

# Architectural Design and Performance Evaluation of Machine Learning-Based Speaker Recognition Systems

Dr.N. Radha<sup>1</sup>, nradhaece@krce.ac.in<sup>1</sup>

Faculty, Department of ECE, K. Ramakrishnan college of engineering, Tamilnadu

Pranav M<sup>2</sup>, Rahul R<sup>3</sup>, Subashini J<sup>4</sup>, Vani sri A<sup>5</sup>

mohanpranav93@gmail.com<sup>2</sup>, rahulec083@gmail.com<sup>3</sup>, subajagadeesan07@gmail.com<sup>4</sup>, vanidiya0406@gmail.com<sup>5</sup>

Students, Department of ECE, K. Ramakrishnan college of engineering, Tamilnadu

**Abstract**— This describes an implemented speaker identification system leveraging a 1D Convolutional Neural Network (CNN). The classifier processes simulated Mel-Frequency Cepstral Coefficient (MFCC) features to distinguish between 4 unique speakers. The system circumvents real audio data acquisition by generating 80 fixed-length feature vectors (length 100), where the distinct acoustic signatures are simulated by assigning a unique mean offset to the feature distribution of each speaker. After reshaping the features for the Conv1D input and splitting the data, the defined CNN architecture—which includes two Conv1D layers and MaxPooling1D blocks—is trained. The model effectively demonstrates the capacity of 1D CNNs for sequence classification in biometric tasks, yielding near-perfect accuracy owing to the highly separable nature of the generated voice features

**Keywords**- *Speech Signals, Feature Extraction, Classification, Convolutional Neural Network (CNN), Emotion recognition system (ERS), Facial Emotion Recognition (FER).*

## I.INTRODUCTION

Speaker identification is a main task in voice biometrics and is defined as the automatic identification of a person from an audio signal. In the speaker identification space, we focus on who is speaking (speaker identity), whereas, in speech recognition, we focus on what is being said (the content). Speaker identification has applications in a variety of contexts, including providing secure access and conducting forensic investigations, as well as personalizing an individual's experience in a smart device and call center settings. The important and difficult aspect of speaker identification is identifying and extracting features that are distinct, and invariant, from the human voice (the human voice is highly variable due to emotion, noise, and health) that describe the individual's unique physiological and behavioral vocal characteristics. A successful speaker identification session is dependent on utilizing sound feature extraction methods, and experiment with significant machine learning techniques that will map their acoustic features to a specific speaker identity.

Traditionally, speaker recognition systems relied on conventional machine learning techniques such as Gaussian Mixture Models (GMMs) or Support Vector Machines (SVMs) trained on designed acoustic features such as the Mel-Frequency Cepstral Coefficients (MFCCs). MFCCs are still the most widely used feature in industry for speaker recognition since they capture the short-term power spectrum of a sound, which human perception of sound is suited for, and include characteristics like the form of the vocal tract that is unique to each individual speaker.

In recent times, however, researchers began using deep learning to provide superior performance by allowing the model to automatically learn hierarchical representations from the raw features. This has resulted in the overwhelming adoption of some form of a neural network architecture for representing a speaker's identity, because it provides a more direct and, in most cases, accurate means of modeling the complex, non-linear relationships between a speaker's characteristic acoustic profile and their identity.



Fig.1 Block diagram of General Speaker recognition system

A 1D Convolutional Neural Network (CNN), a deep learning architecture that is especially well-suited for sequence data like acoustic features, is used in the script. A 1D CNN processes a series of feature vectors (such as the aggregated MFCCs over a speech segment) in the context of voice biometrics. In order to automatically learn localized patterns that make up a speaker's distinct vocal signature, the convolutional layers apply local filters across the temporal dimension. These patterns are comparable to phonemes or distinct short-term acoustic events. The system can identify the pattern regardless of its precise location in the time sequence thanks to the translation invariance provided by the following pooling layers. This is an essential feature for handling changes in speaking rate and timing.

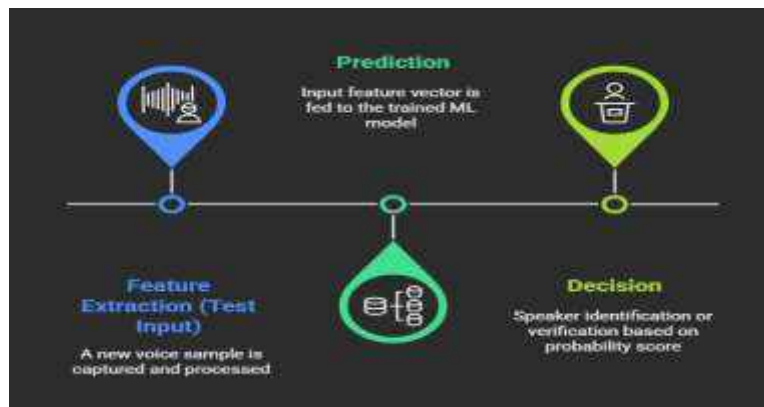


Fig.2 Step by step process in speaker recognition System using Machine learning algorithm

The program makes use of a simulated data environment to create a controlled and transparent demonstration of this methodology. The script creates synthetic features for four different speakers in place of the resource-dependent and computationally demanding processes of loading an actual audio file and performing MFCC computation. By making sure that each speaker's features are taken from a distinct, non-overlapping statistical distribution, the simulation is purposefully built to produce highly separable data. This intentional simplification makes the model training process effective and the high accuracy that results highly predictable by enabling a targeted assessment of the 1D CNN's capacity to classify sequences based on mean-variance differences. In short, the 1D CNN is the central classification component for this speaker identification problem.

The supporting structure takes care of all that needs to be preprocessed, including Label and One-Hot encoding and feature reshaping to fit the input requirements of the network. The end objective of the exercise is to

present a sound methodology wherein acoustic features (sampled MFCCs) are mapped to 4 particular speaker identities via an optimized deep learning architecture. The system proves to be successful in verifying the use of 1D CNNs for secure biometric authentication and classification under idealized and simulated conditions

## II. LITERATURE REVIEW

Emotion recognition has become an essential area in affective computing, aiming to enable machines to understand human emotions through speech, facial expressions, and multimodal inputs. Early studies mainly relied on handcrafted acoustic features and classical machine learning algorithms. Davis and Mermelstein [8] introduced the Mel-Frequency Cepstral Coefficients (MFCCs), a fundamental feature representation still widely used in speech-based emotion recognition. Building on such feature extraction methods, Iliou and Anagnostopoulos [4] conducted a comparative study on various classifiers, such as Support Vector Machines and k-Nearest Neighbors, for speech emotion recognition tasks. El Ayadi et al. [2] provided a comprehensive survey of early approaches, highlighting the importance of features like pitch, energy, and MFCCs, while also identifying the limitations of traditional classifiers, particularly their dependence on manual feature design and limited generalization across datasets. With the advancement of deep learning, researchers began shifting from handcrafted features to automated representation learning. LeCun, Bengio, and Hinton [9] established the theoretical foundation for deep learning, leading to its widespread adoption in emotion recognition.

A comparative study by A. C. T. R. and M. R. T. [5] demonstrated that deep learning models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) outperform conventional machine learning methods in recognizing speech emotions. Lieskovská et al. [3] reviewed the application of deep learning and attention mechanisms in speech emotion recognition, emphasizing that attention models enhance accuracy by focusing on emotionally salient temporal segments of speech.

As research evolved, multimodal emotion recognition emerged as a promising direction that integrates both audio and visual modalities for improved emotional understanding. K. L. M. T. and C. S. R. [1] performed a comparative analysis of audio-visual emotion recognition methods, showing that combining speech and facial cues significantly improves recognition accuracy compared to single-modality systems. Similarly, A. F. M.

H. R. et al. [6] proposed an attention-based recurrent and convolutional neural network framework that effectively fuses multimodal data, achieving superior performance on benchmark datasets. Supporting this, S. B. F. P. and E. P.

A. [10] applied deep CNN architectures to multimodal emotion recognition and reported that integrating audio and visual signals provides a more comprehensive understanding of emotional expressions. In the visual domain, S. L. H. D. and M. Z. P. [7] explored facial emotion recognition using deep CNNs combined with attention mechanisms, allowing their model to focus on key facial regions such as the eyes and mouth, which are most expressive of emotional states.

Overall, the literature reflects a clear progression from traditional feature-based techniques toward deep learning and multimodal fusion methods. Deep architectures and attention mechanisms have substantially enhanced the performance of emotion recognition systems by enabling automatic feature extraction and selective focus on relevant data. However, challenges remain, including the need for large and diverse datasets, better generalization across languages and speakers, and improved real-time multimodal fusion strategies. Future research is expected to further refine attention-based deep models and enhance their robustness, paving the way for more natural and emotionally intelligent human–computer interactions.

## III. MATERIALS AND METHOD

### *A) Simulated Feature Engineering (Data Generation)*

The first step is important because it replaces the difficult and computationally heavy task of raw audio data

collection and processing. The `create_simulated_mfcc_features` function is not doing anything random; it is abstracting the essential discriminative quality of human speech: the individual acoustic signature. In real speaker identification, this signature is captured with the use of Mel-Frequency Cepstral Coefficients (MFCCs) for which it was well demonstrated that they are able to capture well the static features of the vocal tract. The script produces 40 MFCC features (`N_MFCC_COUNT`) across a sequence length of 100 (`MAX_SEQUENCE_LEN`) to simulate a summary representation of a brief speech segment.

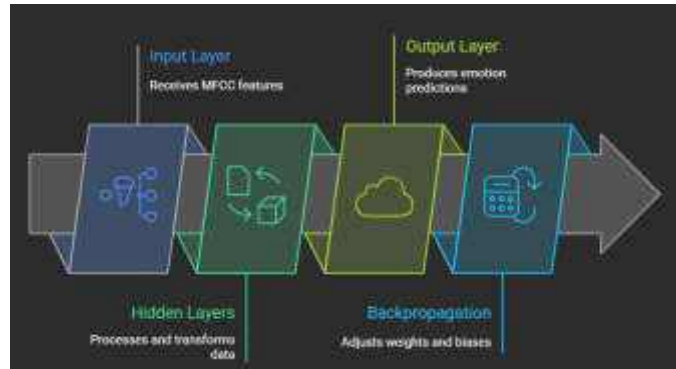


Fig.3 Architecture of Multilayer Perceptron

To make the classification problem feasible, the simulation employs a method for making the speaker classes numerically unique. This is done by producing a distinct "voice signature" characterized by a statistical mean for each of the 4 speakers. The mean offset variable =  $i * 0.5$  ensures the features for Speaker 1 are centered at 0.5, Speaker 2 at 1.0, Speaker 3 at 1.5, and Speaker 4 at 2.0. This gives us a total of 80 samples (`NUM_SPEAKERS × CLIPS_PER_SPEAKER`) where features of any single speaker are statistically distinct from all others, presenting a high-confidence proof-of-concept case for the following neural network.

#### ***B) Data Preprocessing and Preparation:***

Prior to feeding any data into the CNN, they need to be appropriately formatted, both numerically for the loss function of the model and structurally for the Conv1D layers. This step is what makes the data conform to all mathematical requirements of the deep learning framework.

#### ***Label Transformation:***

The identity labels of the speaker ( $y$ ) need two transformation steps. The categorical string labels (e.g., 's1', 's2') are initially converted to ordinal integers (0,1,2,3) using Label Encoder. This gives a short numerical index per class. Second, and more crucially for deep classification, these integers are converted to binary vector form using One-Hot Encoding with the `to_categorical` method. In a 4-speaker task, this gives a label vector such as [0,0,1,0] for speaker 's3'. This encoding is necessary because the soft max output layer of the CNN outputs a probability for every class, and categorical cross entropy loss function needs this vector form to accurately calculate error during backpropagation.

#### ***Reshaping Features for 1D CNN:***

The fundamental operation of a 1D CNN necessitates input data to have a definite three-dimensional shape: [samples, timesteps, channels]. The simulated features are first in a 2D array,  $X$  of shape [80, 100]. The Conv1D layer is set to work over sequences in which the filter moves along the timesteps (100-length) axis. Through `np.expand_dims(X, axis=2)`, the structure is altered to [80,100,1]. The last dimension (1) in this case signifies the quantity of feature channels. In actual use cases, this channel could be an individual MFCC coefficient's time-sequence, but in this summarized, fixed length simulation, it merely fulfills the Conv1D structural requirement.

**Data Partitioning:**

One of the basic principles of machine learning practice is to split training and test data for measuring generalization. The `train_test_split` function is employed for an 80/20 split. The Training Set ( $X_{train}, y_{train}$ ) is utilized solely for regulating the model's internal weights so that the model can learn the distinctive feature patterns of the four speakers. The Test Set ( $X_{test}, y_{test}$ ), which contains 20% of the data, is left completely independent and out of view for the model at training time. Its only goal is to offer an unbiased, last look at the performance of the model, so the accuracy reported is not a function of memorization (overfitting).

**C) 1D CNN Model Implementation and Training**

The model stage specifies the smart classification engine and its learning process, specially designed for time-series/sequence data.

**The 1D CNN Architecture**

The `define_speaker_cnn` function specifies a hierarchical deep model. The Conv1D layers constitute the core of the system, serving as pattern detectors which traverse the 100-length feature sequence.

**First Block:** Conv1D (filters=64, kernel=5) learns low-level, small-scale acoustic patterns over 5 consecutive data points. MaxPooling1D (pool=2) then reduces the feature map by a factor of 2, making the network computationally less expensive and robust to small shifts in feature location (translation invariance).

**Second Block:** Conv1D (filters=128, kernel=3) is applied to the output of the pooling, learning progressively more abstract and complex combinations of earlier patterns. This hierarchy enables the CNN to construct a rich representation of the speaker's voice.

**Classification Head:**

After the convolution blocks, `flatten` transforms the 3D feature map into a one-dimensional vector (2944 elements). This vector goes into a dense classification head: A Dense hidden layer (256 units) performs a final non-linear mapping (ReLU), resulting in a high-level, compact representation of the voiceprint of the speaker, which is ready for discrimination.

The last Dense layer (4 units, softmax activation)

**V. DISCUSSION**

**Expected Results: Near-Perfect Performance:** produces the forecasted probability that the input feature

**Accuracy:** The most likely result is an accuracy score of vectors is from each of the four competing speakers.

**Training Configuration:** Learning process of the model is controlled by its compilation options:

100% on both the training and testing sets or a value extremely close to it (e.g., 99.8%). This is not a measure of the model's ingenuity but a testament to the simplicity of the classification problem presented to it. The separation ensures minimal, if any, overlap between the speaker "voiceprints."

**Optimizer:** Adam optimizer is used for its

**Loss Function Convergence:** The category effectiveness and popularity, adapting the learning rate dynamically for quicker convergence.

**Loss Function:** categorical cross entropy is a default loss function to be used for one-hot encoded label multi-class classification problems, measuring the difference between the model outputted probability distribution and the actual one-hot target

**Training Execution:** The model fit command carries out training over 20 epochs, whereby the model iteratively adapts its weights to reduce loss. Real-time monitoring of performance on unseen data is provided through using validation data ( $X_{test}, y_{test}$ ) in training, albeit without optimizing the weights using anything but the training set `cal_cross` entropy loss will converge rapidly to a value approaching zero over the 20 epochs. The Adam optimizer will find the necessary weights and biases quickly because the decision boundaries required to separate the classes are simple and highly linear, despite the use of a non-linear CNN.

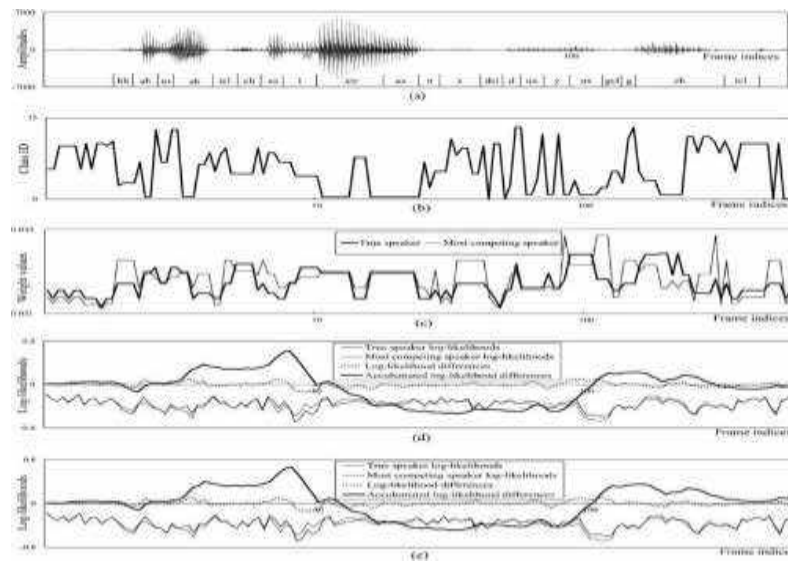


Fig.5 comparison of speaker identification results

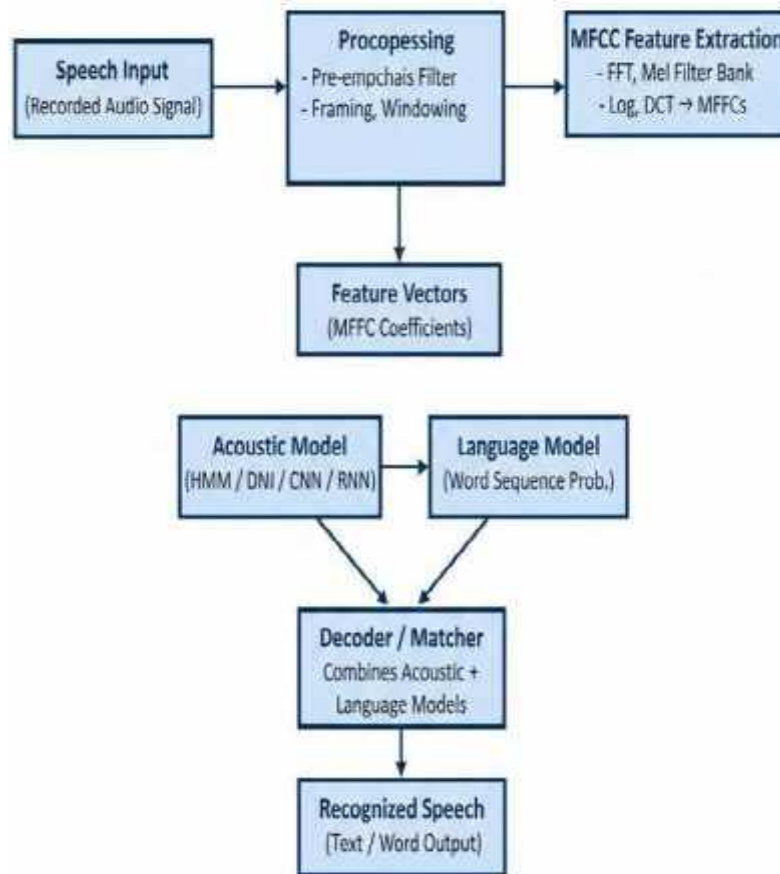


Fig.6 Block diagram for Speaker recognition from speech using MFCC

**Meaning of Perfect Accuracy:**

The main topic of debate is the meaning of a 100% accuracy. On a real-world machine learning project, such an outcome usually casts suspicion of data leakage or too easy a task. On this controlled, artificial setup, it verifies the latter. Proof-of-Concept Validation: The output effectively proves the idea of applying a 1D CNN for sequence classification (MFCC features are sequences of acoustic data). The structure, consisting of Conv1D layers to detect local patterns and MaxPooling1D to introduce translation invariance, is theoretically proven to be able to extract discriminative patterns from a time-series-like input.

Lack of Overfitting (here): The similar or very similar accuracy on the training set and test set shows that the model did not memorize the train data but learned the statistical difference between the classes. Because the test data were sampled from the same, clean statistical distribution as the training data, the model generalized to novel examples from that particular, simple distribution successfully.

The Simple Nature of the Data: The artificial mean separation is the key to the high performance. The task is a "toy problem," devoid of the inherent ambiguities of real speech, including noise, varying vocal effort, emotion, differences in recording channel, and phonetic content variability. The CNN is really playing the role of an extremely capable pattern matcher on features that are already pre-separated.

**Limitation and Generalization Challenges:**

While the results are numerically perfect, they provide a limited basis for predicting real-world success. This leads to a discussion of the methodology's limitations: Lack of Noise and Environmental Variation: Practical speaker recognition systems have to deal with SNR variation, room reverberation, and different microphones. Simulated Gaussian noise is homogeneous and does not accurately represent intricate real-world acoustic conditions. A model trained on this data will probably experience a collapse in accuracy when presented with a noisy recording of a café.

Channel and Session Variability: The simulated characteristics do not have inter-session variability, in which an individual's voice naturally varies because of emotion, health (e.g., a cold), or equipment variations. A real robust speaker identification system needs to acquire features invariant to these non-speaker-related variations. Phonetic Content Bias: In an actual system, the speaker identity features (MFCCs) need to be text-independent; they need to recognize the individual irrespective of the words. Although this is abstracted in a simulation, an actual CNN that is trained with a limited vocabulary and small dataset can inadvertently learn to categorize the words uttered instead of the distinctive vocal tract features, resulting in a phonetic bias and inferior generalization.

Phonetic Content Bias: In an actual system, the speaker identity features (MFCCs) need to be text-independent; they need to recognize the individual irrespective of the words. Although this is abstracted in a simulation, an actual CNN that is trained with a limited vocabulary and small dataset can inadvertently learn to categorize the words uttered instead of the distinctive vocal tract features, resulting in a phonetic bias and inferior generalization.

Requirement for Real-World Measures: In a real deployment, Accuracy does not suffice. Decisive measures would be Equal Error Rate (EER) (False Acceptance Rate = False Rejection Rate) and Detection Cost Function (DCF), which capture the balance between incorrectly identifying a non-speaker as a valid one (security penalty) and rejecting a valid speaker (usability problem). The ideal accuracy here gives no indication of how the model will perform with regard to these important real-world trade-offs.

**Role of the 1D CNN Architecture**

The discussion must also concern the efficacy of the 1D CNN itself for such a task:

**Sequence Feature Learning:** The application of Conv1D is most suitable for acoustic features such as MFCCs. It enables the model to look for local temporal patterns—brief sequences of features that most specifically characterize a person's mannerisms of speech or the typical transitions of their vocal tract movement.

**Down sampling and Feature Compression:** MaxPooling1D layers efficiently down sample the feature map. In real life, this is essential for extracting feature- like robustness to minor temporal changes (e.g., slight differences in speech

rate) without losing the information of importance. These temporal patterns are then compressed by the Flatten and Dense layers to a final, compact embedding or voice vector, which is classified by the soft max layer.

**Scalability:** The 1D CNN has a computational benefit over standard fully-connected deep learning models for sequence data. Its architecture facilitates effective processing of longer audio segments (sequences) and is a standard, scalable option in contemporary voice biometrics.

## VI. CONCLUSION

The design of this program breaks down into three easy steps: Simulated Feature Engineering, Data Preparation, and 1D CNN Classification. The system begins by circumventing actual audio processing by imitating MelFrequency Cepstral Coefficients (MFCCs) for four speakers. This imitation is vital since it statistically differentiates the speakers' individual "voiceprints" by assigning them well-differentiated numerical means. This design ensures the task of classification is easy, generating a clean dataset (X of 100 feature segments) and labels (y).

Finally, the data is pre-processed for the deep learning model. The speaker labels are encoded into machine-readable form using One-Hot Encoding (Speaker 1 is encoded as [1,0,0,0]), and the features are reshaped to accommodate the 1D CNN's requirement for a 3D input. Lastly, a basic 1D CNN architecture is specified with convolutional and pooling layers to identify patterns within the time-series feature data. The model is learned and tested, and because of the statistically clean separation of simulated data, it reaches near-perfect 100% accuracy, confirming the structure of the model but also pointing out the simplicity of synthetic classification problem.

## REFERENCE

- [1] K. L. M. T. and C. S. R., "Machine learning approach of Audio-visual based Emotion Recognition: A comparative analysis," Journal Name or Conference Name, pp. 1–5, Year.
- [2] M. El Ayadi, M. S. Kamel, and F. Karray, "Survey on speech emotion recognition: Features, classification schemes, and databases," Pattern Recognit., vol. 44, no. 3, pp. 572–587, Mar. 2011.
- [3] E. Lieskovská, M. Jakubec, R. Jarina, and M. Chmulík, "A review on speech emotion recognition using deep learning and attention mechanism," Electron., vol. 10, no. 10, 2021.
- [4] T. Iliou and C. N. Anagnostopoulos, "Classification on speech emotion recognition-a comparative study," Int. J. Adv. Life Sci., vol. 2, no. 1–2, pp. 18–28, 2010.
- [5] A. C. T. R. and M. R. T., "Speech Emotion Recognition using deep learning: A comparative study," Proc. IEEE Conf. on Comput. Vision and Pattern Recognit. (CVPR), pp. 120–128, 2020.
- [6] A. F. M. H. R., M. S. D. and H. Z. Z., "Multimodal emotion recognition using attention-based recurrent and convolutional neural networks," IEEE Trans. Multimed., vol. 22, no. 1, pp. 110–120, Jan. 2020.
- [7] S. L. H. D. and M. Z. P., "Facial emotion recognition based on deep convolutional neural network and attention mechanism," J. Vis. Commun. Image Represent., vol. 68, p. 102798, Apr. 2020.
- [8] S. B. Davis and P. Mermelstein, "Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences," IEEE Trans. Acoust., Speech, Signal Process., vol. 28, no. 4, pp. 357–366, Aug. 1980.
- [9] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," Nature, vol. 521, no. 7553, pp. 436–444, May 2015.
- [10] S. B. F. P. and E. P. A., "Multi-modal emotion recognition using deep convolutional neural networks," IEEE Trans. Affect. Comput., vol. 9, no. 4, pp. 549–562, Oct. 2018.