

Deep Learning for the Recognition of Human Speech

Kandula Akhil Sai¹

Department of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Greenfields, Vaddeswaram, Guntur, 522302, India
190030731cse@gmail.com

Kamineni Gopal³

Department of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Greenfields, Vaddeswaram, Guntur, 522302, India
190030708cse1@gmail.com

Ms. G. Lakshmi Sowjanya²

Department of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Greenfields, Vaddeswaram, Guntur, 522302, India
chlks1709@kluniversity.in

Abstract—Humans primarily communicate through speech, and language is their main means of doing so. Emotion is crucial to social connection. Knowing how to recognise the emotions in a speech is important and challenging given that we are dealing with human-machine interaction. Depending on the mood a person is going through, they will express themselves differently. When someone expresses their emotions, each one has a distinctive energy, pitch, and tone fluctuation that is categorised according to the situation. The identification of spoken emotions is thus a future goal for computer vision. The objective of our project is to develop intelligent speech that utilises a convolutional neural network to recognise emotions, which uses a variety of modules for emotion recognition and classifier.

Keywords— Human speech, Emotion, recognition of speech, GMM, Mel-Frequency-Cepstrum Coefficient (MFCC), linear prediction Cepstral coefficients (LPCC), (LSTM)

I. INTRODUCTION

Understanding emotions in human speech has increased recently in order to enhance the effectiveness and naturalness of interactions between humans and machines. Due of the ambiguity in classifying performed and natural emotions, recognising human emotions is a challenging endeavour in and of itself. The spectral and prosodic components that might lead to an accurate assessment of emotions have been the subject of numerous studies. We explored how to classify emotions using human speech and computerised data. We all discussed the process of categorising emotions using the estimated pitch of human speech. We have looked at how to classify emotions and extract speech's acoustic features. Science may be used to analyse and comprehend emotions expressed through language, facial expressions, and physical cues. In order to establish more spontaneous and transparent interactions between humans and robots [1], feelings communicated through audio signals must be regularly detected and recorded. The discrepancy between acoustic features and human emotions, which primarily relies on the discriminative acoustic characteristics gathered for a specific

identification job, makes the Automated Speech Emotion Recognition process challenging. Emotions vary from person to person and are expressed in many ways. Speech emotion has a range of intensities when discussing a variety of topics, and pitch changes are accentuated [4]. Therefore, voice emotion identification is a challenging task in computer vision. The Convolutional Neural Network (CNN) algorithm is used in this study to identify spoken emotions. It has numerous modules for emotion detection and classifiers to distinguish between emotions including happiness, surprise, anger, neutral mood, grief, and so on.

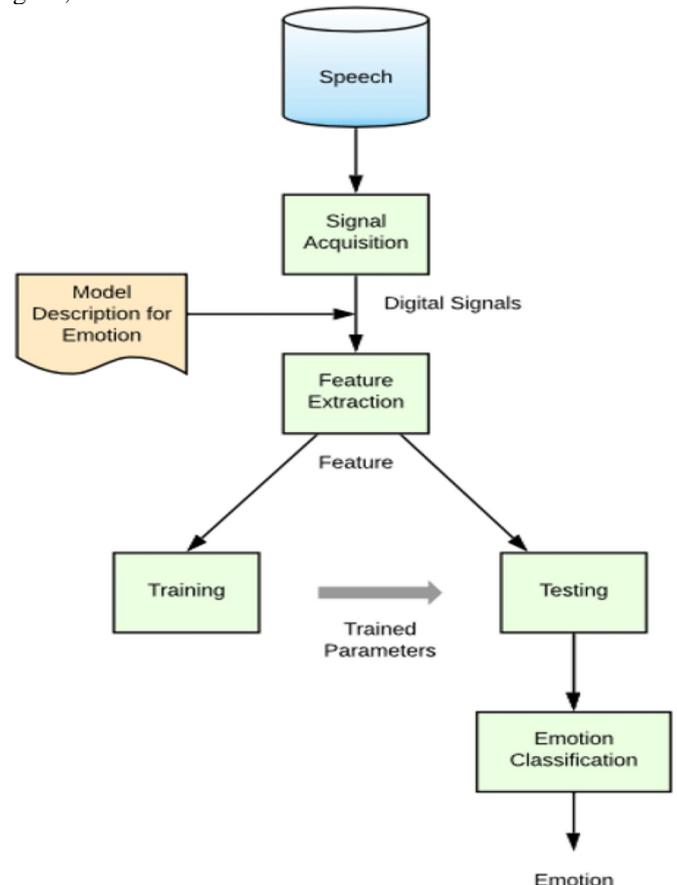


Figure.1 Display the Examples for speech Processing and Feature Analysis

II. LITERATURE SURVEY

In recent years, a significant effort has been made to identify emotions from voice data. We provided a ranking SVM method for synthesising information related emotion recognition to overcome the binary classification difficulty. With each utterer's input treated as a separate query, this ranking approach directs SVM algorithms for certain emotions before aggregating all rankers' predictions to employ multi-class prediction. There are two benefits to ranking SVM. For the purposes of testing and training, it first obtains data specific to speakers [3]. To determine the dominant emotion, it also considers the possibility that each speaker may display a variety of emotions. In two open datasets of performed emotive speech, Berlin and LDC, ranking approaches exceed standard SVM in terms of accuracy. They worked on two things, namely the detection of dangerous situations and recognition of these using their detection algorithm and discovering the specificity of these weapons. Their proposed system also deals with lower quality images. We investigated to find the false signs alerts occurring when a situation is happening on the women.

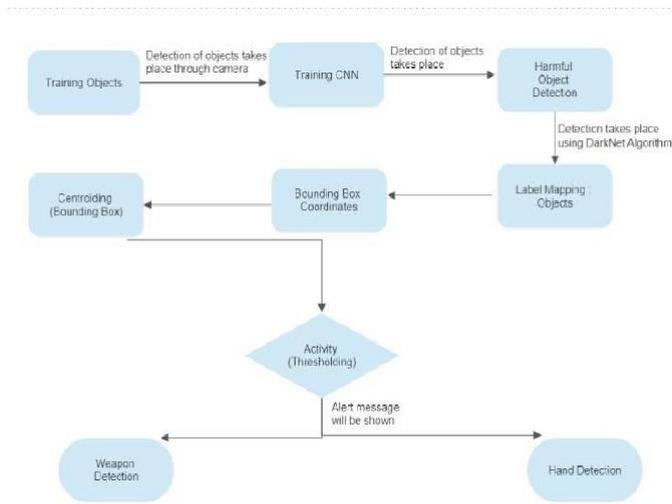


Figure.2 Display the Architectural Diagram

We aimed to improve speaker-independent speech emotion recognition by implementing a three-level speech emotion recognition system. This method organises different emotions into coarse, medium, and fine categories before choosing the appropriate feature based on the Fisher rate. An input parameter for a multi-level GMM-based is the Fisher rate output. We offered a cutting-edge method for identifying emotions in speech signals. The system uses discrete GMM and Mel-frequency cepstral coefficients to describe the speech signals and classifier. Using the private data, this technique split the emotions into six different groups before training and testing the new system. Are contrasted with LFPC to see how effective the proposed method is (LPCC).

The average and highest levels of categorization accuracy reached were, respectively, 80% and 92%. The outcomes also demonstrate that LFPC outperforms traditional characteristics in the classification of emotions. We proposed a multi-dimensional model based on emotion primitives for voice emotion identification [6]. Three distinct emotion primitive values—valence, activation, and dominance—called 020105-4 were combined to produce three dimensions. These factors' values are regarded as falling between [-1, +1]. A text-free, image-based technique was created to analyse emotion primitives, and it produces the highest levels of inter-evaluator agreement. In order to extract acoustic properties like as energy, pitch, and spectral parameters, rule-based estimators as well as fuzzy logic are used. The method is validated using two EMA and VAM datasets: performance emotion and spontaneous spoken emotion [2]. Both datasets were acquired from an Indian-language chat programme.

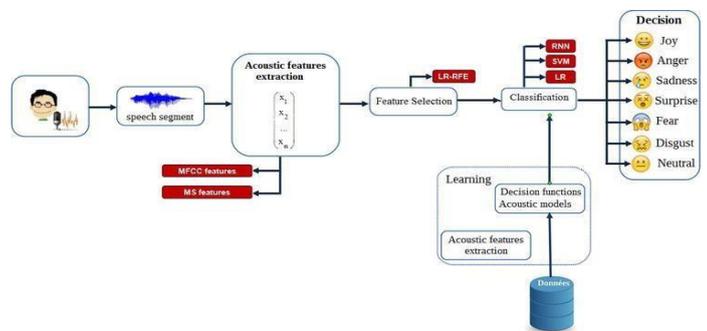


Figure.3 Display the Speech Emotion Feature Extraction

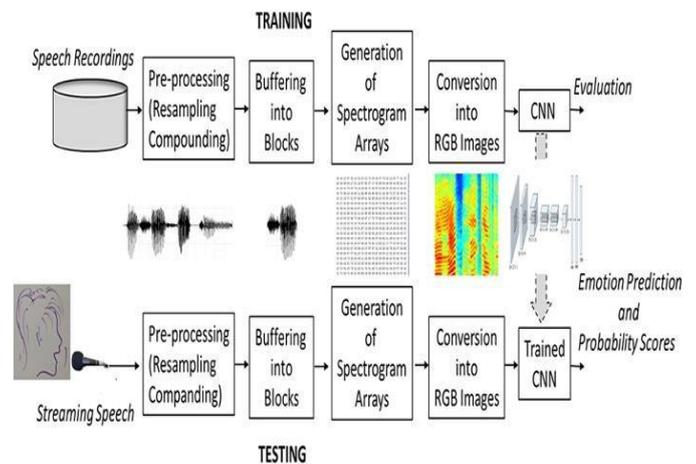


Figure.4 Display the Flow diagram of the process.

III. THEORETICAL ANALYSIS

Emotions that take place in a context of natural speech is, nevertheless, a difficult and little-studied issue. One of the most significant areas of application for digital signal processing is emerging as speech processing. To extract, define, and recognise emotional information from speakers is

the aim of automatic emotion recognition. The initial stage of speech recognition is feature extraction. Numerous feature extraction techniques are proposed or created by researchers. The Mel-Frequency-Cepstrum Coefficient (MFCC) function was utilised in this report to create an automatic Some modifications[2]. An approach for identifying emotions from voice signals is presented in this study. This framework is used to extract speech signal elements that can be utilised to gauge the speaker's emotional state. The automatic recognition and classification of the emotions that are most likely to be present in the speech input is a crucial step in producing expressive speech synthesis. Keywords for speech recognition, MFCC, and GMM. EMOTIONS are important in normal human interactions. Recent research reveals that emotions play a crucial role in our ability to make logical decisions. By letting us express our emotions and receiving feedback, it aids in the formation of connections. When creating human-machine interfaces, this crucial component of human contact must be considered (HMI). A close connection exists between the emotional expression [3]. It is generally known that when the brain's emotional control regions are injured, even though reasoning capacity is intact, the brain is unable to make rational decisions. Most emotional computing applications, including natural language interfaces, e-learning environments, instructional or entertaining games, opinion extraction and mood analysis, etc., humour detection, and computer security, emphasise the critical necessity of persuasive emotion analysis. For instance, emotion recognition is a crucial technique for keeping an eye out for the use of violence or hate speech. It is extremely desirable to take into account emotional states as a crucial component of human-machine interaction in human-machine communication, both at the input and output levels. Feature Extraction from Phonetics Although there are other techniques to implement feature extraction, Mel-based cepstral coefficients are one of the most often used approaches. Then it is claimed that algorithms for detecting emotions that employ prosodic cues are insufficiently accurate[5]. To differentiate between emotions, the phonetic feature offers less information. In actuality, the prosodic aspects of language have fewer separate components than the phonetic ones. As an effective phonetic characteristic for voice recognition, Mel Frequency Cepstral Coefficients (MFCC) of 12–16 dimensions have been used. By adding more independent phonetic functions, the accuracy of emotion identification can be increased if even a little quantity of meaningful information is incorporated in each one. Consequently, an emotion detection algorithm that emphasises precise MFCC classification is needed. We will assign an emotion label to each picture to create such an accurate classification by using MFCC multi-template clustering. The technique has more precision than the conventional approach and is straightforward enough to produce an immediate response even on a basic machine. The human peripheral hearing system is the foundation of MFCC. Below 1000 Hz, the mel frequency scale uses linear spacing, and above 1 kHz, logarithmic spacing is used. The pitch of a 1 kHz tone that is

40 dB over the threshold of human hearing is defined as 1000 mels. Second, a short-time converter is used to convert the analysis frame to the frequency domain. We updated the LSTM cell states using the attentional mechanism, which focused on intercellular communication and took past cell states into consideration. To make calculation of candidate cell state weights simpler, we developed a linked LSTM that controls the impact of previous states on the present state of the cell using just one gate. In contrast to other research, this study modifies the LSTM forget gate with the self-awareness algorithm and concentrates on the calculation of the cell interior. As a result, the forget gate computation is different from the prior LSTM.

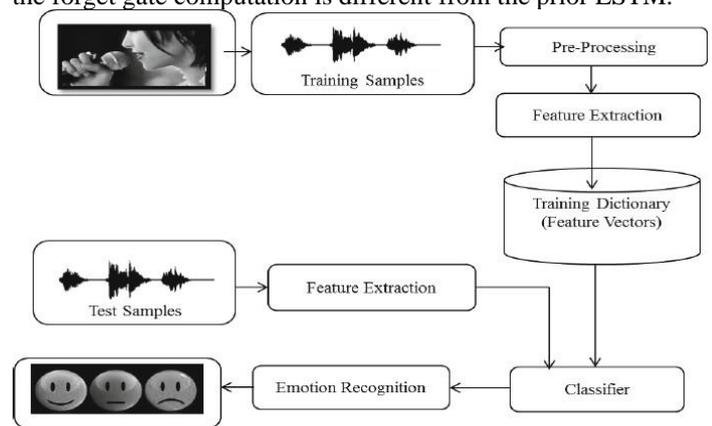


Figure.5 Display the Emotion Recognition

IV. SYSTEM DESIGN

There are various professions that call for knowledge of emotional state. As technology advances, there is more demand on human-computer communication to be more precise and straightforward. Language is increasingly being used as an input-output interface in applications nowadays. Due to the lack of knowledge on the emotional state, this type of encounter may lead to two issues. The first is when someone in a stressed circumstance misunderstands a sentence or a command. Compared to a hearing person, the machine recognises human speech differently. Changes in the spoken signal brought on by vocal tract stress have an impact on accuracy. Lack of emotional state as it relates to the speaker's machine language is the second issue [5]. Such speech has an impact on people and is artificially unreliable. The way a person expresses their emotions depends on their environment and way of life. The way we say also changes as our way of life and our surroundings do, which is another venture in front of the speech emotion reputation gadget. Applications of Speech Emotion Recognition include psychiatric diagnosis, intelligent toys, lie detection, mastering environments, educational software, and detection of the emotional state in phone call middle conversations to provide feedback to an operator or a manager for tracking purposes. The most crucial application for the automatic reputation of feelings from the speech is in vehicle board devices, which keep track of the driver's mental state. Enter speech signal,

pre-processing, characteristic extraction and selection, type, and finally feelings reputation make up the proposed human emotion reputation device [5]. the architecture of the emotional speech reputation mechanism. The degree of naturalness of the database utilised as an input to the speech emotion reputation device determines the accuracy of the emotional speech reputation device. The actual global feelings or the acted-out ones may likewise be included in the database as an input to the speech emotion reputation mechanism.

Applying a database that is compiled from real-world situations makes it more realistic.

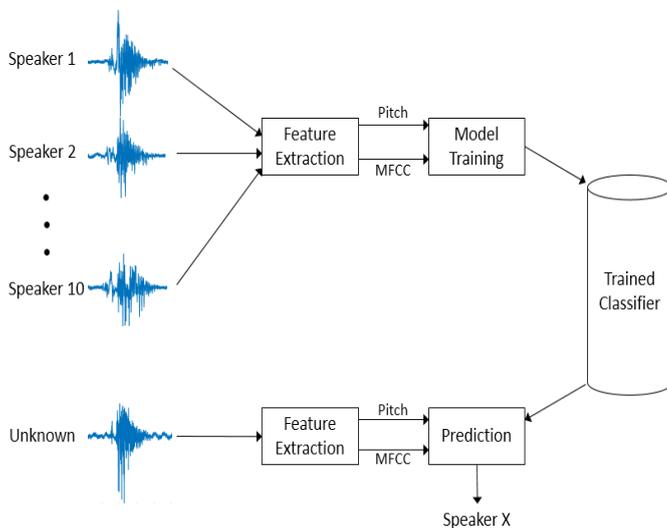


Figure.6 Display the Flow diagram of the process.

V. EXPERIMENTAL INVESTIGATIONS

Body language (posture) and facial mimicry (expressions) are common. Although there may be regional differences in behaviour, it is reasonable to assume that everyone will be able to identify an angry person no matter where they are from. However, this is not clearly the case when it comes to speaking. Any spoken language has two independent components: the language itself, which includes its syntax, vocabulary, and context, and "what" is above the language (how it is expressed) [6]. Language is referred to as "linguistics" in the first half and "paralinguistic" in the second. Even though we may not understand what is being said due to our ignorance of linguistics, we may be able to infer the emotions being expressed, even if paralinguistic elements may vary between languages. Different people may view languages as harsher than others, yet for native speakers of a particular nation, their language has highly special paralinguistic qualities that enable everyone (aside from very specific people) to grasp the emotions being expressed in speech. With this project, we're seeking to allay these worries. We will employ LSTM and CNN to classify conflicting emotions categorise to study the connection between gender and the

emotional content of communication, the speech by speaker gender. Our understanding of which components convey the most emotional information and why is aided by these trials' precision in emotion identification.

VI. RESULTS AND DISCUSSION

```
In [26]: #Calculate Accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: {:.2f}%".format(accuracy*100))
```

Accuracy: 75.97%

Figure.7 Display the Accuracy of the model.

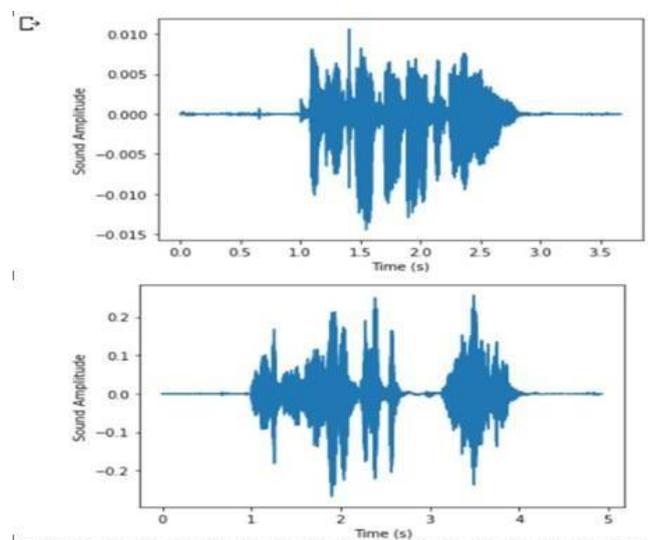


Figure.8 Display the sound amplitude respected to time(s)



Figure.9 Display the Emotion: Calm

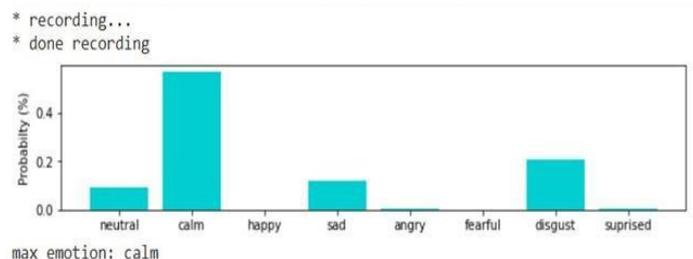


Figure.10 Display the Emotion: Calm

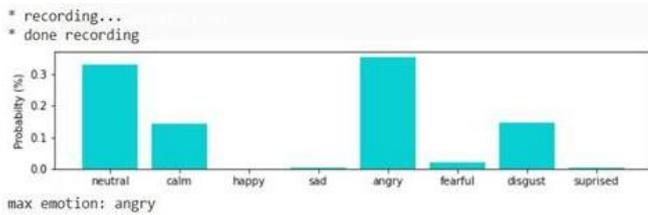


Figure.11 Display the Emotion: Angry

VII. CONCLUSION

We are attempting to address these concerns with this project. To categorize opposing emotions, we will utilize Convolutional Neural Networks and LSTM. To study the association between gender and emotional content of speech, we segregate the speech by speaker gender. Human speech may be used to extract a range of temporal and spectral properties. As inputs to classification algorithms, we employ pitch statistics, Mel Frequency Cepstral Coefficients (MFCCs), and Formants of Speech. These trials' emotion identification accuracy helps us to understand which aspects convey the most emotional information and why. And where emotion recognition has advantages which plays a major role which are improving learning performance, lowering computational complexity, building better generalizable models, decreasing required storage, important for security and healthcare purposes. Which made the model very efficient.

VIII. REFERENCES

- [1] D. Navanith, K. Likhith, M. S. Vardhan and S. Kavitha, "Optimized Sentiment Analysis of Hotel Reviews using Machine Learning Algorithms," 2022 6th International Conference on Electronics, Communication and Aerospace Technology, Coimbatore, India, 2022, pp. 1075-1081, doi: 10.1109/ICECA55336.2022.10009104.
- [2] M.Govindarajan, Sentiment Analysis of Movie Reviews using Hybrid Method of Naive Bayes and Genetic Algorithm, International Journal of Advanced Computer Research (ISSN (print): 2249-7277 ISSN (online): 2277-7970), Volume-3 Number-4 Issue-13 December-2013.
- [3] Apoorv Agarwal, BoyiXie Ilia Vovsha, Owen Rambow, Rebecca Passonneau, Sentiment Analysis of Twitter Data.
- [4] Apoor v Agarwal, Jasneet Singh Sabharwal, End-to-End Sentiment Analysis of Twitter Data, Proceedings of the Workshop on Information Extraction and Entity Analytics on Social Media Data, pages 39–44, COLING 2012, Mumbai, December 2012.
- [5] V.K. Singh, R. Piryani, A. Uddin, P. Waila, Sentiment analysis of movie reviews: A new feature-based heuristic for aspect-level sentiment classification, Conference Paper March 2013, DOI: 10.1109/iMac4s.2013.6526500 .
- [6] Ajinkya Ingle, Anjali Kante, ShriyaSamak, Anita Kumari, Sentiment Analysis of Twitter Data Using Hadoop, International Journal of Engineering Research and General Science Volume 3, Issue 6, November-December 2015, ISSN 2091-2730, www.ijergs.org
- [7] Krishna, P. Venkata, et al. "Learning automata based sentiment analysis for recommender system on cloud." Computer, Information and Telecommunication Systems (CITS), 2013 International Conference on. IEEE, 2013
- [8] Sangeeta, Twitter Data Analysis Using FLUME & HIVE on HadoopFrameWork, Special Issue on International Journal of Recent Advances in in Science, Technology & Management (NCRISTM) ISSN (Online): 2347-2812, Gurgaon Institute of Technology and Management, Gurgaon.
- [9] Sushith, Mishmala. "Semantic Feature Extraction and Deep Convolutional Neural Network-based Face Sentimental Analysis." Journal of Innovative Image Processing 4, no. 3 (2022): 157-164.
- [10] Verma, Jai Prakash, Bankim Patel, and Atul Patel. "Big data analysis: recommendation system with Hadoop framework." Computational Intelligence & Communication Technology (CICT), 2015 IEEE International Conference on. IEEE, 2015.
- [11] Karuppusamy, Dr P. "Artificial Recurrent Neural Network Architecture in Customer Consumption Prediction for Business Development." Journal of Artificial Intelligence and Capsule Networks 2, no. 2 (2020): 111-120.
- [12] Shrote, Khushboo R., and A. V. Deorankar. "Review based service recommendation for big data." Advances in Electrical, Electronics, Information, Communication and Bio-Informatics (AEEICB), 2016 2nd International Conference on. IEEE, 2016.