

# Multimodal Depression Detection and Knowledge Infused Mental Health Therapeutic and Empathetic Response Generation

Anonymous ACL submission

## Abstract

The detection of depression through non-verbal cues has gained significant attention. Previous research predominantly centered on identifying depression within the confines of controlled laboratory environments, often with the supervision of psychologists or counselors. Unfortunately, datasets generated in such controlled settings may struggle to account for individual's behaviors in real-life situations. In response to this limitation, we present the Extended D-vlog dataset, encompassing a collection of 1,261 YouTube vlogs. Additionally, the emergence of large language models (LLMs) like ChatGPT has sparked interest in their potential if they can act like a Mental Health Professionals. yet, the readiness of these LLM models to use it in real life setting is still a concern as they can give wrong responses which can be harmful for the users. We Proposed a Virtual Agent that can act as a first point of contact for mental health patient and provide a empathetic as well as the Therapeutic response alike human therapists. The system is divided into two part 1. Detection of a Disorder 2. Providing Empathetic and Therapeutic response to the user. The use of TVLT model on our Multimodal Extended D-vlog Dataset produced the most promising results, achieving a remarkable **F1-score of 67.8%**

## 1 Introduction

Depression, is a prevalent and significant medical condition. It has a harmful impact on one's emotional state, thought processes, and behavior. It manifests as persistent feelings of sadness and a diminished interest in previously enjoyed activities. This condition can give rise to various emotional and physical challenges, affecting one's ability to perform effectively both at work and in personal life. Depression symptoms range from mild to severe and can include persistent sadness, loss

of interest in once-enjoyable activities, appetite changes, sleep disturbances, fatigue, psychomotor changes, feelings of worthlessness, cognitive challenges, and, in severe cases, suicidal Thoughts. Symptom severity varies, requiring careful clinical evaluation for diagnosis and treatment (Cleveland Clinic).

According to the Statistics of The World Health Organisation (WHO) (World Health Organization) 3.8% of the world's population experience depression, including 5% of adults less than 60 years of age (4% of men and 6% of women) and 5.7% of Adults above 60 years of age. Approximately 280 million people have depression in which depression is 50% more common in women than men. Depression is 10% more in pregnant women and women who have just given birth (Evans-Lacko et al., 2018). If the depression is left untreated can lead to several serious outcomes such as suicide (Ghosh et al., 2022).<sup>1</sup>

The process of clinically diagnosing depression relies on interviews that incorporate **PHQ-8** (Kroenke et al., 2009) <sup>2</sup> Questionnaires which include questions such as "Do you experience feelings of sadness, depression, or hopelessness?" and inquire about the duration of these feelings, ranging from 0-1 day to 1-3 days and so forth", and using **BDI-II Questionnaires** (Smarr and Keefer, 2011) <sup>3</sup> having questions such as *How often do you have Guilty feelings? with four options as 1) dont feel particularly guilty 2) I feel guilty over many things i have or should have done 3) I feel quilty most of the times 4) I feel guilty all the times.* However, these questionnaires has a limitations, as patients may exhibit hesitance expressing their gen-

<sup>1</sup><https://www.who.int/news-room/fact-sheets/detail/depression>

<sup>2</sup><https://www.childrenshospital.org/sites/default/files/2022-03/PHQ-8.pdf>

<sup>3</sup><https://naviauxlab.ucsd.edu/wp-content/uploads/2020/09/BDI21.pdf>

077 uine emotions during these interviews, potentially  
078 resulting in an inaccurate assessment of their de-  
079 pression (Yoon et al., 2022). In contrast, the evolu-  
080 tion of social media platforms like Twitter, Reddit,  
081 and YouTube has provided users with a medium to  
082 openly share their thoughts, perspectives, and cur-  
083 rent life situations. These online platforms contain  
084 valuable and diverse emotional information.

085 Majority of methods utilized for Depression  
086 detection predominantly focus on analyzing text-  
087 ual information gathered from social media to  
088 infer users’ emotional states. However, (Yang  
089 et al., 2017a), (Yang et al., 2017b), (Gong and  
090 Poellabauer, 2017), (Ray et al., 2019) demon-  
091 strates that adding modalities, such as utilization of videos  
092 to extract facial expressions, body gestures, as well  
093 as incorporating audio analysis to detect fluctua-  
094 tions in speech patterns, has the potential to en-  
095 hance the precision of depression identification.

096 We utilize the power of **Vision-Language TVLT**  
097 (Tang et al., 2022) **Transformer model**, which  
098 has demonstrated its power by achieving state-  
099 of-the-art performance on various tasks such as  
100 video-captioning, Multimodal sentiment analysis  
101 and Multimodal emotion recognition. TVLT (Tang  
102 et al., 2022) serves as a Encoder-Decoder Modal  
103 which take raw video, raw audio and text as input  
104 and produce a comprehensive Multimodal repre-  
105 sentation that can be leveraged for downstream  
106 task related to detection and classification. Incorpor-  
107 ating additional wav2vec2 (Baevski et al., 2020)  
108 features with spectrograms for audio along with  
109 video and text, we achieved an impressive accuracy  
110 of 67.8% on the Extended D-vlog dataset (Yoon  
111 et al., 2022).

112 In the realm of mental health support, the no-  
113 tion of therapy chatbots has intrigued both re-  
114 searchers and the public since the introduction of  
115 Eliza (Shum et al., 2018) in the 1960s. Recent ad-  
116 vancements in large language models (LLMs) like  
117 ChatGPT have further fueled this interest. How-  
118 ever, concerns have been raised by mental health  
119 experts regarding the use of LLMs for therapy as  
120 the therapy provided may not be accurate. De-  
121 spite this, many researchers have begun exploring  
122 LLMs as a means of providing mental health sup-  
123 port (Sharma et al., 2023).

124 Yet, our understanding of how LLMs behave in  
125 response to clients seeking mental health support  
126 remains limited. It is unclear under what circum-  
127 stances LLMs prioritize certain behaviors, such as  
128 reflecting on client emotions or problem-solving,

129 and to what extent (Chung et al., 2023), (Ma et al.,  
130 2023). Given the critical nature of mental health  
131 support, it is essential to comprehend LLM be-  
132 havior, as undesirable actions could have severe  
133 consequences for vulnerable clients. Additionally,  
134 identifying desirable and undesirable behaviors can  
135 inform the adoption and improvement of LLMs in  
136 mental health support.

### 137 **Our Contribution are:**

- 138 • Extended D-vlog dataset (**Original no. of**  
139 **videos: 961, Total videos** (after adding 300  
140 videos to the Original dataset): **1261**) which  
141 contains videos of various type such as Ma-  
142 jor depressive disorder, postmortem disorder,  
143 anxiety and videos from different age group  
144 and gender which was lacking in the original  
145 D-vlog dataset.
- 146 • TVLT (Tang et al., 2022) model for depression  
147 detection, which outperforms baseline models  
148 by **4.3%** and establishes a new benchmark, on  
149 Extended D-vlog dataset.
- 150 • Replacing spectrogram with combination of  
151 spectrogram and wav2vec2 (Baevski et al.,  
152 2020) features which captures the vocal cues  
153 associated with depression more effectively  
154 than spectrogram, which further increases the  
155 accuracy by **2.2%** resulting in the final F1-  
156 score of **67.8 %**.
- 157 • To the best of our knowledge, this work is first  
158 to propose Virtual Agent that deliver therapeutic  
159 responses and Empathic responses to the  
160 users using LLM with Domain Knowledge  
161 as an External Knowledge base on Mental  
162 Health.

## 163 **2 Related Work**

164 With the increase in mental health conditions these  
165 days, there is an increase in the popularity of the  
166 detection of depression. However, very little work  
167 is done in creating the dataset to detect depres-  
168 sion. Due to privacy concerns, most datasets are  
169 only used for their research and are not publicly  
170 available. Among the relatively scarce publicly  
171 available datasets suitable for analysis, the DAIC-  
172 WOZ (Gratch et al., 2014) is one of most famous  
173 and used dataset. This dataset encompasses clinical  
174 interviews in various formats, including verbal  
175 (text) and non-verbal (audio and video). Notably,  
176 The DAIC-WOZ (Gratch et al., 2014) dataset re-  
177 lies on self-reporting through the PHQ-8 question-  
178 naire. Another well-known dataset is the Pittsburgh

dataset (Keenan et al., 2010), which comprises clinical interviews primarily in audio and video formats. However, this DAIC-WOZ (Gratch et al., 2014) dataset is relatively small, containing only 189 samples, making it a valuable resource for research purposes. The AViD-Corpus (Audio-Video Depressive Language Corpus) is another prominent dataset, with subsets used in AVEC 2013 (Valstar et al., 2013) and AVEC 2014 (Valstar et al., 2014) competitions. This dataset includes video recordings of participants engaging in various activities, such as singing and delivering speeches. Notably, the self-reporting in this dataset is conducted in the presence of mental health professionals, including psychiatrists, psychologists, and experienced counselors. The questionnaires used in AViD-Corpus gauge the severity of patients’ symptoms over specific periods, with responses recorded on a scale of 0 (not at all), 1 (several days), 2 (more than half the days), and 3 (nearly every day). The overall score is derived by summing the responses. These datasets have been instrumental in gaining insights into depression patterns. Nevertheless, as they are predominantly assembled and curated within controlled laboratory environments, they may not fully encapsulate the typical behaviors exhibited by individuals experiencing depression.

dataset	Modality	# Subjects	# Samples
DAIC-WOZ	A+V+T	189	189
Pittsburg	A+V	49	130
AViD-Corpus	A+V	292	340
D-vlog	A+V	816	961
E-Dvlog	A+V+T	1016	1261

Table 1: Comparison of various Depression datasets with E-Dvlog (Extended D-vlog). Where A: Audio, V: Video, T: Text.

The utilization of social media for depression detection has been increasingly used instead of clinical interviews. Social media datasets can reveal the patient’s unusual and atypical behavior, which cannot be seen in the clinical interview conducted under supervision, where the individual may not authentically express their emotion or actual behavior, which is shown in their daily lives (Yoon et al., 2022). Therefore, many approaches have been taken to detect Depression using data from social media sites such as Twitter, Reddit and Facebook. In recent years, depression detection using text from social media has been focused on (Fa-

tima et al., 2019), (Burdisso et al., 2019), (Chiong et al., 2021). Textual-based features focus on the linguistic features of the social media text, such as words, POS, n-gram, and other linguistic characteristics. In (Wang et al., 2013) uses text and tags from micro-blog (Sinba Weibo) used in China. They extracted content behavioral features from the blogs to detect Depression. In comparison, this method focuses on detecting Depression from Social media using text. However, more attention should be paid to video data and multimodal Fusion.

Multimodal Fusion is to combine multiple modalities to predict output. Work is done to detect Depression using multiple fusion. (Haque et al., 2018) in this paper the Authors have uses 3D Facial Expressions and spoken language as features from the dataset to detect Depression. (Yang et al., 2018) use text and video features and hybridizes deep and shallow models for depression estimation and classification from audio, video and text descriptors. (Ortega et al., 2019) proposed an end-to-end deep neural network (DNN) model that integrated three different modalities of speech features, facial features, and text transcription. Each modality is first encoded independently with fully connected layers and then combined into a single representation for estimating the emotional state of subject. Previous studies have investigated depression detection using multimodalities; however, the combination of Multimodal Transformer with wav2vec2 features alongside spectrograms, which has shown to yield superior results in detecting depression, has not been explored.

**Virtual Agents:** In recent years the with the increase in mental health problems, people have started taking emotional support from the text based platforms such as in (Eysenbach et al., 2004), (De Choudhury and De, 2014), (talkelife.co). there is also a rise in Empathetic virtual agents (?), which impart empathy in their responses by giving motivational responses and responses with hope and reflections which is seen as an important to uplift the spirit of an individual who is seeking support. Additionally, efforts have been made to enhance the therapeutic value of these platforms by incorporating insights (Fitzpatrick et al., 2017), (Xie and Pentina, 2022) encouraging exploration through open-ended questioning, and providing guidance and problem-solving techniques, all aimed at aiding users in their healing process.

### 3 Datasets

The D-vlog dataset (Yoon et al., 2022) is a collection of Depression vlogs of various people posted on YouTube. The D-vlog dataset has around 961 vlogs in total out of which 505 are categorized as depressive vlogs and 465 are categorized as Non-depressive vlogs. However, the D-vlog dataset (Yoon et al., 2022) has some limitations, such as the dataset majorly having Major Depressive Disorder and lacking Other Disorder such as Bipolar Disorder, Postmortem Disorder, and Anxiety with depression. Which will make the dataset more generalized. So, we extended the D-vlog dataset by adding around 300 more vlogs to the D-vlog dataset (Yoon et al., 2022) which now have more vlogs on various depressive disorder from varying age groups and different gender. Figure 1.

#### 3.1 Dataset Collection:

We have collected the dataset vlogs using certain keywords using YouTube API (Yin and Brown, 2018) and downloaded them using the yt-dlp package <sup>4</sup>.

**Depressive vlogs:** ‘depression daily vlog’, ‘depression journey’, ‘depression vlog’, ‘depression episode vlog’, ‘depression video diary’, ‘my depression diary’, and ‘my depression story’, ‘post-partum depression vlogs’, ‘Anxiety vlogs’.

**Non-Depressive vlogs:** ‘daily vlog’, ‘grwm (get ready with me) vlog’, ‘haul vlog’, ‘how to vlog’, ‘day of vlog’, ‘talking vlog’, and etc.

We used the same approach to collect the dataset as used in the D-vlog dataset (Yoon et al., 2022). We focused our analysis on vlogs featuring content creators who have a documented history of depression, currently manifesting symptoms of the condition. We specifically excluded vlogs that solely discussed having a bad day without a deeper connection to depressive experiences.

#### 3.2 Dataset Statistics:

The Extended D-vlog dataset has 1261 vlogs with 680 Depressive vlogs and 590 Non-Depressive vlogs as Shown in below Table 2

<sup>4</sup><https://github.com/yt-dlp/yt-dlp/wiki/Installation>

	Gender	# Samples
Depression	Male	273
	Female	406
Non-Depression	Male	232
	Female	350

Table 2: Extended D-vlog Statistics

Extended D-vlog dataset exhibits more representation of Female vlogs as compared to Male vlogs within Depressed category, reflecting high prevalence of depression among Female. In Non-depressive category similar trend is observed with more female representation than Male vlogs as predominantly "get ready with me vlogs", "Haul vlogs" are uploaded by Females. In Extended D-vlog we added vlogs which have other disorder other than Major Depressive Disorder such as Anxiety disorder, Bipolar Disorder <sup>5</sup>, and Postmortem Disorder. The below Figure 1 show the Distribution of various type Depressive vlogs.

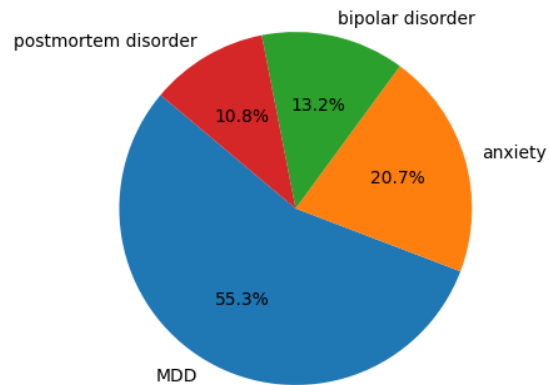


Figure 1: The Above figure shows the distribution of various type Depressive vlogs. where MDD is Major Depressive Disorder, Bipolar Disorder also called as Manic Disorder.

#### 3.3 Datasets for Therapeutic and Empathetic Conversations:

Acquiring datasets of therapy conversations poses a significant challenge as they are typically private and rarely shared. Moreover, potential privacy issues may arise when exposing therapy datasets to public LLM APIs as they may contain sensitive client information. Publicly available therapy conversation datasets are limited. Here, we use Three datasets that carefully preprocess publicly available

<sup>5</sup><https://www.nimh.nih.gov/health/topics/bipolar-disorder>

on therapy. This ensures high-quality transcripts while maintaining the confidentiality of sensitive personal information.

1. **High-and-Low-Quality Therapy Conversation Dataset (High-Low Quality):** The initial dataset, established by (Pérez-Rosas et al., 2018), encompasses 259 therapy dialogues, predominantly centering on evidence-based motivational interviewing (MI) therapy. Assessing the conversations in accordance with MI psychotherapy principles, the authors identify 155 transcripts of high quality and 104 of low quality within the dataset. Both high-quality and low-quality therapy dialogues conducted by human therapists are utilized to examine desirable and undesirable conversational behaviors.
2. **HOPE Dataset:** The second dataset from (Malhotra et al., 2022) was used to study dialogue acts in therapy. This dataset contains 212 therapy transcripts and includes conversations employing different types of therapy techniques (e.g., MI, Cognitive Behavioral Therapy).
3. **MotiVAte Dataset:** The MotiVAte Dataset (Saha et al., 2022) contains 7076 dyadic conversations with support seekers who have one of the four mental disorder: MDD, OCD, Anxiety or PTSD.

## 4 Methodology

We used TVLT (Textless Vision Language Transformer) (Tang et al., 2022), a minimal end-to-end vision and language Multimodal transformer model that takes raw video, raw audio, and text as input to the transformer model. TVLT (Tang et al., 2022) is a Textless Model, which implicitly does not use text, but with the ASR model (Whisper) (Radford et al., 2023), we can extract text from the audio segments. The TVLT model is more effective for Multimodal classification because the TVLT (Tang et al., 2022) model can capture visual and acoustic information, providing a more comprehensive fused representation of video, audio, and text.

For **Textual Feature**, we make use of the powerful BERT (Kenton and Toutanova, 2019) Language model, a pre-trained model described in to capture important features from text. This means we can understand not only the specific details in the text

but also the overall context. These BERT embeddings help us understand text thoroughly, making them perfect for tasks like analyzing sentiment or identifying depression. We apply BERT (Kenton and Toutanova, 2019) to our text, using specific dimensions (dt = 786), and we start with good initial weights using Xavier’s method (Kumar, 2017).

To extract **Audio features**, we employ a combination of techniques. First, we utilize low-level features like spectrograms, which are generated using the librosa Library (McFee et al., 2015). Additionally, we incorporate features from wav2vec2, as described in (Baevski et al., 2020). These wav2vec2 (Baevski et al., 2020) features encompass various acoustic attributes, including MFCC (Hossain et al., 2010), Spectral (Pachet and Roy, 2007), Temporal (Krishnamoorthy and Prasanna, 2011), and Prosody (Olwal and Feiner, 2005) features which help in identifying pitch, Intonation, Tempo of the audio segment. They excel in capturing both local and contextual information from the raw audio waveform. To create our final audio representation, we compute the average across the spectrogram vector and the wav2vec2 vector.

Our video processing pipeline involves several essential steps. First, we load the video file using a tool called VideoReader (Frith et al., 2005). Next, we randomly select a subset of frames from the video clip. These frames are then resized and cropped to focus on the subject’s frontal view. For extracting **visual features**, we rely on the powerful ViT (Vision Transformer) model introduced in (Dosovitskiy et al., 2020). This model helps us create what we call "vision embeddings." It does this by breaking down each video frame into smaller 16x16 patches. We then apply a linear projection layer to these patches, resulting in a 768-dimensional patch embedding. This vision embedding module is a critical component of our model. It takes each video frame or image and transforms it into a sequence of 768-dimensional vectors. These vectors are rich in both spatial and temporal information, making them invaluable for our model to comprehend the visual content within the input data.

We have implemented the architecture illustrated in Figure 1, where our TVLT (Tang et al., 2022) transformer model comprises a 12-layer encoder and an 8-layer decoder. To obtain the fused representation of all three modalities, we exclusively utilize the encoder portion of the model. These fused representations are subsequently fed into the

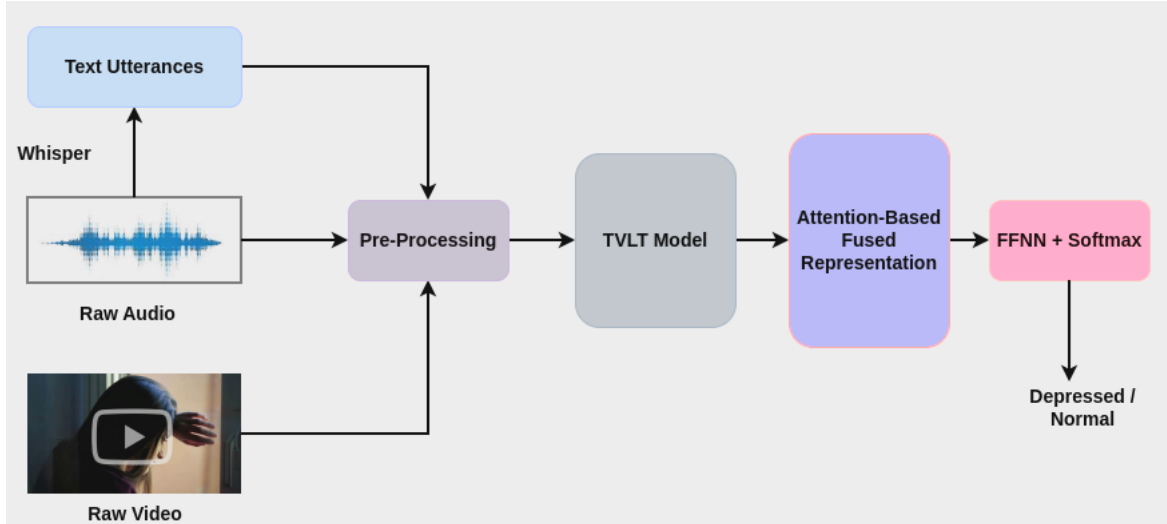


Figure 2: In the Above **Architecture** we leverage three different modalities such as video, audio and text where text is extracted from the audio segment using Whisper ASR Model. we then preprocess all the three modalities and pass to the model where we get the fused representation of all three modalities. These fused representation is then passed to the feed forward Neural Network with sigmoid function to get whether the individual exhibits the sign of Depression or is it in Normal state.

downstream task for depression prediction. Our model’s evaluation is conducted on the Extended D-vlog dataset, which consists of 35,046 video clips collected from 1016 different speakers. For each video clip, we generate text using the ASR model and manually correct any errors to ensure ground-truth transcriptions. In line with previous studies, we employ a 7:1:2 train-valid-test split and evaluate using weighted accuracy (WA) and F1 score metrics. For each downstream task, we introduce a task-specific head (a two-layer MLP) on top of the encoder representation. We train the model jointly using binary cross-entropy loss for these tasks.

$$L(y, \hat{y}) = - [y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})] \quad (1)$$

where  $y$ : True label and  $\hat{y}$ : Predicted label.

We have constructed a RAG pipeline incorporating domain-specific documents as an external knowledge base. This external knowledge is employed to validate responses generated by Fine-tuned Mistral Models on the Mental Health Dataset. In generating responses to user queries, we utilize a system prompt instructing "Act as if you’re a professional therapist. You provide evidence-based therapy, utilizing motivational interviewing techniques, to help clients in making behavioral changes, such as quitting smoking or alcohol consumption. You should maintain your therapist persona while responding. Communicate in a con-

versational style, mirroring the style of previous therapist responses"

## 5 Experiments

To obtain fused representation of audio, video and text modalities, we employ pretrained text-based TVLT (Tang et al., 2022) model on video dataset and subsequently fine-tune on the Extended D-vlog dataset.

We split our dataset into Train, Valid and Test Set in the ratio of 7:1:2. These three sets do not overlap i.e. no dataset is used in more than one set.

Gender	Train	Valid	Test
Male	354	51	100
Female	530	74	152

Table 3: Number of vlogs in Train, Valid and Test Split of Extended D-vlog dataset

For training the model, we utilized Adam’s Optimizer with learning rates ranging from 0.0001 to 0.00001 and batch sizes of 32 and 64. The model underwent four iterations with different seed values, each taking approximately three hours to train on Nvidia RTX A6000. Binary cross-entropy served as our chosen loss function for the depression detection task, and F1-scores were reported based on the test set results.

In the next phase, to generate empathetic and

therapeutic responses, we employed pretrained Mistral and Phi-2 models, finetune them on the Hope and Motivate Dataset. Our approach involved training these models using PEFT QLoRA, a novel technique combining 4-bit quantization with Low-rank Adapters to enhance memory utilization and computational efficiency. Additionally, we implemented an RAG Pipeline to ensure accurate responses without generating false information, utilizing Adam’s Optimizer with a learning rate set to 0.00025, known for its superior results.

## 6 Result and Discussion

To analyse the importance of each modality for depression detection, we trained our model on each modality separately and reported the results in the Table 4 below. We found that audio modality have

Modalities	F1-scores
T	0.57
A	0.60
V	0.56
V + A	0.631
V + T	0.628
A + T	0.634
V + A + T	<b>0.656</b>

Table 4: Results obtained on the Extended D-vlogs dataset via experiments with Video (V), Audio (A), Textual (T).

best F1-Score than the other modalities, implies that audio features are more importance than the visual and textual features for Depression detection. This suggest that people with depression have distinct speech features. Even though the audio features are more important than the visual features but when we combine the two modalities we can see that combining the two modalities yield better score than just using audio modality. Also, combination of audio and text modality prove to be better than just using audio as a single modality. Finally, combining all the three modalities we get much better result on relying on just two modalities which tells that combining audio features, visual features and textual features and their relationship is more effective for the Depression detection.

Modalities	F1-scores
V + A + T	0.656
V + A + T(Mask)	0.663
V + A(W2V2+Spect) + T	<b>0.678</b>
V(Mask) + A + T	0.661

Table 5: Results obtained on the Extended D-vlogs dataset via experiments with Video (V), Audio (A), Textual (T). T(Mask) is text with word-masking, V(Mask) is Video frames with frame-masking and A(W2V2+Spect) is Audio with wav2vec2 +spectrogram features.

The introduction of random word masking in text modalities proves instrumental in enhancing the model’s understanding of textual information. This improvement becomes apparent when analyzing the results, where a subtle yet noteworthy performance boost of 0.007% is observed in comparison to not employing word masking in the text. Moreover, the application of frame masking to video data, when combined with audio and text modalities, also contributes to a slight enhancement of 0.005%. These findings underscore the efficacy of incorporating diverse modalities in the model. Table 6 provides a comprehensive overview of the results. It unmistakably highlights the significance of leveraging all three modalities—text, video, and audio—in conjunction with wav2vec2 (Baevski et al., 2020) features and spectrograms, as opposed to using spectrograms for audio processing. This approach leads to an impressive F1-score of 67.8%.

We have extensively evaluate the TVLT (Tang et al., 2022) model performance on the D-vlog dataset (Yoon et al., 2022) and compared its results with several baseline models. The purpose was to check the effectiveness of the TVLT (Tang et al., 2022) model in evaluating task on Depression detection on D-vlog dataset. The TVLT (Tang et al., 2022) model in isolation was when performed on D-vlog dataset (Yoon et al., 2022), surpasses the Cross Attention State-of-the-Art model by 2.2 %. This improvement established the TVLT (Tang et al., 2022) model as the New Benchmark for the D-vlog dataset (Yoon et al., 2022). The result Obtained by TVLT model on D-vlog dataset (Yoon et al., 2022) states that the TVLT (Tang et al., 2022) model was correctly able to understand the characteristics and handling the complexity of the D-vlog dataset (Yoon et al., 2022). The exceptional performance of the TVLT (Tang et al., 2022) model could serve as a catalyst, motivating researchers to explore and develop more advanced techniques in the realm of

Model Type	Model	Precision	Recall	F1-Score
<b>Fusion Baseline</b>	Concat	62.51	63.21	61.1
	Add	59.11	60.38	58.1
	Multiply	63.48	64.15	63.09
<b>Depression Detector</b>	Cross-Attention	65.4	65.5	63.5
<b>Our Model</b>	TVLT Model	<b>67.3</b>	<b>68.3</b>	<b>67.8</b>

Table 6: Comparison of various baseline models with our model on the D-vlog dataset

multi-modal analysis and deep learning.

## 7 Qualitative Analysis

This section demonstrates how integrating wav2vec2 (Baevski et al., 2020) features significantly improves our TVLT (Tang et al., 2022) model’s ability to detect depression accurately. By capturing vocal cues in audio data, wav2vec2 (Baevski et al., 2020) enhances our model’s performance compared to relying solely on spectrogram data. This inclusion enables our model to make predictions that would otherwise be challenging.

As shown in the Table 7 provided contains selected examples where our model correctly identifies depression. In the first instance, despite consistent facial expressions, the audio analysis reveals monotone tone, low pitch, and crying, supported by distressing textual content. Our TVLT (Tang et al., 2022) model, enhanced with wav2vec2 (Baevski et al., 2020) and spectrogram features, accurately predicts depression in this case, highlighting the crucial role of wav2vec2 (Baevski et al., 2020) features in capturing vocal cues that spectrogram data may miss.

In the second example, we observe a girl expressing both a smile and tears, with audio and textual evidence indicating depression. While the model lacking wav2vec2 (Baevski et al., 2020) features fails to predict accurately, the inclusion of wav2vec2 (Baevski et al., 2020) correctly identifies the depressive nature of the example. This showcases wav2vec2’s (Baevski et al., 2020) ability to extract vital information from audio, enhancing the model’s comprehension of depression cues. In the third example, despite minimal facial expression variation and unremarkable audio, textual analysis reveals indications of depression, challenging our model’s predictive accuracy.

## 8 Summary, Conclusion & Future Work

In this study, we introduced an Extended D-vlog dataset, comprising 1261 videos that include both vlogs by individuals with depression (680 videos) and those without (590 videos). Our objective is to detect depression in non-verbal and non-clinical vlogs. To achieve this, we employed a TVLT model, which is a multimodal transformer (Tang et al., 2022), to create a multimodal representation using text, video, and audio data. We utilized the ViT model for visual embeddings and extracted audio features with wav2vec2 (Baevski et al., 2020) and spectrograms. The TVLT model, incorporating all three modalities, yielded a noteworthy F1-score of 0.656. By introducing text word-masking, we improved the F1-score to 0.663, representing a 0.007% enhancement over the absence of word-level masking. Additionally, with frame-masking, the F1-score reached 0.661.

Our TVLT model, combined with the supplementary wav2vec2 and spectrogram features, outperformed all baseline models on the D-vlog dataset (Yoon et al., 2022) and established a new benchmark on the Extended D-vlog dataset. We believe that our introduced dataset and the multi-modal depression detection model have the potential to play a significant role in the early identification of individuals experiencing depression through their social media presence. This proactive approach aims to ensure timely access to essential clinical interventions for those in need.

In future work, we plan to extend the scope of our research to detect various type of mental health diseases. Additionally, we aim to investigate and incorporate multilingual capabilities into our work frame so that it can be scaled and utilized for various language settings. we are also building virtual agent which will provide Empathetic and Therapeutic responses.

## 9 Limitation:

We have curated our dataset to exclusively feature vlogs from individuals who have experienced or are currently experiencing depression, acknowledging the potential for bias inherent in this selection process. However, we have taken measures to mitigate this bias to the best of our ability. Our model encounters challenges in accurately predicting certain depressed classes, notably in cases such as 'smiling depression,' where individuals conceal genuine emotions behind a facade of cheerfulness and high functionality, making detection of this issue particularly challenging. Creation of Virtual Agents is still under process and we will working on the results, Qualitative and Quantitative analysis.

## 10 Ethics Statement:

All vlogs included in the dataset were voluntarily uploaded by individuals onto the YouTube platform, and each vlog is authored by its respective uploader. None of the vlogs in the dataset have been sourced from other platforms without explicit consent. All the data in the dataset has been sourced from open-access platforms, and none of the videos or text within it contain any offensive or derogatory language aimed at any particular team or entity.

## References

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.

Sergio G Burdisso, Marcelo Errecalde, and Manuel Montes-y Gómez. 2019. A text classification framework for simple and effective early depression detection over social media streams. *Expert Systems with Applications*, 133:182–197.

Raymond Chiong, Gregorious Satia Budhi, and Sandeep Dhakal. 2021. Combining sentiment lexicons and content-based features for depression detection. *IEEE Intelligent Systems*, 36(6):99–105.

Neo Christopher Chung, George Dyer, and Lennart Brocki. 2023. Challenges of large language models for mental health counseling. *arXiv preprint arXiv:2311.13857*.

Cleveland Clinic. [Depression](#).

Munmun De Choudhury and Sushovan De. 2014. Mental health discourse on reddit: Self-disclosure, social support, and anonymity. In *Proceedings of the international AAAI conference on web and social media*, volume 8, pages 71–80.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.

Sara Evans-Lacko, Sergio Aguilar-Gaxiola, Ali Al-Hamzawi, and et al. 2018. Socio-economic variations in the mental health treatment gap for people with anxiety, mood, and substance use disorders: Results from the who world mental health (wmh) surveys. *Psychological Medicine*, 48(9):1560–1571.

Gunther Eysenbach, John Powell, Marina Englesakis, Carlos Rizo, and Anita Stern. 2004. Health related virtual communities and electronic support groups: systematic review of the effects of online peer to peer interactions. *Bmj*, 328(7449):1166.

Iram Fatima, Burhan Ud Din Abbasi, Sharifullah Khan, Majed Al-Saeed, Hafiz Farooq Ahmad, and Rafia Mumtaz. 2019. Prediction of postpartum depression using machine learning techniques from social media text. *Expert Systems*, 36(4):e12409.

Kathleen Kara Fitzpatrick, Alison Darcy, and Molly Vierhile. 2017. Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): a randomized controlled trial. *JMIR mental health*, 4(2):e7785.

Simon Frith, Andrew Goodwin, and Lawrence Grossberg. 2005. *Sound and vision: The music video reader*. Routledge.

Saptarshi Ghosh, Asif Ekbal, and Pushpak Bhattacharyya. 2022. A multitask framework to detect depression, sentiment, and multi-label emotion from suicide notes. *Cognitive Computation*, pages 1–20.

Yuan Gong and Christian Poellabauer. 2017. Topic modeling based multi-modal depression detection. In *Proceedings of the 7th annual workshop on Audio/Visual emotion challenge*, pages 69–76.

Jonathan Gratch, Ron Artstein, Gale M Lucas, Giota Stratou, Stefan Scherer, Angela Nazarian, Rachel Wood, Jill Boberg, David DeVault, Stacy Marsella, et al. 2014. The distress analysis interview corpus of human and computer interviews. In *LREC*, pages 3123–3128. Reykjavik.

Albert Haque, Michelle Guo, Adam S Miner, and Li Fei-Fei. 2018. Measuring depression symptom severity from spoken language and 3d facial expressions. *arXiv preprint arXiv:1811.08592*.

Md. Afzal Hossan, Sheeraz Memon, and Mark A Gregory. 2010. A novel approach for mfcc feature extraction. In *2010 4th International Conference on Signal Processing and Communication Systems*, pages 1–5.

737	Kate Keenan, Alison Hipwell, Tammy Chung,	Verónica Pérez-Rosas, Xueting Sun, Christy Li, Yuchen	791
738	Stephanie Stepp, Magda Stouthamer-Loeber, Rolf	Wang, Kenneth Resnicow, and Rada Mihalcea. 2018.	792
739	Loeber, and Kathleen McTigue. 2010. The pittsburgh	Analyzing the quality of counseling conversations:	793
740	girls study: overview and initial findings. <i>Journal of</i>	the tell-tale signs of high-quality counseling. In <i>Pro-</i>	794
741	<i>Clinical Child &amp; Adolescent Psychology</i> , 39(4):506–	<i>ceedings of the Eleventh International Conference on</i>	795
742	521.	<i>Language Resources and Evaluation (LREC 2018)</i> .	796
743	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	797
744	Toutanova. 2019. Bert: Pre-training of deep bidirec-	man, Christine McLeavey, and Ilya Sutskever. 2023.	798
745	tional transformers for language understanding. In	Robust speech recognition via large-scale weak su-	799
746	<i>Proceedings of naacL-HLT</i> , volume 1, page 2.	pervision. In <i>International Conference on Machine</i>	800
747	P Krishnamoorthy and SR Mahadeva Prasanna. 2011.	<i>Learning</i> , pages 28492–28518. PMLR.	801
748	Enhancement of noisy speech by temporal and spec-	Anupama Ray, Siddharth Kumar, Rutvik Reddy, Pre-	802
749	tral processing. <i>Speech Communication</i> , 53(2):154–	rana Mukherjee, and Ritu Garg. 2019. Multi-level	803
750	174.	attention network using text, audio and video for	804
751	Kurt Kroenke, Tara W. Strine, Robert L. Spitzer,	depression prediction. In <i>Proceedings of the 9th in-</i>	805
752	Janet B.W. Williams, Joseph T. Berry, and Ali H.	<i>ternational on audio/visual emotion challenge and</i>	806
753	Mokdad. 2009. The phq-8 as a measure of current	<i>workshop</i> , pages 81–88.	807
754	depression in the general population. <i>Journal of Af-</i>	Tulika Saha, Saichethan Reddy, Anindya Das, Sriparna	808
755	<i>fective Disorders</i> , 114(1-3):163–173.	Saha, and Pushpak Bhattacharyya. 2022. A shoulder	809
756	Siddharth Krishna Kumar. 2017. On weight initial-	to cry on: towards a motivational virtual assistant	810
757	ization in deep neural networks. <i>arXiv preprint</i>	for assuaging mental agony. In <i>Proceedings of the</i>	811
758	<i>arXiv:1704.08863</i> .	<i>2022 conference of the North American chapter of</i>	812
759	Zilin Ma, Yiyang Mei, and Zhaoyuan Su. 2023. Under-	<i>the association for computational linguistics: Human</i>	813
760	<i>language technologies</i> , pages 2436–2449.	814	
761	language model-based conversational agents for men-	Ashish Sharma, Inna W Lin, Adam S Miner, David C	815
762	tal well-being support. In <i>AMIA Annual Symposium</i>	Atkins, and Tim Althoff. 2023. Human–ai collabora-	816
763	<i>Proceedings</i> , volume 2023, page 1105. American	tion enables more empathic conversations in text-	817
764	Medical Informatics Association.	based peer-to-peer mental health support. <i>Nature</i>	818
765	Ganeshan Malhotra, Abdul Waheed, Aseem Srivastava,	<i>Machine Intelligence</i> , 5(1):46–57.	819
766	Md Shad Akhtar, and Tanmoy Chakraborty. 2022.	Heung-Yeung Shum, Xiao-dong He, and Di Li. 2018.	820
767	Speaker and time-aware joint contextual learning for	From eliza to xiaoice: challenges and opportunities	821
768	dialogue-act classification in counselling conversa-	with social chatbots. <i>Frontiers of Information Tech-</i>	822
769	tions. In <i>Proceedings of the fifteenth ACM interna-</i>	<i>nology &amp; Electronic Engineering</i> , 19:10–26.	823
770	<i>tional conference on web search and data mining</i> ,	Karen L Smarr and Autumn L Keefer. 2011. Measures	824
771	pages 735–745.	of depression and depressive symptoms: Beck depres-	825
772	Brian McFee, Colin Raffel, Dawen Liang, Daniel P	sion inventory-ii (bdi-ii), center for epidemiologic	826
773	Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto.	studies depression scale (ces-d), geriatric depression	827
774	2015. librosa: Audio and music signal analysis in	scale (gds), hospital anxiety and depression scale	828
775	python. In <i>Proceedings of the 14th python in science</i>	(hads), and patient health questionnaire-9 (phq-9).	829
776	<i>conference</i> , volume 8, pages 18–25.	<i>Arthritis care &amp; research</i> , 63(S11):S454–S466.	830
777	Alex Olwal and Steven Feiner. 2005. Interaction tech-	Zineng Tang, Jaemin Cho, Yixin Nie, and Mohit Bansal.	831
778	niques using prosodic features of speech and audio	2022. Tvlt: Textless vision-language transformer.	832
779	localization. In <i>Proceedings of the 10th international</i>	<i>Advances in Neural Information Processing Systems</i> ,	833
780	<i>conference on Intelligent user interfaces</i> , pages 284–	35:9617–9632.	834
781	286.	Michel Valstar, Björn Schuller, Kirsty Smith, Timur Al-	835
782	Juan DS Ortega, Mohammed Senoussaoui, Eric	maev, Florian Eyben, Jarek Krajewski, Roddy Cowie,	836
783	Granger, Marco Pedersoli, Patrick Cardinal, and	and Maja Pantic. 2014. Avec 2014: 3d dimensional	837
784	Alessandro L Koerich. 2019. Multimodal fusion with	affect and depression recognition challenge. In <i>Pro-</i>	838
785	deep neural networks for audio-video emotion recog-	<i>ceedings of the 4th international workshop on au-</i>	839
786	nition. <i>arXiv preprint arXiv:1907.03196</i> .	<i>dio/visual emotion challenge</i> , pages 3–10.	840
787	François Pachet and Pierre Roy. 2007. Exploring bil-	Michel Valstar, Björn Schuller, Kirsty Smith, Florian	841
788	lions of audio features. In <i>2007 international work-</i>	Eyben, Bihan Jiang, Sanjay Bilakhia, Sebastian	842
789	<i>shop on content-based multimedia indexing</i> , pages	Schnieder, Roddy Cowie, and Maja Pantic. 2013.	843
790	227–235. IEEE.	Avec 2013: the continuous audio/visual emotion and	844
		depression recognition challenge. In <i>Proceedings of</i>	845
		<i>the 3rd ACM international workshop on Audio/visual</i>	846
		<i>emotion challenge</i> , pages 3–10.	847

848 Xinyu Wang, Chunhong Zhang, Yang Ji, Li Sun, Leijia  
849 Wu, and Zhana Bao. 2013. A depression detection  
850 model based on sentiment analysis in micro-blog  
851 social network. In *Trends and Applications in Knowl-  
852 edge Discovery and Data Mining: PAKDD 2013  
853 International Workshops: DMApps, DANTh, QIMIE,  
854 BDM, CDA, CloudSD, Gold Coast, QLD, Australia,  
855 April 14-17, 2013, Revised Selected Papers 17*, pages  
856 201–213. Springer.

857 World Health Organization. [Depression](#).

858 Tianling Xie and Iryna Pentina. 2022. Attachment the-  
859 ory as a framework to understand relationships with  
860 social chatbots: a case study of replika.

861 Le Yang, Dongmei Jiang, and Hichem Sahli. 2018. Inte-  
862 grating deep and shallow models for multi-modal de-  
863 pression analysis—hybrid architectures. *IEEE Trans-  
864 actions on Affective Computing*, 12(1):239–253.

865 Le Yang, Dongmei Jiang, Xiaohan Xia, Ercheng Pei,  
866 Meshia Cédric Oveneke, and Hichem Sahli. 2017a.  
867 Multimodal measurement of depression using deep  
868 learning models. In *Proceedings of the 7th Annual  
869 Workshop on Audio/Visual Emotion Challenge*, pages  
870 53–59.

871 Le Yang, Hichem Sahli, Xiaohan Xia, Ercheng Pei,  
872 Meshia Cédric Oveneke, and Dongmei Jiang. 2017b.  
873 Hybrid depression classification and estimation from  
874 audio video and text information. In *Proceedings  
875 of the 7th annual workshop on audio/visual emotion  
876 challenge*, pages 45–51.

877 Leon Yin and Megan Brown. 2018. [Smappnyu/youtube-  
878 data-api](#).

879 Jeewoo Yoon, Chaewon Kang, Seungbae Kim, and Jiny-  
880 oung Han. 2022. D-vlog: Multimodal vlog dataset  
881 for depression detection. In *Proceedings of the AAAI  
882 Conference on Artificial Intelligence*, volume 36,  
883 pages 12226–12234.

## A Appendix

### A.1 TVLT Model:

Textless Vision-Language Transformer (TVLT), a model designed for vision-and-language representation learning using raw visual and audio inputs. Unlike traditional approaches, TVLT employs homogeneous transformer blocks with minimal modality-specific design and does not rely on text-specific modules such as tokenization or automatic speech recognition (ASR). TVLT is trained using masked autoencoding to reconstruct masked patches of continuous video frames and audio spectrograms, as well as contrastive modeling to align video and audio. Experiments demonstrate that TVLT achieves comparable performance to text-based models across various multimodal tasks, including visual question answering, image retrieval, video retrieval, and multimodal sentiment analysis. Additionally, TVLT offers significantly faster inference speed (28x) and requires only one-third of the parameters. These results suggest the feasibility of learning compact and efficient visual-linguistic representations directly from low-level visual and audio signals, without relying on pre-existing text data.

### A.2 Pretaining details:

- **HowTo100M:** We used HowTo100M, a dataset containing 136M video clips of a total of 134,472 hours from 1.22M YouTube videos to pretrain our model. Our vanilla TVLT is pretrained directly using the frame and audio stream of the video clips. Our text-based TVLT is trained using the frame and caption stream of the video. The captions are automatically generated ASR provided in the dataset. We used 0.92M videos for pretraining, as some links to the videos were invalid to download.
- **YTTemporal180M:** YTTemporal180M includes 180M video segments from 6M YouTube videos that spans multiple domains, and topics, including instructional videos from HowTo100M, lifestyle vlogs of everyday events from the VLOG dataset, and YouTube’s auto-suggested videos for popular topics like ‘science’ or ‘home improvement’.

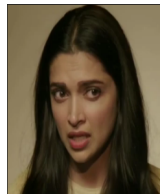
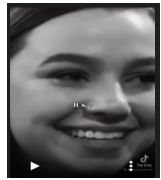
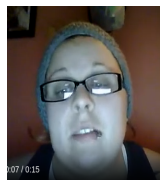
Utterance	Ground Truth	Prediction (w2v2 + spect)	Prediction w/o (w2v2 + spect)	Video frames
I knew what I was feeling, but I don't think I was able to communicate entirely what I was feeling. Like I knew I had this pittish feeling in my stomach. I knew that I'd be scared to wake up. I didn't want to wake up. Yeah, I think waking up was tough because I didn't want to face a day.	Depression	Depression	Normal	
Some days, it's really, really hard to just move. It's... I like it. I, yeah, it's hard to get out of bed. It's hard to even go downstairs to get something to eat.	Depression	Depression	Normal	
No concept of time, no sense of feeling. Have I become cold, dead to the world, where I once mattered? I can't even remember when I was important to someone last. Everything has escaped me. Deeper I fall into a void.	Depression	Normal	Normal	

Table 7: A **Qualitative analysis**, In the given instances, the model, equipped with both wav2vec2 and spectrogram features, effectively detects depression through audio analysis. In the first example, despite seemingly normal facial expressions, the model accurately detects depression. In the second case, the model succeeds in identifying depression even when the individual smiles while crying, whereas the model relying solely on spectrogram data falls short in these situations. In the third scenario, the woman's facial expressions and audio do not exhibit evident signs of depression, while text analysis reveals potential indicators that challenge our model's accuracy, resulting in an incorrect prediction.