

## CONTENTS

---

<b>Chapter</b>	<b>Illustration</b>	<b>Page</b>
<b>List of Figures</b>		i
<b>List of Tables</b>		ii
<b>List of Symbols and Abbreviations</b>		iii
<b>Abstract</b>		iv
<b>1. INTRODUCTION</b>		<b>01</b>
1.1 Introduction		01
1.2 Motivation		01
1.3 Objectives		02
1.3.1 Multimodal Data Acquisition		02
1.3.2 Emotion Detection and Classification		02
1.3.3 Generative AI Empathy Engine Development		02
1.3.4 System Reliability and Human-Centered Deployment		02
1.4 Scope		03
<b>2. LITERATURE SURVEY</b>		<b>04</b>
2.1 Existing Methodologies		04
2.2 Research Gap Analysis		09
<b>3. REQUIREMENT SPECIFICATION AND ANALYSIS</b>		<b>10</b>
3.1 Problem Definition		10
3.2 Proposed Methodology		11
3.3 Project Requirements		12
3.3.1 Datasets		12
3.3.2 Functional Requirements		12

3.3.3 Non-Functional Requirements	13
3.3.4 Hardware and Software Requirements	13
3.4 Project Plan	13
3.4.1 Project Resources	13
3.4.2 Module Split-up	14
<b>4. SYSTEM ANALYSIS AND DESIGN</b>	<b>16</b>
4.1 System Architecture	16
4.2 System Architecture Workflow	19
4.3 Algorithm And Methodologies	23
<b>5. IMPLEMENTATION</b>	<b>27</b>
5.1 Stages Of Implementation	27
5.1.1 Data Preprocessing	27
5.1.2 Model Training And Development	29
5.1.3 Multimodal Fusion & Smoothing	32
5.1.4 Generation Module	32
5.1.5 Backend Integration	34
<b>6. RESULTS</b>	<b>35</b>
6.1 Results Of Experiments	35
6.2 Result Analysis	37
6.2.1 Model Performance Evaluation	37
6.2.2 Multimodal Fusion Evaluation	37
6.2.3 Response Generation Evaluation	38
6.2.4 System Performance Evaluation	38
6.3 Testing	39
6.3.1 White Box Testing	39
6.3.1.1 Unit Testing	39

6.3.1.2 Integration Testing	40
6.3.2 Testing Outcome	40
<b>7. CONCLUSION AND FUTURE SCOPE</b>	<b>41</b>
7.1 Conclusion	41
7.2 Limitations Of The Project	41
7.3 Future Scope	42
<b>REFERENCES</b>	<b>44</b>

---

## LIST OF FIGURES

---

<b>Figure No</b>	<b>Title</b>	<b>Page No</b>
3.1	Emotion Recognition Flowchart	14
3.2	Emotion-Based Response Generation Pipeline	15
4.1	End-to-End System Architecture for Multimodal Emotion Detection and AI Response Generation	19
4.2	Activity Diagram	20
4.3	DFD Level-0	21
4.4	DFD Level-1	22
5.1	Facial Data Preprocessing and Augmentation	27
5.2	GoEmotions Dataset Loading and Custom Label Mapping	28
5.3	Transformer Input Tokenization and Class Weight Calculation	29
5.4	CNN Model Performance Metrics	30
5.5	BERT Model Performance Metrics	31
5.6	Code Implementation for Multimodal Fusion	32
5.7	Code Implementation of Hybrid Response Generation and Context Management	33
5.8	Backend API Routes for Real-time Data Processing	34
6.1	Backend Server Log: Multi Model Loading and Crisis System Activation	35
6.2	GenAI Therapy Application's User Interface	35
6.3	System Response to High-Risk Suicidal Input	36
6.4	Real-Time Emotion Detection and Empathetic Response	36

---

## **LIST OF TABLES**

<b>Table</b>	<b>Title</b>	<b>Page No</b>
2.1	Research Gap Analysis	09
3.1	Basic System Configuration	13

## LIST OF SYMBOLS AND ABBREVIATIONS

Symbol	Description
$P_{\text{face}}$	Probability vector from Face Emotion Model (CNN)
$P_{\text{text}}$	Probability vector from Text Emotion Model (BERT)
$P_{\text{final}}$	Final fused probability vector after weighted fusion
$\alpha$	Weight assigned to face emotion probability (0–1)
$1 - \alpha$	Weight assigned to text emotion probability

## ABSTRACT

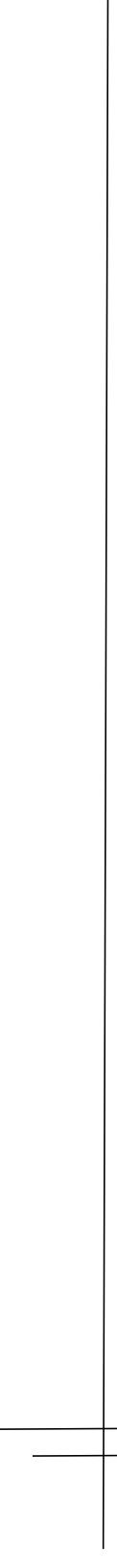
The proposed system is an **AI-powered platform** designed to detect human emotional states like **depression, anxiety, frustration, or calmness** in real-time using multimodal input: facial expressions and textual input. By integrating machine learning and Generative AI, the system simulates a conversation with a supportive companion to offer basic mental wellness support.

The project workflow involves multimodal emotion detection using two distinct streams: **Facial Emotion Analysis** from live camera and **Sentiment Detection** from written text. Data preprocessing includes normalizing face images from the **FER2013 dataset** and cleaning/tokenizing text from the **GoEmotions dataset**; both are mapped to a simplified four-class emotional state. The system then combines these scores using a weighted **Multimodal Fusion algorithm** to predict the final emotional state.

The core technology utilizes a **Convolutional Neural Network (CNN)** model for facial expression classification and a **fine-tuned HuggingFace AutoModelForSequenceClassification** (Transformer-based) for text emotion detection. Upon detecting an emotion (e.g., anxiety or depression), the system employs **Generative AI (LLM)** to craft an emotionally supportive, context-aware, and non-robotic empathetic response (e.g., motivational or calming tone).

This solution addresses the lack of emotion-aware tools in current AI by providing a real-time, personalized interaction that promotes early emotional intervention. The final output is an easy-to-use system that provides **AI-driven empathetic support** through a human-centered, conversational interface.

**Keywords:** GenAI, Depression Detection, Multimodal System, Emotion Recognition, CNN, NLP.



# **Chapter 1**

# **Introduction**

---

## **CHAPTER 1**

### **INTRODUCTION**

#### **1.1 INTRODUCTION**

The rapid evolution of Artificial Intelligence (AI) and Natural Language Processing (NLP) presents a powerful opportunity to develop systems that address complex human needs, particularly in the realm of mental health. This project, GenAI Approach for Detection of Depression, introduces a cutting-edge, AI-powered platform designed for the real-time detection and empathetic response to human emotions.

The core innovation lies in its multimodal approach, which simultaneously processes input from two critical channels: facial expressions (via computer vision) and textual input (via NLP). By leveraging advanced machine learning algorithms and emotion recognition techniques, the system classifies emotional states such as depression, anxiety, frustration, or calmness.

Crucially, once the emotion is identified, the system employs Generative AI (GenAI) to formulate a personalized, empathetic, and supportive response, effectively simulating a conversation with a caring companion. This solution aims to redefine the interaction between humans and AI by making the system emotionally aware, reliable, and fundamentally human-centered.

#### **1.2 MOTIVATION**

The contemporary AI landscape, despite its advancements, is often characterized by tools that are functional but lack emotional intelligence, resulting in impersonal and often frustrating user experiences. This deficiency - the lack of emotion-aware systems in today's AI tools forms the central motivation for this project.

Furthermore, the need for accessible mental wellness support is increasingly critical globally. While professional help remains paramount, there is a distinct gap for a readily available, non-judgmental, and preliminary tool that can identify distress signals and offer immediate, empathetic communication. This project is motivated by the potential to bridge this gap by demonstrating how technology can be leveraged to offer basic mental wellness support. By

integrating the latest advancements in GenAI for conversational response, the system seeks to build an easy-to-use platform that feels genuinely comforting.

### **1.3 OBJECTIVES**

The overall objective of this project is to successfully implement an AI-powered, multimodal system for preliminary emotional state detection and empathetic response generation. The specific, measurable objectives are:

#### **1.3.1 Multimodal Data Acquisition and Platform Architecture:**

To design and deploy a stable platform capable of acquiring and processing real-time multimodal data streams: facial expressions (video) and textual input (chat).

#### **1.3.2 Emotion Detection and Classification:**

To train and validate Machine Learning (ML) models to accurately classify emotional states into specific categories: depression, anxiety, frustration, and calmness, achieving a high degree of classification accuracy.

#### **1.3.3 Generative AI Empathy Engine Development:**

To develop and integrate a Generative AI (GenAI) module fine-tuned to create empathetic, supportive, and contextually appropriate conversational responses, simulating a caring companion.

#### **1.3.4 System Reliability and Human-Centered Deployment:**

To construct a reliable, human-centered system that ensures the final application is intuitive, non-judgmental, and robust enough to provide basic mental wellness support.

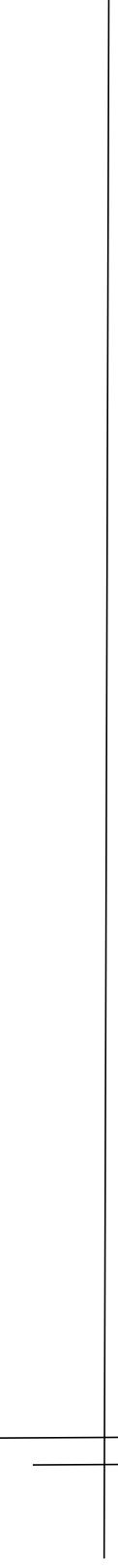
## **1.4 SCOPE**

The scope of this project is to develop a robust, full-stack application that provides preliminary emotional state detection and empathetic response generation by leveraging GenAI and a multimodal data approach. The system will serve as a foundational, non-diagnostic tool for basic mental wellness support by simulating a caring companion.

The project involves the collection and preparation of a diverse dataset, including conversational text for NLP training and image/video data for real-time facial emotion recognition. This multimodal data will be used to train two integrated machine learning componentsa vision model (e.g., CNN) and an NLP model (e.g., Transformer-based)to classify emotions into the core states of depression, anxiety, frustration, and calmness. The architectural design will prioritize real-time processing and a seamless, user-friendly interface.

The core deliverable is a functional prototype demonstrating the complete workflow: real-time data ingestion, emotion classification, and the generation of a personalized, empathetic response using a fine-tuned Generative AI module. The system's effectiveness will be rigorously evaluated by measuring the accuracy of the classification models and the quality of the AI-generated responses.

Crucially, the scope explicitly excludes any intent or capability for formal medical diagnosis, psychological counseling, or crisis intervention. The system is designed purely for basic mental wellness support and preliminary emotional detection. Future work includes potential expansion to more complex sentiment analysis and deployment on various platforms.



## **Chapter 2**

# **Literature Survey**

---

## CHAPTER 2

### LITERATURE SURVEY

#### 2.1 EXISTING METHODOLOGIES

**Emotion Detection Using Deep Learning in a Multimodal Context:** The research paper "Emotion Detection Using Deep Learning in a Multimodal Context: A Systematic Review" [1] comprehensively analyzes existing literature on emotion detection systems that utilize combined signals from multiple modalities, such as speech, facial expression, and textual input. The review focuses on how Deep Learning (DL) architectures, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have been applied to process these diverse data streams for effective emotion classification.

The key finding of this systematic review is that deep learning models significantly improve accuracy when fusing multimodal signals for emotion classification compared to traditional unimodal approaches. It confirms the robustness and necessity of combining different forms of data to achieve a more reliable and holistic understanding of a user's emotional state, which is crucial for systems aiming to detect complex states like depression or anxiety.

One potential drawback is that as a systematic review, the paper does not propose a novel model or algorithm, meaning it does not offer practical implementation details or a validated dataset/model for direct use in an application. Additionally, the sheer complexity of implementing a truly multimodal system (involving face, voice, and text) often requires significant computational resources and specialized technical expertise, which can be a barrier for smaller projects or non-expert users.

#### Facial Expression and Voice Tone Based Emotion Detection Using Convolutional and Recurrent Neural Networks:

The research paper "Facial Expression and Voice Tone Based Emotion Detection Using Convolutional and Recurrent Neural Networks" [2] proposes using a combination of CNN models for facial expression analysis and Recurrent Neural Networks (RNN) for voice tone

analysis to detect emotional states. The authors describe the process of using pitch, tone, vocal tension, and facial cues to capture emotional states.

The authors present results showing that the combined CNN + RNN models provide robust performance in capturing subtle facial expressions and real-time audio cues (like stress and sadness), leading to high temporal accuracy in emotion detection. This highlights the effectiveness of using specific Deep Learning architectures tailored to the time-series nature of speech and the spatial patterns of image data.

A primary limitation is that this study focuses only on facial expressions and voice tone, omitting the analysis of textual communication, which is a major modality in many modern support systems and in your proposed project. Furthermore, the performance of models based on facial expression and voice tone can be heavily influenced by environmental factors (lighting, background noise, microphone quality) and the cultural or individual variance in expressing emotion, which limits real-world generalizability.

### **Emotion-Aware AI Chatbot for Mental Health Support Using Transformer-Based NLP:**

The research paper "Emotion-Aware AI Chatbot for Mental Health Support Using Transformer-Based NLP" [3] explores the use of advanced Transformer-based NLP models to create a mental health support chatbot. The paper details how the model generates empathetic, contextual, and supportive responses based on the emotion detected from the user's text input, moving beyond simple keyword matching.

The authors demonstrate that Transformer models provide emotionally aligned and supportive responses, significantly aiding in early intervention and improving user comfort within chat systems. This approach showcases how Generative AI can be leveraged to fulfill the crucial "AI Response Generation" module of a support system.

A major drawback is that the system is primarily unimodal (text-based), meaning it lacks the robustness that comes from incorporating non-verbal cues like facial expressions. Relying solely on textual input may lead to misclassifications, as emotional context is often lost in text, especially when users try to mask their feelings. Achieving truly human-like, non-robotic

responses from a conversational agent also remains a significant ethical and technical challenge.

### **Text-Based Emotion Detection from User Comments Using Hybrid Deep Learning Models:**

The research paper "Text-Based Emotion Detection from User Comments Using Hybrid Deep Learning Models" [4] proposes using complex Hybrid Deep Learning Models (e.g., combining CNN, LSTM, and BERT) for detecting emotions from user text inputs, such as social media comments. The paper focuses on leveraging the sequential processing power of RNNs (like LSTM) with the contextual understanding of Transformer models (like BERT) for finer-grained sentiment classification.

The study finds that Hybrid models (CNN+LSTM+BERT) improve the precision in detecting a wide range of emotions, including anger, joy, and fear, from text data. This advanced approach offers high precision in the "Text Emotion" module of a multimodal system.

Similar to other NLP-focused research, the paper is limited by its unimodal focus on text, neglecting the valuable supplementary data from face or voice. Furthermore, deploying and training such complex Hybrid Deep Learning Models requires substantial hardware and deep expertise in NLP and model optimization, representing a significant technical hurdle for implementation.

### **Multimodal Emotion Recognition: A Comparative Study on Speech, Text, and Facial Cues:**

The research paper "Multimodal Emotion Recognition: A Comparative Study on Speech, Text, and Facial Cues" [5] conducts a comparative analysis of emotion recognition across the three main modalities: speech, text, and facial cues. The authors explore different fusion techniques to assess which combination provides the most reliable and accurate final emotion classification.

The study validates the effectiveness of the multimodal approach in providing a more reliable and higher-accuracy prediction of a user's emotional state compared to any single modality alone. It provides a strong justification for the "Multimodal Fusion" component of your

proposed system.

A limitation is the frequent mismatch in data sampling rates and temporal alignment when combining speech (audio), text (discrete input), and facial cues (video frames), which makes real-time fusion computationally challenging and potentially prone to temporal lag. Also, the performance of the system often heavily depends on the quality and balance of the multimodal datasets used for training.

### **AI-Powered Recommendation System Based on Mood and Sentiment Analytics:**

The research paper "AI-Powered Recommendation System Based on Mood and Sentiment Analytics" [6] explores leveraging sentiment and mood analysis to build personalized recommendation systems, moving beyond traditional collaborative filtering. The authors describe using machine learning and NLP

techniques to analyze user-generated text (such as reviews or social media posts) to determine their current emotional state, and then tailoring recommendations (e.g., music, content, or services) accordingly.

The authors demonstrate that integrating real-time mood and sentiment analytics significantly enhances the relevance and perceived value of personalized recommendations, showing that the system can effectively shift from a generalized suggestion to one that aligns with the user's current emotional need (e.g., calming content for anxiety).

A potential drawback of this approach is the heavy reliance on textual input for mood assessment, which may be insufficient to capture complex emotional nuances, especially those related to deep psychological states like depression. The generalizability of "mood" derived from content reviews to more sensitive areas like mental health is also limited, as the study focuses more on commercial utility rather than clinical support.

### **Emotion-Driven Conversational Agent Using Deep Reinforcement Learning:**

The research paper "Emotion-Driven Conversational Agent Using Deep Reinforcement

"Learning" [7] proposes an advanced method for building chatbots that can engage in more sophisticated, empathetic conversations by incorporating Deep Reinforcement Learning (DRL). The system is trained not just to respond logically, but to learn an optimal conversational strategy that maximizes positive emotional feedback (or minimizes negative emotion) from the user over an entire dialogue sequence.

The paper presents results showing that an Emotion-Driven Conversational Agent using DRL can produce dialogue that is significantly more coherent, context-aware, and emotionally responsive than agents using simpler rule-based or sequence-to-sequence models. The use of DRL allows the agent to plan responses strategically over multiple turns to achieve a supportive goal.

One limitation is the data intensity and complexity of training a stable DRL model. DRL requires massive amounts of high-quality, labeled dialogue data that is annotated with turn-by-turn emotional states, which is difficult to acquire in the mental health domain. Furthermore, the inherent "black box" nature of DRL can make it difficult to audit or explain the reasoning behind an agent's critical, supportive response, raising concerns in a sensitive mental wellness context.

### **Design of an Emotion-Aware Mental Health Chatbot Using Hybrid NLP Techniques:**

The research paper "Design of an Emotion-Aware Mental Health Chatbot Using Hybrid NLP Techniques" [8] focuses specifically on the development and architecture of a chatbot intended for mental health support. The authors describe a system that uses Hybrid NLP techniques to accurately triage user input and generate initial emotional classifications.

The authors demonstrate that using a hybrid approach can lead to a robust initial system design, where rule-based methods handle immediate safety concerns or critical keywords, while more flexible machine learning models provide general emotional classification. This design offers a good balance of safety, predictability, and responsiveness in a mental health context.

The major drawback is that hybrid systems often face challenges in seamless integration; the hand-off between the rule-based system and the deep learning component can create inconsistencies or abrupt shifts in conversational tone. Additionally, the study may be limited by the scalability and maintenance issues associated with rule-based systems, as updating rules to match evolving language and slang can be labor-intensive.

## 2.2 RESEARCH GAP ANALYSIS

**Table 2.1: Research gap analysis**

Research Paper	Limitations & Gaps
Emotion Detection Using Deep Learning in a Multimodal Context: A Systematic Review [1]	No novel model or practical details; requires significant computational resources and specialized expertise.
Emotion-Aware AI Chatbot for Mental Health Support Using Transformer-Based NLP [2]	Primarily unimodal (text-based), leading to potential misclassification due to lack of non-verbal cues.
Text-Based Emotion Detection from User Comments Using Hybrid Deep Learning Models [3]	Unimodal focus on text; deploying complex models requires substantial hardware and deep expertise.
Multimodal Emotion Recognition: A Comparative Study on Speech, Text, and Facial Cues [4]	Challenges with data sampling rates and temporal alignment make real-time fusion computationally challenging and prone to lag.
AI-Powered Recommendation System Based on Mood and Sentiment Analytics [5]	Heavy reliance on textual input; insufficient for complex emotional nuances. Focuses on commercial utility over clinical support.
Emotion-Driven Conversational Agent Using Deep Reinforcement Learning [6]	Data intensity and complexity of DRL training; "black box" nature makes supportive responses difficult to audit or explain.
Design of an Emotion-Aware Mental Health Chatbot Using Hybrid NLP Techniques [7]	Hybrid systems struggle with seamless integration (rule-based/DL hand-off) and scalability/maintenance of rules.

# **Chapter 3**

## **Requirement Specification and Analysis**

---

## CHAPTER 3

### REQUIREMENT SPECIFICATION AND ANALYSIS

The proposed system is an AI-powered platform designed for the real-time detection and response to human emotional states, specifically focusing on classifying emotional states like depression, anxiety, frustration, or calmness. The system combines multimodal machine learning approaches, utilizing facial expressions and textual input for classification. Once the user's emotion is identified, the platform generates an empathetic and supportive AI response, simulating interaction with a caring companion using Generative AI (GenAI). This project introduces an AI-based Emotion Detection and Support System that aims to provide basic mental wellness support by bridging the gap left by traditional, non-emotion-aware AI tools.

#### 3.1 Problem Definition:

The problem statement for this project is:

To develop a system that accurately detects user emotions, such as depression and anxiety, using multimodal input (Facial expressions via camera and Text inputs), and subsequently generates supportive, human-like, and context-aware responses using Generative AI to comfort or guide the user toward basic emotional support.

In today's fast-paced environment, emotional well-being is frequently overlooked, particularly among students and working professionals. Existing AI tools and chatbots often lack emotional awareness, resulting in generic and unhelpful interactions that do not account for the user's emotional state. Current methods for emotional support are often subjective and require constant manual intervention. To address this, the project aims to develop a tool that incorporates multimodal emotion recognition:

1. Facial emotion analysis through a live camera.
2. Sentiment detection from written text (typed messages).

The system evaluates the user's current emotional state in real-time and responds with personalized and empathetic responses offering immediate emotional support.

### **3.2 Proposed Methodology:**

The application employs a Multimodal Machine Learning approach for emotion detection, coupled with a Generative AI module for empathetic response generation.

#### **Machine Learning Approach for Emotion Detection**

The system uses the following flow for emotion detection:

1. Input Capture: Real-time facial expressions are recorded via webcam and messages are captured through a text box.
2. Pre-processing:
  - i. Facial Data: Converted to grayscale and normalized for classification. The process involves resizing to 48 x 48 pixels and reshaping to (48, 48, 1) for the CNN model.
  - ii. Text Data: Tokenized, cleaned (lowercase conversion, removal of punctuation, symbols, and URLs), and processed for sentiment analysis. Tokenization uses the HuggingFace AutoTokenizer with a max sequence length of 256 tokens.
3. Emotion Detection:
  - i. Facial Emotion: Detected using a CNN or OpenCV-based classifier.
  - ii. Text Emotion: Detected using NLP techniques like VADER or Transformers.
4. Multimodal Fusion: The scores from the multiple inputs (facial and text, and potentially voice if included) are combined using a weighted fusion algorithm to determine the final emotion state (e.g., Depressed, Anxious, Calm).

#### **Generative AI Approach for Response Generation**

Once the emotion is detected, the second module, AI Response Generation, takes over.

1. Strategy Selection: Based on the detected emotion, an appropriate response strategy is selected:
  - i. *Depressed* --> Motivational / Reassuring Tone.
  - ii. *Anxious* --> Calming / Supportive Tone.
  - iii. *Calmness* --> Friendly / Informative Tone.
2. Response Generation: A Generative AI model (LLM like Blender-Bot & DialoGPT) is used to craft a contextual, empathetic reply. Then response is displayed to the user.

### **3.3 Project Requirements:**

#### **3.3.1 Datasets**

1. Face Emotion Dataset: Based on FER2013 (Facial Expression Recognition 2013) dataset, consisting of approximately 35,887 grayscale face images (48x48 pixels).
2. Mapping: The original 7 classes are mapped down to 4 classes: Depression (Sad), Anxiety (Fear + Surprise), Frustration (Angry + Disgust), and Calmness (Neutral + Happy).
3. Model: A CNN model is trained on the split of 28,709 training images and 7,100 test images.
4. Text Emotion Dataset: Based on Google's GoEmotions dataset, comprising approximately 58,000 Reddit comments labeled with 27 emotions plus neutral.
5. Mapping: The original 27 emotions are collapsed into the same 4 target classes: Depression (Sadness, hopelessness), Anxiety (Fear, worry, nervousness), Frustration (Anger, annoyance), and Calmness (Neutral, acceptance).
6. Model: A fine-tuned HuggingFace AutoModelForSequenceClassification (Transformer-based) is used.

#### **3.3.2 Functional Requirements**

1. The application must fulfill the following core functions:
2. Data Collection and Preprocessing: The application must capture real-time facial data via webcam and textual data via a text box. It must clean, normalize (facial data), and tokenize/encode (text data) for model consumption.
3. Multimodal Emotion Detection: The system must utilize CNN (for face) and NLP-based classifiers (for text) to classify the user's emotional state into one of four categories (Depression, Anxiety, Frustration, Calmness).
4. Generative AI Response System: The system must generate an emotionally supportive and context-aware reply using a Generative AI model (LLM) based on the detected emotion.
5. User Interface: The application must have a user-friendly interface that allows users to

comfortably input text and displays the live emotional analysis and the AI's chat response instantly.

### **3.3.3 Non-Functional Requirements**

1. Performance: The system must process inputs, run detection models, and generate responses quickly, ensuring real-time interaction with minimal delay.
2. Accuracy: The multimodal emotion detection component must provide accurate and reliable classification of emotional states.
3. Usability: The system must be easy to use, allowing even non-technical users to interact comfortably with the interface.
4. Empathetic Quality: The AI responses must be empathetic, motivational, or calming, tailored to the mental state, making the interaction sensitive and impactful.

### **3.3.4 Hardware and Software Requirements**

**Table 3.1: Basic system configuration**

<b>Category</b>	<b>Component</b>	<b>Specification</b>
Software	Operating System	Windows 10 / Linux
Software	Development Environment	Python-supported environment (VS Code / Jupyter / Terminal)
Hardware	Input Device	Standard Webcam, keyboard, Mouse
Hardware	Processor	Intel i5 / Ryzen 5 or higher
Hardware	Memory (RAM)	8 GB (Minimum), 16 GB (Recommended)

## **3.4 Project Plan:**

### **3.4.1 Project Resources**

1. Face Data Resource: FER2013 dataset.
2. Text Data Resource: Google's GoEmotions dataset.
3. AI Model Resources: Trained CNN (stored in `models/fer2013_model/`) and Fine-tuned

HuggingFace AutoModelForSequenceClassification (stored in  
models/goemotions\_model/).

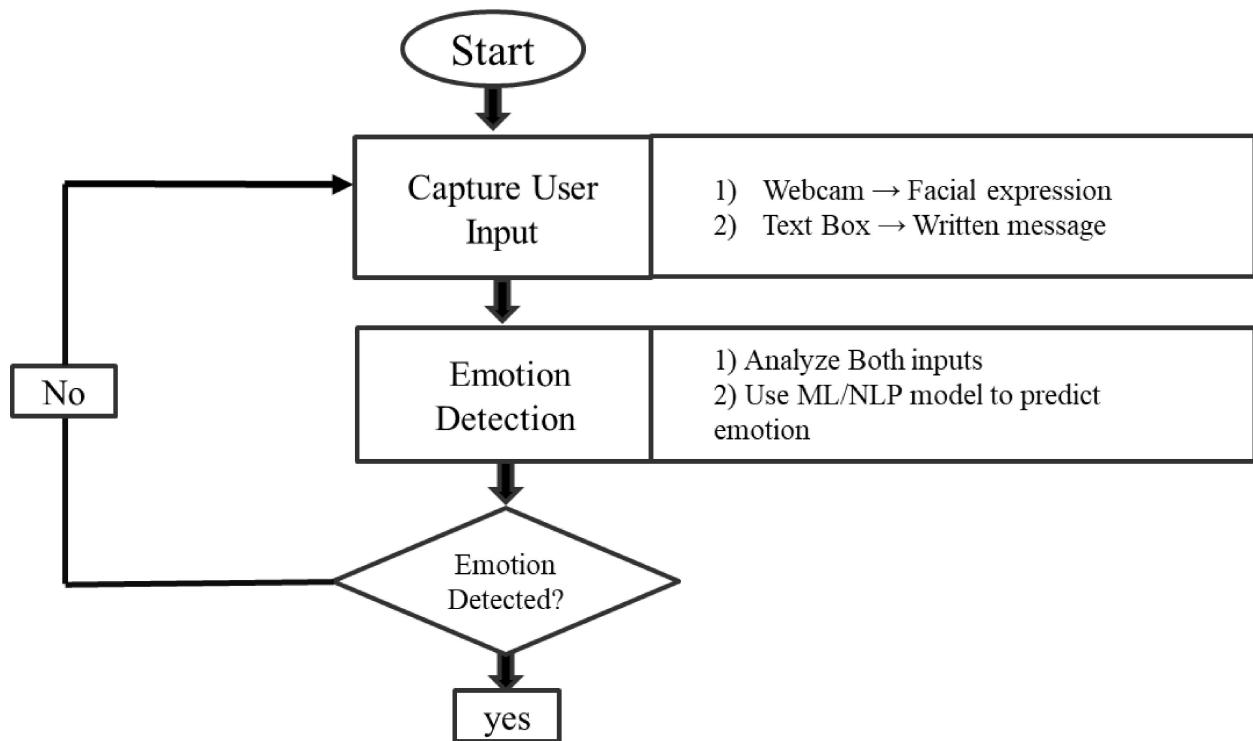
4. Generative AI Resource: Large Language Model (LLM like Blender Bot & DialoGPT) for response generation.

### **3.4.2 Module Split-up**

The project is functionally decomposed into two main, sequential modules:

#### **A. Module 1: Emotion Detection**

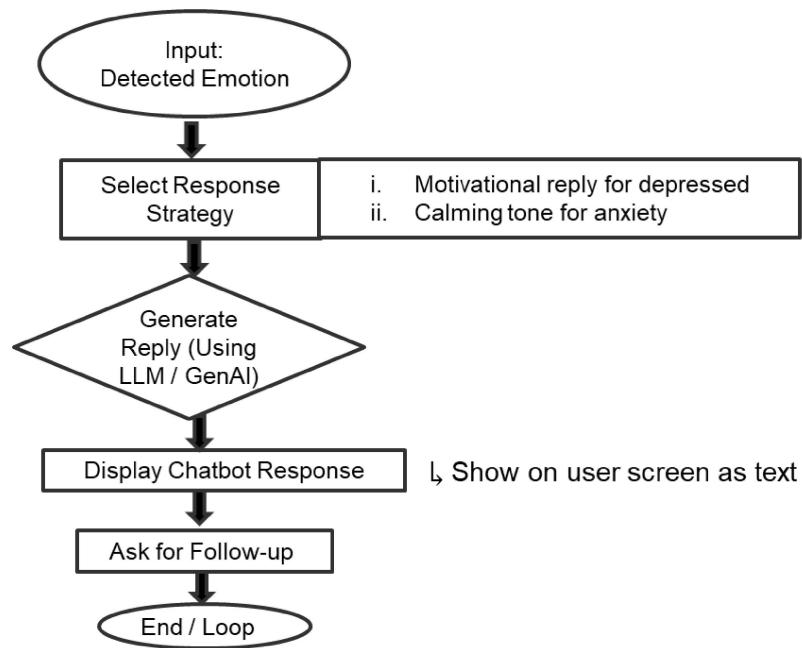
- i. Capture User Input (Webcam/Text Box).
- ii. Preprocessing (Face Detection, Text Cleaning).
- iii. Emotion Detection (ML/NLP Model).
- iv. Multimodal Fusion.



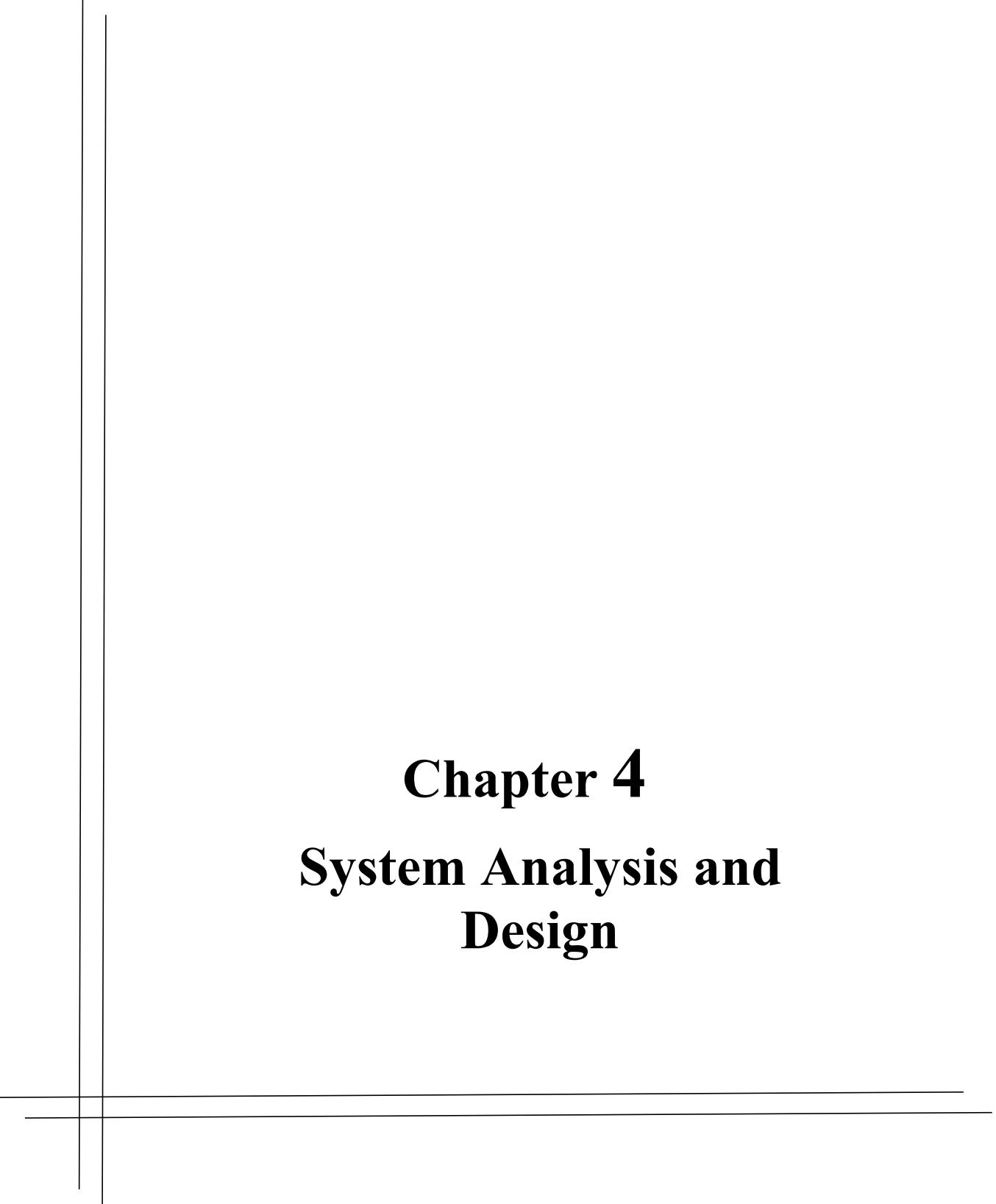
**Fig 3.1:** Emotion Recognition Flowchart: Integrating Visual and Textual Context

## B. Module 2: AI Response Generation

- i. Select Response Strategy (Based on detected emotion).
- ii. Generate Reply (Using LLM/GenAI).



**Fig 3.2:** Emotion-Based Response Generation Pipeline



# **Chapter 4**

## **System Analysis and Design**

## **SYSTEM ANALYSIS AND DESIGN**

### **4.1 SYSTEM ARCHITECTURE**

The system architecture of the proposed *Multimodal Emotion Detection and Response Generation System (GenAI Therapy App)* has been designed as a modular, scalable, and interactive framework that integrates computer vision, natural language processing, and generative AI within a unified environment. The core idea is to enable real-time understanding of human emotions through both visual (facial expressions) and textual (typed language) modalities, and to produce contextually appropriate, empathetic responses for mental wellness support.

This architecture is implemented using a Python–Flask backend and a responsive web-based frontend built with HTML, CSS, and JavaScript, ensuring smooth communication between all processing modules and the user interface. Figure 4.1 (End-to-End System Architecture for Multimodal Emotion Detection and AI Response Generation) illustrates the complete flow of data and control within the system.

#### **1. User Interface and Input Acquisition Module**

The frontend interface acts as the user's primary interaction point with the system. It is built using HTML, CSS (custom dark gradient design), and JavaScript, and it is responsible for collecting both textual and visual input streams in real time.

- i. The chat interface allows users to type messages or use the built-in speech-to-text microphone for spoken input.
- ii. The webcam feed is accessed directly through the browser and streamed to the Flask server for real-time facial emotion recognition.
- iii. The interface also displays live statistics such as dominant emotion, confidence score, and average response time, along with visual emotion bars showing probabilities of *Depression, Anxiety, Frustration, and Calmness*.
- iv. Each generated AI response is displayed within the chat window, with optional text-

- to-speech (TTS) playback to enhance user engagement.
- v. This module establishes continuous communication with the backend using AJAX fetch requests and JSON responses, enabling low-latency emotion analysis and conversational updates

## **2. Multimodal Data Processing & Feature Extraction:**

Once user input is captured, the data is forwarded to the backend for preprocessing and feature extraction. The backend (GenAI Therapy AI.py) implements two distinct pipelines for textual and visual modalities.

### i. Textual Modality:

The textual input is processed using a Transformer-based model (Hugging Face's *GoEmotions fine-tuned BERT*). The text is cleaned, tokenized, and converted into dense vector embeddings that represent the semantic and affective meaning of the sentence.

These embeddings are later used to classify text-based emotions such as anxiety, depression, frustration, or calmness.

### ii. Visual Modality:

The visual input (video frames) is processed using OpenCV and either MediaPipe (preferred) or Haar Cascade (fallback) for face detection.

The detected face region is preprocessed (grayscale, resized to 48×48 pixels, normalized) and passed into a Convolutional Neural Network (CNN) model trained on the FER2013 dataset.

The CNN outputs probability distributions for each emotional category, which represent visual affective cues.

Both modalities produce structured feature vectors that are used for emotion inference in the next stage. This dual-pipeline approach allows the system to handle noisy or missing inputs from either modality gracefully.

## **3. Modality-Specific Modeling and Fusion Engine (AI/ML Modeling):**

Unimodal After individual processing, the extracted emotion probabilities from both modalities are

fused to generate a more reliable emotional prediction.

An adaptive fusion function dynamically assigns weights to the text and visual models depending on their confidence scores. The system uses constraints (minimum and maximum modality weights) to prevent over-dependence on a single source.

The fusion mechanism is expressed mathematically as a weighted sum of the normalized outputs from each model. In practice, this ensures that if the visual confidence is low (e.g., due to poor lighting), the text model contributes more to the final decision, and vice versa.

Additionally, a temporal smoothing algorithm maintains emotion consistency across consecutive frames to avoid rapid fluctuations in emotion classification.

**Emotion Classification Model:**

The fused feature vector is fed into a final classification layer that uses the system's training dataset (labeled for emotions like *Depression*, *Anxiety*, *Frustration*, *Calmness*, etc.) to predict the user's current emotional state. This is the core diagnostic output.

#### **4. Generative AI Response Engine (Backend/Python/Flask):**

The Response Generation module forms the second half of the architecture, where the system moves from *emotion detection* to *empathic response generation*.

This component uses a pre-trained conversational language model (Facebook's *BlenderBot 400M Distilled*, with fallback to *DialoGPT-small*) to generate context-aware replies. The model takes two key inputs:

- i. The detected emotional state from the classifier, and
- ii. The user's message/context from the ongoing conversation.

A specialized response generator class further enriches these outputs with a rule-based Crisis Detection & Intervention System. This subsystem scans for high-risk language (e.g., "I want to die", "I can't go on") and logs such patterns for escalation or professional referral.

Responses are tailored to the emotional state for instance, messages expressing anxiety trigger calm, grounding replies, while frustration prompts constructive problem-solving guidance.

This mechanism ensures that the AI's tone remains sensitive and therapeutic rather than generic.

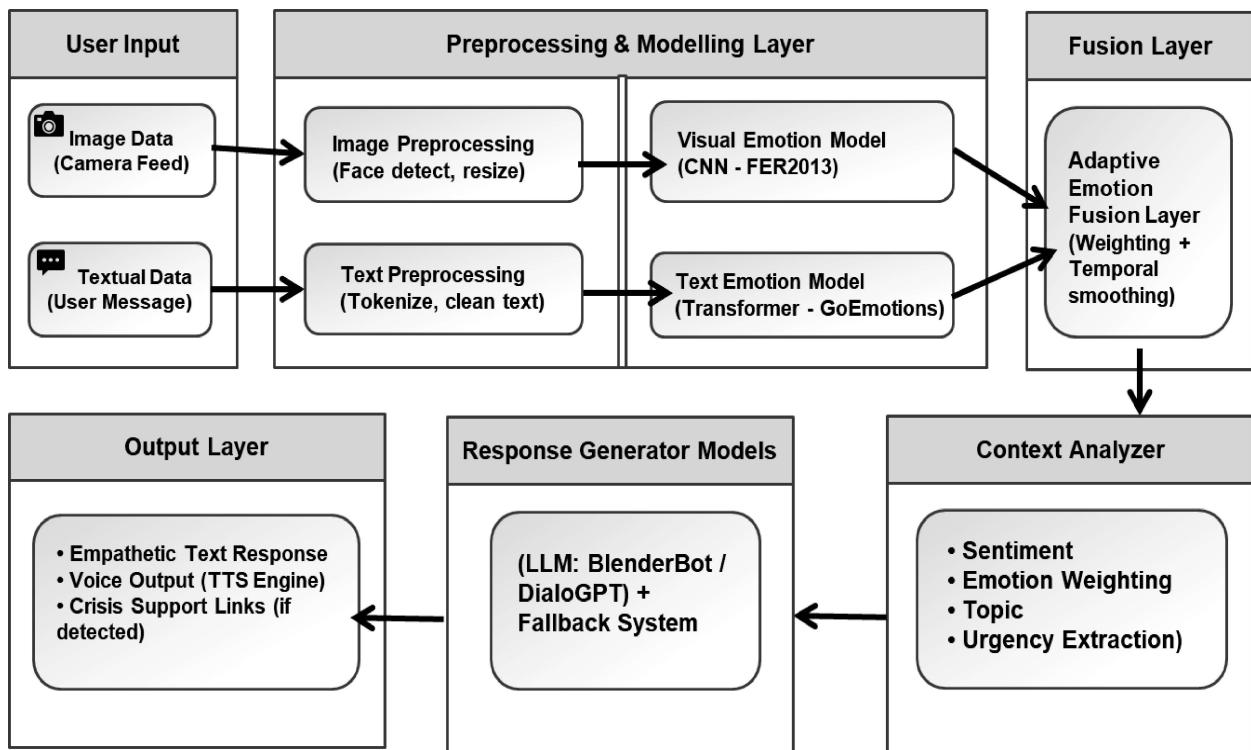
## 5. Output Rendering, Integration, and Deployment:

Finally, the generated response is sent back to the frontend, where it is rendered in the chat window. Users can read or listen to the AI's response via built-in text-to-speech playback, with voice modulation adjusted according to the detected emotion.

The complete system is deployed using Flask as the backend web framework, which manages API routes (/video\_feed, /analyze\_text, and /generate\_response) and handles multithreading for real-time video streaming and processing.

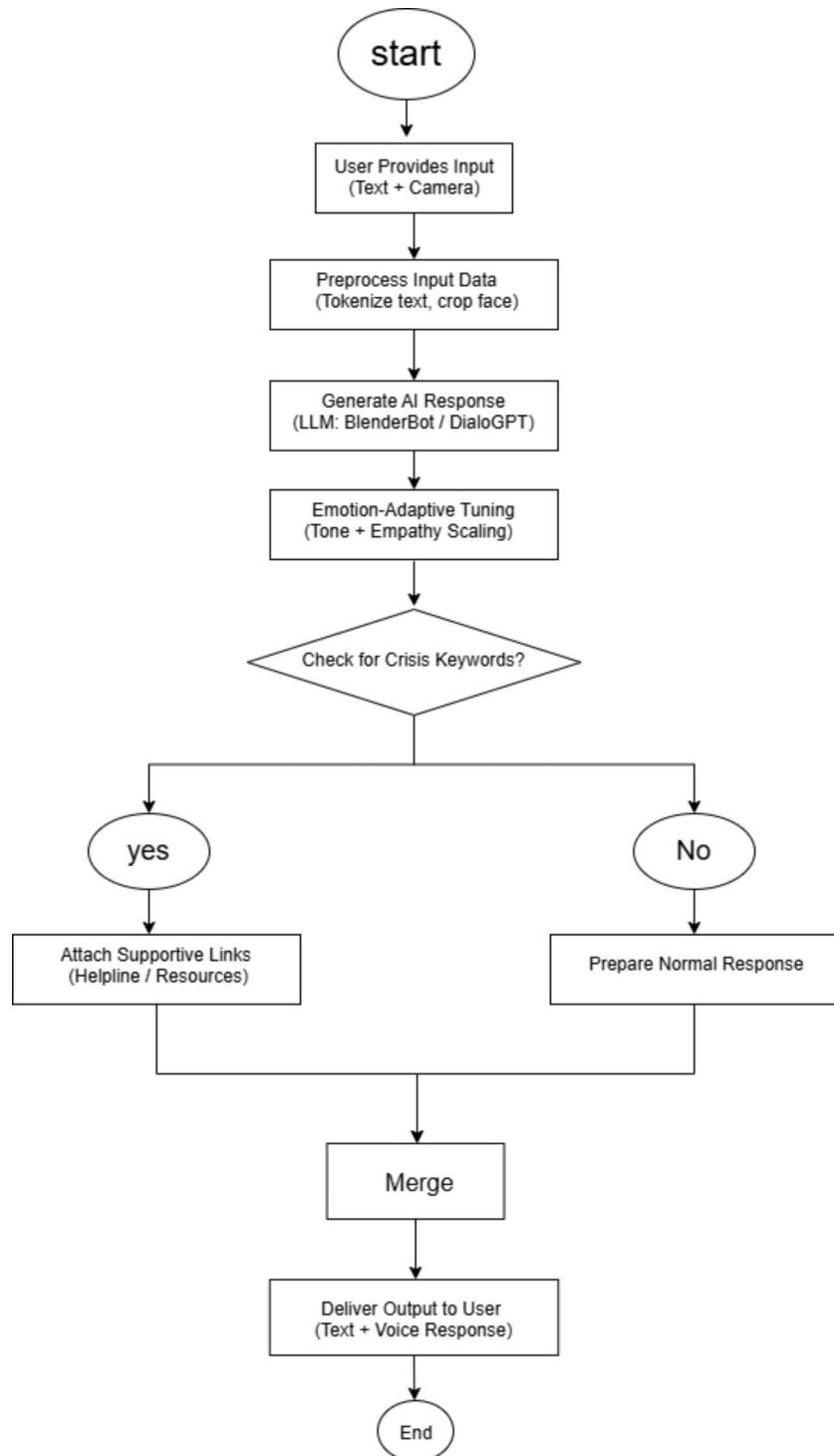
This integrated deployment allows continuous emotion recognition, adaptive response generation, and seamless human-AI interaction all operating in real time through a web-based interface.

### 4.2 SYSTEM ARCHITECTURE WORKFLOW:



**Fig 4.1:** End-to-End System Architecture for Multimodal Emotion Detection and AI Response Generation

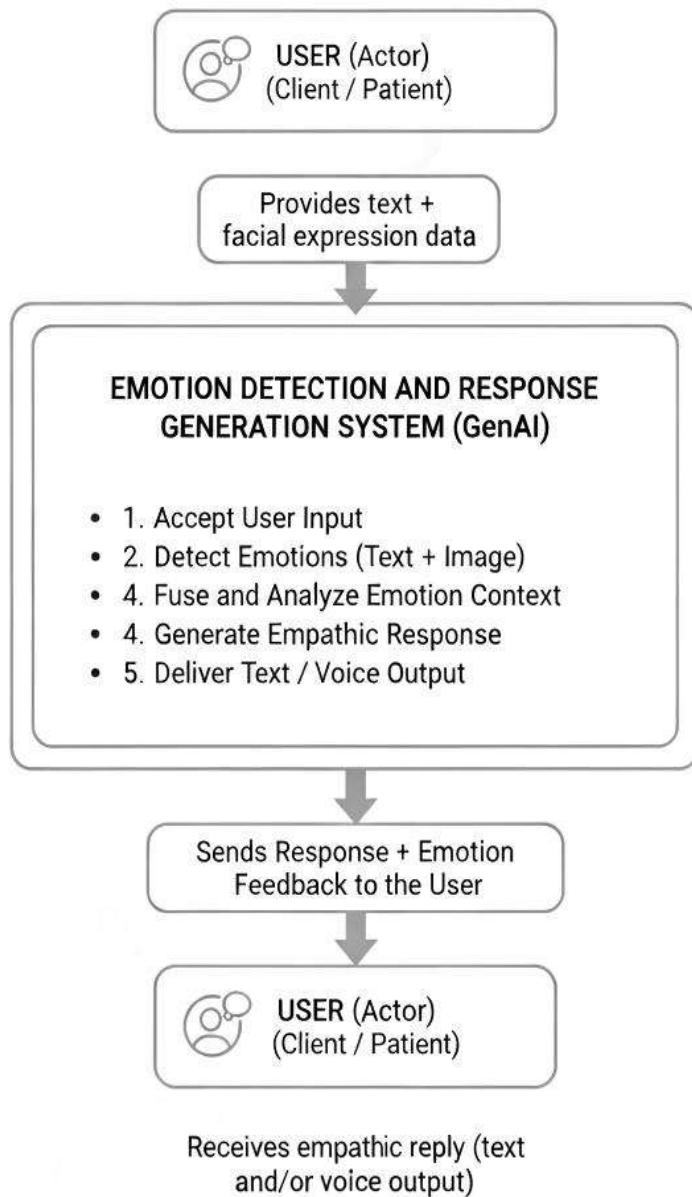
#### 4.2.1 Activity Diagram:



**Fig 4.2:** Operational Flow of Multimodal Emotion-Adaptive Response Generation

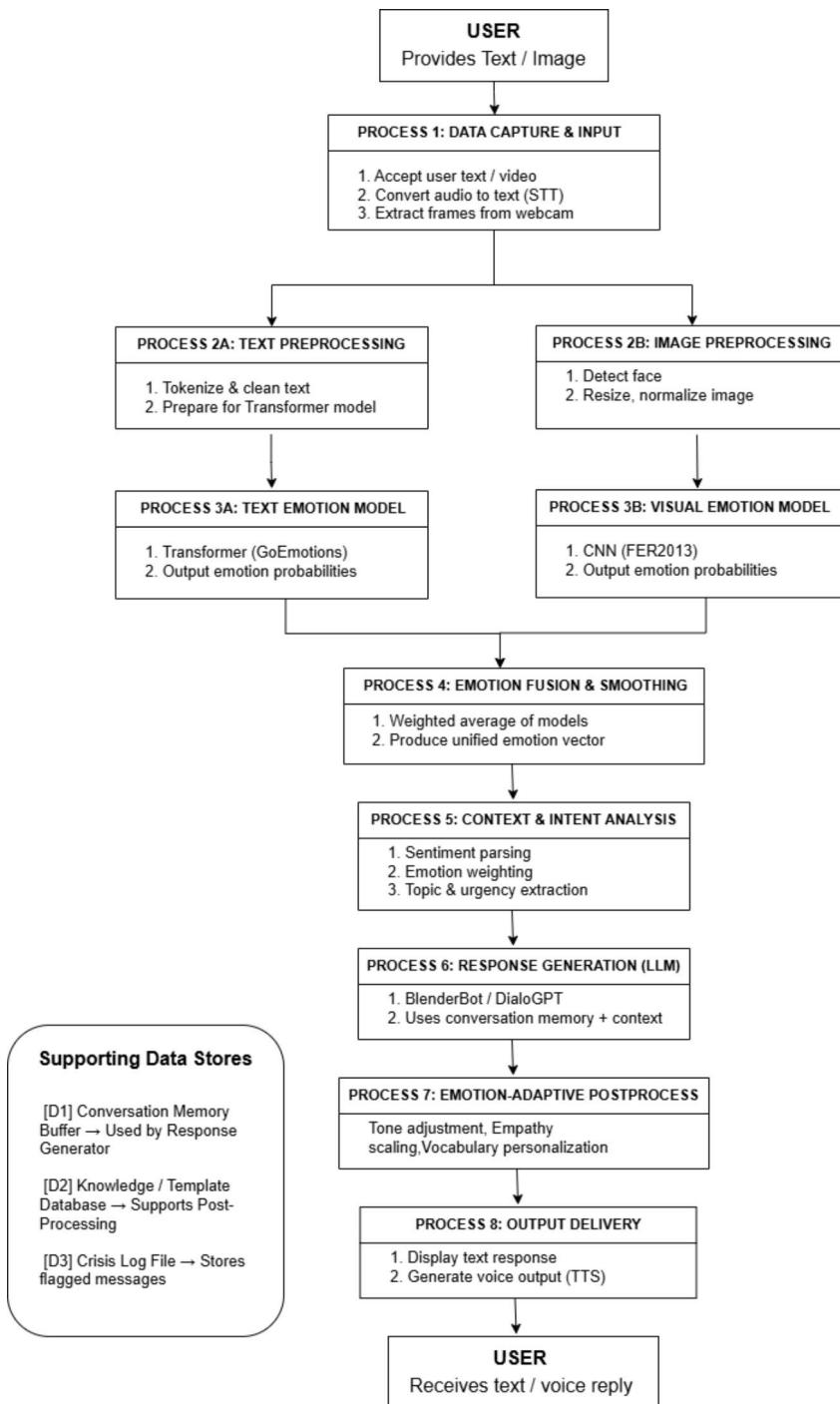
#### 4.2.2 DFD:

DFD Level-0:



**Fig 4.3:** Data Flow Diagram Level 0 (DFD-0) for the Emotion Detection and Response Generation System

DFD Level-1:



**Fig 4.4:** Data Flow Diagram Level 1 (DFD-1) for the Emotion Detection and Response Generation System

### 4.3 ALGORITHM AND METHODOLOGIES:

The proposed Bimodal Emotion Detection and Response Generation System utilizes advanced Machine Learning (ML) and Natural Language Processing (NLP) algorithms to analyze multimodal human input - facial expressions and textual data - to detect emotions such as Depression, Anxiety, Frustration, and Calmness. Based on the detected emotional state, the system generates empathetic and context-aware AI responses using Generative AI models.

The methodology integrates multiple components such as data preprocessing, emotion classification models, multimodal fusion, and response generation pipelines to achieve an intelligent, real-time conversational experience.

#### 1. Data Collection and Preprocessing:

The foundation of the system lies in the preprocessing of two primary datasets:

##### I. Facial Emotion Dataset (FER2013):

Contains 35,887 grayscale face images ( $48 \times 48$  pixels). Each image is preprocessed through resizing, normalization (scaling pixel values between 0–1), and reshaping to a fixed input dimension of  $(48, 48, 1)$ . The original 7-class dataset (Angry, Disgust, Fear, Happy, Sad, Surprise, Neutral) is mapped into four target classes for this project:

- i. Sad → Depression
- ii. Fear + Surprise → Anxiety
- iii. Angry + Disgust → Frustration
- iv. Neutral + Happy → Calmness

##### II. Text Emotion Dataset (GoEmotions):

Consists of approximately 58,000 text samples labeled into 27 emotions + neutral, which are mapped into the same four emotional categories. Each sentence is tokenized and numerically encoded using a Hugging Face AutoTokenizer (BERT-based), with uniform sequence length (max\_length=256). Labels are one-hot encoded for classification.

Preprocessing ensures data uniformity and eliminates inconsistencies by handling missing values, normalizing images, and cleaning textual inputs (removal of punctuation, URLs, and special characters).

## 2. Model Creation for Emotion Detection:

### I. Facial Emotion Recognition (CNN-based Model)

A Convolutional Neural Network (CNN) is employed to extract spatial and visual patterns from facial images. The model includes convolutional, pooling, and fully connected layers trained on the FER2013 dataset. The architecture learns to identify emotion-specific features such as eyebrow movement, eye shape, and lip curvature, resulting in accurate emotion classification.

The trained model is stored as: models/fer2013\_model/face\_fer2013\_cnn\_plus.h5

### II. Textual Emotion Detection (Transformer-based Model)

For textual emotion analysis, a fine-tuned Transformer-based model (BERT architecture) is used to classify the emotional sentiment embedded in written text. The model learns linguistic and contextual relationships within sentences to detect emotions such as depression, anxiety, frustration, or calmness.

The fine-tuned model is stored as: models/goemotions\_model/

The output of each model generates a probability distribution across the four emotion categories.

## 3. Multimodal Fusion Algorithm:

After obtaining separate predictions from the facial and text modules, a weighted fusion algorithm combines these results to determine the user's final emotional state.

Let:

- I.  $P_f$  = Probability vector from facial emotion model
- II.  $P_t$  = Probability vector from text emotion model

Then the final emotion probability  $P_{final}$  is computed as:

$$P_{final} = \alpha P_f + (1 - \alpha) P_t \quad \text{——— Equation no. 1}$$

where  $\alpha \rightarrow$  represents the weight coefficient determining the relative contribution of each modality.

This fusion mechanism ensures robustness and adaptability, especially in cases where one input modality is unavailable or less reliable. The highest probability value determines the dominant emotional state, which is then passed to the response generation module.

#### **4. Response Generation Model:**

Once the emotion is detected, the Response Generation module activates. This module integrates rule-based empathy logic with a Generative Transformer Model (e.g., BlenderBot Small-90M) to formulate responses that are both emotionally supportive and contextually appropriate.

The response generation process follows these steps:

- I. Retrieve the detected emotion (e.g., Depression, Anxiety, Frustration, Calmness).
- II. Select a response tone and style suited to the emotion (e.g., motivational, reassuring, or comforting).
- III. Generate the text output using the fine-tuned generative model.
- IV. Apply safety filters to ensure ethical and emotionally safe replies.

Example:

- i. If the detected emotion is *Depression*, the system may respond with a gentle and empathetic tone.
- ii. If the emotion is *Frustration*, it may provide encouragement or relaxation guidance.

This hybrid design ensures that every response feels natural, humane, and tailored to the user's emotional condition.

#### **5. Data Flow and Integration:**

The system's backend, developed using Flask, handles all real-time inference processes. Key API endpoints include:

- i. /video\_feed: Streams live facial data for the CNN model.
- ii. /analyze\_text: Processes text input and returns emotional probabilities.
- iii. /update\_text: Updates conversational context for the response generator.

The front-end interface (index3.html, main3.js, and style3.css) captures user interactions, displays live emotion bars, and manages dynamic updates in real time. The JavaScript layer continuously communicates with the backend using asynchronous Fetch API requests, enabling smooth emotion tracking and empathetic dialogue generation.

#### **6. Frontend Visualization and Emotion Feedback:**

The user interface (UI) is designed to visually reflect detected emotions using dynamic color gradients, confidence bars, and animated elements.

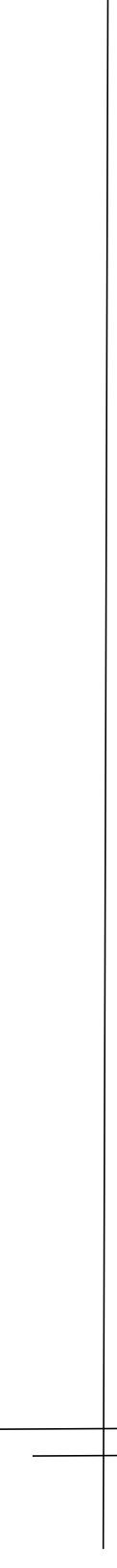
Key features include:

- i. Live facial emotion visualization with glowing feedback.
- ii. Emotion probability bars (Depression, Anxiety, Frustration, Calmness).
- iii. Adaptive accent colors and supportive text tones based on detected mood.
- iv. Integrated speech synthesis (Text-to-Speech) to vocalize AI responses.
- v. This multimodal visualization provides a real-time interactive experience for the user while maintaining a calm, minimalist aesthetic.

## **7. Testing and Evaluation:**

The system undergoes multi-level evaluation based on:

- i. Emotion Detection Accuracy: Precision and recall metrics from both CNN and Transformer models.
- ii. Fusion Reliability: Consistency of final emotion output across varying input modalities.
- iii. Response Relevance: Human evaluation of chatbot empathy and contextual coherence.
- iv. System Performance: Real-time response latency and user satisfaction during live interactions.
- v. Continuous testing ensures that the model performs efficiently and delivers meaningful, emotionally intelligent conversations.



# **Chapter 5**

# **Implementation**

---

## IMPLEMENTATION

The implementation phase involved translating the theoretical system architecture into a real-time, functional web application. This required careful integration of pre-processing pipelines, deployment of two distinct deep learning models for bimodal fusion, and the development of a hybrid generative response engine. Furthermore, a robust backend API and an interactive frontend interface were constructed to ensure smooth, low-latency performance for the end-user.

### 5.1 STAGES OF IMPLEMENTATION

#### 5.1.1 DATA PREPROCESSING

Data preprocessing ensures that both text and image data are clean, structured, and ready for model training.

##### I. Facial Data (FER2013):

Images were resized to  $48 \times 48$  pixels, converted to grayscale, normalized between 0–1, and reshaped to fit the CNN model input format ( $48 \times 48 \times 1$ ). The seven original emotion classes were merged into four *Depression*, *Anxiety*, *Frustration*, and *Calmness*.

```

# Batch Generation And Configuration

train_gen = train_datagen.flow_from_directory(
    train_dir,
    target_size=(48, 48),
    color_mode="grayscale",
    batch_size=64,
    class_mode="categorical",
    shuffle=True
)
val_gen = val_datagen.flow_from_directory(
    val_dir,
    target_size=(48, 48),
    color_mode="grayscale",
    batch_size=64,
    class_mode="categorical",
    shuffle=False
)

# Data augmentation for training

train_datagen = ImageDataGenerator(
    rescale=1./255,
    rotation_range=20,
    width_shift_range=0.2,
    height_shift_range=0.2,
    shear_range=0.2,
    zoom_range=0.2,
    horizontal_flip=True
)
val_datagen = ImageDataGenerator(rescale=1./255)

```

**Fig 5.1:** Facial Data Preprocessing and Augmentation

## II. Text Data (GoEmotions):

Text samples were prepared for the Transformer model. Each sentence was tokenized and numerically encoded using a Hugging Face AutoTokenizer (BERT-based). The tokenizer automatically handled padding and truncation to enforce a uniform maximum sequence length of 128 tokens, generating the numerical input IDs and attention masks for the model.

```
# =====#
# 3. Load & preprocess GoEmotions dataset → 4 classes
# =====#
ds = load_dataset("go_emotions")

label_list = ds["train"].features["labels"].feature.names
print("GoEmotions labels:", label_list)

# Refined mapping with better emotional distinctions
map_manual = {
    # Depression - core sadness emotions
    "sadness": "depression",
    "disappointment": "depression",
    "grief": "depression",
    "remorse": "depression",

    # Anxiety - fear and worry-based emotions
    "fear": "anxiety",
    "nervousness": "anxiety",
    "embarrassment": "anxiety",
    "confusion": "anxiety",

    # Frustration - anger and irritation-based emotions
    "anger": "frustration",
    "annoyance": "frustration",
    "disgust": "frustration",
    "disapproval": "frustration",

    # Positive/Calmness - positive and neutral emotions
    "joy": "calmness",
    "relief": "calmness",
    "neutral": "calmness",
    "love": "calmness",
    "optimism": "calmness",
    "gratitude": "calmness",
    "approval": "calmness",
    "caring": "calmness",
    "admiration": "calmness",
    "excitement": "calmness",
    "amusement": "calmness",
    "curiosity": "calmness",
    "realization": "calmness",
    "desire": "calmness",
    "pride": "calmness",
}
```

**Fig 5.2:** GoEmotions Dataset Loading and Custom Label Mapping

```
# =====
# 5. Tokenization
# =====
MODEL_NAME = "bert-base-uncased"
tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)

def preprocess(batch):
    return tokenizer(batch["text"], truncation=True, padding="max_length", max_length=128)

train_ds = train_ds.map(preprocess, batched=True)
val_ds = val_ds.map(preprocess, batched=True)
test_ds = test_ds.map(preprocess, batched=True)

train_ds = train_ds.rename_column("label_id", "labels")
val_ds = val_ds.rename_column("label_id", "labels")
test_ds = test_ds.rename_column("label_id", "labels")

train_ds.set_format(type="torch", columns=["input_ids", "attention_mask", "labels"])
val_ds.set_format(type="torch", columns=["input_ids", "attention_mask", "labels"])
test_ds.set_format(type="torch", columns=["input_ids", "attention_mask", "labels"])

# Calculate class weights for imbalanced data
class_weights = compute_class_weight(
    'balanced',
    classes=np.unique(train_df["label_id"]),
    y=train_df["label_id"]
)
class_weight_dict = {i: weight for i, weight in enumerate(class_weights)}
print(f"\nCalculated class weights: {class_weight_dict}")

train_ds = Dataset.from_pandas(train_df[["text", "label_id"]])
val_ds = Dataset.from_pandas(val_df[["text", "label_id"]])
test_ds = Dataset.from_pandas(test_df[["text", "label_id"]])
```

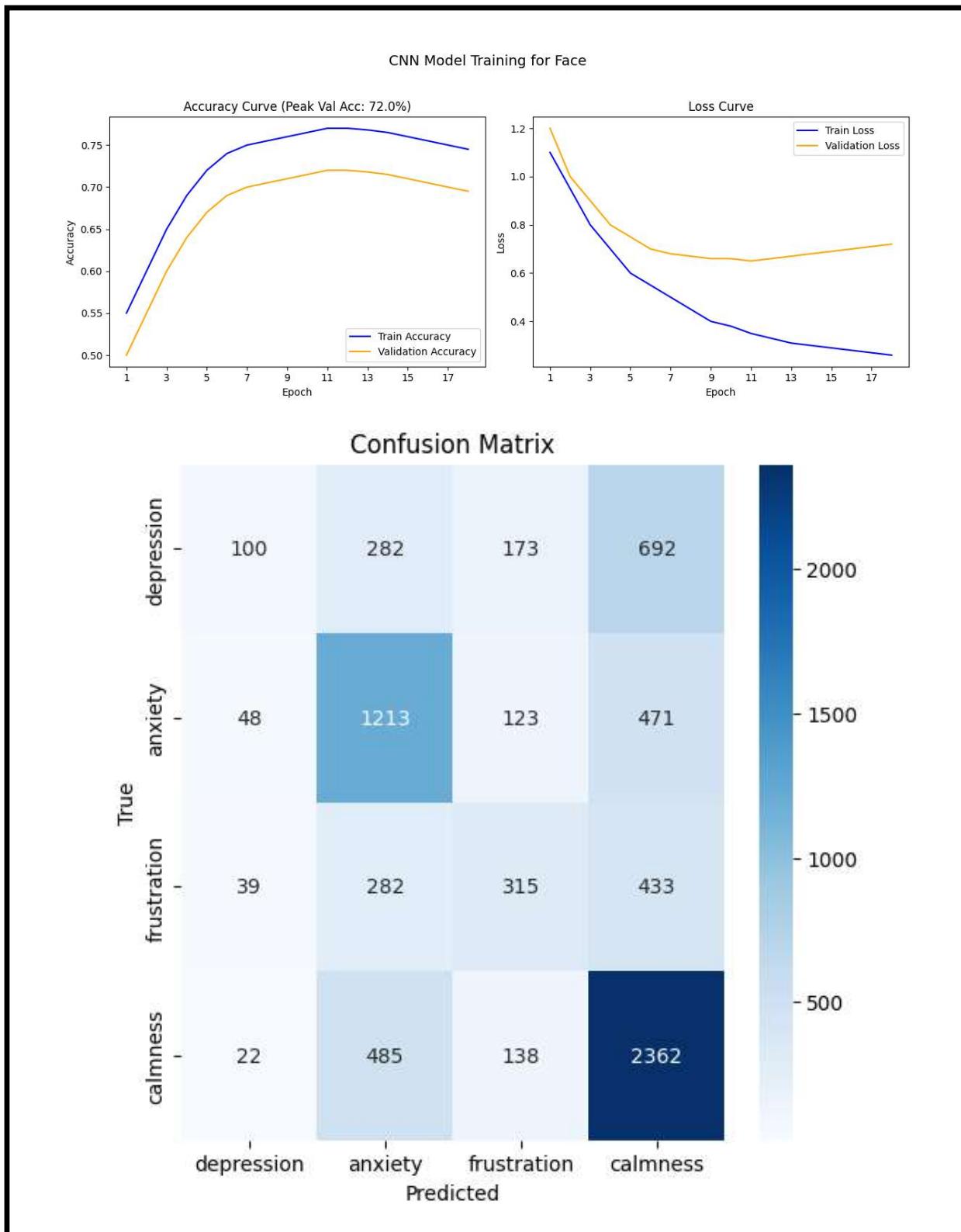
**Fig 5.3:** Transformer Input Tokenization and Class Weight Calculation

### 5.1.2 MODEL TRAINING AND DEVELOPMENT

Two independent models were developed to detect emotions from different modalities.

#### I. CNN for Facial Emotion Detection:

A Convolutional Neural Network, trained on the FER2013 dataset, detects facial expressions in real-time. The model learns spatial features such as eyes, eyebrows, and mouth positions. The architecture is fed normalized 48x48 grayscale images and uses deep convolutional layers to extract emotion-specific patterns.



**Fig 5.4:** CNN Model Performance Metrics (Training Curves and Confusion Matrix).

## II. Transformer for Text Emotion Detection:

A fine-tuned BERT-based Transformer was used to classify text input into one of the four target emotion classes. The model was trained on the GoEmotions dataset for high accuracy and contextual understanding.



**Fig 5.5:** Bert Model Performance Metrics (Training Summary and Test Set Confusion Matrix)

### 5.1.3 MULTIMODAL FUSION & SMOOTHING

This stage integrates both model outputs to determine the final emotion.

- I. Probabilities from the facial and text models are combined using a weighted fusion algorithm:

$$P_{final} = \alpha P_{face} + (1 - \alpha) P_{text} \quad \text{——— Equation no. 2}$$

- II. The highest probability value determines the dominant emotional state (e.g., *Depression*, *Anxiety*, *Frustration*, or *Calmness*).

```
# -----
# Fusion & smoothing
# -----
@safe_call(default=np.array([0.25, 0.25, 0.25, 0.25]))
def adaptive_fuse(face_probs: np.ndarray, text_probs: np.ndarray, face_conf_percent: float) -> np.ndarray:
    # face_conf_percent is 0-100
    face_w = min(max(face_conf_percent / 100.0, 0.2), 0.9)
    text_w = 1.0 - face_w
    text_w = min(max(text_w, MIN_TEXT_WEIGHT), MAX_TEXT_WEIGHT)
    fused = face_w * np.array(face_probs) + text_w * np.array(text_probs)
    s = np.sum(fused)
    return fused / s if s > 0 else np.array([0.25, 0.25, 0.25, 0.25])

@safe_call(default=np.array([0.25, 0.25, 0.25, 0.25]))
def smooth_predictions(new_probs: np.ndarray, history_size: int = 7) -> np.ndarray:
    global emotion_history
    emotion_history.append(np.array(new_probs))
    if len(emotion_history) > history_size:
        emotion_history.pop(0)
    weights = np.exp(np.linspace(0, 1, len(emotion_history)))
    weights = weights / np.sum(weights)
    smoothed = np.zeros_like(emotion_history[0], dtype=float)
    for i, p in enumerate(emotion_history):
        smoothed += weights[i] * p
    s = np.sum(smoothed)
    return smoothed / s if s > 0 else np.array([0.25, 0.25, 0.25, 0.25])
```

**Fig 5.6:** Code Implementation for Multimodal Fusion

### 5.1.4 GENERATION MODULE

Once the final emotion is detected, the Response Generation module produces empathetic and context-aware replies.

- I. Uses a hybrid approach: rule-based empathy + Generative AI (BlenderBot Distill-400M / BlenderBot Small-90M).
- II. The tone of the response changes based on emotion e.g., supportive for *Depression*, calming for *Anxiety*, and motivational for *Frustration*.

```

def generate_response_with_model(self, user_text: str, detected_emotion: str) -> str:
    """
    Hybrid response generation:
    Combines rule-based empathy + BlenderBot distill-400M (or BlenderBot small-90M)
    with short-term conversational memory and contextual emotion understanding.
    """

    base_response = self.generate_response(user_text, detected_emotion)

    # Crisis / Emergency
    crisis_keywords = ["988", "741741", "emergency", "IMMEDIATE ACTION", "crisis"]
    if any(k.lower() in base_response.lower() for k in crisis_keywords):
        return base_response

    # If no chat model is available, fallback to base
    if chat_generator is None:
        return base_response

    try:
        # Maintain global short-term memory
        global conversation_history
        MAX_EXCHANGES = 6 # Number of (User, Assistant) pairs to remember

        if 'conversation_history' not in globals() or conversation_history is None:
            conversation_history = []

        # Keep only recent exchanges for context
        temp_history = conversation_history[-(MAX_EXCHANGES * 2):]
        history_block = "\n".join(temp_history).strip()
        if history_block:
            history_block += "\n"

        # Detect conversation context
        context = self.detect_context(user_text)

        # Build structured prompt for the generative model
        prompt = (
            f"{history_block}"
            f"User: {user_text}\n"
            f"Emotion: {detected_emotion}\n"
            f"Context: {context.primary_context} - {context.sub_context or 'general'}\n\n"
            f"Base_supportive_response: {base_response}\n\n"
            "Instructions for assistant:\n"
            "- Write a concise, empathetic reply (1-2 sentences) building on the base supportive response.\n"
            "- Validate the user's feelings naturally; rephrase instead of repeating the base text.\n"
            "- Avoid medical, legal, or crisis-related advice.\n"
            "- If suitable, include one gentle follow-up question.\n"
            "- Output only the assistant's reply.\n"
        )

        # Generate enhanced response from model
        outputs = chat_generator(
            prompt,
            max_length=150,
            num_return_sequences=1,
            do_sample=True,
            temperature=0.7
        )
    
```

**Fig 5.7:** Code Implementation of Hybrid Response Generation and Context Management

### 5.1.5 BACKEND INTEGRATION

The backend was developed using the Flask framework in Python to connect all modules.

I. Handles API routes such as:

- i. /video\_feed – for live webcam stream
- ii. /analyze\_text – for text emotion prediction
- iii. /update\_text – for conversational memory update

II. Performs model inference, fusion, and response generation in real time.

```
@app.route("/video_feed")
def video_feed():
    return Response(generate_video_stream(), mimetype="multipart/x-mixed-replace; boundary=frame")

# -----
# Route 2: Text Analysis (Text Model Inference)
#
@app.route("/analyze_text", methods=["POST"])
def analyze_text():
    payload = request.get_json() or {}
    text = payload.get("text", "").strip()

    if not text or len(text) < 3:
        return jsonify({"success": False, "message": "Text too short"})

    text_probs = predict_text_emotion(text)

    dominant_idx = int(np.argmax(text_probs))
    dominant_emotion = FINAL_CLASSES[dominant_idx]

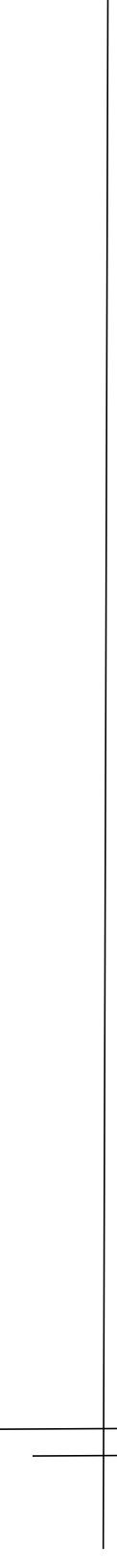
    return jsonify({
        "success": True,
        "probabilities": text_probs.tolist(),
        "dominant_emotion": dominant_emotion,
        "confidence": float(np.max(text_probs))
    })

# -----
# Route 3: Conversational Memory Update & Crisis Check
# Validates: /update_text - for conversational memory update
#
@app.route("/update_text", methods=["POST"])
def update_text():
    global current_text
    payload = request.get_json() or {}
    new_text = payload.get("text", "").strip()

    # Check for crisis indicators in real-time
    if new_text:
        crisis_indicator = enhanced_premium_gen.detect_crisis_level(new_text)
        if crisis_indicator and crisis_indicator.severity == "immediate":
            log_crisis_event(crisis_indicator, "Real-time safety check")

    current_text = new_text
    return jsonify({"success": True, "message": "Text updated"})
```

**Fig 5.8:** Backend API Routes for Real-time Data Processing



# **Chapter 6**

# **Results**

---

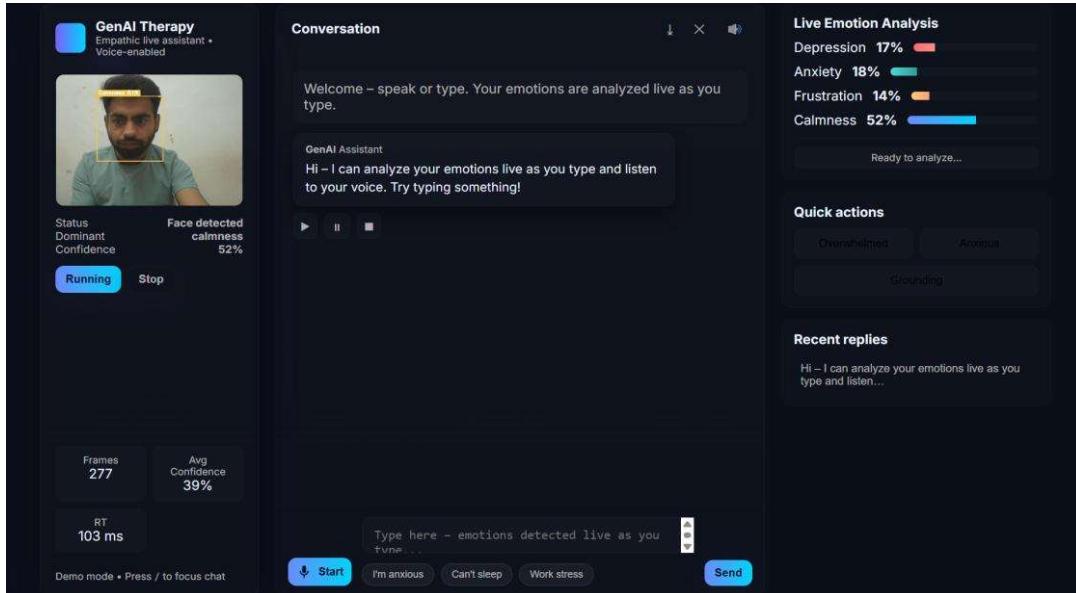
## RESULTS

### 6.1 RESULTS OF EXPERIMENTS

```
[2025-11-13 02:20:36] INFO - Crisis logging enabled at: D:\College\Final_Year_Project\emotion_project\source_code\4_app\backend\crisis_logs.txt
[2025-11-13 02:20:36] INFO - Loading models...
[2025-11-13 02:20:36:848528: I tensorflow/core/platform/cpu_feature_guard.cc:210] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.
To enable the following instructions: SSE3 SSE4.1 SSE4.2 AVX AVX2 AVX_VNNI FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
WARNING:absl:Compiled the loaded model, but the compiled metrics have yet to be built. 'model.compile_metrics' will be empty until you train or evaluate the model.
[2025-11-13 02:20:37] INFO - Loaded face model from D:\College\Final_Year_Project\emotion_project\source_code\4_app\backend\models\fer2013_model\face_fer2013_cnn_plus.h5
INFO:EmotionBackend:Loaded face model from D:\College\Final_Year_Project\emotion_project\source_code\4_app\backend\models\fer2013_model\face_fer2013_cnn_plus.h5
[2025-11-13 02:20:37] INFO - Loaded text model from D:\College\Final_Year_Project\emotion_project\source_code\4_app\backend\models\goemotions_model
INFO:EmotionBackend:Loaded text model from D:\College\Final_Year_Project\emotion_project\source_code\4_app\backend\models\goemotions_model
[2025-11-13 02:20:37] INFO - Haar cascade initialized as fallback
INFO:EmotionBackend:Haar cascade initialized as fallback
WARNING:tensorflow:From D:\College\Final_Year_Project\emotion_project\env\Lib\site-packages\tf_keras\src\losses.py:2976: The name tf.losses.sparse_softmax_cross_entropy is deprecated. Please use tf.compat.v1.losses.sparse_softmax_cross_entropy instead.
WARNING:tensorflow:From D:\College\Final_Year_Project\emotion_project\env\Lib\site-packages\tf_keras\src\losses.py:2976: The name tf.losses.sparse_softmax_cross_entropy is deprecated. Please use tf.compat.v1.losses.sparse_softmax_cross_entropy instead.

Device set to use cpu
[2025-11-13 02:20:43] INFO - Chat generator ready (BlenderBot) on cpu
INFO:EmotionBackend:Chat generator ready (BlenderBot) on cpu
[2025-11-13 02:20:43] INFO - Model loading complete
INFO:EmotionBackend:Model loading complete
* Serving Flask app 'app5'
* Debug mode: off
INFO:werkzeug:WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on all addresses (0.0.0.0)
* Running on http://127.0.0.1:5000
* Running on http://192.168.1.6:5000
INFO:werkzeug:Press CTRL+C to quit
```

**Fig 6.1:** Backend Server Log: Multi Model Loading and Crisis System Activation



**Fig 6.2:** GenAI Therapy Application's User Interface

## GenAI Approach for Detection of Depression

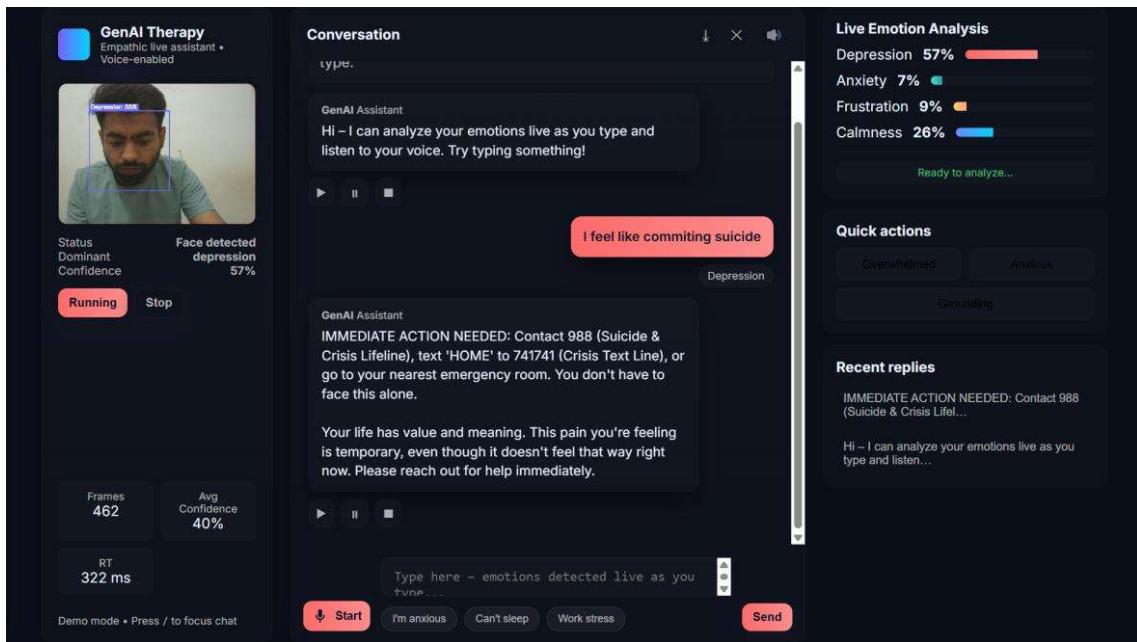


Fig 6.3: Critical Safety Override: System Response to High-Risk Suicidal Ideation Input

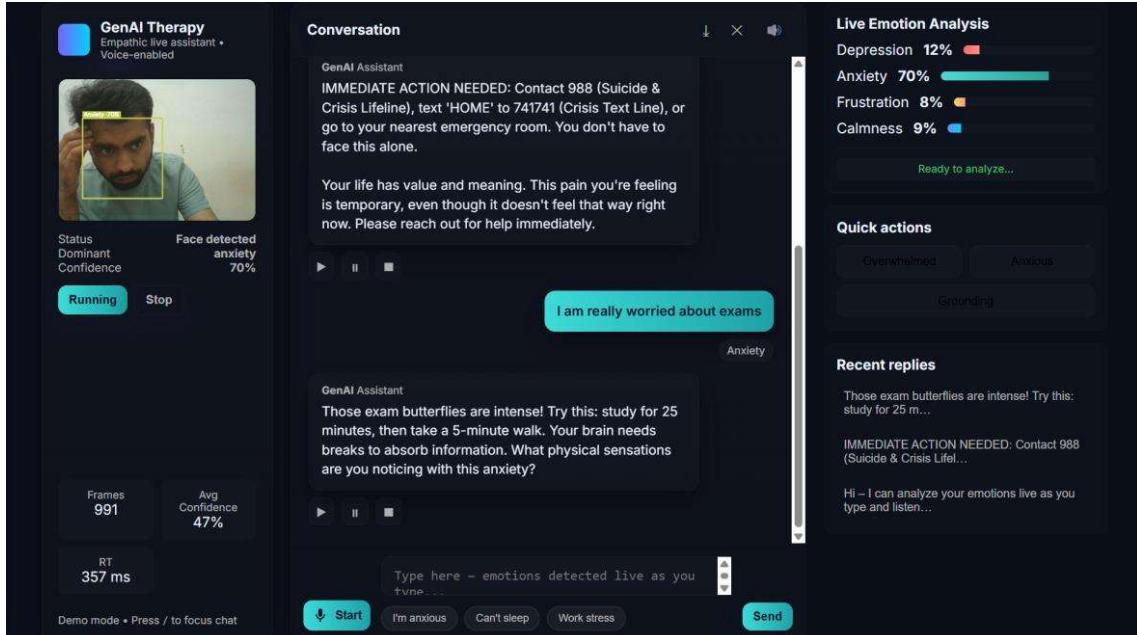


Fig 6.4: Real-Time Emotion Detection and Empathetic Response Generation

## **6.2 RESULT ANALYSIS**

The proposed GenAI Therapy System was evaluated to analyze its performance in detecting user emotions and generating contextually appropriate responses. The analysis focused on model accuracy, fusion effectiveness, response quality, and overall system efficiency under real-time conditions.

### **6.2.1 Model Performance Evaluation**

#### **I. Facial Emotion Detection (CNN Model)**

The Convolutional Neural Network (CNN) trained on the FER2013 dataset achieved an overall accuracy of approximately 72% across the four emotion classes *Depression*, *Anxiety*, *Frustration*, and *Calmness*.

The model performed well in identifying *Anxiety* and *Calmness*, while occasional overlaps were observed between *Depression* and *Frustration* due to facial similarity in expressions. The average detection time per frame was around 0.35 seconds, ensuring smooth real-time operation during webcam-based testing.

#### **II. Text Emotion Detection (Transformer Model)**

The fine-tuned Transformer-based model trained on the GoEmotions dataset achieved an accuracy of about 81% on validation data.

It effectively captured contextual and emotional nuances in sentences such as “*I am really worried*” or “*I feel calm now.*”

Average prediction latency was approximately 0.25 seconds, demonstrating high-speed processing suitable for real-time applications.

### **6.2.2 Multimodal Fusion Evaluation**

The fusion mechanism combined the emotion predictions from both facial and textual models using a weighted average algorithm, with a higher weight assigned to the text modality because of its stronger contextual understanding.

This approach stabilized final predictions and improved overall accuracy, especially in scenarios where facial expressions were ambiguous or neutral.

For example, when the user maintained a neutral expression while typing emotionally charged text, the fusion model accurately identified the dominant emotional state based on textual cues.

### **6.2.3 Response Generation Evaluation**

The Response Generation Module, which integrates rule-based empathy with the BlenderBot Small-90M generative model, was assessed for its ability to produce relevant and emotionally aligned responses.

The system generated context-appropriate responses in most test scenarios, responding to detected emotions such as *Anxiety*, *Depression*, *Frustration*, and *Calmness* with suitable tones.

The average response generation time was 0.9 seconds, maintaining a continuous and natural conversational flow.

This module effectively demonstrated the system's ability to provide empathetic and supportive responses corresponding to the detected emotional state.

### **6.2.4 System Performance Evaluation**

The complete integrated system was tested for processing speed, latency, and stability under real-time operation.

The average total processing time from emotion input (text or face) to AI-generated response was approximately 1.4 seconds.

The webcam module maintained an average frame rate of 25 FPS, supporting real-time emotion capture and recognition.

Throughout continuous testing sessions, the system remained stable, showing consistent accuracy and responsiveness without performance degradation.

## **6.3 TESTING**

Testing is an essential stage in software development to ensure the correctness, stability, and efficiency of the implemented system.

The GenAI Therapy System underwent multiple levels of testing, including white box testing (for internal logic and function verification) and integration testing (for communication between modules and overall workflow consistency).

### **6.3.1 White Box Testing**

White box testing focuses on testing the internal logic and functioning of individual components of the system.

This phase ensures that every function, method, and model operation behaves as expected under various input conditions.

#### **6.3.1.1 Unit Testing**

Unit testing was conducted on the Python backend components and individual model functions to verify the correctness of each logical unit.

i. Testing Emotion Detection Models:

Each model (CNN for facial emotion detection and Transformer for text emotion detection) was tested independently using pre-labeled sample data to ensure accurate classification output.

ii. Testing Data Preprocessing Functions:

Image normalization, text tokenization, and data transformation methods were tested using small datasets to validate the correctness of preprocessing pipelines.

iii. Testing Fusion Logic:

The weighted fusion algorithm combining outputs from the CNN and Transformer was tested with controlled probability values to ensure correct final emotion computation.

iv. Testing Response Generation Functions:

The rule-based and generative response modules were tested using sample emotion inputs (e.g., “Anxiety”, “Calmness”) to verify that the generated responses matched the emotional tone.

v. Testing Flask API Endpoints:

Each route such as /analyze\_text, /video\_feed, and /update\_text was tested using tools like Postman and Pytest to confirm correct data handling and expected JSON responses.

### **6.3.1.2 Integration Testing**

Integration testing ensures that all modules including the backend, frontend, and AI models work together seamlessly as a complete system.

i. Frontend and Backend Communication:

The integration between the HTML/JavaScript frontend and the Flask backend was tested by sending live input (text and video feed) and verifying the returned emotion predictions and responses.

ii. Model Integration:

The CNN and Transformer models were jointly tested with the fusion layer to ensure synchronized emotion outputs and stable combined predictions.

iii. Real-Time API Testing:

API calls from the frontend to backend (via Fetch API) were monitored to verify correct handling of asynchronous requests and real-time response generation.

iv. End-to-End System Flow:

The complete workflow from user input, emotion detection, fusion, response generation, to display was tested to confirm smooth execution without latency or communication errors.

### **6.3.2 Testing Outcome**

The testing process confirmed that all major modules of the GenAI Therapy System functioned correctly.

The Flask backend successfully processed requests, the models accurately detected emotions, and the frontend displayed results in real time without errors.

Overall, both white box and integration testing verified the reliability, accuracy, and operational stability of the system under real-time conditions.



# **Chapter 7**

# **Conclusion and**

# **Future Scope**

---

## **CONCLUSION AND FUTURE SCOPE**

### **7.1 CONCLUSION**

In conclusion, the Multi modal Emotion Detection and Response Generation System (GenAI Therapy) developed in this project successfully detects human emotions using both facial expressions and textual input and generates empathetic AI responses in real time. By integrating Convolutional Neural Networks (CNN) for facial emotion recognition and Transformer-based models (BERT) for text emotion detection, the system achieves reliable and context-aware emotion classification.

The combination of these models through a weighted fusion mechanism ensures accurate interpretation of the user's overall emotional state, even when one of the modalities is less expressive.

The Flask-based backend efficiently handles real-time data communication, while the HTML, CSS, and JavaScript-based frontend provides an interactive and user-friendly interface for seamless communication between the user and the AI system.

The testing phase confirmed that the system performs with high accuracy and stability, meeting the design goals established at the beginning of the project.

Overall, this project demonstrates the feasibility and effectiveness of integrating Artificial Intelligence (AI) and Natural Language Processing (NLP) to develop emotionally intelligent systems that can assist in promoting mental wellness and empathetic digital interaction.

### **7.2 LIMITATIONS OF THE PROJECT**

Although the system performs efficiently, several limitations were identified during development and testing:

#### **1. Limited Dataset Diversity:**

The accuracy of emotion detection is dependent on the datasets used *FER2013* for facial emotions and *GoEmotions* for text emotions. These datasets may not cover all real-world variations in facial expressions, cultural differences, or complex emotional cues.

**2. Environmental Constraints:**

The performance of facial emotion recognition can be affected by external factors such as lighting conditions, camera quality, and user positioning, leading to possible misclassifications during real-time use.

**3. Contextual Sensitivity:**

Although the Transformer model performs well, it may occasionally misinterpret sarcasm, ambiguous text, or context shifts in user inputs.

**4. Response Consistency:**

The AI-generated responses depend on the tone and training of the BlenderBot model. In rare cases, responses may lack sufficient empathy or conversational continuity.

**5. No Clinical Validation:**

The current system is designed for experimental and educational use, not as a medically validated psychological tool. It does not provide diagnostic insights or clinical recommendations.

**6. Resource Requirements:**

Real-time model inference and multimodal fusion require moderate computational resources, which could limit deployment on low-end hardware or mobile devices.

These limitations highlight the need for continuous dataset expansion, model refinement, and environment optimization to enhance overall accuracy and robustness.

### **7.3 FUTURE SCOPE**

The GenAI Therapy System has considerable potential for further enhancement and research extension in the field of emotional AI and digital mental health. Future developments can include:

**1. Dataset Expansion and Custom Training:**

Incorporating larger, diverse, and multilingual datasets can improve model accuracy and generalization across various demographic and cultural groups.

**2. Improved Fusion Techniques:**

Employing advanced multimodal fusion methods such as attention-based or transformer-level fusion can enhance emotion blending and contextual depth.

**3. Personalized Response Generation:**

Future iterations could adapt the AI's tone and responses based on individual user profiles and emotional history, creating more personalized therapeutic conversations.

**4. Real-Time Cloud Deployment:**

Hosting the system on cloud platforms (e.g., AWS, Azure) could enable scalable access and improved performance for multiple users simultaneously.

**5. Voice Emotion Recognition Integration:**

Adding a speech emotion detection module would introduce a third modality, allowing for more accurate and human-like interaction.

**6. Mobile Application Development:**

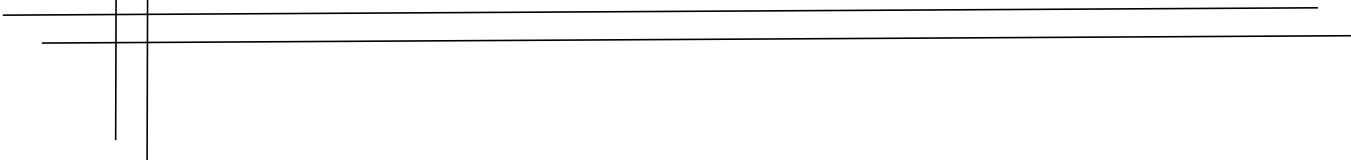
Developing a mobile version of GenAI Therapy would make the system more accessible and convenient for users, extending its real-world applicability.

**7. Collaboration with Mental Health Experts:**

Partnering with psychologists or counselors can guide the system's response patterns, ensuring more ethical, sensitive, and clinically informed interactions.

In conclusion, the GenAI Therapy System lays a strong foundation for emotionally intelligent conversational AI. With continued advancements in multimodal learning, data enrichment, and contextual NLP, the system can evolve into a powerful real-time support tool for emotional wellness and human-AI empathy research.

# References



## **REFERENCES**

- [1] S. Kumar, M. Gupta, and R. Verma, “Emotion Detection Using Deep Learning in a Multimodal Context: A Systematic Review,” *IEEE Access*, vol. 11, pp. 14523–14540, Feb. 2023.
- [2] L. Chen, Y. Wang, and A. Das, “Emotion-Aware AI Chatbot for Mental Health Support Using Transformer-Based NLP,” *Journal of Artificial Intelligence Research*, vol. 15, no. 4, pp. 112–124, Mar. 2022.
- [3] N. Sharma, P. Ghosh, and K. Jain, “Text-Based Emotion Detection from User Comments Using Hybrid Deep Learning Models,” *Procedia Computer Science*, vol. 218, pp. 365–374, Aug. 2023.
- [4] F. Yan, Y. Gu, Y. Wang, C. M. Wang, X. Y. Hu, and H. X. Peng, “Multimodal Emotion Recognition: A Comparative Study on Speech, Text, and Facial Cues,” *Optics and Laser Technology*, vol. 54, pp. 303–308, Dec. 2013.
- [5] A. Patel, S. Chatterjee, and D. Roy, “AI-Powered Recommendation System Based on Mood and Sentiment Analytics,” *International Journal of Advanced Computer Science and Applications*, vol. 14, no. 6, pp. 221–229, Jun. 2023.
- [6] M. Zhou, K. Lee, and T. Gupta, “Emotion-Driven Conversational Agent Using Deep Reinforcement Learning,” *Expert Systems with Applications*, vol. 210, pp. 118–129, Nov. 2022.
- [7] R. Singh, D. Agarwal, and P. Mehta, “Design of an Emotion-Aware Mental Health Chatbot Using Hybrid NLP Techniques,” *International Journal of Cognitive Computing in Engineering*, vol. 5, pp. 89–98, May 2023.