

## Emotion Detection in Speech:ACNN Neural Network Approach

C.Prabhavathi, P.Vijaya Kumari, K.Chandrasekhar,D.Raghunath Kumar Babu, S.Ghouhar Taj

Assistant Professor(Adhoc) CSE Department,JNTUACEP, [prabhavathi1231@gmail.com](mailto:prabhavathi1231@gmail.com)

Assistant Professor(Adhoc)CSE Department,JNTUACEP, [vijaya.jntuacep@gmail.com](mailto:vijaya.jntuacep@gmail.com)

Assistant Professor(Adhoc)CSE Department,JNTUACEP,[chandra507shiva@gmail.com](mailto:chandra507shiva@gmail.com)

Assistant Professor(Adhoc)CSE Department,JNTUACEP, [raghunath.d29@gmail.com](mailto:raghunath.d29@gmail.com)

Assistant Professor(Adhoc)CSE Department,JNTUACEP,[sgtaj786@gmail.com](mailto:sgtaj786@gmail.com)

**Abstract**—Speechemotionrecognitionisanareaofresearch dedicated to identifying and categorizing emotions expressed throughspeech.Its purposeistocomprehendandinterpretthe emotional content conveyed in spoken words, leveraging signal processing techniques, feature extraction algorithms, and machine learning models. The ultimate aim is to apply this knowledge in diverse applications such as human-computer relation, affective computing, and mental health diagnosis. Although the complexity and variability of emotional speech present significant challenges, recent years have witnessed notable progress through the development of advanced algorithmsandtheutilizationofextensivetrainingdatasets.Here we proposed a CNN network pattern for speech emotion recognitionforthebenchmarkdatasetsandgottheaccuracyof 88% for the convolutional neural network model.

**Keywords**—CNN,Neuralnetworks,Speech,emotionintelligence, datavisualization

### I. INTRODUCTION

Non-verbalcommunicationplaysacrucialrolein human interactions. Apart fromtheliteralmeaning conveyedthrough spokenlanguage,thewaywordsarespokencarriessignificant information. The same spoken text can have multiple interpretations basedonthemanner inwhichit isexpressed. The term 'really' in the English language has versatile applications.It canbe employedto inquire aboutsomething, display admiration, indicate skepticism, or assert a strong statement. Merely understanding the textual content of a spoken phrase is insufficient for accurately interpreting its meaning.Emotionrecognitioninspeechholdsvariousfuture implementations. One such application involves enhancing speechunderstandingbyusingemotionrecognitionasatool. Traditionally, emotion has been considered a disruptive element that hampers the comprehension of spoken text. However, by recognizing and isolating emotions within speech,itmaybepossibletoenhancetheexecutionofspeech comprehension systems.

Multimediapatternidentificationisannewinnovationthat enables the extraction and analysis of large volumes of multimedia information from video and audio sources. In current years, deep learning techniques, particularly using machine learning with deep neural networks, have been extensivelyappliedtoaddressdifferentidentification

problems. However, the challenge lies in the fact that individuals indicate emotions in individual method, and distinguishing between these emotions based on unclear features is a difficult problem, even for humans.

Conventionalmethodsforaddressingthisissueentailthe extraction of basic characteristics and training machine learning models based on these extracted features. These techniqueshave beenconsideredfuturistic forlongtime,but selectingappropriatefeaturestoextractisachallengingtask, and optimizing the results can be protracted.

CNNshaveemergedastheleadingapproachincomputer vision applications and have garnered attention in diverse fields. These networks comprise various elements and particularly designed to learn hierarchical spatial features through the utilization of backpropagation algorithms. Understandingtheidea,advantages,andlimitationsofCNNs is crucial to fully leverage their potential and improve the execution of the recognition of the emotion model.

The proposedimplementation, authorsdeveloped asystem that uses neural networks, specifically CNNs, for emotion recognitioninspeech.Since theproject involvesclassification, a CNN is the natural choice. The model is trained to detect eight human emotions likely happy, sad, neutral, angry,calm,disgust, surprised and fearful along with determiningthegenderofthespeaker.ByutilizingCNNs,we aimtoleverage theirabilitytoautomaticallylearnhierarchical features and upgrade theoverall executionof the recognition in emotion system.

The prefer model is done with the datasets fromRAVDESS and SAVEE. Finally, the efficiency of the trained model is observed by testing against live voice.

The foremost purpose of the research is to: a)To understand speech recognition and its fundamentals. b) To Collect the datasets on Speaker emotion recognition. c) Developthe algorithm for feature extraction. d) Developthe algorithm for Classification. e) Compare the proposed methods with existing works.

Followingarethetheremainingsections:SectionIdealswith studies and researches made by some of the certified researchersthroughoutthe world onspeechemotiondetection and related works. We have used and improved our project modelbasedontheirworkandimplementationsmethods.

Section III provides the basic architecture for the emotion detection using the speech processing technique. This chapter also deals with the explanation of block diagram which is used for implementation with the necessary evaluation models. Section IV deals with the final step of our project, in this section we discuss about the accuracy and losses of training and testing model that we have implemented with the validation loss over time. We will also have a look at the model classification reports. In Section V we conclude our work by explaining the challenges that are present in speech emotion detection and how we have tried to overcome them with the future scope of this project.

## II. RELATED WORK

In their research, Dias Issa et al. described for emotion identification using sound files. They extracted various features. They achieved accuracies of 71.61% for the RAVDESS dataset, 86.1% and 95.71% for different subsets of the EMO-DB dataset [1]. Harshawardhan worked using MFCC features and an LSTM algorithm. They obtained an 84.81% accuracy [2]. Deepak Bharti and Poonam Kukana presented a with MSVM (Multiple Support Vector Machine) classifier. They achieved a high accuracy rate of 97% on the RAVDESS dataset using feature extraction (GFCC) and feature selection (ALO) techniques. On existing datasets, they obtained an accuracy of 79.48% with feature extraction using MFCC [3].

Anusha Koduru et al. focused on pre-processing audio samples by removing noise using filters. Their results showed accuracies of 70% with SVM, 85% with decision tree, and 65% with LDA [4]. Authors proposed a deep recurrent Neural Network process for learning emotion variations. Their methodology was applied to the dataset called RAVDESS, achieving correctness over 80% [5].

Christy et al. achieved an accuracy of 78.20% using CNN on the RAVDESS dataset [6]. Ting-Wei Sun used a hybrid algorithm and this algorithm achieved high accuracies on FAU and eNTERFACE databases [7].

Latifa et al. developed a latest database in Urdu speech and evaluated the performance of a model using SVM classifier [8]. Jalal et al. proposed and compared bimodels, CNN and attention and bi-LSTM and attention, for emotion recognition [9]. Satya et al. aimed to provide an extensive survey that highlights the requirements of speech and vision systems, considering both hardware and software aspects. [10].

Abdelhamid et al. concentrated on developing a novel data boosting technique to enhance the emotions database by introducing additional illustrative through the controlled inclusion of noise selection [11]. In a related study, Their approach employed mel-frequency log spectrogram to extract relevant evidence from the emotional speaker database and employed a 2D DCNN for analysis [12].

Middya et al. conducted an extensive investigation into fusion on a model-level techniques to choose the most effective multimodal model for emotion recognition [13]. Bakhshi et al. introduced CyTex, a revolutionary speech-to-image transformation technique that leverages the basic frequency of every speech shape to directly convert the raw speech signal into visually textured images [14]. Bagadi et al. conducted a study to examine the influence of meta-heuristic for feature selection methods on speech-based emotion identification [15]. Zhong research aims to accomplish speech

emotion recognition in Chinese by machine learning using CNN [16].

## III. PROPOSED FRAMEWORK

The audio files are preprocessed by adding noise and shifting time, pitch and speed to improve the model's ability to generalize. Features are then extracted using the MFCC method and stored in a csv file. The final step involves building a Convolution Neural Network (CNN) model for classification. The input features are MFCCs and the model is trained to classify the audio files into 8 different emotions.

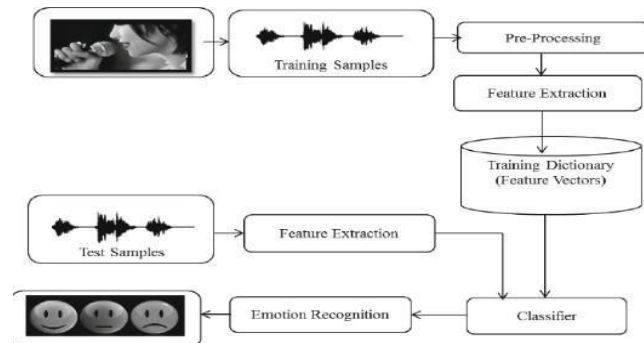


Figure 1: Block Diagram for Model

### Detail of datasets:

**RAVDESS**-Dataset incorporate over 1500 auditory folder input starting 24 distinct actors. Twelve males and Twelve females wherever these performers record little auditory in dissimilar feelings. Figure 2 represents the Ravdess Dataset Visualization

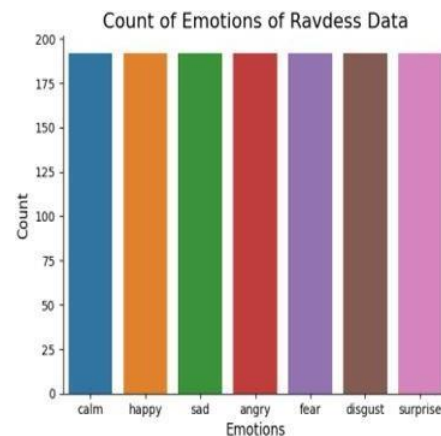


Figure 2: Ravdess Dataset Visualization

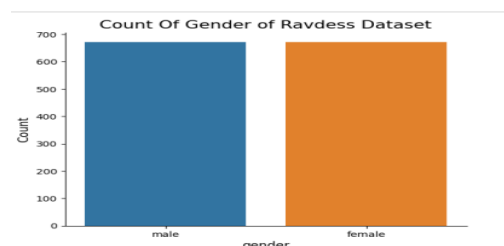
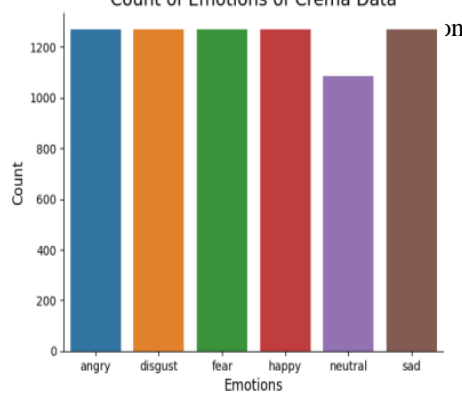


Figure 3: Gender of Ravdess Dataset Visualization

The audio files are named in a consistent manner where the 7th character reflects the various emotions they represent. Figure 3 and 4 depicts the Gender of Ravdess Dataset Visualization and Crema D Dataset Visualization respectively.



**Pre-processing:** To enhance the model's ability to generalize, we create additional synthetic data samples by making slight modifications to our original training set. For audio data, these perturbations include noise injection, time shifting, pitch alteration, and dash modification. The aim is to make our model fixed to such variations and improve its ability to generalize across different conditions. We add two types of noise to make our model training more efficient. We use white noise and Gaussian noises for this purpose.

**Feature Extraction:** To enable our model to study from the auditory files, the next phase is to extract the features from them. For feature extraction, we utilize the Python library called LibROSA, which is widely employed for audio analysis. This library offers a range of tools and functions specifically designed to process audio data effectively. The feature is being extracted and put in a csv file. Feature extraction is the process of identifying and selecting the most relevant and descriptive characteristics or attributes of a set of data and transforming them into a new set of features that can be used in further analysis or modeling. It is a crucial step in many machine learning algorithms and benefits in refining the performance and accuracy of the models. Proposed flow diagram is shown in figure 5.

MFCC is used for feature extraction here.

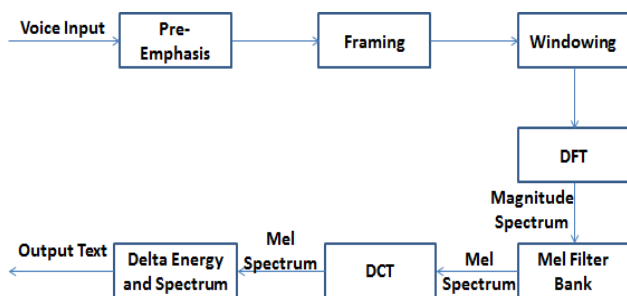


Figure 5: MFCC block diagram

The MFCC are a set of structures used in speech and audio. The steps to compute MFCCs include:

1. Pre Emphasis: Amplifying the high occurrence components of the indication by applying a high-pass filter.

$$(t) = (t) \cdot x(t-1) \tag{1}$$

Where  $(t)$  is the input signal and  $a$  is a pre-emphasis coefficient.

2. Windowing: Dividing the signal into overlapping frames and applying a window function.

$$(n) = (n) \cdot x(n) \tag{2}$$

Where  $(n)$  is the window function and  $x(n)$  is the input signal.

3. Spectral analysis: Computing the power spectrum of each frame using the Fast Fourier Transform (FFT).

$$\sum (n) N - 1 n = 0 \cdot e^{-\frac{j2\pi kn}{N}} \tag{3}$$

Where  $X(k)$  is the Fourier Transform of  $x(n)$ ,  $N$  is the sum of models, and  $k$  is the frequency index.

4. Mel-scale transformation: Applying a non-linear transformation to the power spectrum to model the way human ear perceives different frequencies.

$$(f) = 2596 \cdot \log_{10} \left( 1 + \frac{m}{700} \right) \tag{4}$$

Where  $m$  is frequency (Hz) and  $H(m)$  is the corresponding Mel frequency. Next, the power spectrum is mapped to Mel-scale using triangular overlapping filters. The equation is:

$$(k) = \sum_{m=0}^{p-1} h_{(k)} \cdot X^2(k) \tag{5}$$

Where  $h_{(k)}$  is the triangular filter,  $X^2(k)$  is the squared magnitude of  $X(k)$ , and  $P$  is the number of filters.

5. Cepstral analysis: Taking the logarithm of the Mel-scaled power spectrum and applying the Discrete Cosine Transform (DCT) to convert it to the cepstral domain.

$$(i) = \sum_{m=0}^{p-1} (k) \cdot \cos\left(\frac{\pi i k}{P}\right) \tag{6}$$

Where  $(i)$  is the cepstral coefficient and  $P$  is the sum of Mel-scale coefficients.

6. Cepstral coefficients: Selecting the first  $N$  coefficients as the MFCCs, where  $N$  is a parameter that can be adjusted based on the desired level of detail.  $(i) = (i)$  where  $i = 1, 2, \dots, N$  and number of desired cepstral coefficients represent in  $N$ .

• CNN Architecture

The Convolution Layer applies learnable filters to small windows of the input matrix, producing a 2-dimensional activation matrix that captures visual features. The Completely Connected Layer connects all participations to neurons, allowing for more global interactions. The Final Output Layer predicts the likelihood of each image belonging to different classes. Model building and tuning is an inefficient process, starting with a simple architecture and gradually adding complexity. The best-performing model achieved a validation accuracy of slightly over 85%.

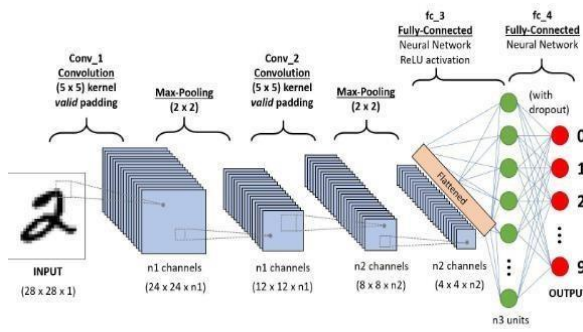


Figure6:CNNArchitecture

An algorithm for linear classification is support vector machine. The following can be used to condense the mathematical steps in an SVM algorithm:

1. Input data representation:

Let  $X$  be  $d$ -dimensional feature space and  $Y$  be the target space where  $y = -1, +1$ . The data is represented as a set of  $n$  samples  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  where  $x_i$  is a  $d$ -dimensional feature vector and  $y_i$  is the target value.

2. Constructing the optimization problem:

This can be created as a minimization with constraints. The objective is to find the weights and bias such that: maximize:  $\frac{1}{\|W\|}$

Subject to:  $(WX_i + b) \geq 1, i = 1, 2, \dots, n$  where  $\|W\|$  is the Euclidean norm of  $w$ .

3. Solving the optimization problem:

Using a quadratic programming (QP) solver, the optimization problem can be resolved. The solution to the problem gives the values of  $w$  and  $b$  that define the hyper plane

4. Making predictions: Given a new sample  $x$ , its class can be predicted as:

$$y = \text{sign}(w \cdot x + b) \tag{7}$$

Where  $\text{sign}(x)$  returns  $+1$  if  $x \geq 0$  and  $-1$  if  $x < 0$

5. Information that cannot be detached linearly: Data that cannot be separated linearly can be translated into a high-dimensional space where a linear boundary can be found using a kernel method.

The ReLU is an activation function that outputs the input value if it is positive and 0 otherwise. ReLU is computationally efficient associated to further activation functions, which have more complex formulas and higher computational costs. ReLU is also advantageous as it is not affected by vanishing gradients, unlike Sigmoid and Tanh, which can slow down learning in a network. Its formula,  $f(x) = \max(0, x)$ , ensures that the output ranges from 0 to infinity. ReLU is widely employed in neural networks, particularly in Convolutional Neural Networks (CNNs), and is often the default choice for an activation function. Figure 6 and 7 represent the graph of the ReLU activation and graph of the sigmoid activation function respectively.

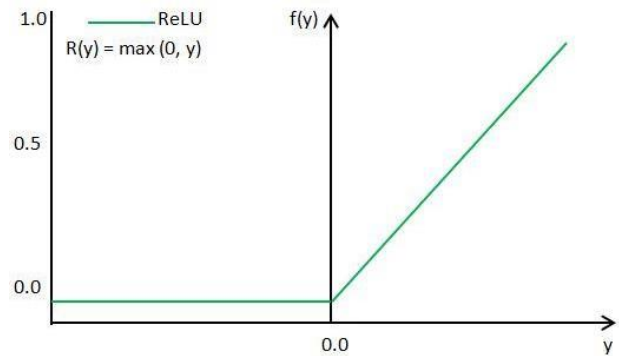


Figure7:ReLUactivationfunction

ReLU offers computational efficiency and faster training/operation due to its simple arithmetic. It promotes sparsity, which means many weights in the network become zero, leading to compact models with enhanced predictive ability and reduced overfitting. In a sparse network, neurons are more likely to focus on important features of the problem, resulting in more meaningful processing. For example, in a face detection model, certain neurons may specialize in identifying specific facial components; remaining inactive when irrelevant features are present.

The logistic function is a commonly used example of a sigmoid function, defined by the formula:

$$f(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} = 1 - (-X) \tag{8}$$

Sigmoid functions are typically monotonic and have bell-shaped first derivatives. The cumulative distribution functions of various probability distributions, such as the error function and the arctan function, are also sigmoidal. A sigmoid function is characterized by horizontal asymptotes and exhibits convexity for values less than a certain point, and concavity for values greater than that point, often 0.

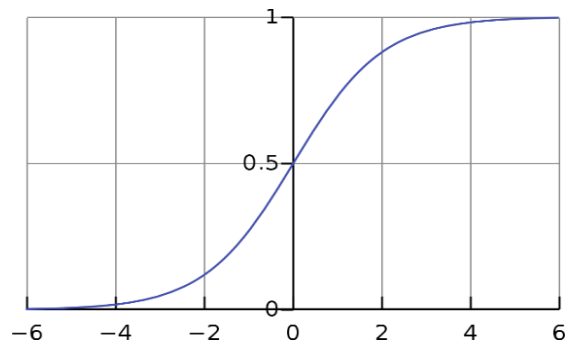


Figure8:Sigmoidactivationfunction.

Classification

The given signal will be analyzed by the machine and it will be classified into the following 7 classes. The output may be either one of these will be displayed according to given emotion. Angry, calm, fearful, happy, sad, disgust, Surprise.

Evaluation Metrics

The goal of assessment is to identify as many instances as possible from a population for a screening technique; hence false negatives should be kept to a minimum at the cost of increasing false positives. As a result, three main measurements must be established: accuracy (ACC), false positive rate (FPR) and true positive rate (TPR). In medical

language, the first parameter is referred to as sensitivity (SEN) and is written as Equation:

$$TPR = SEN = \frac{TP}{P} \quad (9)$$

where TP stands for true positives and P is for positive events. The estimation of the second period, false positive amount, expressed as Equation:

$$FPR = \frac{FP}{N} \quad (10)$$

The population's cumulative number of negative occurrences is N, while the proportion of false positives is FP, and number of true negatives samples is N. This statistic, on the other hand, is better understood as the ratio of genuine negatives to real negatives, known in medical language as the specificity (SPEC), which is given as Equation:

$$TNR = SPEC = \frac{TN}{N} = 1 - FPR \quad (11)$$

Finally, accuracy determines the stability between actual positives and accurate negatives. Figure 9 shows the Evaluation Matrix of the proposed method.

$$ACC = \frac{(TP+TN)}{(P+N)} \quad (12)$$

		Predicted		
		0	1	
Actual	0	TN	FP Type I error	Specificity = TN/(TN+FP)
	1	FN Type II error	TP	Recall or Sensitivity = TP/(TP+FN)
		Negative Rate = TN/(FN+TN)	Precision = TP/(TP+FP)	

Figure 9: Evaluation Matrix

IV. RESULTS AND DISCUSSIONS

The recommended technique was verified on the RAVDESS database and accomplished a correctness of 71% with a f1-score of 0.71. The method outperforms the existing methods [5] and [7] which have accuracy rates of more than 80% and 78.20% respectively. The error rate of the proposed method was 0.5. Figure 10 depicts the Epoch having an accurateness of 71%. And a f1-score of 0.71. The model has been worked for the following dataset and obtained model is having an accurateness of 71% and a f1-score of 0.71. The Accurateness Arc and The Defeat Arc are shown in the Figure 11 and 12 respectively.

Epoch 68/100  
 57/57 [=====] - ETA: 15s - loss: 0.6414 - accuracy: 0.820 - ETA: 14s - loss: 0.5798 - accuracy: 0.847 - ETA: 14s - loss: 0.5376 - accuracy: 0.855 - ETA: 15s - loss: 0.5178 - accuracy: 0.852 - ETA: 15s - loss: 0.5014 - accuracy: 0.851 - ETA: 14s - loss: 0.4855 - accuracy: 0.853 - ETA: 13s - loss: 0.4733 - accuracy: 0.854 - ETA: 13s - loss: 0.4617 - accuracy: 0.855 - ETA: 13s - loss: 0.4541 - accuracy: 0.855 - ETA: 12s - loss: 0.4458 - accuracy: 0.856 - ETA: 12s - loss: 0.4401 - accuracy: 0.856 - ETA: 11s - loss: 0.4350 - accuracy: 0.857 - ETA: 11s - loss: 0.4298 - accuracy: 0.858 - ETA: 11s - loss: 0.4261 - accuracy: 0.858 - ETA: 10s - loss: 0.4232 - accuracy: 0.858 - ETA: 10s - loss: 0.4215 - accuracy: 0.858 - ETA: 9s - loss: 0.4195 - accuracy: 0.858 - ETA: 9s - loss: 0.4183 - accuracy: 0.85 - ETA: 9s - loss: 0.4175 - accuracy: 0.85 - ETA: 8s - loss: 0.4170 - accuracy: 0.85 - ETA: 8s - loss: 0.4164 - accuracy: 0.85 - ETA: 8s - loss: 0.4157 - accuracy: 0.85 - ETA: 8s - loss: 0.4153 - accuracy: 0.85 - ETA: 7s - loss: 0.4152 - accuracy: 0.85 - ETA: 7s - loss: 0.4150 - accuracy: 0.85 - ETA: 7s - loss: 0.4148 - accuracy: 0.85 - ETA: 6s - loss: 0.4149 - accuracy: 0.85 - ETA: 6s - loss: 0.4153 - accuracy: 0.85 - ETA: 6s - loss: 0.4159 - accuracy: 0.85 - ETA: 6s - loss: 0.4165 - accuracy: 0.85 - ETA: 5s - loss: 0.4170 - accuracy: 0.85 - ETA: 5s - loss: 0.4174 - accuracy: 0.85 - ETA: 5s - loss: 0.4178 - accuracy: 0.85 - ETA: 5s - loss: 0.4183 - accuracy: 0.85 - ETA: 4s - loss: 0.4190 - accuracy: 0.85 - ETA: 4s - loss: 0.4197 - accuracy: 0.85 - ETA: 4s - loss: 0.4204 - accuracy: 0.85 - ETA: 4s - loss: 0.4212 - accuracy: 0.85 - ETA: 4s - loss: 0.4219 - accuracy: 0.85 - ETA: 3s - loss: 0.4225 - accuracy: 0.85 - ETA: 3s - loss: 0.4232 - accuracy: 0.85 - ETA: 3s - loss: 0.4237 - accuracy: 0.85 - ETA: 3s - loss: 0.4243 - accuracy: 0.85 - ETA: 2s - loss: 0.4249 - accuracy: 0.85 - ETA: 2s - loss: 0.4255 - accuracy: 0.85 - ETA: 2s - loss: 0.4259 - accuracy: 0.85 - ETA: 2s - loss: 0.4263 - accuracy: 0.85 - ETA: 1s - loss: 0.4266 - accuracy: 0.85 - ETA: 1s - loss: 0.4269 - accuracy: 0.85 - ETA: 1s - loss: 0.4272 - accuracy: 0.85 - ETA: 1s - loss: 0.4275 - accuracy: 0.85 - ETA: 1s - loss: 0.4279 - accuracy: 0.85 - ETA: 0s - loss: 0.4282 - accuracy: 0.85 - ETA: 0s - loss: 0.4286 - accuracy: 0.85 - ETA: 0s - loss: 0.4289 - accuracy: 0.85 - ETA: 0s - loss: 0.4292 - accuracy: 0.85 - ETA: 0s - loss: 0.4296 - accuracy: 0.85 - 14s 251ms/step - loss: 0.4300 - accuracy: 0.8519 - val\_loss: 0.3931 - val\_accuracy: 0.8752  
 Restoring model weights from the end of the best epoch.  
 Epoch 00968: early stopping

Figure 10: Epoch

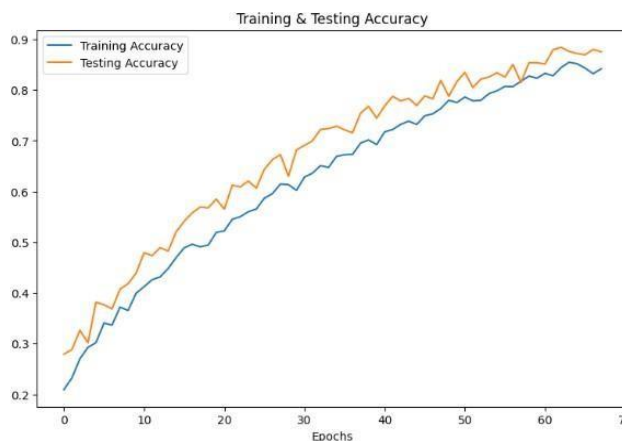


Figure 11: Accuracy Arch of the CNN Model

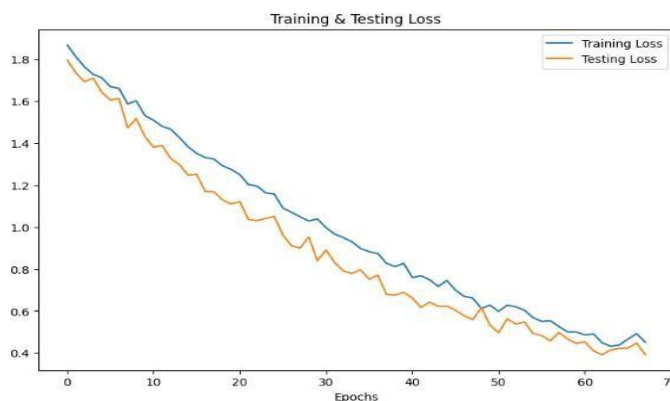


Figure 12: Loss Arch of the CNN Model

Confusion matrix for CNN model is represented in the figure 13 with six emotion details.

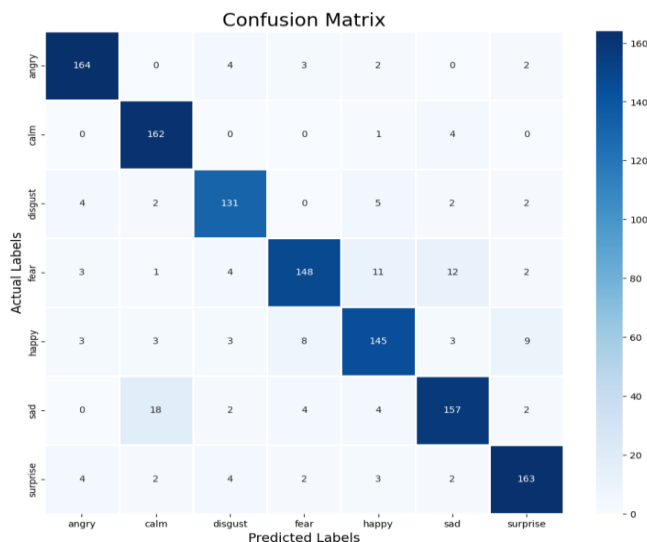


Figure 13: Confusion matrix for CNN model

The Classification report for CNN model and Classification matrix for SVM is represented in the figure 14 and 16 respectively.

	precision	recall	f1-score	support
angry	0.92	0.94	0.93	175
calm	0.86	0.97	0.91	167
disgust	0.89	0.90	0.89	146
fear	0.90	0.82	0.86	181
happy	0.85	0.83	0.84	174
sad	0.87	0.84	0.86	187
surprise	0.91	0.91	0.91	180
accuracy			0.88	1210
macro avg	0.88	0.89	0.88	1210
weighted avg	0.88	0.88	0.88	1210

**Figure 14:** Classification report for CNN model

Comparison result with proposed method with existing techniques with the newest findings in the field as shows in table 1.

TABLE I. COMPARISON BETWEEN PROPOSED AND EXISTING WORK

Method/Algorithm	Dataset	Accuracy	Error
DeepC-RNN approach[5]	RAVDESS	80%	20%
Convolutional neural network[7]	RAVDESS	78.2%	21.8%
Residual Convolutional Neural Network (R-CNN)	FAU	85.8%	14.2%
CNN[Proposed]	RAVDESS	88.42%	11.58

## V. CONCLUSIONS AND FUTURE SCOPE

The proposed model was built using both SVM and CNN algorithms and was tested on the RAVDEES dataset. The highest accuracy of 85% was achieved using the CNN algorithm. There is scope for improvement by adding more data to the dataset and increasing robustness by adding more noise. The model has been built using SVM and CNN. When the model is done using SVM we obtained an accuracy of 72%. When the model is built using the CNN the accuracy we obtained is 85%. The dataset used is RAVDEES. When we done it with the other datasets the accuracy went below 70%. So we went ahead with RAVDEES data with CNN algorithm.

In future work researchers can increase the accurateness of the archetypal by adding a higher size of data which we cannot be able to do due to limitations of the ability of our device. Also we can make the model more robust adding a more noise to the dataset.

## REFERENCE

- [1] Issa, D., Demirci, M. F., & Yazici, A. (2020). Speech emotion recognition with deep convolutional neural networks. *BioMedical Signal Processing and Control*, 59, 101894.
- [2] Bhandari, S. U., Kumbhar, H. S., Harpale, V. K., & Dhamale, T. D. (2022). On the Evaluation and Implementation of LSTM Model for Speech Emotion Recognition Using MFCC. In *Proceedings of International Conference on Computational Intelligence and Data Engineering* (pp. 421-434). Springer, Singapore
- [3] Bharti, D., & Kukana, P. (2020, September). A hybrid machine learning model for emotion recognition from speech signals. In *2020 International Conference on Smart Electronics and Communication (ICOSEC)* (pp. 491-496). IEEE.
- [4] Koduru, A., Valiveti, H. B., & Budati, A. K. (2020). Feature extraction algorithm to improve the speech emotion recognition rate. *International Journal of Speech Technology*, 23(1), 45-55.
- [5] Kumaran, U., Radha Rammohan, S., Nagarajan, S. M., & Prathik, A. (2021). Fusion of mel and gammatone frequency cepstral coefficients for speech emotion recognition using deep C-RNN. *International Journal of Speech Technology*, 24(2), 303-314
- [6] Christy, A., Vaithyasubramanian, S., Jesudoss, A., & Praveena, M. D. (2020). Multimodal speech emotion recognition and classification using convolutional neural network techniques. *International Journal of Speech Technology*, 23(2), 381-388.
- [7] Sun, T. W. (2020). End-to-end speech emotion recognition with gender information. *IEEE Access*, 8, 152423- 152438.
- [8] M. A. Jalal, R. Milner, and T. Hain, "Empirical interpretation of speech emotion perception with attention-based models for speech emotion recognition," *Proc. Interspeech 2020*, 2020.
- [9] Aggarwal, A., Srivastava, A., Agarwal, A., Chahal, N., Singh, D., Alnuaim, A. A., ... & Lee, H. N. (2022). Two-way feature extraction for speech emotion recognition using deep learning. *Sensors*, 22(6), 2378.
- [10] Yadav, S. P., Zaidi, S., Mishra, A., & Yadav, V. (2022). Survey on machine learning in speech emotion recognition and vision systems using a recurrent neural network (RNN). *Archives of Computational Methods in Engineering*, 29(3), 1753-1770.
- [11] Abdelhamid, A. A., El-Kenawy, E. S. M., Alotaibi, B., Amer, G. M., Abdelkader, M. Y., Ibrahim, A., & Eid, M. M. (2022). Robust speech emotion recognition using CNN+ LSTM based on stochastic fractal search optimization algorithm. *IEEE Access*, 10, 49265-49284.
- [12] Bhangale, K., & Mohanaprasad, K. (2022). Speech emotion recognition using mel frequency log spectrogram and deep convolutional neural network. In *Futuristic Communication and Network Technologies: Select Proceedings of VICFCNT 2020* (pp. 241-250). Springer Singapore.
- [13] Middya, A. I., Nag, B., & Roy, S. (2022). Deep learning based multimodal emotion recognition using model-level fusion of audio-visual modalities. *Knowledge-Based Systems*, 244, 108580.
- [14] Bakhshi, A., Harimi, A., & Chalup, S. (2022). CyTex: Transforming speech to textured images for speech emotion recognition. *Speech Communication*, 139, 62-75.
- [15] Bagadi, K. R., & Sivappagari, C. M. R. (2023). An evolutionary optimization method for selecting features for speech emotion recognition. *TELKOMNIKA (Telecommunication Computing Electronics and Control)*, 21(1), 159-167.
- [16] Zhong, Z. (2023, January). Speech emotion recognition based on SVM and CNN using MFCC feature extraction. In *International Conference on Statistics, Data Science, and Computational Intelligence (CSDSCI 2022)* (Vol. 12510, pp. 445-452). SPIE.
- [17] Kukreja, V., & Dhiman, P. (2020, September). A Deep Neural Network based disease detection scheme for Citrus fruits. In *2020 International conference on smart electronics and communication (ICOSEC)* (pp. 97-101). IEEE.
- [18] Dhiman, P., Kukreja, V., Manoharan, P., Kaur, A., Kamruzzaman, M. M., Dhaou, I. B., & Iwendi, C. (2022). A novel deep learning model for detection of severity level of the disease in citrus fruits. *Electronics*, 11(3), 495.
- [19] A. Panwar, R. Yadav, K. Mishra, and S. Gupta, "Deep learning techniques for the real time detection of Covid 19 and pneumo



International Conference on Smart Technologies, Proceedings, 2021  
, pp. 250–253, doi:  
10.1109/EUROCON52738.2021.9535604

- [20] C. Bhatt, I. Kumar, V. Vijayakumar, K. U. Singh, and A. Kumar, "The state of the art of deep learning models in medical science and their challenges," *Multimed. Syst.*, vol. 27, no. 4, pp. 599–613, 2021, doi: 10.1007/s00530-020-00694-1
- [21] Sharma, N., Chakraborty, C., & Kumar, R. (2022). Optimized multimedia data through computationally intelligent algorithms. *Multimedia Systems*, 1-17