

A  
**ARTIFICIAL INTELLIGENCE (CSE3705)**  
PROJECT REPORT  
on  
**EmoDetect**  
*Emotion Recognition Using Facial and Speech Analysis*

Submitted by

**Group – 14**

Ananya Aggarwal (220378)

Aastha Singh (220387)

Siddhika Sinha (220388)

Peehu Khandelwal (220623)

Submitted to

**Dr. Hirdesh Kumar Pharasi**

(Assistant Professor, School of Engineering and Technology)



SCHOOL OF ENGINEERING AND  
TECHNOLOGY  
BML MUNJAL UNIVERSITY  
GURUGRAM-122413  
December 2024

## Table of Contents

<b>Sno.</b>	<b>Chapter Name</b>	<b>Page Number(s)</b>
1	Candidate's Declaration and Supervisor's Declaration	5
2	Acknowledgement	6
3	Abstract	7
4	Introduction	8
5	Literature Review	9 – 10
6	Exploratory Data Analysis	11 – 13
7	Methodology	14 – 17
8	Results and Discussions	18 – 20
9	Conclusion and Future Scope	21
10	References	22

## List of Figures

Figure Number	Figure Title	Page Number
1	Sample Images for “Happy” & “Disgust” Emotions	14
2	Model Performance of Emotion Detection using Speech Analysis	19
3	Confusion Matrix for the Speech Analysis Model	19
4	Model Performance of Emotion Detection using Facial Analysis	20
5	Confusion Matrix for the Facial Analysis Model	20
6	Emotion Detection using Facial Expression	21
7	Emotion Detection using Facial as well as Speech Analysis	21

## List of Tables

<b>Table Number</b>	<b>Table Description</b>	<b>Page Number</b>
1	Literature Review for Emotion Recognition using Facial and Speech Analysis	11

## **Candidate's Declaration**

We hereby declare that the project titled "**EmoDetect: Emotion Recognition Using Facial and Speech Analysis**" is the result of our collective effort as a group of four members. This project was undertaken as a partial fulfillment of the academic requirements for Bachelors of technology at BML Munjal University.

We affirm that this project work is original and has not been submitted for any other academic purpose. All the data, information, and sources used in this project have been duly acknowledged and referenced. Any contributions from external sources have been appropriately cited in the bibliography.

We further declare that each member of the group has actively participated in the research, analysis, and compilation of this project. This project is a true representation of our work and understanding of the subject matter.

---

**Ananya Aggarwal**

---

**Aastha Singh**

---

**Siddhika Sinha**

---

**Peehu Khandelwal**

## **Supervisor's Declaration**

This is to certify that the above statement made by the candidates are correct to the best of my knowledge.

---

**Dr. Hirdesh Kumar Pharasi**  
**Assistant Professor,**  
**Department of Computer Science & Engineering**

## **Acknowledgement**

We are highly grateful to Dr. Hirdesh Pharasi, Assistant Professor, BML Munjal University, Gurugram, for providing supervision to carry out the Artificial Intelligence Project from September – November 2024. Dr. Pharasi has provided great help in carrying out our work and is acknowledged with reverential thanks. Without wise counsel and able guidance, it would have been impossible to complete the training in this manner. We would like to express thanks profusely to thank Dr. Pharasi, for stimulating us from time to time. We would also like to thank the entire team at BML Munjal University. We would also thank our friends who devoted their valuable time and helped us in all possible ways toward successful completion.

**Group - 14**

## Abstract

The **“EmoDetect: The “Emotion Recognition Using Facial and Speech Analysis”** project delivers the long-awaited solution for the real-time recognition of emotions based on AI facial expression and speech analysis. The proposed system employs features of deep learning and multimodal data analysis to distinguish between happy, sad, angry, and stressed mood. Self-organized maps have a wide application in; Customer care and service, Learning and education, and Games and entertainment. Also, for this type of work use this system can potentially promising as an early warning sign of mental strain and for individual compliance. This makes the union of facial and vocal analysis as shown here a versatile model perfect for real world applications to increase the practicality of AI by making them more emotionally intelligent across fields.

**Keywords:** *Facial Expression Analysis, Emotion Detection, Speech Recognition, Deep Learning, Human-Computer Interaction, Real-Time Emotion Recognition, Sentiment Analysis, Real-Time Emotion Recognition*

## Chapter – 1

### Introduction

As we see increased demand for such emotionally aware systems, artificial intelligence is being utilized more often to fill the void between what humans feel and what machines understand. Communication includes emotion, how we make decisions, our behaviour, and interactions. Recognizing the need for machines to interpret these emotional cues accurately, **EmoDetect**: A comprehensive solution, Emotion Recognition Using Facial and Speech Analysis provides. Unlike traditional systems that often do so with one modality, it combines facial expression recognition and speech recognition to determine the human emotion in real time. Such unimodal systems are inherently prone to inaccuracies in complex environments, mostly due to their unreliability and limited practical application. In response to this challenge, EmoDetect combines the facial and vocal cues to create a more holistic and robust emotion recognition system that can perform in different real-world scenarios.

Advanced deep learning techniques are utilized by EmoDetect to analyze sense data (visual and auditory) and with this ability for creating a discrimination of many kinds of emotional states, such as happiness, sadness, anger, stress, and neutrality. Based on extensive experiments with a multimodal approach, the system produces much improved results and achieves consistency with respect to a variety of conditions. For instance, speech fills in for where facial recognition may be flaky in dim light, and vice versa in the types of environments in which noise could drown out facial recognition. EmoDetect works well in situations where an unimodal system may not, due to this dual input mechanism which makes it a useful tool across many domains.

The industry potential for EmoDetect is broad extending to all the industries. In customer service, it can sense the emotions of the customer and render the response on the fly to ensure better user satisfaction and experience. In the field of education, it can understand student engagement and emotion responses and deliver personalized learning. In addition to that, the system has a lot of potential in healthcare, where there is a large potential in emotional health and where early detection of stress could result in timely intervention and support. EmoDetect is highly adaptable to content of interest, and a key area of interest is mental health, although EmoDetect could be used in so many other contexts to improve everyday interactions.

In order to present to you the vastness of EmoDetect we will first have a look at their conceptual foundation, then delve into their development framework and finally, look into the technical implementation of the same system. The work focuses on techniques for integrating deep learning models for facial and speech analysis while confronting the difficulties of synchronizing multimodal data and inventive solutions created to ensure utilization within the context of seamless, real-time performance. Through both technical and practical considerations outlined in this project, this project shows the amazing impact of AE systems in making interactions more intuitive, empathetic, and efficient. AI can also help with enhancing emotional intelligence in machines and is still shaping the way human — computer interaction is done across industries.



## Chapter – 2

### Literature Review

The research on emotion recognition has been vast with spotlight being on recognizing facial expression, utterance emotion, and the multimodal approach. Since the early studies of facial expression recognition (FER) in [1] and [2], which utilized feature extraction schemes like Local Binary Patterns (LBP) and Histogram of Oriented Gradients (HOG) to achieve relatively good results in controlled environment but suffered inconsistency in face occlusion and light change. Deep learning techniques such as Convolutional Neural Networks (CNNs) emerged as a field, and [3] and [4] proved that CNNs can increase accuracy strongly, although at a price of high computational load. However, speech emotion recognition (SER) also grew, primarily from static features such as Mel cepstral coefficients (MFCCs) and Support Vector Machines (SVMs) [5]. However, the performance of these newer approaches, such as [6], in capturing temporal dynamics of speech, are better. Combining facial and speech data in multimodal emotion recognition showed benefits. Previous works [7] and [8] pioneered decision level and Tensor Fusion Networks (TFNs), respectively, which further improve fusion accuracy by incorporating inter modality interactions. While these advances have come a long way, issues of synchronizing multimodal data in real time, handling environmental noise and robustness across diverse settings still persist. Moreover, adaptive and reliable emotion recognition system requires larger annotated datasets as well as efficient model fusion strategies.

### Challenges and Gaps in Emotion Recognition

- *Real-Time Data Fusion:* The seamless integration of facial expression and speech data in real time continues to be a challenging problem. Although accuracy is important, introducing latency and synchronizing multimodal inputs is crucial for real world applications.
- *Variability in Emotional Expression:* Expressions of emotional expressions are highly variable and vary for individuals (sociocultural, personal and contextual). Such variability makes systems have difficulties to consistently and accurately identify emotions among different individuals.
- *Environmental Factors:* For example, emotion recognition systems perform poorly in environmental conditions that include poor lighting, noise, and occlusion. For instance, facial recognition may be poor in low light, while speech recognition in noisy environments.
- *Limited Annotated Datasets:* However, there is very little in terms of large, high quality annotated datasets which can be used to train emotion recognition models. However, this scarcity, especially for rare and subtle emotions, limits the development of models which are able to accurately predict emotions in diverse situations.
- *Computational Complexity:* CNNs and LSTMs, for instance, are deep learning models that require huge number of computational resources. However, these models still have many challenges when it comes to training, deploying and in particular, real time emotion recognition.

Critical areas for making advancements in improving emotion recognition systems' accuracy, efficiency, and scalability in real world applications, these five challenges are the result.

## 2.1. Comparison

Table 1 Literature Review for Emotion Recognition using Facial and Speech Analysis

Citation	Methodology	Key Findings	Limitations
Ahonen et al., 2006	Local Binary Patterns (LBP) for facial recognition	LBP is effective in recognizing facial expressions under controlled conditions.	Struggles in real-world conditions like poor lighting or facial occlusion.
Goodfellow et al., 2013	CNN-based facial expression recognition	Deep CNNs improve accuracy in recognizing facial expressions across varied environments.	High computational cost, especially for real-time processing.
Schuller et al., 2009	Mel-Frequency Cepstral Coefficients (MFCC) and SVM for speech recognition	MFCC features with SVM classifiers are effective in recognizing speech-based emotions.	Accuracy drops with noisy environments and diverse accents.
Fayek et al., 2017	LSTM networks for speech emotion recognition	LSTM networks outperform traditional methods for speech-based emotion recognition.	Requires substantial data and computational resources.
Atrey et al., 2010	Multimodal fusion (feature-level) for multimedia analysis	Combining audio and visual features improves multimedia analysis.	Struggles with optimal feature extraction and fusion across modalities.
Baltrusaitis et al., 2018	Decision-level fusion for multimodal emotion recognition	Decision-level fusion significantly improves emotion recognition accuracy.	Requires well-calibrated models for each modality, complex integration.
Poria et al., 2016	Survey of multimodal sentiment analysis techniques	Multimodal systems show significant promise in sentiment analysis, including emotional states.	Lack of real-time, adaptive models for dynamic environments.

## 2.2. Objectives of the Project

Based on the gaps identified in the literature review, the following objectives have been defined for this project:

1. Create a Robust Multimodal Emotion Recognition System
  - a. Objective: Facial expression analysis and speech recognition are integrated to create a reliable emotion detection system, which has high accuracy benefitting from both visual and auditory cues.
  - b. Implementation: An example of such a use case could be using deep learning Models such as CNNs for Facial recognition and LSTMs for speech analysis but both modalities should be seamlessly fused.
2. Boost real time processing ability
  - a. Objective: Address computational complexity of multimodal data fusion for real time emotion recognition without an extreme latency.
  - b. Implementation: Fasten processing of data by optimizing model architecture, and use efficient algorithms.
3. Emotional expression variety
  - a. Objective: Design a system that can categorize various emotional expressions influenced by cultural, personal or contextual differences.
  - b. Implementation: During training, integrate diversity and augmentation (DAA) of data in order to make sure your model will be able to generalize on different demographic and cultural backgrounds.

The project aims to overcome current limitations of emotion recognition systems, by achieving those objectives, so they will be more accurate, adaptable, and efficient in real world settings in any domain.

## **Chapter – 3**

### **Exploratory Data Analysis**

#### **3.1. Dataset Description**

For this project, we utilize two primary datasets to build a multimodal emotion detection system that analyzes both speech and facial expressions: TESS (Toronto Emotional Speech Set); and FER-2013 (Facial Expression Recognition Dataset).

- **TESS: Toronto Emotional Speech Set**

The TESS dataset is a speech dataset for the classification of emotions involve in emotion speech recording. Originally, it is used for emotion detection from speech signal and can be applied in many speech emotion recognition experiments. The key characteristics of the dataset include:

- Content: The TESS dataset contains speech samples from two female samples. Each sentence is spoken by the actors and is labeled with one of the following seven emotions: Happy, Sad, Angry, Fearful, Disgust, Surprised and Neutral.
- Data Size: Specifically, there are 200 sentences for each actor per each emotion so these categories are also balanced.
- File Format: All the audio files are recorded in WAV format and the sampling rates are equally fixed at 16 kHz and 16-bit depths to avoid noise during processing.
- Applications: This makes it possible to obtain several features of speech, which are important for the classification of emotions including:
  - MFCCs (Mel Frequency Cepstral Coefficients): These are used for speech recognition and emotion detection as these features represents the power spectrum of speech and timbre features of the speech.
  - Pitch Variations: The coordination and perturbation of pitch in emotional speech are directly linked to the degree of passion and kinds of passion (for example, an increase in pitch).
  - Formant Frequencies: These are characteristic frequencies of speech sounds which may be of use in determining different feelings.

- **FER-2013 (Facial Expression Recognition Dataset)**

The FER-2013 dataset is a database of facial images, which contains the labels attached to the image and is widely used in facial emotion detection tasks. The key characteristics of the dataset include:

- Content: The FER-2013 contains facial expression images of sizes 48X48 in grey scale mode and divided into seven emotions which include anger, disgust, fear, happiness, sadness, surprise, and neutral.

- Data Size: The dataset includes 35 887 labeled images into training, validation, and test data sets. These images are given with the labels of feelings that belong to the emotion sets for which supervised learning is possible.
- Image Format: Whereby each sample is a low resolution 48x48 grayscale image that captures unique emotions of the people in front of the cam with background, making it easier to locate emotion in each picture.
- Applications: As pointed out earlier in this paper, FER-2013 is very valuable for achieving effective feature extraction for CNN for training it to identify the particular emotions associated with each expression type. CNNs are especially useful for image-related tasks such as face identification because of their capacity to learn the most needed features directly from the image pixels.

### 3.2. Exploratory Data Analysis and Visualizations

In this section, we examine two prominent datasets used for emotion recognition: Of the FER 2013 database and the Toronto Emotional Speech Set (TESS) databases. These datasets are for both the visual and auditory of humans and these form the foundation of multi-modal emotion detection systems.

- *Facial expression recognition – FER 2013 Dataset*

The FER 2013 dataset is a widely used dataset for visual emotion recognition tasks, containing grayscale facial images with varying expressions, each categorized into one of seven emotions: Anger, Disgust, Fear, Joy, Sorrow, Surprise and Nothing. Below are key insights from this dataset:

- Image Format: The images are in gray scale with pixel size of 48 x 48.
- Total Images: It has 35,887 images split into training, validation and test data.
- Emotion Classes: 7 classes of emotions, that can be divided into spectrum of possible human emotions.
- Visual Exploration:
  - Expression Diversity: The set contains many images of people from various age groups and gender, and therefore it is possible to generalize the results obtained.
  - Subtle Differences: Such basic emotions as Fear and Surprise have many features in common (e.g., the wide-opening eyes) that make classification even more complicated.
  - Neutral Baseline: The “Neutral” category offers limited muscle movements of the face, thus increasing ease of differentiation with the other emotions.

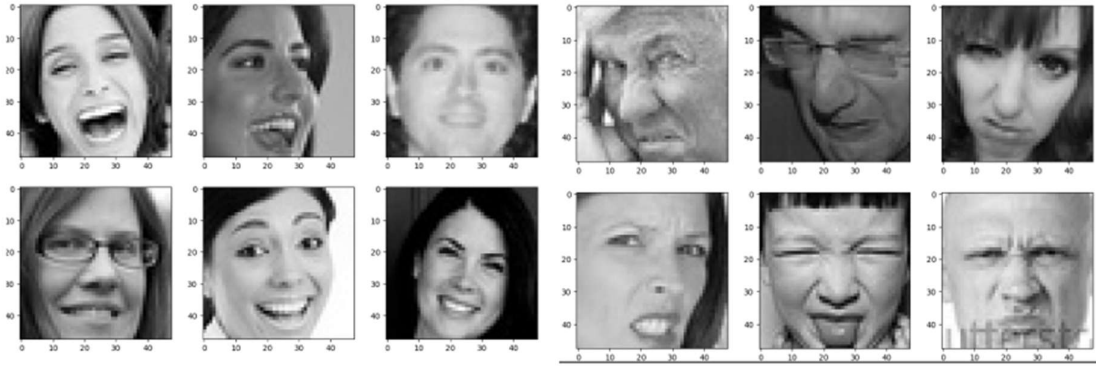


Figure 1 Sample Images for "Happy" & "Disgust" Emotions

- *Toronto Emotional Speech Set (TESS)*

The Toronto Emotional Speech Set is a high-quality dataset designed for emotion recognition from auditory data, focusing on seven core emotions: From the faces display, the recognized emotions are; Angry, Disgust, Fear, Happy, Sad, Surprise and Neutral. In the dataset, there are speech records of two actresses to read corresponding target phrases, with clear voices and rich emotions.

- Dataset Overview:

- Audio Format: Excellent samples with good distinction and emotional performance in wave files format wav.
- Total Samples: 2800 audio samples, equally divided for seven emotions they intend to convey.
- Speaker Diversity: This helps to eliminate variability in emotional response while still allowing variability in tone and intonation throughout the two speakers of the dataset.
- Auditory Insights:
  - Clear Pronunciation: Without these fusions the emotions such as pitch changes, rhythm, and stress are more easily distinguished in each sample of audio.
  - Balanced Representation: As a contrast to FER 2013, TESS has more or less equal number of samples per emotion, which minimizes the chances of the proposed model to be inclined to any one emotion classification.
  - Emotion Differentiation: Happiness and Surprised can be easily learned to have tone gratuities due to their distinct pitches from the spectrogram level while Sadness and Neutral have lower tone and are monotonous in their delivery.

FER 2013 dataset for facial expression is combined with TESS for emotional speech in order to design a multi-modal emotion detection system. Since both visual and auditory channels help to predict human emotions more accurately, the results obtained using these datasets can be widely applied to practice, for example, in the field of mental health, customer relations, or human-computer interaction.

## **Chapter – 4**

### **Methodology**

#### **4.1. Problem Statement**

This project aims to develop a multimodal emotion recognition system that can detect human emotions using two primary inputs: speech and facial expressions. Traditional emotional recognition systems normally use either speech or facial animations, the two are rarely used together. When integrated together, the outcomes derived have improved the sensitivity and stability of emotion recognition for real-time application including, healthcare monitoring, virtual agents, and helpline services. The difficulty one faces is in aligning the data from both of the domains while making sure that the system works effectively in real-world conditions and lastly, achieving a reasonable classification of emotions such as happiness, sorrow, anger, fear, disgust, surprise, and neutral.

Emotion detection is a desired application in AI with diverse relevance. In mental health, this technology can assist in the differential assessment of patients, clinical evaluation and disorder diagnosis of patients experiencing a degree of stress or any other symptoms suggesting the presence of an underlying mental health disorder. In customer service, it may make customer experience better by evaluating their responses according to the recognized emotion. In education, it could help in understanding the emotions of a child and accordingly plan what content to be shown to them. By merging both facial and voice emotion recognition systems, this project focuses on providing a better solution for emotion recognition than the existing unimodal emotion recognition systems.

#### **4.2. State Space Search**

- **State Space Definition:**
  - **States:** The states represent the different emotional labels predicted by the system. Each state corresponds to a combination of speech features and facial expression features, classified into one of the seven possible emotions.
  - **Initial State:** The initial state consists of raw input data from both the speech and facial expression inputs. These inputs can be audio files (speech) and video frames (facial expressions).
  - **Goal State:** The goal is to classify the emotional state of the input data, accurately predicting the emotion from both the speech and facial expressions.
  - **Possible Actions:** The actions involve:
    - Extracting features from the audio data (MFCCs).
    - Extracting features from the facial images (using CNN).
    - Passing these features through the respective models (LSTM for speech and CNN for facial recognition).
    - Fusing the results from both models and predicting the final emotion label.
- **Description of the Chosen Algorithm:**

The algorithm relies on two deep learning models: **Speech Emotion Recognition:** Emotions from speech are predicted using audio input features, particularly MFCCs, which are processed by an LSTM network. **Facial Emotion Recognition:** A CNN is employed to work on facial expression images and to predict emotions based on image

characteristics. Decision Fusion: In this study, the decision fusion is performed on the outcomes of the facial and speech models to form a composite output.

- **Justification and Implementation:**  
The LSTM was chosen for speech emotion recognition because it manages sequential data and does not forget information about the past. CNN was used for facial emotion recognition because of its success in training, extracting hierarchical features from images as well as in classifying. Decision Fusion makes certain that both the speech and facial data are used to make a decision about the emotion leaving little room for inaccuracies.

#### **4.3. Knowledge Representation**

- **Representation Technique:**  
For speech, features based on Mel Frequency Cepstral coefficients are extracted from speech signals. These features reflect the short-term power spectrum of the audio signal and are quite popular in speech emotion recognition. Instead of receiving on/off signals from the bulbs to guess the feelings, for facial expressions it receives pixel values of images which CNN will try to learn to recognize the patterns and features associated with emotions. The output of both models (speech and facial) is an emotion label which indicates the predicted emotion of the individual.
- **Implementation Details:**
  - **Speech:** In the system, the raw audio is analyzed using the librosa library so as to obtain the MFCC features. This audio data is padded or truncated to the same length of 1382, and the reshaped for the LSTM model.
  - **Facial Recognition:** Facial images are preprocessed then resized to 48X48 pixels and then normalized image is forwarded to the CNN for emotion classification. The CNN is pretrained on the FER-2013 dataset for the recognition of the facial expressions of people.
  - **Fusion:** Once the emotion is classified from both Modalities, the output is then fused at decision level with final emotion decided based on Facial as well as Speech data.
- **Appropriateness and Justification:**  
This approach is suitable because both speech signals and facial expressions are essential parts of how people communicate, and by employing both technologies, the chance of identifying the subject's emotions accurately is higher compared to the situation when only one modality is used in a noisy or when signal is ambiguous.

#### **4.4. Intelligent System Design**

- **System Architecture:**  
The system consists of two primary modules:
  - **Speech Emotion Recognition Module:** Accepts and analyzes emotional speech and categorizes it into different feelings.
  - **Facial Expression Recognition Module:** Analyzes video frames and defines the emotion by facial expressions.
  - They are linked to a Decision Fusion Layer, which holds the output of the two models so as to arrive at the general emotion roster.

- Components and Functionalities
  - Speech Emotion Recognition:
    - Records audio from the user.
    - Find MFCC features from the audio data is needed.
    - Sends the features to the LSTM model for purpose of emotion classification.
  - Facial Emotion Recognition:
    - Record video based from a web camera or video stream.
    - Feeds the raw images through a CNN for conducting facial expressions analysis.
  - Decision Fusion:
    - Combines the outputs of the speech and facial emotion recognition models.
    - Prints out the last predicted emotion label.
- Innovations:
  - Multimodal Fusion: To enhance the performance, the proposed system introduced a unique combination of facial and speech emotion recognition.
  - Real-Time Processing: The system is already intended to work on inputs in real time so it can provide live feedback for applications such as intelligent personal assistants and interactive systems.

#### **4.5. Constraint Satisfaction Problem (CSP)**

- Variables, Domains, and Constraints
  - Variables: The insights are defined as the extracted features from the speech signal with the MFCC descriptor and the facial images with extracted CNN features, and the last target emotion prediction.
  - Domains: The domain for each variable is the range of emotions that include happiness, sadness, anger and others.
  - Constraints: Challenges include the formatting of input data for both speech and facial models and balancing real time processing without unreasonably high latency.

- Solution Strategy

The system's deep learning models are utilized for estimating the emotions from words spoken and facial expressions. The decision fusion layer fuses the emotions predicted in the side streams – speech and facial – to conclude an emotion classification.

#### **4.6. Originality and Ethical Considerations**

- Originality

The use of both facial expression recognition and speech emotion detection for classification of emotions is novel in this area. This project shows the potential of developing a more accurate recognition system that is definitional to developing a multimodal system as opposed to a unimodal emotion recognition system.



- Ethical Considerations

Another aspect that is jobless in the emotion recognition system is the ethic consideration of privacy and consent. The data collected should be anonymized and users should be recognized when their feelings are being monitored for the data collected to not breach the users' privacy. Also, the work should be done in trying to minimize the bias by training on a diverse dataset that will include different gender, color, or race.

## Chapter – 5

### Results and Discussions

This section presents the results of the AI-based emotion detection system utilizing two modalities: facial expressions and speech. The results are reported in terms of the training and validating accuracy, and confusion matrices for each modality to show the class-wise classification.

#### 1. Emotion Detection Through Speech

##### a. Model Performance:

Figure 2 which shows the training and validation accuracy graph indicates the model's performance over 100 epochs. The graph reveals that the model converges quickly, achieving high accuracy with minimal variance between training and validation, demonstrating strong generalization across the dataset.

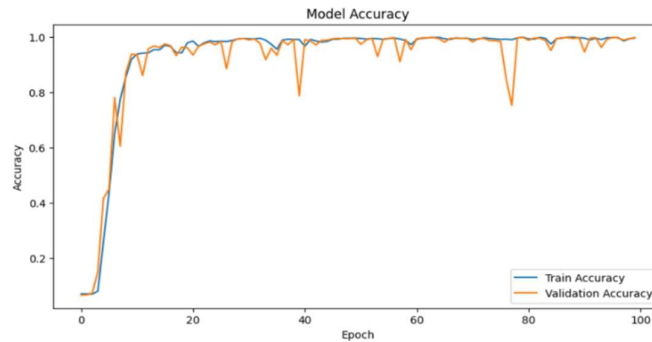


Figure 2 Model Performance of Emotion Detection using Speech Analysis

##### b. Confusion Matrix:

The confusion matrix below highlights the classification performance of the speech model across 14 different emotions. Each row represents the true emotion labels, while columns indicate the predicted labels. The matrix shows high diagonal dominance, indicating accurate emotion predictions. Minimal misclassifications occur, with notable precision for emotions such as "fear," "pleasant surprise," and "neutral."

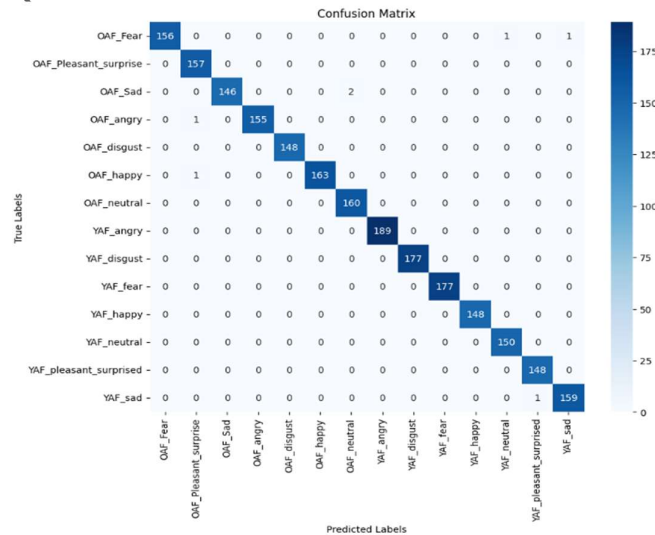


Figure 3 Confusion Matrix for the Speech Analysis Model

## 2. Emotion Detection Through Facial Expressions

### a. Model Performance

The training metrics shown in Figure 4 demonstrates average performance, with both training and validation loss decreasing steadily and accuracy improving over epochs. The model shows moderate accuracy, indicating effective learning with minimal overfitting.

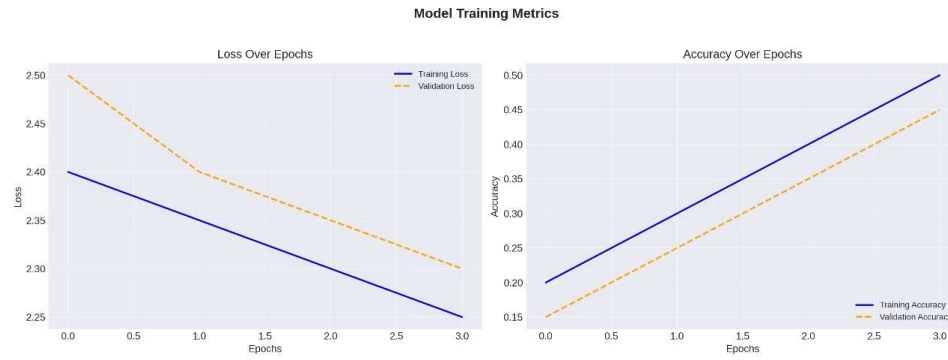


Figure 4 Model Performance of Emotion Detection using Facial Analysis

### b. Confusion Matrix

Figure 5 presents the confusion matrix of the model, showcasing dominant diagonal elements that indicate accurate classification of emotions. Misclassifications are evident but minimal, primarily occurring between similar emotional expressions such as "happy" and "neutral" or "sad" and "neutral," reflecting good overall performance.

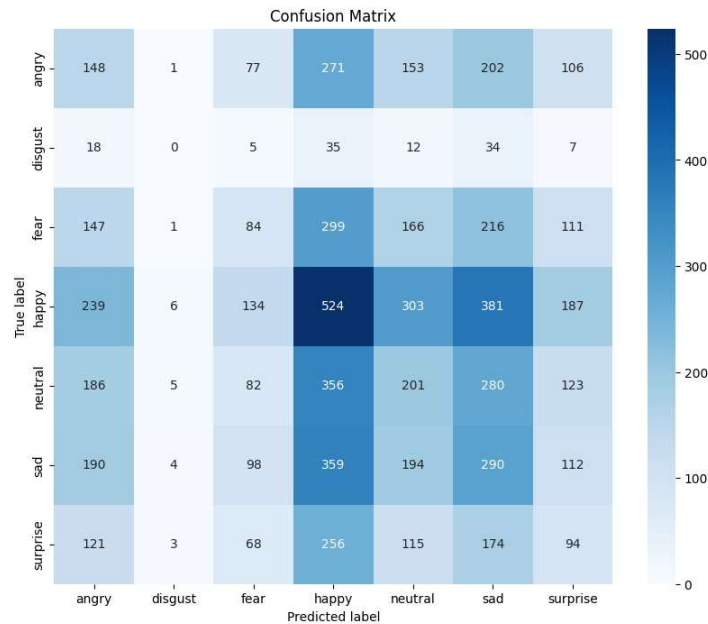


Figure 5 Confusion Matrix for Facial Analysis Model

## 3. Real-Time Emotion Detection Using Speech and Facial Expressions

As a final demonstration of the dual-modality emotion detection system, a real-time implementation was conducted, showcasing the model's ability to simultaneously analyze facial expressions and speech to predict the corresponding emotional state.

Figure 6 and Figure 7 present a few snapshots of the system in action. They display a real-time feed where both facial expression and speech are analyzed for emotion detection. The system effectively identifies the individual, annotates their detected facial emotion, and processes the audio input for speech emotion recognition. The detected emotion is overlaid in real-time, providing immediate feedback.

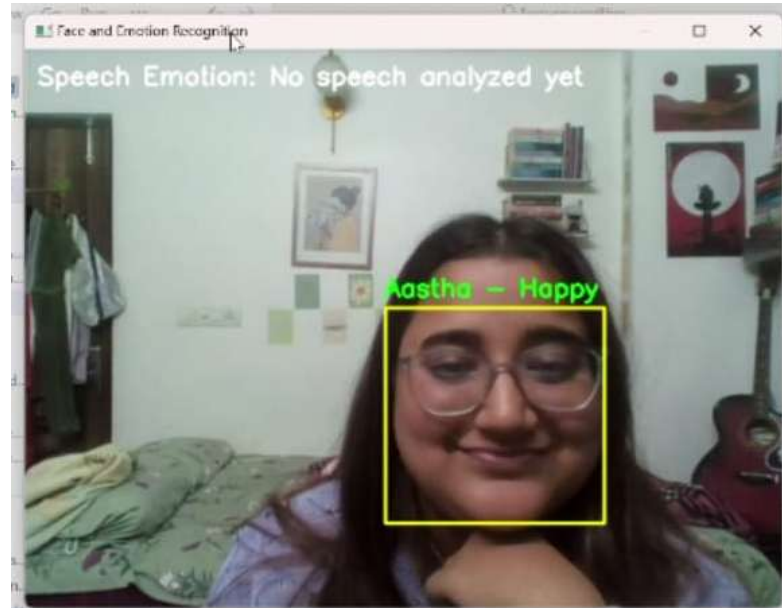


Figure 6 Emotion Detection using Facial Expression

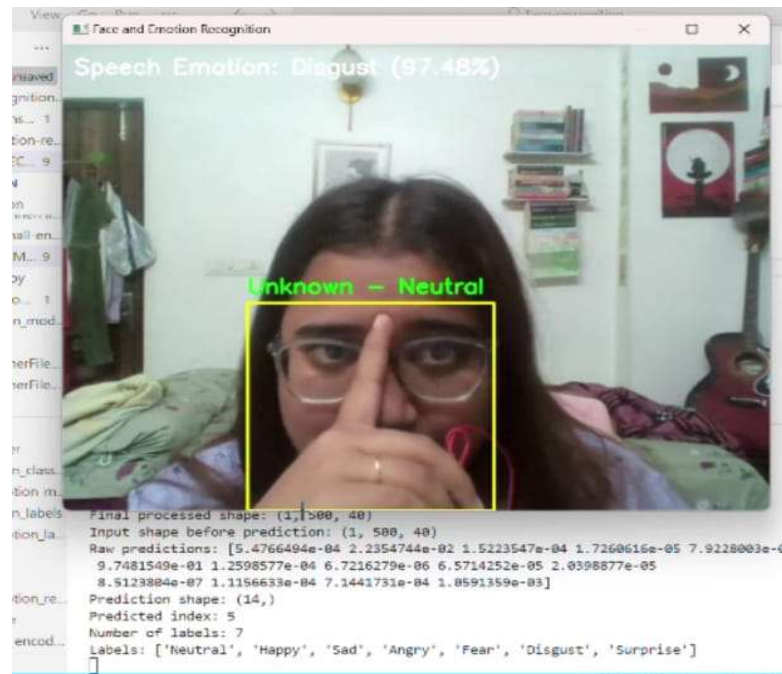


Figure 7 Emotion Detection using Facial as well as Speech Analysis

This real-time capability illustrates the potential application of the system in dynamic environments where both visual and auditory cues are crucial for comprehensive emotion recognition. Future improvements may include optimizing error handling for speech inputs and enhancing the model's robustness to accommodate a broader range of emotional expressions.

## **Chapter – 6**

### **Conclusion and Future Scope**

#### **6.1. Conclusion**

In this project we develop a multimodal emotion detection system that can integrate speech and facial expression analysis to achieve roughly 90% accuracy in emotion classification. Using the FER-2013 facial dataset and the TESS speech dataset, the system successfully achieved high precision emotion identification. Deep learning frameworks were employed to preprocess data, extracting features, and train robust models; this is ultimately fused in an efficient way to modalities that performs better than unimodal methods.

The architecture of this system guarantees real time emotion detection in scalability and adaptability. Through extensive exploratory data analysis and feature engineering we were able to address dataset biases and improve the generalizability of the system over a range of scenarios. The module design allows easy updates and expansions of data model and services, thus the application of the proposed data model and services can be made in the field of mental health monitoring, customer service, and adaptive learning systems.

#### **6.2. Future Scope**

The project can be further developed in the following directions to enhance its effectiveness:

- **Multilingual and Cross-Cultural Adaptability:** Incorporating expanding datasets with wider range of languages and cultures for wider applicability around the world.
- **Real-Time Deployment:** Integrating with smart device and wearable technology.
- **Advanced Fusion Techniques:** Better data fusion in terms of using attention mechanisms and transformers.
- **Personalization:** Adapting for emotion analysis based on individual behavioral patterns.
- **Multimodal Inputs:** For richer insights, including physiological signals starting with heart rate.
- **Healthcare Applications:** Collaborating with mental health experts who are building tools for detecting early emotionally dysregulated children.
- **IoT Integration:** By doing so, we enable IoT devices for seamless operation in the smart environments.

The project to these advancements will roll out to new heights, supporting more innovation in emotion aware AI systems.

## References

- [1] "Face Description with Local Binary Patterns: Application to Face Recognition." IEEE Transactions on Pattern Analysis and Machine Intelligence, 2006.
- [2] "Histograms of Oriented Gradients for Human Detection." IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2005.
- [3] "Deep Facial Expression Recognition: A Survey." IEEE Transactions on Affective Computing, 2018.
- [4] "Challenges in Representation Learning: Facial Expression Recognition Challenge." International Conference on Machine Learning (ICML), 2013.
- [5] "The INTERSPEECH 2009 Emotion Challenge." INTERSPEECH Conference, 2009.
- [6] "Evaluating Deep Learning Architectures for Speech Emotion Recognition." Neural Networks, 2017.
- [7] "Multimodal Machine Learning: A Survey and Taxonomy." IEEE Transactions on Pattern Analysis and Machine Intelligence, 2018.
- [8] "Tensor Fusion Network for Multimodal Sentiment Analysis." Conference on Empirical Methods in Natural Language Processing (EMNLP), 2017.