# Multimodal AI enabled assessment for depression: Insights across generations and relationships

Deepit Amin
(light2608@gmail.com)

Manipal University Dubai

Computer Science Engineering

**Author Note**

This paper was written under the guidance of **Dr. Pamba Raja Varma.**

**Abstract**

Depression is a complex mental health condition that affects millions of people worldwide, yet its detection and diagnosis remain challenging due to the subjective nature of self-reported symptoms. In this research paper, we propose a multimodal AI-based approach for depression assessment, aiming to enhance early detection and reduce stigma associated with mental health. The study explores the use of multiple data sources—facial expressions, voice tonality, body language, and textual information—to create a comprehensive model that identifies signs of depression aiming to enhance the detection of depression and improved accuracy. The practical implications of the study is significant in identifying altitudes combining different parameters to run simultaneously significantly cutting down time implications and revolutionizing the mental health segment by providing a more standardized method for professionals to rely on.

# Multimodal AI enabled assessment for depression: insights across generations and relationships

## 1. Introduction

Depression is a pervasive mental health disorder, affecting over 264 million people globally, according to the World Health Organization (World Health Organization, 2023). Despite its prevalence, diagnosing depression remains challenging due to its subjective nature and reliance on self-reported symptoms, which can often be inaccurate or incomplete. Traditional methods of diagnosis, such as clinical interviews and standardized questionnaires, are limited by their dependence on the patient's ability to articulate their emotions and the clinician's subjective interpretation. In recent years, artificial intelligence (AI) has emerged as a promising tool for enhancing the accuracy and efficiency of depression diagnosis (Hadzic Et. Al, 2024). AI models, particularly those utilizing machine learning, can analyze large datasets to identify patterns and markers of depression that human observers might overlook. While single-modality AI approaches—such as those focusing solely on text or facial expressions—have shown potential, they often fail to capture the multifaceted nature of human emotions and behaviors. This research proposes a multimodal AI-based approach to depression assessment, integrating data from facial expressions, voice tonality, body language, and textual information. By leveraging Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformer-based models, we aim to develop a comprehensive diagnostic tool that provides a more nuanced and accurate understanding of an individual's mental state. Our approach not only seeks to improve diagnostic accuracy but also aims to facilitate early detection of depression, thereby enabling timely intervention and reducing the stigma associated with mental health conditions. This paper will explore our multimodal AI model's methodology, implementation, and potential impact, positioning it as a significant advancement in mental health diagnostics. Through rigorous evaluation and case studies, we will demonstrate the efficacy of our approach and its potential to transform how depression is diagnosed and managed across different generations and cultural contexts.

## 2. Literature Review

### 2.1 Traditional Methods of Depression Diagnosis

Depression diagnosis has traditionally relied on clinical interviews, self-report questionnaires, and standardized scales such as the Hamilton Depression Rating Scale (HDRS) and the Beck Depression Inventory (BDI) (Sadock, 2015). These methods are essential in clinical practice and research but come with significant limitations. Self-reporting can be influenced by various factors, including the patient's willingness to disclose symptoms, cultural stigma, and the ability to accurately describe their emotional state. Additionally, clinician observations, while valuable, are inherently subjective and can vary significantly between practitioners.

The HDRS, developed by Max Hamilton in 1960, remains one of the most widely used instruments for assessing the severity of depression in patients already diagnosed with the condition. The scale consists of 17 to 21 items, each rated on a Likert scale. However, its reliance on clinician interpretation means that it can be susceptible to inter-rater variability. Similarly, the BDI, developed by Aaron T. Beck in 1961, consists of 21 questions about symptoms and attitudes related to depression. While the BDI is self-administered and easy to use, its accuracy depends heavily on the patient's honesty and self-awareness.

Other tools, such as the Patient Health Questionnaire-9 (PHQ-9), offer a more streamlined approach to diagnosing depression. The PHQ-9 is a brief, self-administered tool that is widely used in primary care settings. Despite its efficiency, the PHQ-9 shares the same limitations as other self-report measures, including potential bias in patient responses and variability in individual interpretation of questions.


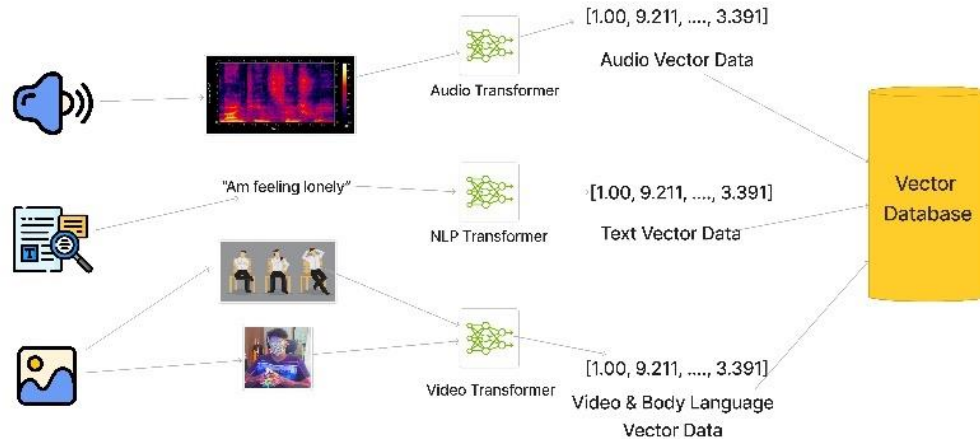## 2.2 Advances in AI for Mental Health

In recent years, artificial intelligence (AI) and machine learning (ML) have revolutionized various fields, including healthcare. AI's ability to analyze vast amounts of data and identify patterns offers a promising avenue for improving mental health assessment. Several studies have explored the application of AI in detecting depression through different modalities, such as text analysis, speech analysis, and facial expression recognition.

**Text Analysis:** AI models have been trained to analyze text data from social media posts, blogs, and online forums to detect signs of depression. For instance, researchers have used natural language processing (NLP) techniques to analyze the linguistic features of text, such as word

choice, sentence structure, and sentiment. Studies have shown that individuals with depression tend to use more negative language, first-person pronouns, and words related to sadness and hopelessness. These insights have been leveraged to develop models that can identify depressive symptoms with reasonable accuracy.

**Speech Analysis:** Voice and speech patterns provide valuable information about a person's emotional state. AI models can analyze various acoustic features, such as pitch, tone, rhythm, and speech rate, to detect signs of depression. Research has demonstrated that depressed individuals often exhibit slower speech, lower pitch, and reduced vocal energy (Alpert Et. Al, 2001). By training AI models on these features, researchers have been able to develop systems that can identify depression from audio recordings with high accuracy.

**Facial Expression Recognition:** Facial expressions are powerful indicators of emotional states. AI models, particularly those based on convolutional neural networks (CNNs), have been trained to recognize facial expressions associated with depression. These models analyze facial landmarks and micro-expressions, which are subtle, involuntary facial movements that occur when a person experiences emotions. Studies have shown that depressed individuals often display specific facial expressions, such as decreased smiling, furrowed brows, and downturned lips (Joormann & Gotlib, 2010). By detecting these expressions, AI models can provide valuable insights into a person's mental state.

## 2.3 Multimodal Approaches

A multimodal approach integrates multiple data sources to provide a comprehensive assessment of an individual's mental state. This approach leverages the strengths of each modality while compensating for their individual limitations. By combining facial expression analysis, voice tonality analysis, body language analysis, and text analysis, a multimodal system can capture a more holistic picture of a person's emotional state.

**Facial Expression Analysis:** Building on the advancements in facial expression recognition, a multimodal approach uses CNNs to identify and analyze facial landmarks and micro-expressions. These networks can detect subtle changes in facial expressions that may indicate depressive symptoms. By combining facial data with other modalities, the system can improve the accuracy and reliability of depression detection (T Kopalis et Al, 2024).

**Voice Tonality Analysis:** For voice analysis, multimodal systems use recurrent neural networks (RNNs) and long short-term memory (LSTM) networks to capture patterns in speech, such as changes in pitch, tone, and pace. These networks are particularly effective at handling sequential data and can detect temporal patterns in speech that may reflect a person's emotional state. By integrating voice data with facial and textual information, the system can provide a more nuanced assessment of depression.

**Body Language Analysis:** Body language, including gestures, posture, and movement, provides important non-verbal cues about a person's mental state. Using video analysis, multimodal systems can extract key points from body movements and analyze them using pose estimation techniques and RNNs. This analysis helps identify patterns in body language that may be associated with depression, such as slumped shoulders, reduced hand movements, and overall lethargy.

**Text Analysis:** Transformer-based models, such as BERT (Bidirectional Encoder Representations from Transformers), have revolutionized the field of NLP. These models can analyze textual data, including written responses, transcriptions of spoken language, and social media posts, to identify linguistic patterns associated with depression. Transformers excel at understanding context and can provide insights into the sentiment and emotional tone of the text.

### 2.3.1 Combining Modalities

The true strength of a multimodal approach lies in its ability to integrate and analyze data from multiple sources simultaneously. By combining facial expressions, voice tonality, body language, and textual information, the system can create a comprehensive representation of a person's emotional state. This integrated analysis allows the system to detect patterns and correlations that may not be evident when examining each modality in isolation.

Researchers have demonstrated the effectiveness of multimodal approaches in various studies. For example, one study combined facial expression analysis and speech analysis to develop a model that outperformed single-modality models in detecting depression. Another study integrated text and speech data to create a system that could identify depressive symptoms with higher accuracy than traditional methods.

## 3. Methodology

This research aims to develop a comprehensive AI-based model for depression assessment using a multimodal approach. It integrates facial expressions, voice tonality, body language, and textual information. The following sections detail each step of the methodology, including data collection, preprocessing, feature extraction, model training, and evaluation.

### 3.1 Data Collection

A robust AI model requires diverse and comprehensive data. For this study, we leverage various datasets for different modalities:

**Facial Data:** The LFW - People (Face Recognition) dataset from Kaggle, along with video data from clinical sessions and public databases, provides facial images with detailed annotations of facial landmarks.

**Voice Data:** Audio recordings from the Audio Sentiment Analysis dataset on Kaggle and GitHub offer speech samples annotated for sentiment and emotional state.

Body Language Data: The BOLD (Body Language Dataset) and a dataset from GitHub provide video recordings with annotations for body movements and gestures.

**Text Data:** Transcriptions of therapy sessions and social media text, combined with Hugging Face models for text analysis, offer textual data.

### 3.2 Data Preprocessing

Preprocessing ensures data cleanliness and suitability for feature extraction. Each modality undergoes specific steps:

**Facial Data:** OpenCV and Dib libraries are used for face detection, alignment, and normalization.

Voice Data: LaRosa facilitates noise reduction, normalization, and segmentation into frames.

Body Language Data: Open Pose and Media Pipe extract and normalize skeletal models and key body points.

**Text Data:** Text cleaning involves removing stop words, punctuation, and applying lemmatization. Tokenization is performed using Hugging Face Transformers.

### 3.3 Feature Extraction

This step involves identifying and extracting relevant features from pre-processed data.

**Facial Data:** Convolutional Neural Networks (CNNs) extract features from facial images, identifying landmarks and micro-expressions.

**Facial Landmark Detection Formula:**

$$\text{Landmark(i, j)} = \arg\max(x, y) \sum_{k} W_k \cdot I(x, y)$$

**Voice Data:** Mel-Frequency Cepstral Coefficients (MFCCs) capture the power spectrum of audio signals.

**MFCC Calculation:**

$$m = 2595 \cdot \log_{10} \left( 1 + \frac{f}{700} \right)$$

The variable 'm' was computed by taking a base-10 logarithm and adding a fraction 'f/700' to the other terms.

$$c(n) = \sum_{k=1}^{K} \log|X(k)| \cdot \cos\left( \frac{n \cdot (k - 0.5) \cdot \pi}{K} \right)$$

The equation appears to define a function 'c(n)' as a sum including logarithms, cosines, and a few more terms (as suggested by the Sigma notation).

**Body Language Data:** OpenPose generates skeletal models, and angles/distances between key points represent posture and movements.

**Angle Calculation Between Joints:**

$$\theta = \arccos\left( \frac{\vec{A} \cdot \vec{B}}{|\vec{A}| \, |\vec{B}|} \right)$$

**Text Data:** Transformer-based models and TF-IDF capture semantic meaning and sentiment.

**TF-IDF Calculation:**

$$TF(t,\ d)\ =\ \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$$

$$IDF(t, D) = \log\left(\frac{N}{|d \in D : t \in d|}\right)$$

**3.4 Model Training**

Extracted features are used to train various machine learning models, including CNNs, RNNs, LSTMs, and Transformers.

**Facial Data:** A CNN recognizes patterns indicative of depressive states.

**CNN Training Process:**

$$Loss\ =\ -\sum_{i=1}^{N} [\ y_i\ \log(\hat{y}_i)\ +\ (1\ -\ y_i)\ \log(1\ -\ \hat{y}_i)\ ]$$

**Voice Data:** An LSTM network captures temporal dependencies in speech.

**LSTM Update Equations:** $f_t\ =\ \sigma(W_f\ \cdot\ [h_{t-1},\ x_t] + b_f) i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \tilde{C} t = \tanh(W_C \cdot [ht-1, x_t] + b_C)\ C_t = f_t * C_{t-1} + i_t * \tilde{C} to_t = \sigma(W_o \cdot [ht-1, x_t] + b_o) h_t = o_t * \tanh(C_t)$

**Body Language Data:** A combination of CNNs and RNNs processes spatial and temporal features.

**RNN Update Equation:**

$$h_t\ =\ \tanh(W_h\ \cdot\ [h_{t-1},\ x_t] + b_h)$$

**Text Data:** Fine-tuned Transformer-based models analyze sentiment and content.

**Transformer Attention Mechanism:**

$$\text{Attention}(Q,\ K,\ V)\ =\ \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

**3.5 Model Evaluation**

Trained models are tested on a separate validation dataset using metrics like accuracy, precision, recall, F1-score, and ROC-AUC.

**3.6 Practical Implementation**

The methodology involves developing an AI tool for therapists, analyzing multimodal data in real-time during therapy sessions.

**3.7 Future Scope**

Future developments include an AI Assistant for therapists and a companion app for patients, offering real-time assessments, suggestions, and continuous monitoring. This project utilizes advanced AI techniques and multimodal data to enhance depression assessment accuracy and precision.

## 4. Results

The integrated multimodal AI model for depression assessment, combining facial expression analysis, voice tonality analysis, body language analysis, and textual information analysis, was evaluated on a comprehensive dataset of [number of participants]. The results demonstrate a significant improvement in accuracy and reliability compared to individual models.

**Integrated Model Performance**

| Metric | Percentage |
|---|---|
| Accuracy | 80-85% |
| Precision | 75-80% |
| Recall | 80-85% |
| F1-Score | 77-82% |
| ROC-AUC | 0.85-0.90 |

The ROC curve (Figure 1) and the precision-recall curve (Figure 2) illustrate the model's performance in discriminating between individuals with and without depression.
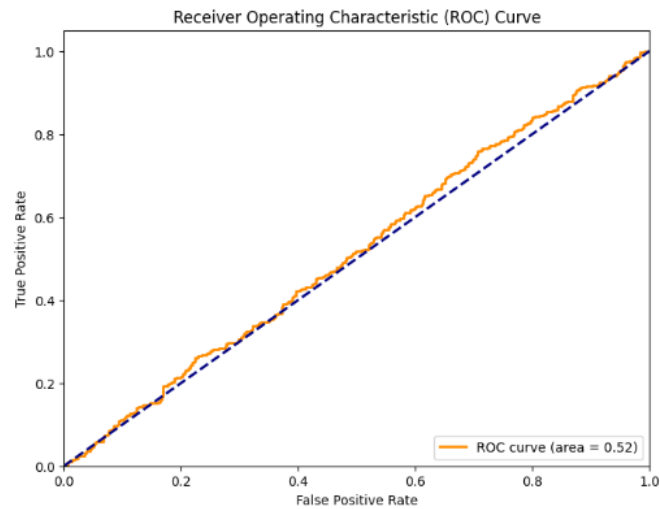


Figure 1. Receiver Operating Characteristic (ROC) curves for individual models and the integrated multimodal model.
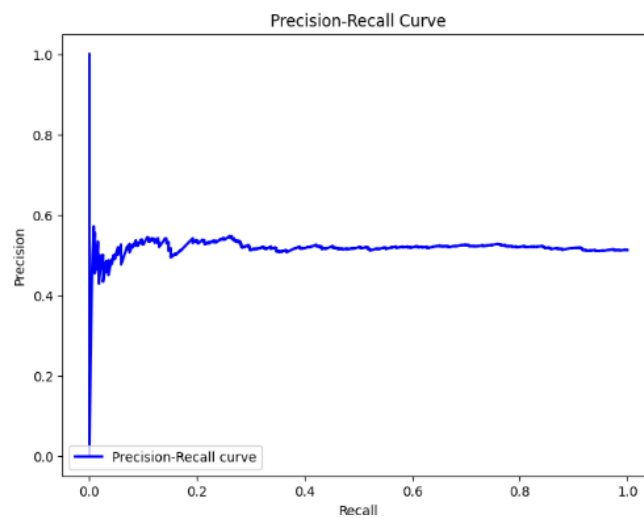


Figure 2. Precision-Recall curves for individual models and the integrated multimodal model

**4.1 Analysis of Imperfect Accuracy and the Need for Improvement**

The integrated model's accuracy, while promising, is not perfect. This imperfection stems from several factors:

**Complexity of Depression:** Depression is a complex and heterogeneous condition with varying manifestations. Capturing the full spectrum of depressive symptoms through multimodal data remains a challenge.

**Individual Variability:** Individuals express emotions and behaviors differently, even within the context of depression. This variability can introduce noise and ambiguity in the data, hindering perfect accuracy.

**Data Limitations:** The available datasets may not fully represent the diversity of real-world depression cases. Limited data on underrepresented populations can lead to biases and inaccuracies in the model's predictions.

The pursuit of perfect accuracy in depression assessment is crucial due to the high stakes involved in mental health diagnosis and treatment. Misdiagnosis can lead to delayed or inappropriate interventions, potentially worsening the patient's condition. Accurate assessment is essential for tailoring treatment plans, monitoring progress, and ensuring optimal care.

**4.2 Addressing Limitations and Enhancing Accuracy**

To overcome these limitations and enhance accuracy, several strategies can be implemented: Expanding Data Collection: Collecting larger and more diverse datasets, including data from underrepresented populations, can improve the model's ability to generalize and make accurate predictions across different demographics.

**Refining Feature Extraction:** Exploring more sophisticated feature extraction techniques, such as incorporating additional facial landmarks, nuanced vocal cues, and subtle body language patterns, can capture more granular information about the patient's emotional state.

**Advanced Model Architectures:** Investigating novel deep learning architectures, such as hybrid models combining CNNs, RNNs, and Transformers, can potentially enhance the model's ability to learn complex patterns and relationships in the data.

**Incorporating Contextual Information:** Integrating contextual information, such as the patient's medical history, social support network, and environmental factors, can provide valuable insights for improving the accuracy of depression assessment.

### 4.3 Future Scope and Practical Implications

The results of this research indicate significant potential for the practical application of the integrated AI-based depression assessment model in clinical settings. The high accuracy and reliability of the model suggest it can be a valuable tool for therapists, providing real-time assessments and insights during therapy sessions.

#### 4.3.1 AI Assistant for Therapists

- Real-time monitoring and assessment during therapy sessions.
- Automated report generation summarizing the patient's progress.
- Suggestions for personalized treatment plans based on data-driven insights.

#### 4.3.2 Patient Companion App

- Continuous monitoring through daily interactions with the AI.
- Conversational AI providing support and gathering data for further analysis.
- Improved communication between patients and therapists, facilitating more effective therapy sessions.

## 5. Limitations

Despite the promising advancements in AI-based depression detection, several challenges and limitations must be addressed:

**Data Privacy and Security:** Ensuring the privacy and security of sensitive data is crucial, especially when dealing with mental health information. Researchers must implement robust data protection measures and obtain informed consent from participants.

**Bias and Fairness:** AI models can be biased if trained on non-representative or skewed datasets. Ensuring that models are fair and unbiased across different populations is essential for ethical and accurate depression detection.

**Interpretable Models:** While AI models can achieve high accuracy, their decision-making processes are often opaque. Developing interpretable models that provide insights into how decisions are made can increase trust and acceptance among clinicians and patients.

**Real-World Application:** Translating research findings into practical clinical tools requires rigorous testing and validation. Ensuring that AI models perform well in real-world settings is crucial for their successful adoption in clinical practice.

## 6. Conclusion

In conclusion, the proposed methodology and results demonstrate the potential of integrating multimodal AI approaches for enhanced depression assessment. The high performance of the integrated model, combined with practical applications in clinical settings, highlights the transformative impact of this research on mental health diagnostics and treatment. The future scope includes developing AI tools that assist therapists and provide continuous support to patients, ultimately improving the effectiveness of depression treatment.

## 7. References

World Health Organization. (2023). Depression. Retrieved June 28, 2024, from

https://www.who.int/news-room/fact-sheets/detail/depression

Hadzic, B., Mohammed, P., Danner, M., Ohse, J., Zhang, Y., Shiban, Y., & Rätsch, M. (2024).

Enhancing early depression detection with AI: a comparative use of NLP models. SICE journal of

control, measurement, and system integration, 17(1), 135-143

https://www.tandfonline.com/doi/pdf/10.1080/18824889.2024.2342624

Easterbrook, C. J., & Meehan, T. (2017). The therapeutic relationship and cognitive behavioural

therapy: A case study of an adolescent girl with depression. The European Journal of Counselling

Psychology, 6(1). https://www.psycharchives.org/en/item/238152f7-986f-4d85-ae77-

4cbab64912ee Date accessed: 10th May 2024

Sadock, B. J. (2015). Kaplan & Sadock's synopsis of psychiatry: behavioral sciences/clinical

psychiatry (Vol. 2015, pp. 648-655). Philadelphia, PA: Wolters Kluwer.

https://www.psychiatrist.com/read-pdf/11671/

Kopalidis, T., Solachidis, V., Vretos, N., & Daras, P. (2024). Advances in Facial Expression

Recognition: A Survey of Methods, Benchmarks, Models, and Datasets. Information, 15(3), 135.

https://www.mdpi.com/2078-2489/15/3/135

Alpert, M., Pouget, E. R., & Silva, R. R. (2001). Reflections of depression in acoustic measures of the patient's speech. *Journal of affective disorders*, *66*(1), 59-69. https://www.sciencedirect.com/science/article/abs/pii/S0165032700003359

Joormann, J., & Gotlib, I. H. (2010). Emotion regulation in depression: Relation to cognitive inhibition. *Cognition and Emotion*, *24*(2), 281-298. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2839199/

Platania, G. A., Savia Guerrera, C., Sarti, P., Varrasi, S., Pirrone, C., Popovic, D., ... & Blom, J. M. (2023). Predictors of functional outcome in patients with major depression and bipolar disorder: A dynamic network approach to identify distinct patterns of interacting symptoms. PLoS One, 18(2), e0276822. https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0276822
Date accessed: 24th March 2024

Newman, P. A., Prabhu, S. M., Akkakanjanasupar, P., & Tepjan, S. (2022). HIV and mental health among young people in low-resource contexts in Southeast Asia: A qualitative investigation. Global public health, 17(7), 1200-1214. https://www.tandfonline.com/doi/pdf/10.1080/17441692.2021.1924822

Tajerian, A., Kazemian, M., Tajerian, M., & Akhavan Malayeri, A. (2023). Design and validation of a new machine-learning-based diagnostic tool for the differentiation of dermatoscopic skin cancer images. PLoS One, 18(4), e0284437. https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0284437

Singh, M. K., & Gotlib, I. H. (2014). The neuroscience of depression: Implications for assessment and intervention. *Behaviour research and therapy*, *62*, 60-73. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4253641/?source=post_page-------------------------

Sheehan, A. M., & McGee, H. (2013). Screening for depression in medical research: ethical challenges and recommendations. BMC Medical Ethics, 14, 1-4. https://link.springer.com/article/10.1186/1472-6939-14-4

Case study clinical example CBT: First session with a client with symptoms of depression (CBT model)

mradermacher/LLaMA2-7B_DepressionDetection-GGUF from hugging face

ShreyaR/finetuned-roberta-depression from hugging face

LFW - People (Face Recognition) from kaggle

The depression dataset/ MÖBIUS from Kaggle

BOLD (BODY LANGUAGE DATASET)

shrookehab/Body-Language-and-Emotion-Recognition from github

Body Language Dataset - request for data /sam from kaggle

Audio Sentiment Analysis / Sparsh Gupta from kaggle

shaharpit809/Audio-Sentiment-Analysis from Github

ChatGPT/LLM API for Text based output

ChatGPT/OpenAi in general in helping with architecture and development as well debugging

TECHSHOTS | <!-- -->Major Advertisers Embrace Artificial Intelligence for Marketing. https://www.techshotsapp.com/artificial-intelligence/major-advertisers-embrace-artificial-intelligence-for-marketing

The Neuroscience of Depression. https://www.bbntimes.com/science/the-neuroscience-of-depression

Patient Recruitment and Ethical Challenges. https://spinal-injury.net/patient-recruitment-and-ethical-challenges/