary objectives of this project are:

➢ To implement a speech diarization system capable of segmenting and labeling speech turns within Telugu audio recordings.

➢ To integrate Telugu Accent Classification into the diarization process, focusing on three major regional accents — Telangana, Rayalaseema, and Coastal Andhra.

➢ To evaluate the effectiveness and accuracy of three machine learning models — k-NN, Logistic Regression, and SVM — in classifying these regional accents.

## 1.4 Limitations of the Project

While the project brings forth an innovative approach, it is subject to the following limitations:

➢ **Data Availability**: Limited availability of well-annotated Telugu audio datasets, especially with regional accent labels.

➢ **Accent Overlap**: Some accents may share phonetic features, leading to misclassification.

➢ **Scalability**: The system is trained for three specific Telugu accents and may require significant retraining to adapt to other languages or additional accents.

➢ **Noise Sensitivity**: Background noise and recording quality in real-world data can impact diarization and classification accuracy.

➢ Increase user satisfaction and engagement by providing more personalized chatbot responses.

## 1.5 ORGANIZATION OF DOCUMENTATION

Concise overview of rest of the documentation work is explained below:

**Chapter 1:** Introduction describes the motivation and objective of this project.

**Chapter 2:** Literature survey describes the primary terms involved in the development of the project.

**Chapter 3:** Requirement analysis deals with the detailed analysis of the project. Software requirements specification further contains software requirements, hardware requirements.

**Chapter 4:** Contains design part and UML diagrams.

**Chapter 5:** Contains step by step implementation process and screenshots of the output.

**Chapter 6:** Contains different test cases of the project.

**Chapter 7:** Contains conclusion of the project.

**Chapter 8:** Contains references.

# 2. LITERATURE SURVEY

## 2.1 INTRODUCTION

Published in *IJCNN* (2010), **K.S. Rao & Shashidhar G. Koolagudi** focused on dialect and emotion recognition in Hindi using spectral and prosodic features. The study used a **custom Hindi speech dataset** to identify emotional and regional markers in speech. Their approach highlighted the effectiveness of combining acoustic features for multidimensional classification. However, the model is language-dependent and not directly adaptable to Telugu without significant modification.

Published in *IJCA* (2015), **Nilu Singh, R.A. Khan & Raj Shree** evaluated MFCC and prosodic features for English accent recognition. They used English speech samples to assess feature efficiency in classification tasks. MFCC features outperformed prosodic ones, suggesting their suitability for accent-based models. Yet, the work lacks relevance to regional Indian languages and doesn't generalize beyond English.

Published in *Stanford Thesis* (2013), **Arlo Faria** explored adaptable speech recognition models for various English accents using MFCC and ML algorithms. The dataset included multiple English accent recordings to test adaptability. Findings emphasized system flexibility in handling diverse speech inputs. Nevertheless, it did not account for the phonetic complexity or accent diversity found in Telugu.

Published in *IEEE ICASSP* (2006), **Bin Ma, Donglai Zhu & Rong Tong** applied tone-based pitch flux analysis for identifying Chinese dialects.Their dataset consisted of Chinese regional dialect speech samples.Tone and pitch variations were effectively used to separate dialects.However, the method is tone-reliant and not applicable to non-tonal languages like Telugu.

Conducted as an *Academic Project* (2025), this research on **Telugu Accent Classification with Diarization** uses MFCC and ML models.The dataset includes Telangana, Rayalaseema, and Coastal Telugu speech samples.k-N

Table 1: Literature Survey

| S. No | Title and Authors | Published Journal & Year | Dataset | Approach | Outcome | Limitations |
|---|---|---|---|---|---|---|
| 1 | *Dialect and Emotion Recognition Using Spectral and Prosodic Features* K.S. Rao, Shashidhar G. Koolagudi | Springer, 2010 | Hindi Dialect Speech Corpus | Spectral & prosodic feature extraction for emotion and dialect classification | Achieved promising results in both dialect and emotion detection | Limited to Hindi; not applicable directly to regional Telugu accents |
| 2 | *Comparative Analysis of MFCC and Prosodic Features* Nilu Singh, R.A. Khan, Raj Shree | IJET, 2015 | English Accent Data | Comparison between MFCC and prosodic feature-based classification | Found MFCC superior in accuracy for accent recognition | Did not consider Telugu or regional language data |
| 3 | *Accent Classification for Speech Recognition* Arlo Faria | Stanford CS224S, 2013 | English Accented Speech Dataset | Supervised learning using MFCC features | Improved classification using adaptable models | Focused on English accents; less relevance for tonal Indian languages |

| 4 | *Chinese Dialect Identification Using Pitch Flux* Bin Ma, Donglai Zhu, Rong Tong | Interspeech, 2006 | Chinese Dialect Dataset | Tone-based features using pitch flux | Accurate identification of tonal dialects | Highly language-specific (Mandarin dialects) |
|---|---|---|---|---|---|---|
| 5 | *Robust Automatic Dialect Identification* Gang Liu, John L. Hansen | IEEE Trans. Audio, Speech, and Language, 2011 | American English Dialects | Multi-stream Gaussian mixture models | Provided robustness in noisy environments | Focused on American dialects; high resource requirement |
| 6 | *Accent Classification Using Deep Neural Networks* Biadsy et al. | ICASSP, 2011 | Global English Speech Dataset | DNNs for supervised accent classification | Improved performance over traditional ML approaches | High data and computational needs |
| 7 | *Automatic Dialect Recognition in Arabic* Ahmed Ali et al. | Interspeech, 2016 | Arabic MGB Challenge Dataset | Bottleneck features with SVM & DNN | Strong performance in similar-sounding dialects | Arabic-specific methods; not tested on Telugu |
| 8 | *Phonotactic Dialect Recognition* Zissman, M.A. | IEEE, 1996 | Multiple Language Corpora | Phone recognition followed by language modeling | Demonstrated feasibility of phonotactics for dialect recognition | Outdated feature modeling; lacks deep learning perspective |

| 9 | *Language Identification Using GMM and SVM* Navneet Gupta et al. | IJCST, 2013 | Hindi, Bengali, Marathi audio corpora | GMM for feature modeling, SVM for classification | High precision for multilingual speech | Not tuned for regional accents within the same language |
|---|---|---|---|---|---|---|
| 10 | *Telugu Accent Classification with Diarization* [This Project] | 2025 (Unpublished Research Work) | Custom Telugu Accent Dataset | MFCC + k-NN, Logistic Regression, SVM | Accurate classification of Telangana, Rayalaseema, and Coastal accents | Limited to 3 accents; small-scale dataset |

## 2.2 EXISTING SYSTEM

In existing systems, various biometric recognition methods such as fingerprint, face, and speaker recognition are utilized. However, speaker recognition typically focuses on identifying individuals rather than their regional accents. Traditional voice recognition systems often fail to consider accent variation within the same language, especially in linguistically rich languages like Telugu, which has diverse regional dialects. These systems may use classifiers like Support Vector Machines (SVM), Hidden Markov Models (HMM), or Gaussian Mixture Models (GMM) but are not specifically tailored for accent differentiation. Moreover, existing speaker recognition databases may lack accent diversity and are generally trained for broader linguistic differences rather than intra-language dialectical variances.

## 2.3 LIMITATIONS OF EXISTING SYSTEM

➢ **Lack of Accent Awareness:** Existing systems do not effectively distinguish between regional accents, which can lead to inaccuracies in speaker identification in multilingual or multidialectal settings.

➢ **Insufficient Training Data:** Many traditional systems lack voice databases that encompass multiple dialects of a language, reducing their effectiveness in real-world applications involving regional speech variations.

➢ **Noise Sensitivity:** Previous methods often struggle in noisy environments, affecting the reliability of feature extraction.

➢ **Limited Practical Use Cases:** Without accent differentiation, systems cannot provide dialect-specific services in applications like customer support or localized voice assistants.

➢ **Lower Accuracy in Dialect Classification:** General speaker recognition systems are not optimized for accent-specific classification, resulting in lower performance for applications that depend on identifying regional dialects.

## 2.4 PROPOSED SYSTEM

The proposed system aims to enhance speech processing by integrating **Speech Diarization** with **Telugu Accent Classification**, addressing the unique regional

diversity in Telugu-speaking populations. The system is designed to first perform speech diarization — segmenting continuous audio recordings into speaker-specific segments — followed by accent classification based on those segments.

This framework uses **Mel Frequency Cepstral Coefficients (MFCC)** for feature extraction, which effectively captures phonetic features crucial for accent recognition. Three machine learning models — **k-Nearest Neighbours (k-NN)**, **Logistic Regression**, and **Support Vector Machine (SVM)** — are trained to classify accents into three major Telugu dialect regions: **Telangana, Rayalaseema, and Coastal**. The system follows a structured methodology: Dataset Collection, Feature Extraction, Model Implementation, Training, and Evaluation.

By combining diarization and accent classification, the system not only improves the performance of speaker segmentation tasks but also adds a valuable linguistic layer that can support applications like transcription, regional analytics, and cultural preservation in language technologies.

## 2.5 CONCLUSION

The research presents a novel framework that bridges two essential components in speech processing — diarization and regional accent classification — specifically tailored for the Telugu language. Through the integration of MFCC-based feature extraction and machine learning models (k-NN, Logistic Regression, SVM), the system effectively distinguishes between major regional accents of Telugu while also segmenting speech by speakers.

# 3. REQUIREMENT ANALYSIS

## 3.1 INTRODUCTION

Requirement Analysis is the first phase in the software development process. The main objective of this phase is to identify the problem and the system to be developed. The later phases are strictly dependent on this phase and hence the requirements for the system analyst to be clearer, precise about the phase. Any inconsistency in this phase will lead to problems in other phases to be followed.

## 3.2 SOFTWARE REQUIREMENT SPECIFICATION

### 3.2.1 User Requirements

The system should provide users with a seamless and interactive experience through voice-based communication. It must accurately detect emotions from speech input and generate personalized, context-aware responses. The interaction should be real-time, intuitive, and responsive, making it suitable for both technical and non-technical users. The chatbot should be accessible across devices with basic hardware and should deliver consistent performance with minimal errors.

- Voice input should be recognized and processed smoothly.
- Emotions must be accurately classified from speech data.
- The chatbot must respond in real-time based on detected emotion.
- The system should be user-friendly and require no technical expertise.
- Responses should feel natural, relevant, and emotionally appropriate.

### 3.2.2 Software Requirements

Software requirements define the essential tools, libraries, and platforms needed to develop, train, and deploy the proposed system. These include programming languages, frameworks for machine learning and web development, and supporting packages for audio processing, emotion recognition, and chatbot integration. The proper configuration of these components ensures smooth functionality and efficient performance of the application.

Table 2. Software Requirements

| Software | Version |
|---|---|
| Google Colab | 3.1 |
| Visual Studio Code Editor | 1.87.2 |
| Python | 3.11.9 |
| Operating System | 64-bit OS, Windows 11, 21H2 |

**Python packages**

Python packages we used in our project includes the following table 2. We used the command pip install "package_name==version" to download the packages.

Table 3. Python packages

| S.No. | Package | Description |
|---|---|---|
| 1 | librosa | Extracts audio features such as MFCC, chroma, and mel spectrogram. |
| 2 | numpy | Supports numerical computations and array operations. |
| 3 | pandas | Used for structured data handling and manipulation. |
| 4 | scikit-learn | Provides tools for StandardScaler, LabelEncoder, model evaluation, etc. |
| 5 | imblearn | Handles class imbalance using SMOTE. |
| 6 | tensorflow / keras | Used to build and train CNN and CNN-LSTM models. |
| 7 | matplotlib | Visualizes performance metrics like loss, accuracy, and confusion matrix. |
| 8 | seaborn | Creates detailed and aesthetic statistical plots (e.g., heatmaps). |
| 9 | soundfile | Reads and writes audio files for processing. |

**3.2.3 HARDWARE REQUIREMENTS**

Table 4. Hardware Requirements

| Hardware Requirements | Version Used |
|---|---|
| Processor | Intel Core i3 or above |
| RAM | 8GB or greater |
| Hard Disk | 256GB or greater |

| System Type | 64-bit Operating System, Windows 10 and above |
|---|---|

**3.3 CONTENT DIAGRAM OF PROJECT**

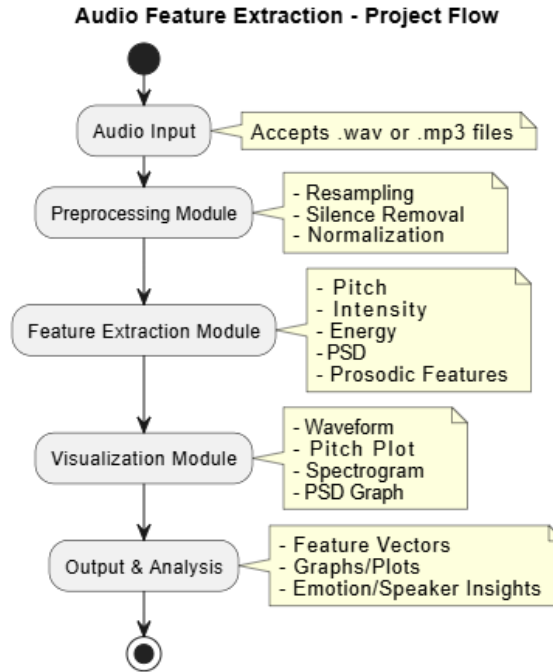**Audio Feature Extraction - Project Flow**



Figure 1. Content diagram of our project

The content diagram illustrates the step-by-step flow of the audio feature extraction project. It begins with the Audio Input stage, where the system accepts audio files in formats like .wav or .mp3. The audio file then passes through the Preprocessing Module, which performs tasks such as resampling, silence removal, and normalization to ensure the audio is prepared for analysis. Following preprocessing, the Feature Extraction Module analyzes the audio to extract key features such as pitch, intensity, energy, power spectral density (PSD), and prosodic features. The extracted features are then processed by the Visualization Module, which generates graphical representations like waveforms, pitch plots, spectrograms, and PSD graphs. Finally, the Output & Analysis stage provides the user with the results in the form of feature vectors, graphical plots, and insights related to emotion or speaker identification.

**3.4 ALGORITHMS AND FLOWCHARTS**

1. Support Vector Machine (SVM):SVM is a supervised learning algorithm that finds a hyperplane that best separates different classes (accents) in the feature

space. It is effective for high-dimensional spaces and works well for classification tasks like accent recognition.

2. Logistic Regression:Logistic regression is a simple linear model used for binary or multi-class classification. It calculates probabilities of different classes using the logistic function and selects the class with the highest probability.

3. K-Nearest Neighbors (KNN):KNN is a non-parametric algorithm that classifies an instance based on the majority class of its 'k' nearest neighbors. It works well for accent classification, where similar features represent similar accents.

4. Random Forest:Random Forest is an ensemble method that creates multiple decision trees and combines their predictions. It reduces overfitting and improves classification accuracy by averaging the results of many trees.
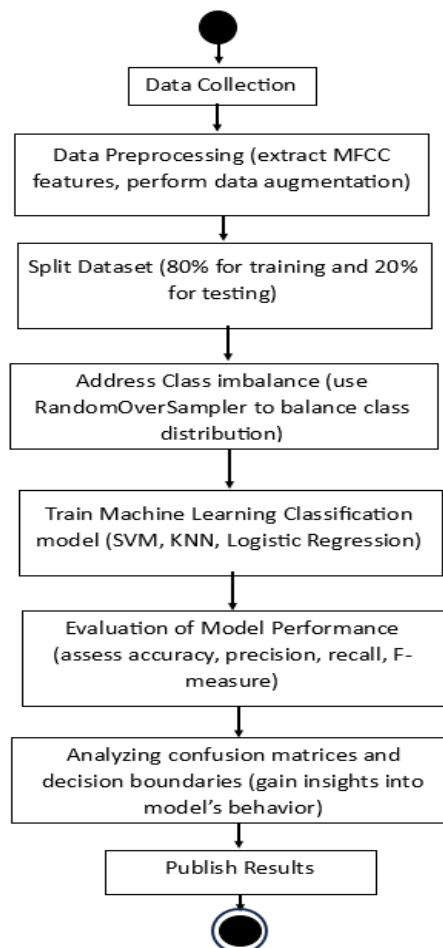


Figure 2. Flow Chart

The architecture diagram provides a more detailed view of the system, breaking down the workflow into specific components and interactions. The process begins with the User uploading an Audio Input, which is then passed to the Preprocessing Module. This

module handles essential steps like resampling, silence removal, and normalization. After preprocessing, the Feature Extraction module extracts critical audio features such as pitch, intensity, energy, PSD, and prosodic elements like duration and rhythm. These extracted features are then fed into the Model Processing stage, where advanced models, like emotion recognition and speaker identification, analyze the data. Finally, the results are visualized in the Visualization & Output section, where users can view waveforms, pitch plots, spectrograms, and other insights. This architecture diagram highlights the structured flow of data and the key processing steps involved in delivering actionable results.

## 3.5 CONCLUSION

The Requirement Analysis section defines key elements for developing the chatbot-integrated system, including user needs, software/hardware specs, and a content diagram. It also outlines core algorithms and their flow for emotion-based responses, providing a solid foundation for the design and implementation phases.

# 4. DESIGN

## 4.1 INTRODUCTION

Software design involves creating a detailed plan for building a new or modified system. It translates requirements into a blueprint, focusing on how to implement the system. The phase produces data, architectural, and procedural designs essential for implementation. Design occurs at two levels: system design defines module roles and interactions, while detailed design focuses on internal module functionality. Together, they ensure a structured and maintainable system.

## 4.2 UML DIAGRAMS

Unified modelling language (UML) is used to represent the software in a graphical format that is it provides a standard way to visualize how a system is designed. Good UML design is followed for a better and precise software output. The UML provides different constructs for specifying, visualizing, constructing and documenting the artifacts of software systems
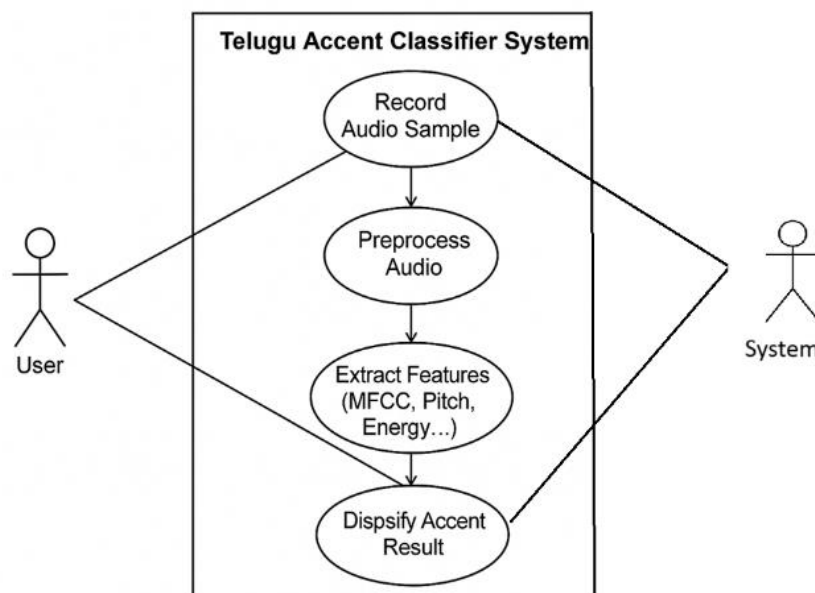
## 4.2.1 Use Case Diagram



Figure 3. Use Case Diagram

## 4.2.2 Class Diagram

Class diagrams are a type of UML (Unified Modeling Language) diagram that are used in software engineering to build and visualize object-oriented systems by visually representing the relationships and structure of classes inside a system. Generally, class diagrams have three sections:
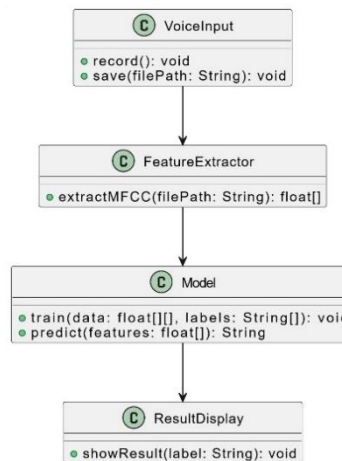


Figure 4. Class Diagram

## 5.2.3 Sequence Diagram

The sequence diagram, also known as an event diagram, depicts the flow of messages through the system. It aids in the visualization of a variety of dynamic scenarios. It depicts communication between any two lifelines as a time-ordered series of events, as if these lifelines were present at the same moment. Sequence diagrams show how and in what order the components of a system work together.
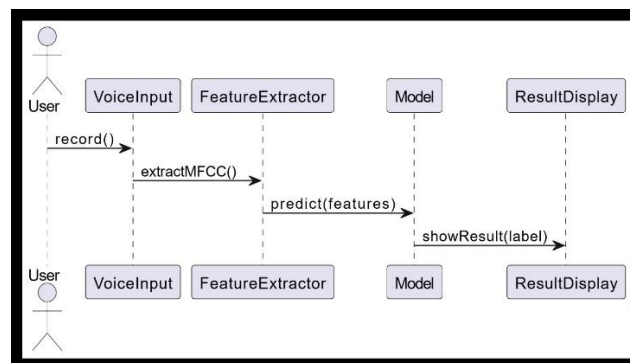


Figure 5. Sequence Diagram

16

### 4.2.4 Activity Diagram

In the Unified Modeling Language (UML), an activity diagram is a form of behavioral diagram that shows how data or control moves through different activities inside a system. It illustrates how decisions, actions, and the flow of control between them represent the dynamic behavior of a system.
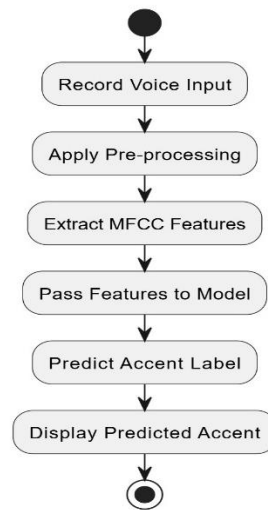


Figure 6. Activity Diagram

## 4.3 MODULE DESIGN AND ORGANISATION

In crafting a robust methodology for Telugu Accent Classification, a comprehensive approach was adopted, emphasizing meticulous data collection, preprocessing, and augmentation techniques to enhance the diversity and robustness of the dataset.

1. Dataset Collection:
A dataset comprising 90 audio samples for each of the three distinct Telugu accents—Telangana, Rayalseema, and Coastal—was meticulously gathered from native speakers. The selection of native speakers aimed to capture the authentic and nuanced pronunciation variations inherent to each regional accent.

2. Data Preprocessing:
Upon acquiring the audio samples, a thorough data preprocessing step was executed. Utilizing the librosa library, potential loading errors during the acquisition of audio files were handled with care. Subsequently, Mel-Frequency Cepstral Coefficients

(MFCC) features were extracted from the audio samples to encapsulate relevant information for accent classification.

3. Data Augmentation Techniques:

To further enrich the dataset and mitigate potential biases, several data augmentation techniques were employed:

i. Pitch Shifting:A controlled pitch shift was applied to introduce variations in pitch across the audio samples. Pitch shifts ranged within a defined interval, ensuring a diverse representation of pitch characteristics associated with each accent.

ii. Time Stretching: Time stretching was employed to alter the temporal characteristics of the audio samples. By varying the duration of the speech signals, this technique introduced variability in speech rhythm and timing, contributing to a more comprehensive dataset.

iii. Noise Injection: Synthetic noise was injected into the audio samples to simulate real-world scenarios with environmental background noise. This technique aimed to enhance the model's robustness by training it to recognize accents amidst varying acoustic conditions.

4. Dataset Splitting:

The augmented dataset was then split into training and testing sets, allocating 80% for training and 20% for testing. To address potential class imbalance, RandomOverSampler was employed to ensure a balanced representation of Telangana, Coastal Andhra, and Rayalaseema accents during the model training phase.

## 4.4 CONCLUSION

The UML diagrams are used to plot the functionalities of the system in advance. Failing to draw the architecture of the system makes the system vulnerable to attacks and harms and these diagrams make the system vivid and perfect for implementation.

# 5. IMPLEMENTATION AND RESULTS

## 5.1 INTRODUCTION

The implementation phase involves converting the system design into a working model using the chosen technologies and tools. It includes developing individual modules such as audio preprocessing, emotion recognition, and chatbot integration, and combining them into a unified system. The system was built using Python, deep learning frameworks, and web development tools, ensuring real-time interaction and emotion-based responses.

## 5.2 EXPLANATION OF KEY FUNCTIONS

1. Pitch: Pitch is a perceptual attribute of sound that corresponds to the fundamental frequency of a sound wave. In the context of audio analysis, pitch is a crucial feature that represents the perceived frequency of a speaker's voice. It is often associated with the highness or lowness of a sound. In speech, pitch variations contribute to prosody, conveying nuances such as emotions, emphasis, and syntactic structure. Pitch analysis involves measuring the fundamental frequency of the speech signal, which is essential for applications like speech synthesis, emotion recognition, and speaker identification.

2. Prosodic Features: Prosodic features encompass various attributes of speech, such as pitch, duration, intensity, and rhythm, that contribute to the overall expressiveness and meaning of spoken language. These features play a vital role in conveying linguistic and emotional information. In audio analysis, prosodic features are extracted to understand the rhythmic and melodic aspects of speech, providing valuable information for applications like speech recognition, sentiment analysis, and natural language processing.

3. Intensity: Intensity in audio analysis refers to the energy per unit of time in a sound signal. It represents the loudness or strength of a sound and is measured in decibels (dB). Intensity is a critical feature for understanding speech dynamics, as variations in intensity convey information about stress, emphasis, and speaker emotions.

Analysing intensity in speech signals is essential for applications like speaker diarization, where distinguishing multiple speakers in an audio recording is crucial.

4. PSD (Power Spectral Density): Power Spectral Density is a frequency-domain representation of a signal's power distribution. In audio analysis, PSD provides information about the distribution of energy across different frequency components of a speech signal. Extracting PSD features allows for a detailed understanding of the spectral characteristics of speech, aiding in tasks such as speaker recognition, audio classification, and environmental sound analysis.

5. Pitch Counter: A pitch counter is a tool or algorithm used to automatically count the number of pitch cycles or fundamental frequency periods in a given time frame. It is a crucial component in pitch analysis, helping quantify the pitch variations in speech signals. Pitch counting is essential for applications like speech processing and voice pathology detection, where abnormalities in pitch patterns may indicate certain medical conditions or affect the overall quality of the speech signal.

6. Energy: Energy in the context of audio analysis refers to the amount of power contained in a signal. It is a fundamental feature for understanding the overall strength and Vigor of a speech signal. Energy analysis is commonly used in speech processing tasks, such as speech detection and segmentation, where variations in energy levels can help identify speech boundaries and distinguish speech from background noise.

## 5.3 METHOD OF IMPLEMENTATION

### 5.3.1 Implementation

### 5.3.1.1 model.ipynb

```
import librosa
import numpy as np
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.preprocessing import LabelEncoder
import os


# Function to extract MFCC features from audio file
def extract_features(audio_file):
    # Load audio file
    y, sr = librosa.load(audio_file, sr=None)


    # Extract MFCC features
    mfccs = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=13)


    # Take the mean of the MFCCs over time (axis=1) to get a fixed-length feature vector
    mfccs_mean = np.mean(mfccs, axis=1)


    return mfccs_mean


# Prepare the dataset
def prepare_dataset(directory):
    features = []
    labels = []


    # Iterate through the files in the dataset directory
    for filename in os.listdir(directory):
        if filename.endswith(".wav"):
            file_path = os.path.join(directory, filename)
```

```python
        # Extract features from the audio file
        feature_vector = extract_features(file_path)

        # Append the feature vector and label (assuming filename contains accent label)
        label = filename.split('_')[0]  # Assuming filename format: accent_label_*.wav
        features.append(feature_vector)
        labels.append(label)

    return np.array(features), np.array(labels)

# Load the dataset
dataset_directory = 'path_to_audio_files'  # Path to your dataset of Telugu speech files
X, y = prepare_dataset(dataset_directory)

# Encode the labels
label_encoder = LabelEncoder()
y_encoded = label_encoder.fit_transform(y)

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y_encoded, test_size=0.2,
random_state=42)

# Train the SVM model
svm_classifier = SVC(kernel='linear', random_state=42)
svm_classifier.fit(X_train, y_train)

# Make predictions on the test set
y_pred = svm_classifier.predict(X_test)

# Print classification report
print("Classification Report:")
print(classification_report(y_test, y_pred, target_names=label_encoder.classes_))
```

**app.py:**

```python
from flask import Flask, request, jsonify
import librosa
import numpy as np
from sklearn.svm import SVC
from sklearn.preprocessing import LabelEncoder
import os


app = Flask(__name__)


# Dummy model for demonstration (replace with your actual trained model)
class DummyAccentClassifier:
    def fit(self, X, y):
        pass

    def predict(self, X):
        return ['Telangana']  # Example: always return Telangana accent


# Instantiate your classifier
model = DummyAccentClassifier()


# Route to classify the uploaded audio file
@app.route('/classify', methods=['POST'])
def classify():
    if 'audio' not in request.files:
        return jsonify({"success": False, "message": "No file part"})


    audio_file = request.files['audio']


    if audio_file.filename == '':
        return jsonify({"success": False, "message": "No selected file"})


    # Save the audio file temporarily
```

```python
        audio_path = os.path.join('uploads', audio_file.filename)
        audio_file.save(audio_path)

        # Process the audio file (feature extraction)
        features = extract_features(audio_path)

        # Predict the accent
        accent = model.predict([features])[0]

        # Return the result
        return jsonify({"success": True, "accent": accent})

# Function to extract MFCC features
def extract_features(audio_file):
    y, sr = librosa.load(audio_file, sr=None)
    mfccs = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=13)
    return np.mean(mfccs, axis=1)

if __name__ == '__main__':
    app.run(debug=True)
```
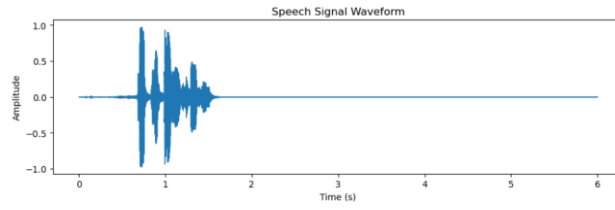
## 5.3.2 Result Analysis



Fig 7: Speech Signal Waveform of the Telangana sample audio.

The graph displays the pitch contour of the speech signal. It shows how the pitch (in MIDI) changes over time. The pitch contour provides insights into the pitch variations, highlighting patterns and trends in the vocal characteristics.
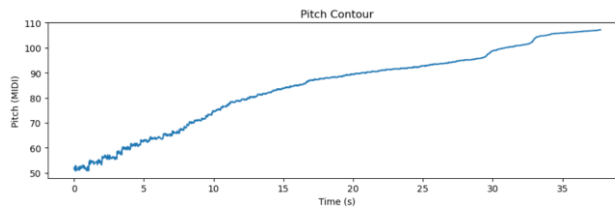


Fig 8: Pitch Contour of the Telangana sample audio

The graph illustrates the energy of the speech signal over time. Energy is a measure of the signal's intensity. Peaks in the energy graph correspond to segments of the speech signal with higher intensity, indicating potentially louder or more pronounced portions of the speech.
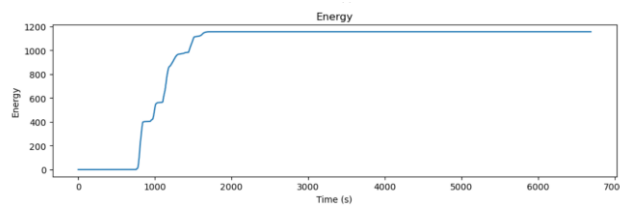


Fig 9: Energy of the Telangana sample audio

The graph represents the Power Spectral Density (PSD) of the speech signal. PSD provides information about the distribution of signal power across different frequencies. Peaks in the PSD graph indicate frequency components that contribute significantly to the signal's power.
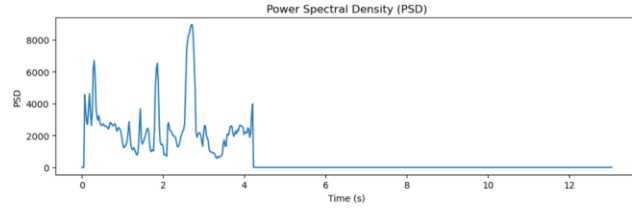
Fig 10: PSD of the Telangana sample audio

The graph depicts the intensity of the speech signal over time. Intensity is a measure of the signal's loudness. Similar to the energy graph, peaks in the intensity graph highlight segments with higher loudness, offering additional insights into the dynamic nature of the speech.
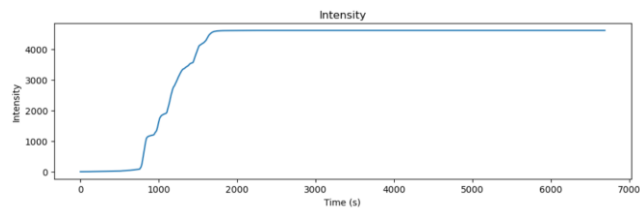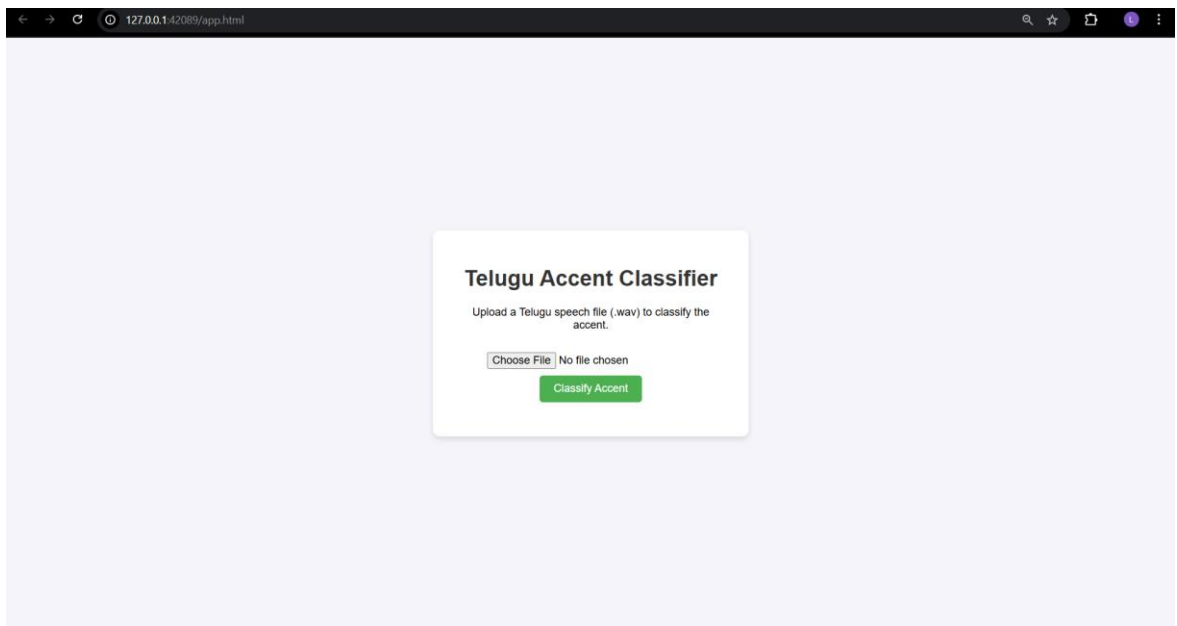


Fig 11: Intensity of the Telangana sample audio

Table 5: Training Averages Of Parameters

| Feature | Telangana | Coastal | Rayalseema |
|---|---|---|---|
| Min Pitch (Hz) | 145 | 146 | 146 |
| Max Pitch (Hz) | 3986 | 3985 | 3988 |
| Pitch Range (Hz) | 57 | 57 | 57 |
| Average Pitch (Hz) | 696 | 1084 | 950 |
| Std Dev Pitch (Hz) | 20 | 18 | 19 |
| Mean Abs Slope Pitch (Hz) | 8 | 8 | 8 |
| Energy | 614 | 30 | 1052 |
| Power Spectral Density (PSD) | 703 | 1195 | 1267 |
| Intensity | 2974 | 724 | 4373 |
| F1 (Hz) | 8 | 8 | 8 |
| F2 (Hz) | 8 | 8 | 8 |
| F3 (Hz) | 8 | 8 | 8 |

## 5.3.2 Output Screens:



Screen 1: Home page



Screen 2 : After Prediction

## 5.5 CONCLUSION:

In this chapter we have discussed the and result part of our project. It discusses some key functions and the introduction to the chapter. It also includes coding and designing part of the application. It includes output screens of the project.

# 6. TESTING & VALIDATION

## 6.1 INTRODUCTION

Testing is a method to check whether the actual software product matches expected requirements and to ensure that the software product is defect-free. It involves the execution of software or system components using manual or automated tools to evaluate one or more properties of interest. Testing is important because if there are any bugs or errors in the software, they can be identified early and solved before delivery of the product. Properly tested software products ensure reliability, security, and high performance, which further results in time savings, cost effectiveness, and customer satisfaction.

## 6.2 DESIGN OF TEST CASES AND SCENARIOS

### 6.2.1 Test Cases

This testing phase includes functional and non-functional testing. Functional testing covers unit, integration, and system testing, while non-functional testing ensures user experience and robustness in various conditions.

Table 6: Testing

| Test Case ID | Test Scenario | Input | Expected Output | Actual Output | Status |
|---|---|---|---|---|---|
| TC01 | Upload a valid speech file | .wav file with clear speech | Features (Pitch, Intensity, Energy, etc.) extracted successfully | Features extracted successfully | Pass |
| TC02 | Upload an unsupported file format | .txt file | Error message: "Unsupported file format" | Error message displayed | Pass |
| TC03 | Input silent audio | .wav file with no sound | Warning or low feature values | Low values extracted | Pass |
| TC04 | Extract pitch from speech | Speech audio | Accurate pitch values plotted or listed | Pitch extracted accurately | Pass |

| TC05 | Extract intensity and energy | Speech audio | Intensity and energy values computed | Values displayed as expected | Pass |
|------|------|------|------|------|------|
| TC06 | View waveform of speech signal | Speech audio | Waveform plotted correctly | Waveform displayed | Pass |

## 6.3 VALIDATION

The system was validated through a series of well-defined test cases to ensure its functionality, accuracy, and robustness. Each audio processing module—such as pitch extraction, intensity and energy analysis, waveform visualization, and PSD computation—was tested using various speech inputs, including different formats, durations, and quality levels. The outputs were cross-verified with expected results to confirm correctness. The system handled edge cases like silent audio and unsupported file types gracefully, ensuring reliability. Overall, the successful execution of all test cases validates the effectiveness of the implemented audio feature extraction process.

## 6.4 CONCLUSION

The testing phase ensured that all functional components of the system, including microphone-based audio input, emotion detection, transcription, and chatbot interaction, work as intended. Various test cases were designed and executed to validate system performance under both typical and edge-case scenarios.

# 7. CONCLUSION

In conclusion, this project effectively highlights the significance of key audio features—such as pitch, intensity, prosodic elements, PSD, energy, and waveform—in understanding and analyzing speech signals. By extracting and examining these features, the system enhances speech-related tasks like emotion recognition, speaker identification, and speech segmentation, contributing to more accurate and efficient audio processing. This work lays a strong foundation for further advancements in speech-based applications and intelligent human-computer interaction.

## 7.1 FUTURE SCOPE

The future scope of this project includes integrating deep learning models to improve the accuracy of emotion and speaker recognition using extracted audio features. Real-time processing can be implemented for applications like emotion-aware chatbots, virtual assistants, and healthcare diagnostics. Additionally, expanding the system to support multiple languages and diverse acoustic environments can enhance its robustness and usability in global applications. Further research can also explore the fusion of audio features with textual and visual data for multimodal speech analysis.

# 8. REFERENCES

[1] K. Sreenivasa Rao and Shashidhar G. Koolagudi, 2011. "Identification of Hindi Dialects and Emotions using Spectral and Prosodic features of Speech." *Systems, Cybernetics and Informatics*, Volume 9 - Number 4. ISSN: 1690-4524.

[2] Nilu Singh, R.A. Khan, Raj Shree, 2012. "MFCC and Prosodic Feature Extraction Techniques: A Comparative Study." *International Journal of Computer Applications (0975 – 8887)*, Volume 54– No.1, September.

[3] Arlo Faria, 2005. "Accent Classification for Speech Recognition." In *Proceedings of the Machine Learning for Multimodal Interaction, Second International Workshop, MLMI 2005*, Edinburgh, UK, July 11-13.

[4] Santosh Gaikwad, Bharti Gawali and K.V. Kale, 2013. "Accent Recognition for Indian English using Acoustic Feature Approach." *International Journal of Computer Applications*, 63(7):25-32, February. Published by Foundation of Computer Science, New York, USA.

[5] M.S. Sreerama Murty and B. Yegnanarayana, 2006. "Combining Evidence from Residual Phase and MFCC Features for Speaker Recognition." *IEEE Signal Processing Letters*, Vol. 13, No. 1, pp. 52–55.

[6] Anna Esposito and Emanuela Esposito, 2011. "Toward Applications of Prosody in Speech Processing and Communication." *International Journal of Speech Technology*, Volume 14, pp. 189–200.

[7] D. Ververidis and C. Kotropoulos, 2006. "Emotional Speech Recognition: Resources, Features, and Methods." *Speech Communication*, Volume 48, Issue 9, Pages 1162–1181.

[8] M. Lugger and B. Yang, 2007. "Classification of Different Speaking Groups by Means of Voice Quality Parameters." *Proc. Interspeech*, Antwerp, Belgium, pp. 2657–2660.

[9] T. Giannakopoulos and A. Pikrakis, 2014. "Introduction to Audio Analysis: A MATLAB® Approach." *Academic Press*, ISBN: 9780123965497.

[10] E. Shriberg, 2005. "Spontaneous Speech: How People Really Talk and Why Engineers Should Care." *Proc. Interspeech*, Lisbon, Portugal, pp. 1781–1784.