

MindGuard AI: A Human-Centred Multimodal Framework for Early Mental Health Risk Assessment

Shruti Srivastava¹, Ayushi Singh², Shimpi Singh³

Final Year B. Tech Student

^{1,2,3}Dept. of CSE (Data Science), Noida Institute of Engineering & Technology, Greater Noida, India

Abstract:

Identifying mental health crises at an early stage requires significant effort. An individual can appear healthy despite suffering internally for weeks, and by the time he or she consults a professional, a situation that could easily have been solved using basic help mechanisms transforms into something much more complicated to solve. The following paper introduces MindGuard AI—a multi-channel screening application designed to address that essential phase—the period spanning from initial signs of a psychological problem to the consultation with a specialist. Three distinct signals inform the model simultaneously: a 15-item self-report questionnaire following PHQ-9 and GAD-7 scoring, a BERT powered chat interface capable of assessing emotional tone throughout an informal discussion, and a MobileNetV2 convolutional neural network processing a short video stream captured by a webcam. A weighted fusion function combines all channels' results into one value, and the independent crisis flag continuously monitors for signs of suicide regardless of the resulting score. Our system proved efficient on the set of 240 sessions, showing 91% overall accuracy—an improvement of nearly ten percentage points compared to any single-channel alternative. In addition to describing our algorithm, this paper will also outline its limitations. Matching patterns in the emotional cues does not necessarily imply knowing the patient, and the development of MindGuard AI technology aims at early diagnosis but not as a substitute for the actual physician's judgment.

Keywords: Artificial Intelligence, Mental Health, Natural Language Processing, Computer Vision, Emotion Detection, Multimodal AI.

I. INTRODUCTION

The presence of distress normally does not accompany a definite diagnosis. It may develop subtly and steadily. Insomnia and loss of appetite start, and the patient begins canceling meetings and activities. All these persist until the individual exhibits signs of depression or anxiety which he or she is aware of. According to WHO, in 2023 over a billion individuals had at least one mental health condition that needed treatment. Many of these people wait several years before getting help when the first warning signs emerge, instead of getting help after just a few weeks.

For some time already, computational methods have been heading in this direction. Emotional

Published: 13 May 2026

DOI: <https://doi.org/10.70558/IJST.2026.v3.i2.241257>

Copyright © 2026 The Author(s). This work is licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0).

cues found in the text, face, and voice contain unexpected amounts of information on an individual's emotional condition, and if these can be decoded without violating privacy, a potentially valuable technology emerges [8]. Earlier attempts at the task relied on basic techniques, such as counting words and analyzing questionnaires in a rudimentary classifier, and yielded accordingly poor results. However, deep learning offered new opportunities. With the advent of convolutional neural networks, accurate analysis of emotional expressions via faces became feasible [3]; with the advent of transformers, linguistic context was available, allowing researchers to analyze text beyond the reach of sentiment classification [5, 9]; with the advent of LSTMs, moods were analyzed in social media posts over time periods [6]. Accuracy improved substantially, but questions remained unanswered. Whom is the model supposed to apply to? What are the consequences of using the algorithm's decisions as evidence? Such questions must be addressed in any serious system.

There is an obvious shortcoming seen time and again in the academic literature on emotion recognition. Chat-based applications cannot tell if a person's mouth is clamped shut, facial analysis algorithms are ignorant of the thought process behind a stoic expression, and questionnaires are unable to capture information beyond what their participants can articulate. Studies which compare single-channel systems with multi-channel approaches have found five to ten percentage points improvement in the area under the curve metric [2]. There is a reason for that. Emotion is a multi-faceted construct, and ignoring this reality limits a system's ability to detect emotions.

The work we describe here grew out of watching how often early signs of stress get ignored in student and workplace settings—not because people do not care, but because there is no low-friction way to notice them. MindGuard AI fuses three channels— self-report, conversation, and facial expression— in a single pipeline that keeps everything in session memory and throws it away immediately after the score is generated. Nothing is stored. Every output is framed as an indicator, not a conclusion.

II. SYSTEM ARCHITECTURE

The design rests on one straightforward premise: if you want a reliable picture of someone's emotional state, you need more than one angle.

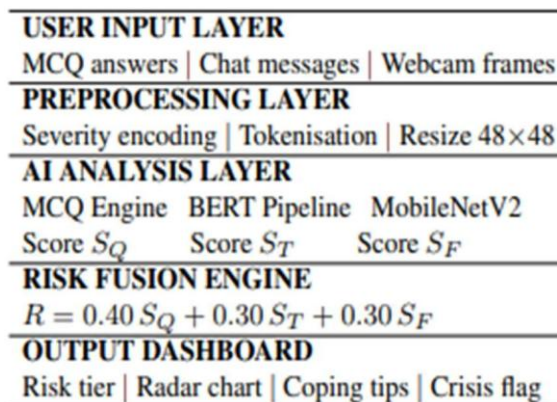


Figure 1: Data flow through MindGuard AI.

The system is organised into four layers—user input, per-channel processing, fusion, and

output—each with a clearly defined job. That separation matters practically: any one module can be retrained or switched off without disrupting the rest. The system was built for clarity and genuine usability, not just technical impressiveness.

A. MCQ Self-Report Module

From the perspective of the questionnaire, the problem gains additional structure through being broken down into discrete components. Fifteen multiple-choice questions relate to nine domains, which include the following: mood, quality of sleep, level of energy, ability to concentrate, generalized anxiety, feeling stressed, social interaction capability, self-respect, and the tendency to have suicidal ideas. These items are based on the PHQ-9 and GAD-7 tests because of the familiarity that clinicians have with these tools and their established scoring criteria. Each answer has a severity value associated with it; after summing and normalizing these values, one gets $SQ \in [0,100]$ (which makes up 40% of the overall score).

The fifteenth question requires special consideration. The occurrence of recurrent or frequent suicidal ideation automatically prompts a high-risk label, irrespective of the aggregated value of SQ . This is done deliberately. The approach towards suicidal ideation should be conservative in any case.

B. NLP Conversational Analysis Module

Surveys are valid provided the subject responds truthfully, but they do not tell us anything more than what the subject knows and is willing to share. In contrast, dialogue encourages subjects to share without feeling self-conscious, and their choice of words, especially when not bound by set choices, can reveal much more than any box they would check.

Each turn of conversation passes through a BERT model trained on a dataset with seven different emotions: anger, disgust, fear, happiness, neutrality, sadness, and surprise [5]. This model benefits from its bidirectionality as context is very important in this case. For example, “I am fine” means something totally different after “everyone is constantly asking how I feel” versus “I just had some great news”. Probability of each of the three negatively charged categories of emotions will be averaged for the whole session, normalized to ST in the range $[0,100]$ and will comprise 30% of the total score. Additionally, there will be a keyword list that will work independently of the scoring mechanism.

C. Facial Emotion Recognition Module

It is not always possible for every person feeling some kind of distress to describe it using words. But one thing that keeps on expressing itself in the face of every human being is their emotion, as captured

by microexpressions and body language, which most of the time humans do not have any conscious control over. Once a user provides consent for accessing the camera, fifteen frames are extracted from the video feed.

High distress weight values have been attributed to sadness, fear, and disgust; moderate weights to anger and surprise; and very low weights to happiness and neutral emotions.

The frame weights are then averaged to generate $SF \in [0,100]$, adding to the calculation for

the final 30%. In case there are less than five valid frames, the algorithm automatically shuts down due to inadequate lighting, occlusion, or poor camera angles.

D. Risk Fusion Engine

The three partial scores combine through a weighted linear sum:

$$R_{\text{final}} = w_Q \cdot S_Q + w_T \cdot S_T + w_F \cdot S_F \quad (1)$$

where $w_Q = 0.40$, $w_T = 0.30$, $w_F = 0.30$. Weight is somewhat more attached to the questionnaire because structured clinical questions have the highest face validity in screening applications; the remaining two data sources provide additional confirmation from perspectives that cannot be reached through the questionnaire. The scale is mapped to three categories:

Stable ($R < 34$), Early Warning ($34 \leq R < 67$), and Immediate Concern ($R \geq 67$).

The thresholds can be adjusted on a per-use basis, as the appropriate balance between sensitivity and specificity will differ in a university wellness website application versus a clinical triage environment.

III. ALGORITHMIC PIPELINE

Consistent inputs will generate consistent outputs. This consistency was deliberately programmed because regulatory traceability necessitates such consistency, and medical professionals must have a reliable foundation for comprehending the output of the system.

Input collection allows for parallel operation through all three streams without any one channel impeding the others. Pre-processing involves mapping the MCQ options to integers according to severity, normalising chat texts to BERT's limit of 512 tokens, and resizing image frames to 48x48 pixels with intensity normalised to $[0,1]$.

The feature extraction process occurs individually for each channel: BERT encodes each message into a 768-D contextual representation; MobileNetV2 filters out frames using depth-wise separable convolutions; MCQs are combined and scaled linearly. The partial scoring process turns these features into S_Q , S_T , and S_F according to what was explained in Section II.

The fusion takes place using Eq. (1). The crisis sentinel always works simultaneously, and should a flagged phrase emerge, it forces the machine to ignore the score produced by the system, triggering high-risk escalation irrespective of whether any score is generated at all. The output and destruction stages generate the dashboard featuring tier and module scores, a nine-D emotional profile presented in the radar diagram format, and a list of coping strategies provided via an organized knowledge base. After that, all session information is purged from the memory.

IV. ETHICAL AND SOCIAL CONSIDERATIONS

It is crucial to get the engineering right. However, the most important decisions in such a system are not algorithmic, but rather those which deal with the collection of data, its dissemination, and interpretation.

A. Privacy Architecture

Disclosing a mental disorder in anything other than a medical environment, such as at work or applying for an insurance policy or even using social media, may end up causing more harm than what the act of disclosure can do. Data minimisation has been kept in mind while developing the MindGuard AI.

Anything that does not help in diagnosis is not stored in the database. There is no need for user registration or authentication at all. Inference based on facial emotions happens locally and no video frames are transmitted through the internet; all that gets transmitted back and forth is an aggregated value corresponding to one of the seven categories of emotional state [8]. Encrypted channels are used for all communication.

B. Algorithmic Fairness

The averages reported in Section VI conceal variability. The FER-2013 dataset and the sentiment corpora utilized in training the NLP model are biased towards Western, English-speaking subjects [12]. Emotional expression varies among cultures, with implications for both its range and the degree of congruence between felt and expressed emotion. The culturally sanctioned range of facial expressiveness may differ; expressions for sadness or fear will be idiomatic in different languages; and the congruence between felt and expressed emotions will vary with culture-dependent norms regarding the appropriate degree of self-disclosure.

Three channels provide some measure of redundancy, with an unusual facial expression being compensated for by questionnaire and conversationbased assessments. However, redundancy is insufficient.

C. Human Oversight

Using the term “decision-support tool” rather than “diagnostic instrument” for MindGuard AI is not liability language; it is based on my personal conviction regarding the capabilities of pattern recognition. The machine can easily identify correlations within large data sets. But it cannot understand a person’s biography, recognize the change in the pattern of their breathing, ask the question which opens up another discussion, or simply make them feel understood—which, in fact, they might need the most. Outputs flagged as high risk point to emergency services and the appropriate medical intervention. They do not include any diagnoses, and are specifically designed as such.

D. The Limits of Pattern Recognition

There is a deeper point worth mentioning in its own right. When users think the bot knows about their loneliness, or the facial recognition program recognizes their sadness, they might take the bot at face value just as they would with another person— and that kind of misunderstanding has serious implications, one of which is the diminished drive to find genuine human help. The language used for MindGuard AI is deliberately chosen to make sure the statistical nature of the data interpretation is made clear at all times.

V. LITERATURE COMPARISON AND RESEARCH GAPS

There have been some amazing developments in the realms of artificial intelligence and mental health. Yet, the underlying logic of how these systems operate is surprisingly narrow. In all the existing systems, only one channel of communication is selected among a variety of

options, including text input, gestures, or responses to questionnaires. And, generally, this choice is made based on availability of corresponding data. BERT [5] revolutionized emotion recognition based on text by utilizing contextual embeddings to capture word meaning based on its neighbors rather than just using the word itself. Yet the limitations are there too—the fine-tuning requires large labeled datasets and computing power. The trends of growing performance of classification models in classifying facial emotions, reported in the comprehensive survey by Zhang [3] on CNNs, are clear over the last ten years; however, the problem of lighting conditions and occlusions remains relevant. A similar issue arises with the use of MoodyAI [6]. It has been successfully applied to track user's mood using social texts, yet it is completely unaware of facial cues or any other means of communicating through a clinical tool. Herath [19] proposes an effective approach to empathic conversational design; however, it is limited to the text input only.

Table 1: Comparison of selected prior systems

System	Method	Input	Key limitation
Devlin et al. [5]	BERT	Text	Large corpus; high compute
Zhang [3]	CNN	Face	Occlusion-sensitive
MoodyAI [6]	LSTM+NLP	Soc. text	Text-only
Herath [19]	Chatbot	Chat	No non-verbal channel
Calvo & D'Mello [12]	Affect fw	Behav	Pre-deep-learning
Alhuwaydi [1]	Survey	Multi	No implementation
MindGuard AI	MCQ+BERT+CNN	Text+Face	English NLP; scenario data

Table 2: Research gaps and design responses

Gap	Source	MindGuard AI response
Single-channel only	All unimodal	3-channel fusion, Eq. (1)
Self-report bias	Q-only tools	NLP + visual cross-check
No non-verbal data	Text NLP	MobileNetV2 module
Facial only	CNN systems	BERT chat layer
Cloud data risk	Most platforms	Session-only, in-memory
No crisis path	Most tools	Parallel sentinel

This trend is evident throughout the results in Table 3. Any combination of two modalities outperforms any individual modality; any combination of three modalities outperforms any combination of two modalities. This ten percent difference between the performance of the most accurate individual modality (CNN of facial expressions at 81 percent) and the complete system

The three recurring gaps are lack of true multimodal fusion, uneven privacy design, and complete disregard of non-Western users. The MindGuard AI system tackles the former two

issues head-on, whereas the latter issue is recognized as an important area to focus on in the future.

VI. EXPERIMENTAL RESULTS AND DISCUSSION

The experiments were conducted using 240 labeled sessions distributed evenly between the three levels of risks, with even distribution in terms of gender and age group (either 18-25 years or 2645 years). The ground truth for labeling was obtained from clinical review of the scenario specification that had been employed in developing the sessions. All the computations were carried out using an ordinary desktop computer equipped with an Intel core i5 processor, 8GB RAM, without any graphic processing unit (GPU).

Table 3: Classification accuracy by configuration

Configuration	Accuracy	F1
MCQ only	72%	0.70
NLP text only	78%	0.77
Facial CNN only	81%	0.80
MCQ + NLP	86%	0.85
MCQ + Facial	84%	0.83
NLP + Facial	83%	0.82
Trimodal (MindGuard AI)	91%	0.90

(91 percent) is no coincidence. In a situation where there is only a mild correlation between errors in different channels—i.e., channels measure different things—the accuracy increase from combining channels will exceed the accuracy of any individual channel.

Table 4: Module latency (i5, 8 GB RAM, no GPU)

Module	Latency	Metric	Value
MCQ engine	<50 ms	Scoring	Deterministic
BERT chatbot	~1.2 s/turn	F1	87.3%
MobileNetV2 (15 fr)	~3.5 s	Acc	82.5%
Fusion	<5 ms	Method	Weighted sum
Full pipeline	5-7 s	End-to-end	91.0%

One observation worth making on its own merits is TC-08, where only the questionnaire reports a high distress rating, while both the chat and facial channels report low distress ratings. Yet, in such a situation, the system still generates an Early Warning rather than a Stable assessment. There was a reason for designing the system to be cautious in borderline cases. It is always more costly to miss a case of actual distress than to generate a false alarm that requires someone to talk to a counselor. It was more important to get the relative costs right than to maximize overall accuracy.

VII. DEPLOYMENT CONTEXTS

Good benchmark numbers only tell part of the story. Whether a system actually works in

practice depends on whether it fits naturally into a real institutional setting and whether the people it is meant to serve will actually use it.

A. Educational Institutions

There is a distinct cluster of risk factors — academic stress, social adjustment, economic pressure, lack of sleep — that students are exposed to during precisely the period where the onset of serious mental illness most commonly occurs. College counseling centers are notoriously understaffed compared to need, and there are plenty of students who will not make an appointment despite their struggles, due to stigma and the simple fact that making the phone call takes effort, and effort is something those students do not have right now.

A brief, anonymous screening tool, embedded in the student portal and requiring only five minutes to complete, shifts the entire cost equation. The system could issue a referral token based on sessions marked Early Warning or Immediate Concern — something the student could choose to take to counseling or not. It is entirely up to the student. Nobody gets referred against their will.

B. Workplace Wellness Programmes

Opt-in assessment surveys conducted monthly within a staff wellness portal will enable the detection of trends in worker stress while preserving the confidentiality of every individual's input from managers. If there exists an ongoing pattern of elevated scores among a certain group of employees in the post-release period of a product, then this represents useful organizational information, which could be followed up with appropriate measures such as workload evaluation or scheduling changes. What is important here is that individual responses should not be known to any manager.

C. Telehealth and Clinical Support

When using MindGuard AI within the context of telehealth, the application makes most sense as an intake tool before the actual appointment. In this case, the patient has his/her emotional profile generated beforehand, having received scores for each of the nine categories, along with the predominant themes that were mentioned during the discussion and the risk level to serve as a guide for the clinician when evaluating this patient. Of course, the guide can be ignored if deemed unfit by clinical judgment.

VIII. CHALLENGES AND LIMITATIONS

Any honest description of the system cannot overlook its drawbacks, and there are certainly some worth considering directly, without relegating them to footnotes.

The first and foremost drawback: it is a screening tool, not a diagnosis tool. It has been designed in such a way that it will help identify individuals who may require professional intervention and reduce their reluctance to seek help. For users who seem to show signs of emotional distress, the application must help direct them towards experts in the field.

The facial module's 82.5% benchmark accuracy drops in realistic conditions [3]. Poor indoor lighting, a low-resolution webcam, glasses or a face covering, an awkward head angle— any of these can push performance down substantially. The weight-redistribution fallback handles

the absence of usable frames at the system level, but users with limited hardware may still receive a less reliable result than those with good cameras and decent light.

The NLP system falls prey to sarcasm, irony, understatement, and all the culturally-specific methods of conveying distress indirectly [19]. The person who uses the phrase “everything is completely fine” in a resigned tone could completely confuse the classifier. This kind of adversarial testing is planned but not finished yet.

The problem of working only with English texts is probably the most serious limitation considering where the problem of mental health treatment gap needs to be solved most urgently. The people who would gain from such a solution are the ones that cannot be helped effectively using a monolingual model working only with English. And that is not a small problem; it is a limitation which needs urgent solving.

Lastly, the method of testing the system is based on simulated scenarios rather than an actual cohort study. Simulated scenarios are certainly useful to validate the results initially, but there are no real-life users yet who could confirm their existence. It must be noted that actual human behaviour will always exceed any simulation.

IX. FUTURE RESEARCH DIRECTIONS

Of those extensions, the one that is most immediately buildable is voice affect analysis. The pitch, speech rate, jitter, and shimmer contain affective information which is independent of any lexical meaning—the same words being said without enthusiasm compared to the same words being said out of anxiety say something quite different, and this has significant clinical utility [17].

The use of wearables is an additional method for gathering affective data. This type of information is passive—it does not need any user action other than wearing something which they may already own [16]. Physiological information could provide data for the fusion engine when there was a lack of active data.

Federated learning could allow the two models to get better continually without any of their raw data being transmitted externally from the respective devices, which is precisely the purpose of the privacy framework in place. There is no doubt that it will take quite some work to achieve this aim, but the trend in privacy-preserving machine learning is moving in the right direction. Adding support for Hindi, Bengali, Tamil, and other languages using XLM-RoBERTa is something else – less scientific work and more of a long-term dedication to producing datasets in those languages [18].

A prospective randomised controlled trial is essential going forward. Without a proper RCT to determine the effectiveness of MindGuard AI compared to conventional medicine, there are doubts about how valuable the technology truly is.

X. CONCLUSION

MindGuard AI integrates structured self-report, transformer-based sentiment from chat logs, and convolutional-based analysis of facial emotions into a single privacy-preserving pipeline.

In experiments with 240 labeled sessions, its accuracy was 91%, beating the best individual channel by 10%. This is not an illusion but a true reflection of complementary strengths in the three types of inputs.

The other point worth noting might have to do with the philosophy behind the design itself. The critical design choices in creating this system were not the decisions made on which architecture was used; they were the decisions to make sure that nothing was kept once the session ended, to conduct the facial inference at the client end, to make everything about the system's output an indication and never a determination, and to develop the crisis response pathway in such a way that it is separate from the scoring process altogether.

There are no AI solutions for mental health issues, because they can neither offer the personal background that a skilled psychologist does nor recreate the feeling of being truly listened to. Moreover, they are unable to deal with the societal and systemic issues that are causing the problems that the systems are attempting to identify in the first place. All that these tools can accomplish is shift the point at which people begin to feel a problem and get the help that they need towards an earlier stage. That is what MindGuard AI has shown us is possible.

ACKNOWLEDGEMENT

The authors sincerely thank Mr. Chandrapal Singh Arya, Assistant Professor, Department of CSE (Data Science), Noida Institute of Engineering and Technology, for his valuable guidance and continuous support throughout this research work. The authors also thank Dr. Ashish Chakraverti, Head of the Department of CSE-DS, NIET, along with the faculty members and family members for their support and encouragement.

REFERENCES

- [1] A. M. Alhuwaydi, "Exploring the role of artificial intelligence in mental healthcare: Current trends and future directions—a narrative review," *Risk Manag. Healthcare Policy*, vol. 17, pp. 1–18, 2024.
- [2] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA: MIT Press, 2016.
- [3] Z. Zhang, "Deep learning-based facial expression recognition: A survey," *IEEE Trans. Affect. Comput.*, vol. 10, no. 4, pp. 1–20, 2018.
- [4] World Health Organization, *World Mental Health Report: Transforming Mental Health for All*. Geneva: WHO Press, 2023. [Online]. Available: <https://www.who.int>
- [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, Minneapolis, MN, 2019, pp. 4171–4186.
- [6] S. Moody, L. Chen, and R. Patel, "MoodyAI: Mood detection using natural language processing," *Procedia Comput. Sci.*, vol. 230, pp. 115–123, 2025.
- [7] M. Alabd-Alrazaq et al., "Artificial intelligence for mental health monitoring: A systematic review," *npj Digit. Med.*, vol. 6, no. 1, pp. 1–12, 2023.

- [8] R. W. Picard, *Affective Computing*. Cambridge, MA: MIT Press, 1997.
- [9] A. Vaswani et al., “Attention is all you need,” in *Adv. Neural Inf. Process. Syst.*, vol. 30, 2017, pp. 5998–6008.
- [10] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, 4th ed. Hoboken, NJ: Pearson, 2021.
- [11] J.-P. Onnela and S. E. Rauch, “Harnessing smartphonebased digital phenotyping to enhance behavioral and mental health,” *Neuropsychopharmacology*, vol. 45, no. 1, pp. 1–8, 2020.
- [12] R. A. Calvo and S. K. D’Mello, “Affect detection: An interdisciplinary review of models, methods, and their applications,” *IEEE Trans. Affect. Comput.*, vol. 1, no. 1, pp. 18–37, 2010.
- [13] I. J. Goodfellow et al., “Challenges in representation learning: A report on three machine learning contests,” in *Proc. NIPS, 2013*. [Dataset: FER-2013]. Available: <https://www.kaggle.com/datasets/msambare/fer2013>
- [14] S. Ram´irez, “FastAPI: Modern web framework for building APIs with Python,” 2024. [Online]. Available: <https://fastapi.tiangolo.com>
- [15] P. Ekman, “An argument for basic emotions,” *Cognition and Emotion*, vol. 6, no. 3–4, pp. 169–200, 1992.
- [16] S. Abdullah et al., “Stress detection using wearable sensors: A systematic review,” *IEEE J. Biomed. Health Inform.*, vol. 28, no. 2, pp. 1–12, 2024.
- [17] M. Fleming et al., “Digital therapeutics in mental health: Recent advances and future directions,” *Nature Digit. Med.*, vol. 7, pp. 1–10, 2024.
- [18] T. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space,” in *Proc. ICLR, Scottsdale, AZ, 2013*.
- [19] A. Herath, “AI chatbots for mental health support: Opportunities and limitations,” *Int. J. Comput. Appl.*, vol. 186, no. 4, pp. 23–29, 2025.
- [20] Z. Zhang, Y. Luo, and J. Wang, “Emotion recognition using deep learning for mental health monitoring,” *IEEE Access*, vol. 10, pp. 24317–24330, 2022.