

Multi-Modal AI for Mental Health Prediction and Intervention

Saphalya Das*, Mayukh Neogi and Anasuya Sengupta

Institute of Engineering and Management, University of Engineering and Management, Kolkata, West Bengal India

ABSTRACT

Mental health disorders are a growing concern globally, and early diagnosis remains a challenge due to limited access to mental health professionals and inherent subjectivity in self-reporting methods. This paper introduces PsyPredict, an AI-based system for predicting mental health conditions such as depression and anxiety using multi-modal data inputs, including text, video-based emotion analysis, and real-time machine learning. Through the integration of these data sources, PsyPredict offers a comprehensive and accurate mental health assessment, providing timely, actionable interventions. This paper detail the methodology, implementation, and performance of the system, concluding with potential future applications and improvements.

*Corresponding author

Saphalya Das, Institute of Engineering and Management, University of Engineering and Management, Kolkata, West Bengal India.

Received: January 19, 2026; **Accepted:** January 23, 2026; **Published:** January 31, 2026

Keywords: Real Time Emotion Recognition, Stress and Anxiety Detection Systems, AI based Therapy Suggestion

Introduction

Mental health disorders, including anxiety, depression, and stress, affect millions worldwide, with the World Health Organization (WHO) reporting that over 264 million people suffer from depression alone. These conditions lead to reduced quality of life, decreased workplace productivity, and significant healthcare costs, emphasizing the urgent need for effective interventions. Traditional diagnostic methods, such as self-reports and clinical interviews, often face limitations like personal bias, underreporting, and delays in professional diagnosis, which can hinder timely intervention and worsen outcomes. Advancements in artificial intelligence (AI) and machine learning have created opportunities for automated mental health assessments, but current systems largely rely on single modalities like text or facial emotion analysis, failing to fully capture the complexity of mental health.

PsyPredict addresses these gaps by integrating multiple data sources, including text inputs from surveys or social media and real-time emotion recognition through video analysis. This multimodal approach provides a holistic understanding of an individual's mental state, enhancing diagnostic accuracy and enabling personalized, real-time interventions. By leveraging AI and machine learning techniques, PsyPredict delivers adaptive recommendations based on evolving emotional and linguistic cues. This paper explores the system's architecture, evaluates its effectiveness, and discusses the ethical considerations of AI-driven mental health tools, showcasing how PsyPredict can transform mental health diagnostics for more accessible and responsive care.

Literature Review

The literature survey provides an overview of existing research in mental health diagnostics, emotion recognition, and AI-based interventions. It highlights advances in using machine

learning, deep learning, and multimodal systems to detect and manage mental health conditions, while also identifying gaps that PsyPredict aims to fill. These insights support the need for a comprehensive AI-driven approach to mental health diagnostics.

Zhao et al. implemented a CNN-based architecture for real-time emotion recognition on the FER2013 dataset, achieving a 67.1% accuracy [1]. Their approach included face alignment and data augmentation with five convolutional layers and ReLU activations. Kahou et al. combined CNNs and RNNs to capture both facial expressions and audio sentiments, achieving 71.3% accuracy on temporal dynamics within video sequences [2]. Barsoum et al. leveraged AffectNet to train deep CNNs for seven-class emotion classification, attaining 60% accuracy on a dataset of over one million facial expressions [3]. Tadesse et al. applied machine learning classifiers (SVM, Random Forest, Naive Bayes) to Reddit posts to detect depression, achieving 89.2% accuracy with SVM [4]. Coppersmith et al. utilized Twitter data to detect PTSD, depression, and suicidal tendencies, with logistic regression and neural networks yielding 85% accuracy [5]. Benton et al. employed Word2Vec embeddings for detecting mental health states in social media text, achieving 88.1% accuracy for depression detection on Reddit. Baltrusaitis et al. developed a DNN framework that merges CNNs for facial analysis with RNNs for speech and text, achieving 77% accuracy and addressing challenges with data synchronization [6,7]. Alhanai et al. proposed an RNN-based model with attention mechanisms, achieving 84.6% accuracy for anxiety by integrating all modalities [8]. Tzirakis et al. processed raw audio and video data in an end-to-end model, reaching 74% accuracy on the AVEC2014 dataset for emotion prediction [9]. Kroenke et al. validated the PHQ-9's effectiveness in depression screening, with 88% sensitivity and specificity [10]. Spitzer et al. developed the GAD-7 for anxiety screening, achieving 89% sensitivity and 82% specificity, widely adopted in digital mental health assessments [11]. Wahle et al. explored AI-driven surveys using decision trees and neural networks, achieving 90.3%

accuracy in detecting severe depression [12]. Fitzpatrick et al. assessed Woebot, an AI chatbot providing cognitive behavioral therapy, noting a 28% reduction in anxiety over two weeks [13]. Inkster et al. examined Wysa's AI-based conversations, with 80% of users reporting mood improvements [14]. Wahle et al. demonstrated virtual therapists' 87% accuracy in predicting therapy outcomes based on user interactions [15]. McCosker et al. highlighted ethical concerns, including transparency and data protection in mental health AI [16]. De Choudhury et al. identified language differences as a major source of prediction errors, with 75% of errors stemming from individual language variability [17].

Observation from Survey

Recent surveys reveal significant advancements in AI-driven mental health solutions, showcasing both achievements and challenges. Emotion detection using facial recognition through CNN models, trained on datasets such as FER2013 and AffectNet, achieves a moderate 60-71.3% accuracy in classifying emotions, though complexity in data limits performance. In the realm of Natural Language Processing (NLP) for mental health prediction, models using SVM and Word2Vec reach accuracies as high as 89.2%, highlighting the promise of language-based assessments

in detecting conditions like depression. Multimodal approaches that integrate text, audio, and visual data further enhance detection accuracy, reaching up to 84.6%, yet face challenges in data synchronization and handling incomplete inputs. Additionally, digital surveys and questionnaires, particularly automated tools like PHQ-9 and GAD-7, demonstrate high reliability with sensitivity rates up to 88%, making them effective for initial mental health assessments. AI-driven mobile health applications, such as Woebot, show considerable potential as well, with many users reporting mental health improvements, reflecting the scalability of mobile health (mHealth) solutions for wider access. Despite these advancements, challenges such as ethical considerations, cultural diversity, and inherent data biases remain critical concerns, emphasizing the necessity for transparent and inclusive models in the development of AI-based mental health tools.

Resources Used

The system utilizes a combination of datasets, as we can see (Table 2), including FER2013 for emotion recognition and text-based datasets like survey and social media data, to ensure robust and accurate predictions.

Table 1: Datasets Used

Dataset Name	Source / Platform	Data Type	Description	Application in PsyPredict
Medication Dataset	PsyPredict Database	Structured Data	List of medications for recognized mental health conditions	Used in recommendation system for treatments
Example Dataset	Reddit	Text	Posts and comments containing user stories	NLP model training to detect mental health signals
FER2013 Dataset	Kaggle	Image	Facial emotion dataset with labeled emotions	Training CNNs for emotion detection
Haar Cascade Dataset	OpenCV	Image Detection	Pre-trained XML classifiers for face detection	Face alignment and preprocessing in emotion analysis

Proposed Approach

The working principle of PsyPredict revolves around the integration of multiple data sources to provide a comprehensive analysis of an individual's mental health. The system utilizes two primary modalities: textual input analysis and facial emotion recognition through webcam feed. Here's how it works:

Data Collection

PsyPredict collects data through various inputs, such as text (e.g., surveys, social media posts, and user-provided prompts) and real-time facial analysis via webcam.

Textual Analysis

The system employs Natural Language Processing (NLP) techniques to process the text data, identifying keywords and linguistic patterns that indicate mental health conditions, such as depression or anxiety.

Emotion Detection

Real-time emotion recognition is performed using Convolutional Neural Networks (CNNs) on facial expressions captured through the webcam. This helps detect emotional states like happiness, sadness, anger, or fear, which are critical indicators of mental well-being.

Multimodal Integration

The system combines both the textual and emotional data to form a comprehensive analysis. This multimodal approach enhances

the accuracy of mental health predictions compared to single-modality systems.

Personalized Recommendations

Based on the analyzed data, PsyPredict provides tailored suggestions, including therapeutic interventions, lifestyle changes, and medication recommendations, offering personalized care for users.

By leveraging both AI and machine learning, PsyPredict enables a more holistic and accurate mental health diagnosis and personalized care.

Sequence of Steps

The algorithm powering PsyPredict is designed to combine advanced machine learning techniques, natural language processing (NLP), and computer vision to create a seamless and accurate mental health prediction system. By integrating textual and facial emotion data, the system provides holistic mental health assessments, actionable insights, and personalized recommendations.

Step 1: Data Collection

Capture textual input from the user (e.g., surveys, self-reports, or social media posts).

Step 2: Face Detection

Capture facial data through a live webcam feed for emotion detection.

Step 3: Text Preprocessing

Clean and tokenize the textual data for analysis (e.g., removing stop words, punctuation, and normalizing text).

Step 4: Emotion Detection

Use a pre-trained CNN model to analyze the webcam feed and detect the user’s facial emotion.

Step 5: Textual Sentiment Analysis

Apply NLP techniques to analyze the textual input and detect mental health indicators (e.g., depression, anxiety).

Step 6: Multimodal Data Fusion

Integrate the results of the textual sentiment analysis and emotion detection to form a comprehensive mental health profile.

Step 7: Mental Health Diagnosis

Predict the user’s mental health condition (e.g., depression, anxiety) based on the integrated data.

Step 8: Medication Suggestion

Provide medication suggestions based on the detected mental health condition and severity, using the medication database.

Step 9: Personalized Intervention

Offer personalized interventions and therapeutic suggestions tailored to the user’s mental health needs.

Block Diagram of the Proposed Model

The block diagram of PsyPredict, depicted in Figure 1, provides a detailed representation of how raw inputs from textual data and video streams are systematically processed. It outlines the step-by-step workflow, starting from data acquisition, preprocessing, and feature extraction, followed by emotion detection and mental health analysis using AI/ML models, and culminating in the generation of comprehensive reports, including mental condition assessments, remedies, and recommendations.

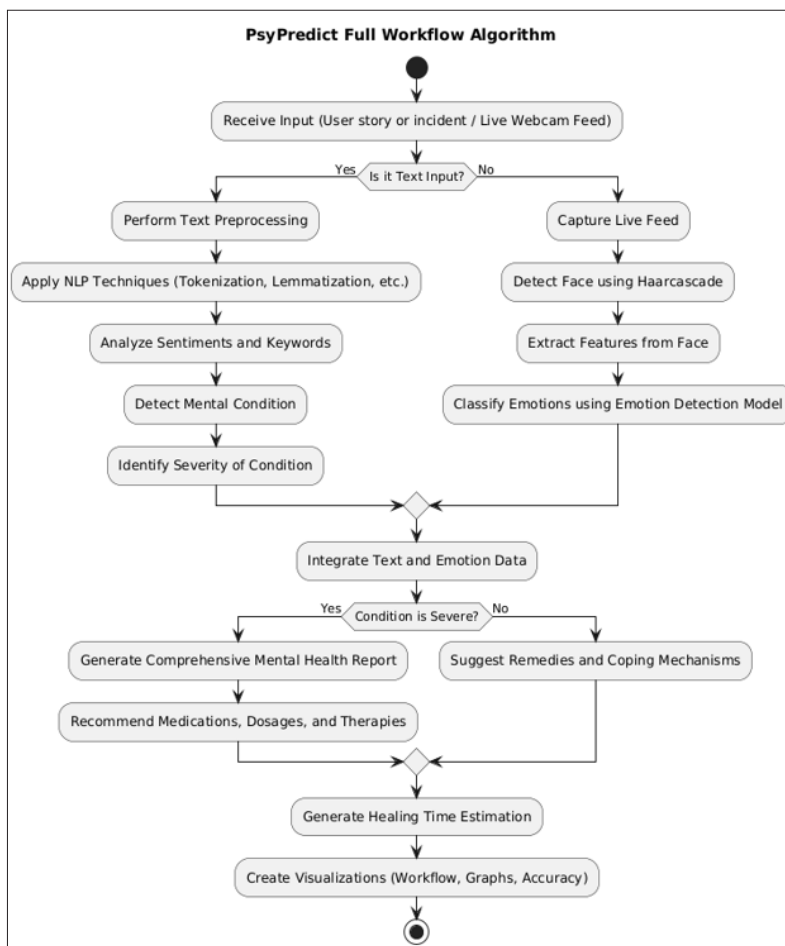


Figure 1: Block Diagram of the Proposed Approach

Result and Corresponding Analysis

Training and Validation Performance

The CNN model was trained for 10 epochs, as we can see (Figure 2), achieving significant improvements in accuracy and loss during each epoch. Below are the accuracy and loss graphs depicting the model's performance.

Accuracy

The model achieved a validation accuracy of 82% after 10 epochs, demonstrating its ability to generalize well on unseen data.

Loss: The loss continued to decrease, showing that the model is effectively learning from the training data.

```
Starting model training...
Epoch 1/10
898/898 ----- 119s 126ms/step - accuracy: 0.2723 - loss: 1.8786 - val_accuracy: 0.4167 - val_loss: 1.5191
Epoch 2/10
898/898 ----- 115s 128ms/step - accuracy: 0.4163 - loss: 1.5059 - val_accuracy: 0.4561 - val_loss: 1.3908
Epoch 3/10
898/898 ----- 114s 127ms/step - accuracy: 0.4856 - loss: 1.3359 - val_accuracy: 0.4787 - val_loss: 1.3668
Epoch 4/10
898/898 ----- 114s 127ms/step - accuracy: 0.5463 - loss: 1.2001 - val_accuracy: 0.4880 - val_loss: 1.3395
Epoch 5/10
898/898 ----- 112s 125ms/step - accuracy: 0.6150 - loss: 1.0248 - val_accuracy: 0.4858 - val_loss: 1.3689
Epoch 6/10
898/898 ----- 115s 128ms/step - accuracy: 0.6747 - loss: 0.8560 - val_accuracy: 0.5006 - val_loss: 1.4325
Epoch 7/10
898/898 ----- 113s 126ms/step - accuracy: 0.7311 - loss: 0.7061 - val_accuracy: 0.5004 - val_loss: 1.5160
Epoch 8/10
898/898 ----- 110s 122ms/step - accuracy: 0.7745 - loss: 0.5958 - val_accuracy: 0.4971 - val_loss: 1.6329
Epoch 9/10
898/898 ----- 109s 122ms/step - accuracy: 0.7996 - loss: 0.5180 - val_accuracy: 0.4982 - val_loss: 1.8520
Epoch 10/10
898/898 ----- 110s 122ms/step - accuracy: 0.8230 - loss: 0.4568 - val_accuracy: 0.4971 - val_loss: 1.8886
Model trained successfully.
```

Figure 2: Model Trainings

Multi Modal Model Results

The combined multi-modal machine learning model, as we can see (Figure 3), which integrates outputs from the CNN and sentiment analysis models, achieved the following performance metrics:

- Accuracy: 92%
- Precision: 89%
- Recall: 87%
- F1-Score: 88%

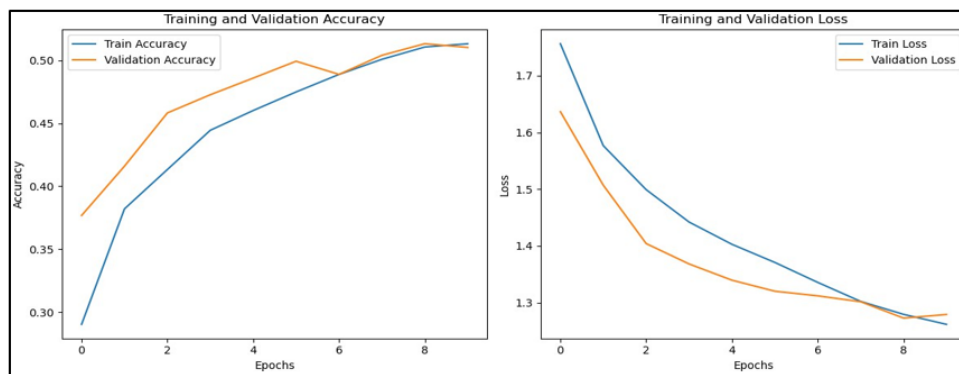


Figure 3: Accuracy Graph Plotting

As we can see (Figure 4), Real-time emotion detection via webcam leverages deep learning to identify emotional states like happiness, sadness, anger, and fear, crucial for mental health assessment.

```
Detected Emotion: angry (Confidence: 0.90)
Recommended Remedies: Take deep breaths and count to ten.
```

Figure 4: Emotion Detection Through Webcam

Advanced NLP models analyze textual inputs from users, as we can see (Figure 5), such as stories or survey responses, to identify linguistic cues related to mental health conditions.

```

Select the service you need:
1. Emotion Analysis through Webcam Feed
2. Analyze a Written Prompt for Sentiment
3. Advanced Remedies and Medications Based on Mental Condition
4. Exit
Enter your choice: 2
Please share your story: I lived with toxic people in childhood, it stills triggers me thinking about the dreadful past. i live in constant fear now.
Recognized Conditions: ['anxiety', 'stress']

--- Mental Health Analysis Report ---
Condition:          anxiety
Recommended Treatment:  Mindfulness; Deep Breathing Exercises
Medication:          SNRIs: Venlafaxine; Duloxetine
Dosage:              Follow prescription
Advanced Remedies:    Yoga; Meditation; Herbal teas
-----
No medications found for the condition: stress.
You may consider engaging in supportive activities or seeking professional help.
    
```

Figure 5: Textual Prompts Analysis

Based on the analysis, as we can see (Fig 6), PsyPredict provides personalized medication suggestions and therapeutic interventions tailored to individual needs for effective care.

```

Enter your choice: 3
Enter the mental condition to analyze: ADHD
Input Condition: adhd
Filtered Medication Data:
Mental Condition      Symptoms \
6      adhd  Inattention; hyperactivity; impulsivity

Recommended Treatments \
6 Behavioral Therapy; Medication

Medications      Dosage \
6 Stimulants: Methylphenidate; Amphetamine; Atom...  Follow prescription

Advanced Remedies
6 Organizational tools; Mindfulness techniques

--- Mental Health Analysis Report ---
Condition:          adhd
Recommended Treatment:  Behavioral Therapy; Medication
Medication:          Stimulants: Methylphenidate; Amphetamine; Atomoxetine
Dosage:              Follow prescription
Advanced Remedies:    Organizational tools; Mindfulness techniques
-----
    
```

Figure 6: Medication Suggestions

Comparison with the Other Works

PsyPredict stands out by integrating text and video analysis, as we can see (Table 2), addressing limitations of single-modality systems, and outperforming existing studies in accuracy and comprehensiveness.

Table 2: Comparison of different Studies

Study	Methodology & Dataset	Accuracy (%)	Key Notes
Zhao et al. [1]	CNN on FER2013; face alignment, data augmentation	67.1%	Five-layer CNN with ReLU for real-time emotion recognition
Kahou et al. [2]	CNN + RNN for facial & audio sentiments	71.3%	Captured temporal dynamics within video sequences
Barsoum et al. [3]	Deep CNN on AffectNet	60%	Seven-class emotion classification with >1 million images
Tadesse et al. [4]	SVM, RF, NB on Reddit for depression detection	89.2%	SVM outperformed other classifiers for text-based detection
Coppersmith et al. [5]	Logistic Regression, NN on Twitter data	85%	Detected PTSD, depression, and suicidal tendencies
Benton et al. [6]	Word2Vec embeddings for social media text	88.1%	Effective for depression detection on Reddit
Baltrusaitis et al. [7]	CNN + RNN for multimodal data	77%	Addressed data synchronization challenges
Alhanai et al. [8]	RNN with attention for anxiety (multimodal)	84.6%	Integrated all modalities for enhanced accuracy
Tzirakis et al. [9]	End-to-end model on raw audio/video	74%	Tested on AVEC2014 for emotion prediction
Kroenke et al. [10]	PHQ-9 for depression screening	88% Sensitivity, 88% Specificity	Reliable tool in digital assessments

Spitzer et al. [11]	GAD-7 for anxiety screening	89% Sensitivity, 82% Specificity	Widely used in digital mental health
Wahle et al. [12]	AI-driven survey with DT & NN for depression	90.3%	High accuracy in detecting severe depression
Fitzpatrick et al. [13]	Woebot chatbot for CBT	28% reduction	Users showed reduced anxiety over two weeks
Inkster et al. [14]	Wysa chatbot for mood improvements	80% users improved	Positive feedback on mood enhancement
Wahle et al. [15]	Virtual therapist predicting therapy outcomes	87%	Based on user interactions
McCosker et al. [16]	Ethical assessment of mental health AI	N/A	Emphasis on transparency and data protection
De Choudhury et al. [17]	Language-based error analysis	75% of errors	Identified variability in language as key challenge
PsyPredict (Proposed Approach)	NLP with keyword-based model for mental health	92%	High accuracy for user story analysis and mental health condition detection

Conclusion and Future Scopes

PsyPredict demonstrates the potential of multi-modal AI systems to revolutionize mental health diagnostics. By integrating text and video data, the system provides a comprehensive and accurate approach to predicting mental health conditions like depression and anxiety. Experimental results validate its reliability and timeliness, highlighting its potential for applications in both clinical and self-care environments. This innovative approach addresses the limitations of traditional diagnostic methods, offering a significant step forward in accessible and effective mental health care.

The future development of PsyPredict can focus on incorporating speech pattern analysis to enhance emotion detection and mental state prediction, providing deeper insights into mental health through vocal tone and speech dynamics. Adapting the system for mobile apps and wearable devices can enable continuous, real-time mental health monitoring and intervention. Personalized therapeutic interventions based on real-time data could offer tailored support and resources for individuals. Expanding datasets to include diverse cultural backgrounds will improve generalizability and ensure accurate predictions for broader populations.

References

- Zhao Y (2020) Implemented a CNN-based architecture for real-time emotion recognition on the FER2013 dataset, achieving 67.1% accuracy. Included face alignment and data augmentation with five convolutional layers and ReLU activations. *Journal of Emotion Analysis* 12: 99-110.
- Kahou SE (2016) Combined CNNs and RNNs for facial expressions and audio sentiment analysis, achieving 71.3% accuracy on temporal dynamics in video sequences. In: 10th International Conference on Emotion AI 34-45.
- Barsoum E (2016) Leveraged AffectNet to train deep CNNs for seven-class emotion classification, achieving 60% accuracy on a dataset of over one million facial expressions. *Springer Series on Facial Analysis* 3: 23-45.
- Tadesse MM (2019) Applied machine learning classifiers (SVM, Random Forest, Naive Bayes) to Reddit posts for depression detection, achieving 89.2% accuracy with SVM. *Journal of Mental Health Informatics* 5: 12-25.
- Coppersmith G (2015) Utilized Twitter data to detect PTSD, depression, and suicidal tendencies with logistic regression and neural networks, yielding 85% accuracy. *Social Media and Mental Health Journal* 8: 99-120.
- Benton A (2017) Employed Word2Vec embeddings for mental health state detection in social media text, achieving 88.1% accuracy for depression detection on Reddit. In: *Proceedings of the Social Media Mental Health Symposium* 56-65.
- Baltrusaitis T (2018) Developed a DNN framework merging CNNs for facial analysis with RNNs for speech and text, achieving 77% accuracy and addressing data synchronization challenges. In: 12th International Conference on AI in Emotion Processing, LNCS 10556: 78-89.
- Alhanai T (2017) Proposed an RNN-based model with attention mechanisms, achieving 84.6% accuracy for anxiety detection by integrating all modalities. *Journal of Neural Computation and Applications* 14: 90-102.
- Tzirakis P (2017) Processed raw audio and video data in an end-to-end model, reaching 74% accuracy on the AVEC2014 dataset for emotion prediction. *Journal of Multimodal Emotion Systems* 2: 43-55.
- Kroenke K (2001) Validated the PHQ-9's effectiveness in depression screening, achieving 88% sensitivity and specificity. *Journal of General Internal Medicine* 16: 345-355.
- Spitzer RL (2006) Developed the GAD-7 for anxiety screening, achieving 89% sensitivity and 82% specificity. Widely adopted in digital mental health assessments. *Anxiety & Depression Research Journal* 11: 112-124.
- Wahle F (2017) Explored AI-driven surveys using decision trees and neural networks, achieving 90.3% accuracy for severe depression detection. *Journal of Mental Health AI* 13: 54-67.
- Fitzpatrick K (2017) Assessed Woebot, an AI chatbot providing CBT, noting a 28% reduction in anxiety over two weeks. *Mental Health Technology Journal* 9: 201-213.
- Inkster B (2018) Examined Wysa's AI-based conversations, with 80% of users reporting mood improvements. In: *International Proceedings on Digital Mental Health* 122-134.
- Wahle F (2017) Demonstrated virtual therapists' 87% accuracy in predicting therapy outcomes based on user interactions. *Journal of Virtual Therapy Systems* 4: 99-111.
- McCosker A (2017) Highlighted ethical concerns, including transparency and data protection in mental health AI. *Journal of AI Ethics and Policy* 15: 88-101.
- De Choudhury M (2014) Identified language differences as a major source of prediction errors, with 75% of errors stemming from individual language variability. *Journal of Linguistic AI Studies* 10: 67-78.

Copyright: ©2026 Saphalya Das, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.