

# Reinforce Trustworthiness in Multimodal Emotional Support System

Huy M. Le<sup>1,3,\*</sup>, Dat Tien Nguyen<sup>1,3,\*</sup>, Ngan T. T. Vo<sup>3</sup>, Tuan D. Q. Nguyen<sup>3</sup>, Nguyen Binh Le<sup>3</sup>, Duy Minh Ho Nguyen<sup>5,6,7</sup>, Daniel Sonntag<sup>5,8</sup>, Lizi Liao<sup>2</sup>, Binh T. Nguyen<sup>4,†</sup>

<sup>1</sup> Mohamed bin Zayed University of Artificial Intelligence (MBZUAI)

<sup>2</sup> Singapore Management University

<sup>3</sup> University of Information Technology, Vietnam National University, Ho Chi Minh City

<sup>4</sup> University of Science, Vietnam National University, Ho Chi Minh City

<sup>5</sup> German Research Center for Artificial Intelligence (DFKI)

<sup>6</sup> Max Planck Research School for Intelligent Systems (IMPRS-IS)

<sup>7</sup> University of Stuttgart

<sup>8</sup> University of Oldenburg

HuyM.Le@mbzuai.ac.ae, lzliao@smu.edu.sg, ngtbinh@hcmus.edu.vn

## Abstract

In today’s world, emotional support is increasingly essential, yet it remains challenging for both those seeking help and those offering it. Multimodal approaches to emotional support show great promise by integrating diverse data sources to provide empathetic, contextually relevant responses, fostering more effective interactions. However, current methods have notable limitations, often relying solely on text-based data, or recognizing emotions only, thus overlooking the full potential of multimodal inputs. Moreover, many studies prioritize response generation without accurately identifying critical emotional support elements or ensuring the reliability of outputs. To overcome these issues, we introduce MULTIMOOD, a new framework that (i) leverages multimodal embeddings from video, audio, and text to predict emotional components and to produce responses aligned with professional therapeutic standards. To improve trustworthiness, we (ii) incorporate novel psychological criteria and apply Reinforcement Learning (RL) to optimize large language models (LLMs) for consistent adherence to these standards. We also (iii) analyze several advanced LLMs to assess their multimodal emotional support capabilities. Experimental results show that MultiMood achieves state-of-the-art on MESC and DFEW datasets while RL-driven trustworthiness improvements are validated through human and LLM evaluations, demonstrating its superior capability in applying a multimodal framework in this domain.

**Code** — <https://github.com/quangtuan-0504/Multimood>

## 1 Introduction

Mental health challenges are an urgent global concern, profoundly affecting individuals and communities. The World Health Organization estimated that 970 million people—one in eight globally—lived with a mental disorder in 2019, primarily anxiety and depression (WHO 2022). In the United

\*These authors contributed equally.

†Corresponding Authors

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

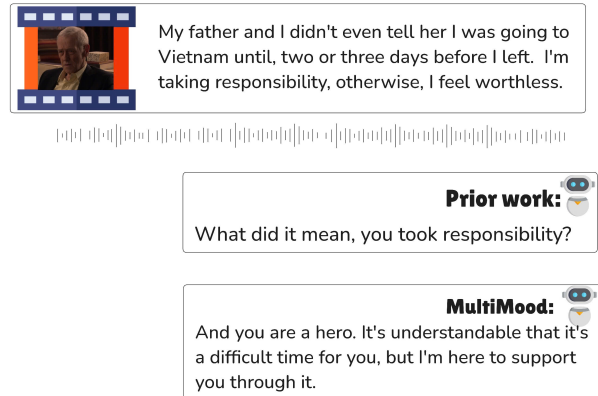


Figure 1: Example conversation illustrating the difference between prior systems and MULTIMOOD. Prior methods respond with factual queries, whereas MULTIMOOD demonstrates emotional awareness and offers empathetic, supportive feedback.

States, nearly 23% of adults experienced mental illness in 2021 (NAMI 2023). These conditions also impose a major economic burden, projected to reach \$6 trillion annually by 2030 (Marquez and Saxena 2016). These figures underscore the need for scalable, innovative tools to support psychological well-being, with artificial intelligence (AI) emerging as a promising aid.

Advances in large language models (LLMs) and vision–language models (VLMs) have transformed text generation, dialogue systems, and multimodal reasoning (Mitsui et al. 2024). These models now support applications in healthcare (Esteva et al. 2019; Nguyen et al. 2024), summarization (Le, Luong, and Luong 2023), and retrieval (Le et al. 2025a,b). Their versatility makes them promising tools for addressing complex social challenges, including mental health care (Margaroli et al. 2025). Systems such as Woebot show that AI-driven dialogue can alleviate de-

pression and anxiety through cognitive behavioral therapy (CBT)–inspired conversations (Fitzpatrick, Darcy, and Vierhile 2017; Rashkin et al. 2019). However, most existing systems remain text-only, overlooking nonverbal cues—tone, facial expression, and gesture—that are essential for empathy and trust (Ekman 2003). Empirical studies indicate that multimodal signals strengthen emotional understanding and engagement in human–computer interaction (Sim, Fortuno, and Choo 2024; Saffaryazdi et al. 2025). Consequently, text-only systems often lack the authenticity and nuance required for effective emotional support.

To address these limitations, we introduce *MultiMood*, a multimodal framework that integrates text, audio, and visual information to enhance emotional understanding in support-oriented dialogue (Figure 1). Unlike prior text-focused approaches, *MultiMood* leverages fine-grained cues—tone, prosody, facial expressions, and dialogue context—to generate empathetic, context-aware responses. Beyond multimodal fusion, *MultiMood* emphasizes *trustworthy alignment*: it employs reinforcement learning with human-defined psychological objectives to guide emotionally appropriate behavior. Specifically, it combines Proximal Policy Optimization (PPO) for stable learning with Group Relative Policy Optimization (GRPO) for fine-grained alignment to therapeutic standards, enabling fluent, safe, and ethically consistent responses suitable for AI-assisted emotional support.

*MultiMood* processes multimodal tokens to infer key emotional-support components—including user and supporter emotions, counseling strategies, and dialogue intent—which then guide response generation aligned with professional psychological frameworks. Our main contributions are as follows:

- (i) Propose the *MultiMood* architecture, integrating multimodal features (text, audio, and vision) for emotional-support dialogue.
- (ii) Design a trustworthiness-alignment framework with reinforcement-learning objectives that promote emotionally appropriate and reliable responses.
- (iii) Evaluate state-of-the-art LLMs on a multimodal emotional-support dataset, demonstrating improvements in empathy, trustworthiness, and contextual accuracy.

Overall, this work advances the development of responsible multimodal emotional-support systems, offering a more holistic and human-centered approach to promoting psychological well-being.

## 2 Background

### 2.1 Dataset

The MESC dataset (Chu et al. 2025), sourced from seasons 1–3 of *In Treatment*, comprises 1,019 dialogues and 28,762 utterances across text, audio, and video, annotated with 7 emotions (e.g., anger, sadness, disgust) and 10 therapeutic strategies (e.g., open questions, interpretation). Initially labeled using GPT-3.5 and refined by experts, it supports tasks like emotion recognition, strategy prediction, and response generation. MESC distinguishes itself from MELD (Poria

et al. 2019) and ESConv (Liu et al. 2021) with its multimodal and therapeutic focus; it advances empathetic AI for mental health. Analysis reveals prevalent neutral therapist emotions, reflecting their neutral stance to build client trust, which is consistent with counseling practices and not affecting model outcomes (Chu et al. 2025). Besides, we also do experiments on the DFEW dataset (Jiang et al. 2020), a dynamic facial expression database. DFEW consists of over 16,000 video clips from movies, which were also annotated with seven emotions. These video clips contain various challenging interferences in practical scenarios such as extreme illumination, occlusions, and capricious pose changes.

### 2.2 Task Definition

Our goal is to emulate a human therapist’s nuanced functions in real-life therapeutic sessions. We decompose the AI-user interaction into four key tasks (Chu et al. 2025) forming an emotionally intelligent support framework. Only Task 1 is referenced in both MESC and DFEW datasets, while Tasks 2–4 are exclusive to MESC:

- (i) **User Emotion Recognition (Task 1)**: Identifies the client’s emotion using multimodal cues (facial expressions, vocal prosody, text), enabling sensitive responses to psychological needs.
- (ii) **System Emotion Prediction (Task 2)**: Predicts the system’s emotional tone (e.g., neutral, angry,...) to align with the chosen strategy, fostering rapport and trust.
- (iii) **System Strategy Prediction (Task 3)**: Selects the optimal therapeutic strategy (e.g., validation, reflection) based on user emotion and dialogue history, mirroring tailored therapist techniques.
- (iv) **System Response Generation (Task 4)**: Generates a natural, contextually appropriate response embodying the predicted tone and strategy, promoting emotional safety and insight.

### 2.3 Related works

**Emotional Support Frameworks** In psychological counseling, several established theoretical frameworks guide practitioners in addressing psychological and emotional difficulties. CBT (Beck and Weishaar 1989) is a structured, evidence-based approach that targets maladaptive thoughts to improve emotions and behaviors, ideal for anxiety and depression. Humanistic Therapy, such as Rogers’ person-centered (Rogers 1957) approach, fosters self-actualization through empathy and unconditional regard, effective for self-esteem and existential concerns.

Recent advancements in AI have enhanced emotional support systems, addressing limitations of smartphone-based conversational agents (Miner et al. 2016), Muffin framework (Sheng et al. 2023) uses model-agnostic AI feedback and contrastive learning to improve response fluency and relevance. Hybrid Empathetic Framework (HEF) (Yang et al. 2024) integrates LLMs with small-scale empathetic models to enhance emotion detection and response generation. The Sequential SMES framework (Chu et al. 2025) leverages multimodal data to simulate therapeutic empathy

Framework	Approach			Training method	Output				
	Visual	Audio	Text		User Emo.	Therapist Emo.	Strategy	Response	Trust. Aware.
<b>InternVideo2.5</b>	✓		✓	SFT+RL					
<b>VideoLLaVA</b>	✓		✓	SFT	✓				
<b>EmotionLLaMA</b>	✓	✓	✓	SFT	✓				
<b>SMES</b>	✓	✓	✓	SFT	✓			✓	
<b>MultiMood (ours)</b>	✓	✓	✓	SFT+RL	✓	✓	✓	✓	✓

Table 1: Comparison between MULTIMOOD and other multi-LLM models for emotion recognition. ‘‘SFT’’ = supervised fine-tuning, ‘‘RL’’ = reinforcement learning, ‘‘Resp.’’ = response generation, ‘‘Trust’’ = trust-awareness. InternVideo2.5 has never been used for emotional tasks before.

and deliver tailored responses. Our MULTIMOOD framework advances these efforts by incorporating trustworthiness through reinforcement learning with PPO and GRPO, ensuring safe, empathetic, and contextually appropriate responses for diverse user needs.

**Multimodal LLMs** Multimodal LLMs like InternVideo2.5 (Wang et al. 2025) and VideoLLaVA (Lin et al. 2024a) advance video understanding. InternVideo2.5 employs a single InternViT encoder with Hierarchical Token Compression (HICO) to merging similar video tokens to reduce computation while preserving quality (Wang et al. 2025). VideoLLaVA aligns images and videos using LanguageBind for unified visual representation via a shared projection layer (Lin et al. 2024a), but lacks audio processing, unlike MULTIMOOD’s modality-specific projectors. These models focus on visual content while missing audio cues (volume, tone, pitch) critical for emotion recognition. Multimodal LLMs also enhance emotional support, overcoming single-modality limitations by capturing nuanced emotional signals for empathetic AI. EmotionLLaMA (Cheng et al. 2024a) uses the MERR dataset (28,618 samples) and specialized encoders for precise emotion recognition. The SMES framework processes multimodal inputs for emotion recognition, strategy prediction, and response generation, improving therapeutic mimicry (Chu et al. 2025). MULTIMOOD stands out with specialized encoders per modality and a reinforcement learning algorithm designed for trustworthiness, as summarized in Table 1.

**Trustworthiness in Responses** Trustworthiness is essential for effective emotional support from therapists and doctors, fostering a safe space for patient vulnerability. Goleman’s emotional intelligence framework (Boyatzis, Goleman, and Rhee 2000) emphasizes empathy, self-regulation, and social skills as key to building trust, enabling clinicians to communicate effectively. Crits-Christoph et al. (Crits-Christoph et al. 2019) highlight that trust, distinct from therapeutic alliance, encourages sharing private information, with racial disparities (e.g., lower trust among Black patients) underscoring equity’s role. Richmond et al. (Richmond et al. 2022) link trustworthiness to communication, fidelity, and fairness, noting that lower trust can delay care. In LLMs, *trustworthiness* is critical for safe, supportive interactions, as outlined in TrustLLM’s eight dimensions (Huang et al. 2024): *truthfulness* ensures accuracy, *safety* fosters healthy dialogue, *fairness* promotes impartiality, and *robustness* ensures reliability. *Privacy* protects autonomy, *machine*

*ethics* ensures moral behavior, *transparency* provides clarity, and *accountability* holds LLMs responsible. In MULTIMOOD, these factors are integrated to train robust LLMs, significantly reducing hallucination. Building on these foundations, we propose a tailored set of trustworthiness dimensions for emotional support systems to improve automatically generated responses to meet therapeutic standards.

### 3 Methodology

#### 3.1 Overview

The MULTIMOOD framework, shown in Figure 2, integrates an audio encoder  $\mathcal{E}^{aud}$ , a visual encoder  $\mathcal{E}^{vis}$ , a conversation compressor  $\mathcal{C}$ , and a large language model  $\phi$ . For an input tuple  $P = \langle \text{Audio, Video, Prompt, History} \rangle$ , the model is defined as:

$$\hat{O} = \Psi(\phi, \mathcal{E}, \Omega, \mathcal{C}, P), \quad (1)$$

where  $\mathcal{E}$  combines audio, vision, and text encoders,  $\Omega$  is the vision pre-processor, and  $\hat{O}$  is the text output. A multi-tower architecture generates modality-specific embeddings: video via vision tower  $f_V$ , audio via audio tower  $f_A$ , and text via  $\phi$  tokenizer  $f_T$ , yielding  $E_T = f_T(\text{Prompt})$ . A compressor distills text histories into concise representations,  $E_H = \text{ConvCompressor}(H)$ , enabling efficient context processing. Embeddings  $[E'_V; E'_A; E_T; E'_H]$  are aligned via modality-specific projectors and fed into the LLM to predict outcomes for four tasks.

#### 3.2 Framework Components

**Modality-Specific Encoder** The vision pre-processor  $\Omega$  uses the input video as a frame sequence, processed by a CLIP-based (Radford et al. 2021) visual encoder  $\mathcal{E}^{vis}$  to extract video features:

$$E_V = \mathcal{E}^{vis}(\Omega(\text{Video})). \quad (2)$$

A Spatial-Temporal Convolution (STC) connector (Cheng et al. 2024b) captures spatial and temporal dynamics:

$$E'_V = \text{STC}(E_V) = P_V(\text{R}_2(\text{Conv3D}(\text{R}_1(E_V)))), \quad (3)$$

where  $\text{STC}(\cdot)$  includes two spatial interaction modules ( $\text{R}_1$ ,  $\text{R}_2$ ) and a 3D convolution ( $\text{Conv3D}$ ), with  $P_V$  projecting features to the language model  $\phi$  space.

For audio, the BEATs model (Chen et al. 2023) serves as the audio encoder  $\mathcal{E}^{aud}$ , extracting features mapped to the language model space via a linear projector  $P_A$ :

$$E_A = \mathcal{E}^{aud}(\text{Audio}), \quad E'_A = P_A(E_A). \quad (4)$$

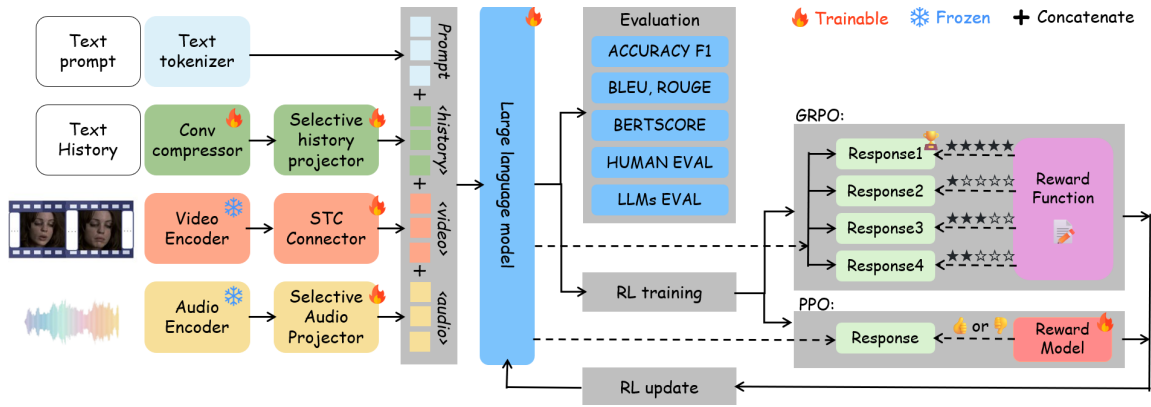


Figure 2: MULTIMOOD overview. Multimodal architecture that processes video, audio, text, and historical conversation data through dedicated encoders. The modality-specific embeddings are fused and passed into an LLM, which is further optimized using reinforcement learning guided by trustworthiness criteria to generate emotionally supportive and responsible responses.

**Cross-Modality Concatenation** Our approach draws from ECoT (Li et al. 2024), a plug-and-play method that boosts LLM performance in emotional generation tasks by aligning with Goleman’s emotional intelligence theory. We design a prompt template to guide the LLM in generating empathetic responses, integrating historical and real-time data. Multimodal features are concatenated into the input using specialized tokens:  $\langle \text{video} \rangle$ ,  $\langle \text{audio} \rangle$ , and  $\langle \text{history} \rangle$ , replaced by processed embeddings  $E'_V$ ,  $E'_A$ , and  $E'_H$ , respectively, forming the input sequence  $X_{LLM}$ . This attention-based fusion enables the model to dynamically prioritize relevant cues (e.g., tone, facial expressions) for safe, context-aware responses, while simultaneously predicting three classifications and generating therapist-like outputs (see Figure 2).

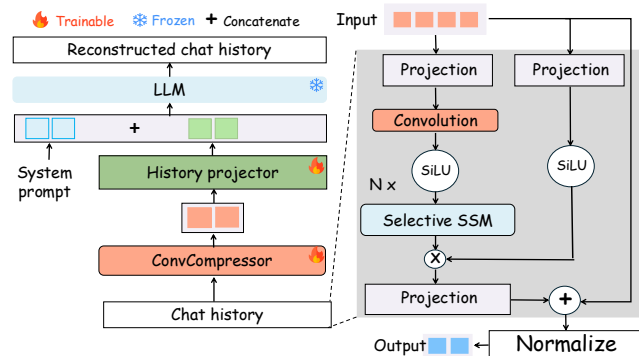


Figure 3: ConvCompressor architecture and pretraining.

**Conversation Compressor** Conversational emotional support systems require effective processing of extensive dialogue histories to deliver contextually appropriate responses. However, long conversation histories pose computational and memory challenges for language models. To address this, we propose Conversation Compressor (ConvCompressor), a lightweight module that distills dialogue histories into compact, semantically rich representations

while retaining critical information. ConvCompressor employs the Mamba state-space model (Gu and Dao 2023) as its core, offering linear computational complexity compared to the quadratic scaling of transformers. It appends a  $\langle \text{MEM} \rangle$  token to each conversational turn  $U_i$  in a history  $H = \text{concat}(\{U_i\}_{i=1}^T)$ , where  $U_i$  includes role information, utterance content, emotional labels, and therapist strategy labels, forming  $H' = U_1 \langle \text{MEM} \rangle U_2 \langle \text{MEM} \rangle \dots U_T \langle \text{MEM} \rangle$ . The Mamba backbone processes  $H'$  to generate hidden representations  $Z$ , from which we extract hidden states only at  $\langle \text{MEM} \rangle$  token positions. The extracted representations then undergo a trainable memory projector  $P_H$  before being fed to LLM (see Figure 3).

ConvCompressor is optimized through a two-stage training process. First, it is pre-trained with a frozen language model on a reconstruction task to regenerate the original dialogue history from compressed  $\langle \text{MEM} \rangle$  representations. Then, it undergoes end-to-end fine-tuning within the multimodal pipeline, adapting its compression strategy to significantly reduce the number of input tokens for the LLM while preserving a comparative overall performance.

### 3.3 Training

The training process comprises two key stages that enhance both robustness and trustworthiness.

**Stage 1: Supervised Fine-Tuning** We fine-tune our framework on the MESC dataset, leveraging multimodal data throughout the training process. To accommodate potential missing modalities (video or audio) during inference, we introduce a random modal selection mechanism. This is defined by a probability vector  $\mathbf{p} = [p_a, p_v, p_{av}]$ , representing the likelihoods of selecting audio or video or both modalities. This approach enhances the framework’s robustness by exposing it to all possible modality combinations during training. For multimodal processing, we employed SigLIP-So400M-Patch14 384 (Zhai et al. 2023) for video, and BEATs (Chen et al. 2023) for audio. The ConvCompressor is built on Mamba-370M (Gu and Dao 2023).

**Stage 2: Trustworthiness-Aware via Reinforcement**

### Trustworthiness Dimensions for Emotional Support

Dimension	Source	Definition
<b>Truthfulness</b>	(Huang et al. 2024) (Richmond et al. 2022) (Boyatzis, Goleman, and Rhee 2000)	The accurate representation of information, facts, and results by the AI system.
<b>Safety</b>	(Huang et al. 2024) (Boyatzis, Goleman, and Rhee 2000)	Promote safe, healthy conversations, avoiding harm, distress, or triggers while supporting user well-being.
<b>Fairness</b>	(Huang et al. 2024) (Richmond et al. 2022) (Boyatzis, Goleman, and Rhee 2000)	The quality of being impartial and equitable, considering multiple perspectives and maintaining a positive, action-oriented tone.
<b>Privacy</b>	(Huang et al. 2024) (Richmond et al. 2022)	Practices that safeguard human autonomy, identity, and data dignity.
<b>Empathy</b>	(Richmond et al. 2022) (Boyatzis, Goleman, and Rhee 2000)	Openness and honesty in expressing sympathy for negative situations or approval for positive ones.
<b>Reliability</b>	(Crits-Christoph et al. 2019) (Richmond et al. 2022) (Boyatzis, Goleman, and Rhee 2000)	Responses foster understanding, connection, and provide encouragement, comfort, or support.
<b>Ethical Guidance</b>	(Huang et al. 2024) (Boyatzis, Goleman, and Rhee 2000)	Ensuring AI behaviors guide emotional health responsibly, avoiding manipulation or harm.

### Inter-annotator Agreement (Top)

Dimension	Flu.	Ide.	Com.	Sug.	Ove.
<b>Fleiss Kappa</b>	0.65	0.61	0.60	0.61	0.67

### Human Evaluation (Middle)

	Flu.	Ide.	Com.	Sug.	Ove.
<b>Qwen2-7B</b>	22%	17%	17%	22%	21%
<b>MM (SFT)</b>	23%	30%	27%	25%	21%
<b>MM (SFT+RL)</b>	<b>55%</b>	<b>53%</b>	<b>56%</b>	<b>53%</b>	<b>58%</b>

### LLMs Evaluation (Bottom)

Model	Judge: GPT-4o							
	Tru.	Saf.	Fai.	Pri.	Emp.	Rel.	Eth.	Avg.
<b>Qwen2-7B</b>	6.0	4.2	5.1	8.0	4.2	4.6	4.3	5.2
<b>MM (SFT)</b>	6.2	4.3	5.8	7.8	4.3	4.8	4.9	5.4
<b>MM (SFT+RL)</b>	<b>7.0</b>	<b>6.3</b>	<b>6.8</b>	<b>8.8</b>	<b>6.2</b>	<b>6.4</b>	<b>6.3</b>	<b>6.8</b>

Model	Judge: Claude 4.0-Sonnet							
	Tru.	Saf.	Fai.	Pri.	Emp.	Rel.	Eth.	Avg.
<b>Qwen2-7B</b>	4.9	7.0	6.0	<b>8.0</b>	5.6	<b>7.0</b>	6.0	6.4
<b>MM (SFT)</b>	5.0	<b>7.0</b>	5.9	<b>8.0</b>	5.9	7.3	6.0	6.5
<b>MM (SFT+RL)</b>	<b>7.2</b>	<b>7.0</b>	<b>6.5</b>	<b>8.0</b>	<b>7.8</b>	<b>7.4</b>	<b>6.6</b>	<b>7.2</b>

Model	Judge: Grok-3							
	Tru.	Saf.	Fai.	Pri.	Emp.	Rel.	Eth.	Avg.
<b>Qwen2-7B</b>	6.1	6.0	7.2	7.5	6.0	6.3	5.8	6.4
<b>MM (SFT)</b>	6.2	6.6	7.6	7.2	6.1	6.9	5.9	6.6
<b>MM (SFT+RL)</b>	<b>7.3</b>	<b>7.8</b>	<b>8.7</b>	<b>9.3</b>	<b>7.5</b>	<b>7.9</b>	<b>7.5</b>	<b>8.0</b>

Table 2: Trustworthiness dimensions for emotional support tasks (left); Inter-annotator agreement, human and evaluation results (right) across different models. Flu., Ide., Com., Sug., Ove., stand for Fluency, Identification, Comfort, Suggestions, and Overall.

**Learning** Initial assessments showed that the system’s responses occasionally lacked the natural flow and trustworthiness needed for effective emotional support. To address this, we defined a set of trustworthiness criteria,  $\mathcal{C} = \{c_1, c_2, \dots, c_k\}$  (see Table 2-left), and employed reinforcement learning to align responses with ethical and therapeutic standards. We used Group Relative Policy Optimization (GRPO) (Shao et al. 2024) and Proximal Policy Optimization (PPO) (Schulman et al. 2017) after supervised fine-tuning to enhance response quality. GRPO optimizes by comparing responses within groups, while PPO stabilizes learning through clipped updates and a Kullback-Leibler divergence penalty. To guide learning, we designed a reward function combining trustworthiness and similarity. The similarity score  $r_{\text{sim}}(y, y^*)$  leverages BGE-M3 embeddings via ColBERT (Khattab and Zaharia 2020), integrating dense, sparse, and ColBERT-specific similarities with weights (1, 0.3, 1), normalized to  $[0, 1]$ :

$$r_{\text{sim}} = \text{scale}(\text{sim}_{\mathcal{C}} + 0.3 \text{sim}_{\text{s}} + \text{sim}_{\text{d}}) \quad (5)$$

Trustworthiness  $r_{\text{trust}}(y)$  is evaluated by GPT-4o per sentence, averaged and scaled to  $[0, 1]$ . GPT-4o is trusted for this task due to its proven capabilities in labeling data (Tan et al. 2024) across various tasks, as well as its robust safety mechanisms (Wei, Haghtalab, and Steinhardt 2023):

$$r_{\text{trust}}(y) = \text{scale}\left(\frac{1}{|y|} \sum_{i=1}^{|y|} \text{GPT-4o}_{\text{trust}}(y_i)\right) \quad (6)$$

The final reward is:

$$r(y|x) = \frac{1}{2}(r_{\text{trust}}(y) + r_{\text{sim}}(y, y^*)) \quad (7)$$

### 3.4 Trustworthiness Dimension Table

To assess response trustworthiness in emotional support, we first developed a domain-specific framework, *Trustworthy Dimensions*. This was built by synthesizing insights

from four key sources: the TrustLLM framework (Huang et al. 2024), which outlines trust principles for LLMs; patient-clinician trust studies (Crits-Christoph et al. 2019; Richmond et al. 2022); and Goleman’s emotional intelligence principles (Boyatzis, Goleman, and Rhee 2000). From TrustLLM, we adopted core technical values such as Truthfulness, Safety, Fairness, Privacy, and Machine Ethics. Clinical trust literature contributed Honesty, Communication, Confidentiality, Fidelity, and Reliability - emphasizing relational trust. Goleman’s work added Empathy and Social Skills, highlighting emotional resonance. These elements were distilled into seven core dimensions, carefully defined to balance technical reliability with emotional sensitivity (Table 2-left).

## 4 Experiments

### 4.1 Experimental Setup

**Metrics** For generation evaluation, we use BLEU-n (B-2), ROUGE-L (R-L), and BERTScore (BS) to evaluate the Therapist’s responses from the models. For classification of the MESC dataset (Chu et al. 2025), we use Accuracy and Weighted-F1 as metrics. These metrics collectively provide a comprehensive overview of model performance across different tasks. For the DFEW dataset (Jiang et al. 2020), we use unweighted average recall (UAR) and weighted average recall (WAR) to compare our method with SOTA methods.

**Baselines** We utilized the pretrained LLM from VideoLLaMA2 (Cheng et al. 2024b) as the multimodal LLM backbone, leveraging its training on multimodal data. We compare MULTIMOOD with API-based LLMs (GPT-4o (OpenAI 2023), Grok3 (xAI 2025), Claude-3.7 (Anthropic 2023), Deepseek-R1 (DeepSeek-AI 2025), LLaMA4 (MetaAI 2025)); Open-source VLMs (Qwen2 (Yang 2024) and Qwen2.5 (Bai 2025), EmotionLLaMA (Cheng et al. 2024a), VideoLLaMA2,

Method	Hap	Sad	Neu	Ang	Sur	Dis	Fea	UAR	WAR
IAL (Li et al. 2023)	87.95	67.21	70.10	76.06	62.22	0.00	26.44	55.71	69.24
VideoMAE (Tong et al. 2022)	93.09	78.78	71.75	78.74	33.44	17.93	41.46	63.60	74.60
S2D (Chen et al. 2025)	93.62	80.25	77.14	81.09	64.53	1.38	34.71	61.82	76.03
EmotionLLaMA (Cheng et al. 2024a)	93.05	79.42	72.47	84.14	72.79	3.45	44.20	64.21	77.06
MultiMood (ours)	<b>96.31</b>	<b>93.68</b>	<b>89.45</b>	<b>88.82</b>	<b>81.68</b>	<b>78.38</b>	<b>85.19</b>	<b>85.94</b>	<b>89.89</b>

Table 3: Comparison of multimodal emotion recognition results on DFEW.

VideoLLaMa3-7B (Cheng et al. 2024b), VideoLLaVA (Lin et al. 2024b) and InternVideo2.5-8B (Wang et al. 2025)) - finetuned on the MESC dataset (Chu et al. 2025) with SFT and PPO; VideoLLaMA2-72B and closed source models are not finetuned due to resource constraints; and SMES-leveraged models (Chu et al. 2025).

**Settings** Experiments were conducted on 4×H100 GPUs including LLMs training, multimodal projectors training, ConvCompressor training and RL training.

## 4.2 Results

In this section, we present a comprehensive evaluation to compare our framework with other frontier models on the MESC and DFEW datasets. The evaluation highlights the strengths and advancements of our framework in handling complex multimodal data in both tasks.

**Overall Performance** Tables 3 and 4 present the primary results of our proposed MultiMood framework compared to baseline models, evaluated across four MESC tasks (Chu et al. 2025) and one DFEW task (Jiang et al. 2020). MultiMood demonstrates exceptional adaptability, achieving robust performance across all tasks, unlike baseline models that often excel in specific domains. It delivers consistent results in emotion recognition, strategy prediction, system emotion prediction, and response generation, surpassing larger models like VideoLLaMA2-72B and specialized classification models like MMGCN. Notably, MultiMood achieves the highest average score (56.45) across the four MESC tasks and a SOTA score on the DFEW dataset. The ConvCompressor module demonstrates remarkable efficiency, achieving 98.6% token reduction while maintaining competitive performance, making our framework significantly more memory-efficient for processing extended dialogue histories. Our framework performance is evaluated from four key perspectives.

**Emotion Recognition:** Our MULTIMOOD framework achieves SOTA performance on the single-labeled DFEW dataset (Jiang et al. 2020), outperforming prior methods in accuracy, unweighted average recall and weighted average recall scores, as shown in Table 3. It achieves the highest UAR of 85.94% and WAR of 89.89%, excelling across all emotion categories, notably Disgust (78.38%), where prior models like IAL (Li et al. 2023), VideoMAE (Tong et al. 2022), S2D (Chen et al. 2025), and EmotionLLaMA (Cheng et al. 2024a) struggled due to under-representation (Jiang et al. 2020). With MESC, the variant utilizing GRPO attains the highest performance, followed closely by the finetuned framework without GRPO. MULTIMOOD surpasses

video understanding models (e.g., VideoLLaMA, InternVideo2.5), the Qwen family, and closed-source models, as well as specialized frameworks like SMES (Chu et al. 2025) and MMDFN (Hu et al. 2022) (shown in Table 4). MULTIMOOD’s robust classification, particularly for nuanced emotions, enhances empathetic response generation, establishing a new benchmark for precise emotion recognition. However, while ConvCompressor improves memory efficiency, it may compromise performance due to information loss.

**Strategy Prediction:** MULTIMOOD achieves a 42.81% accuracy on the Strategy Prediction task, slightly trailing BlenderBot SFT (48%) and SMES (49%) (Chu et al. 2025). This gap reflects MULTIMOOD’s design prioritizing robust, generalized performance across diverse tasks over specialization in strategy prediction. Nonetheless, it delivers a competitive F1 score, surpassing several baselines, though marginally behind SMES in accuracy. Unlike BlenderBot, which benefits from domain-specific retrieval tools, MULTIMOOD faces challenges with class imbalance. However, its instruction-guided framework excels in generating safe, multimodal-aware responses, enhancing generalizability across emotion recognition, strategy planning, and empathetic response generation.

**System Emotion Prediction:** Most fine-tuned models achieve over 90% accuracy in this task, attributed to a data skew where 90% of labels are Neutral. This imbalance is typical in emotional support contexts, as therapists maintain a calm demeanor, enables the system to generate honest, unbiased answers.

**Response Generation:** MULTIMOOD (SFT+GRPO) achieves superior performance across all metrics—BLEU-2 (6.18), ROUGE-L (17.86), and BERTScore (86.80)—demonstrating the efficacy of combining Group Relative Policy Optimization with supervised fine-tuning to produce fluent, contextually aligned responses. It outperforms baselines like VideoLLaMA2-7B (SFT) and Qwen2-7B (SFT + PPO), as well as closed-source models such as GPT-4o (OpenAI 2023) and LLaMA4 (MetaAI 2025), which underperform due to their reliance on textual features alone. MULTIMOOD’s integration of multimodal data enhances its classification and response generation capabilities, surpassing recent SOTA SMES (Chu et al. 2025) and setting a new benchmark for empathetic, high-quality responses.

**Human and LLM Evaluation** We conducted a comprehensive evaluation using both human and LLM assessments to assess the trustworthiness and quality of responses from Qwen2-7B (SFT), MultiMood-MM(SFT), and MultiMood-MM(SFT+RL). Four graduate students served as human annotators, all with expertise in emotional support research

Model	Training method	Modality	Task 1		Task 2		Task 3		Task 4		
			Acc	F1	Acc	F1	Acc	F1	B2	R-L	BScore
MMGCN	SFT	A,V,T	55.80	57.58	-	-	-	-	-	-	-
MMDFN	SFT	A,V,T	58.13	55.86	-	-	-	-	-	-	-
Blenderbot SFT	SFT	A,V,T	-	-	-	-	48.00	<b>46.10</b>	1.31	15.38	86.60
SMES	SFT	A,V,T	54.60	46.80	96.10	64.00	<b>49.00</b>	20.20	5.13	15.42	<b>86.80</b>
VideoLLaMA2-72B	-	A,V,T	55.06	55.68	97.36	98.10	25.77	26.09	3.55	13.77	85.37
VideoLLaVA	SFT	V,T	46.60	47.08	94.18	88.03	27.31	22.28	4.37	9.84	84.23
InternVideo2.5-8B	SFT	V,T	37.22	34.69	98.90	98.79	13.44	4.82	3.92	13.21	85.40
VideoLLaMA3-7B	SFT	A,V,T	45.28	46.23	97.40	72.66	33.96	24.50	3.70	11.55	85.07
EmotionLLaMA	SFT	A,V,T	46.12	41.95	<b>99.11</b>	<b>99.11</b>	37.44	25.41	2.55	10.76	84.28
LLaMA4-Maverick	-	T	23.34	21.16	68.72	81.02	14.53	8.11	3.94	10.03	84.26
Claude-3.7-Sonnet	-	T	32.59	33.33	85.90	91.80	27.97	27.55	2.25	8.45	83.79
Deepseek-R1	-	T	20.48	20.27	59.47	74.03	17.84	15.95	3.22	9.20	83.96
GPT-4o	-	T	38.98	43.56	72.46	83.60	24.88	26.26	2.30	9.20	84.31
Grok-2	-	T	22.46	25.08	65.85	78.80	20.44	18.19	2.29	9.60	84.61
Qwen2-7B	SFT	A,V,T	41.83	37.16	<b>99.33</b>	99.00	37.43	33.52	4.68	13.20	85.61
Qwen2-0.5B	SFT+PPO	A,V,T	44.27	44.99	<b>99.33</b>	99.00	36.34	33.01	4.40	12.31	85.36
Qwen2-7B	SFT+Comp.	A,V,T	44.27	44.03	<b>99.33</b>	99.00	39.42	35.69	4.60	12.90	85.47
Qwen2.5-7B	SFT	A,V,T	53.00	51.13	98.63	98.80	35.14	34.46	4.68	13.81	85.52
<b>MultiMood</b>	<b>SFT+Comp.</b>	A,V,T	53.75	51.75	<b>99.33</b>	99.00	39.29	36.25	5.26	15.34	85.81
<b>MultiMood</b>	<b>SFT</b>	A,V,T	56.38	55.81	99.11	<b>99.11</b>	36.78	34.32	4.58	13.47	85.71
<b>MultiMood</b>	<b>SFT+GRPO</b>	A,V,T	<b>58.60</b>	<b>57.78</b>	<b>99.33</b>	99.00	42.81	39.65	<b>6.18</b>	<b>17.86</b>	<b>86.80</b>
<b>MultiMood</b>	<b>SFT+Comp+GRPO</b>	A,V,T	55.94	55.33	99.11	<b>99.11</b>	38.10	34.58	5.42	15.83	86.00

Table 4: Benchmark of MULTIMOOD against other baselines on MESC. Task 1: User Emotion Recognition, Task 2: System Emotion Recognition, Task 3: Strategy Prediction, Task 4: Response Generation. A=Audio, V=Video, T=Text; B2=BLEU-2; R-L=ROUGE-L; BScore=BERTScore (F1); Comp.=Conversation Compressor.

and advanced English proficiency (IELTS overall  $\geq 7.0$  with reading  $\geq 7.5$ ) to ensure accurate evaluation of text-only outputs. They received training with tutorials and examples, including framework-generated outputs, dialogue contexts, situational details, and responses from a licensed psychologist, followed by a test on 100 MESC dataset validation samples to achieve a Cohen’s kappa inter-annotator agreement above 0.4 (Byrt 1996) (see Table 2-right-middle); retraining was required if unmet. During annotation, two annotators labeled all responses, with discrepancies resolved by a third and persistent disagreements settled by a fourth to establish the majority label, detailed results in Table 2-right-middle.

Human evaluation shows MultiMood (SFT+GRPO) outperforming in Fluency (55%), Comfort (56%), and Overall (58%), highlighting the effectiveness of multimodal fine-tuning and GRPO in enhancing response quality. Simultaneously, LLM evaluation, guided by (Tan et al. 2024), (Beaulieu-Jones et al. 2023), and (Reddy 2023), underscored LLMs’ near-human accuracy in surgical knowledge but noted query inconsistency, stressing stable evaluation needs. LLM scoring pre- and post-application of our trustworthiness dimension table 2-left revealed RL-incorporated frameworks significantly outperformed non-RL frameworks across three LLMs (see Table 2-right-bottom). By aligning with trustworthiness criteria, RL enhances safety, reliability, and ethical soundness, addressing non-RL inconsistencies and boosting utility for critical applications.

## 5 Limitation

Despite the demonstrated effectiveness of our framework, several limitations persist. It underperforms BlenderBot in

strategy prediction (per SMES (Chu et al. 2025)) due to class imbalance and the lack of external retrieval. The inability to fine-tune certain multimodal frameworks, constrained by resource limitations, weakens the robustness of our comparisons. Additionally, although its usage was proved (Tan et al. 2024), using GPT-4o for trustworthiness evaluation may introduce bias, particularly when its reward function influences reinforcement learning training. Furthermore, the experimental datasets, derived from movies and challenges rather than real treatment settings, lack authenticity—a common issue in this field (Kruse et al. 2016; Mudgal et al. 2022), underscoring the need for more realistic emotion support datasets in future research.

## 6 Conclusion

In conclusion, MULTIMOOD leverages multimodal techniques to achieve state-of-the-art results in emotion recognition and response generation, outperforming closed- and open-source models. Enhanced by reinforcement learning, it demonstrates high trustworthiness per human and LLM evaluations, with potential for therapeutic use. However, limitations in strategy prediction, hardware constraints and lack of realistic datasets suggest areas for future enhancement.

## Acknowledgments

This research was supported by the Mohamed bin Zayed University of Artificial Intelligence Travel Grant. The authors also gratefully acknowledge the VNUHCM-University of Information Technology’s Scientific Research Support Fund for their financial assistance. The authors also thank

the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting Duy M. H. Nguyen. Duy M. H. Nguyen and Daniel Sonntag are also supported by the No-IDLE project (BMFTR, 16IW23002), the MAS-TER project (EU, 101093079), and the Endowed Chair of Applied Artificial Intelligence, Oldenburg University. We also thank AISIA Research Lab for supporting us in this paper.

## References

- Anthropic. 2023. Title. <https://claude.ai>.
- Bai, S. 2025. Qwen2.5-VL Technical Report. *CoRR*.
- Beaulieu-Jones, B. R.; Shah, S.; Berrigan, M. T.; Marwaha, J. S.; Lai, S.-L.; and Brat, G. A. 2023. Evaluating Capabilities of Large Language Models: Performance of GPT4 on Surgical Knowledge Assessments. *medRxiv*.
- Beck, A. T.; and Weishaar, M. 1989. *Cognitive Therapy*.
- Boyatzis, R. E.; Goleman, D.; and Rhee, K. S. 2000. Clustering Competence in Emotional Intelligence: Insights from the Emotional Competence Inventory. In Bar-On, R.; and Parker, J. D. A., eds., *The Handbook of Emotional Intelligence: Theory, Development, Assessment, and Application at Home, School, and in the Workplace*, 343–362.
- Byrt, T. 1996. How good is that agreement? *Epidemiology*, 7(5): 561.
- Chen, S.; Wu, Y.; Wang, C.; Liu, S.; Tompkins, D.; Chen, Z.; Che, W.; Yu, X.; and Wei, F. 2023. BEATs: Audio Pre-Training with Acoustic Tokenizers. In *ICML 2023*, volume 202, 5178–5193.
- Chen, Y.; Li, J.; Shan, S.; Wang, M.; and Hong, R. 2025. From Static to Dynamic: Adapting Landmark-Aware Image Models for Facial Expression Recognition in Videos. *IEEE Trans. Affect. Comput.*, 16(2): 624–638.
- Cheng, Z.; Cheng, Z.; He, J.; Wang, K.; Lin, Y.; Lian, Z.; Peng, X.; and Hauptmann, A. G. 2024a. Emotion-LLaMA: Multimodal Emotion Recognition and Reasoning with Instruction Tuning. In *NeurIPS 2024*.
- Cheng, Z.; Leng, S.; Zhang, H.; Xin, Y.; Li, X.; Chen, G.; Zhu, Y.; Zhang, W.; Luo, Z.; Zhao, D.; and Bing, L. 2024b. VideoLLaMA 2: Advancing Spatial-Temporal Modeling and Audio Understanding in Video-LLMs. *CoRR*, abs/2406.07476.
- Chu, Y.; Liao, L.; Zhou, Z.; Ngo, C.-W.; and Hong, R. 2025. Towards Multimodal Emotional Support Conversation Systems. *IEEE Transactions on Multimedia*, 1–12.
- Crits-Christoph, P.; Rieger, A.; Gaines, A.; and Gibbons, M. B. C. 2019. Trust and Respect in the Patient-Clinician Relationship: Preliminary Development of a New Scale. *BMC Psychology*, 7(91).
- DeepSeek-AI. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *CoRR*, abs/2501.12948.
- Ekman, P. 2003. *Emotions Revealed: Recognizing Faces and Feelings to Improve Communication and Emotional Life*. Times Books/Henry Holt and Co.
- Esteva, A.; Robicquet, A.; Ramsundar, B.; Kuleshov, V.; DePristo, M.; Chou, K.; Cui, C.; Corrado, G.; Thrun, S.; and Dean, J. 2019. A Guide to Deep Learning in Healthcare. *Nature Medicine*, 25: 24–29.
- Fitzpatrick, K. K.; Darcy, A.; and Vierhile, M. 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 Ment Health*, 4(2): e19.
- Gu, A.; and Dao, T. 2023. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. *arXiv preprint arXiv:2312.00752*.
- Hu, D.; Hou, X.; Wei, L.; Jiang, L.; and Mo, Y. 2022. MMDFN: Multimodal Dynamic Fusion Network for Emotion Recognition in Conversations. In *ICASSP 2022*, 7037–7041.
- Huang, Y.; Sun, L.; Wang, H.; Wu, S.; Zhang, Q.; and et al. 2024. Position: TrustLLM: Trustworthiness in Large Language Models. In *ICML 2024*.
- Jiang, X.; Zong, Y.; Zheng, W.; Tang, C.; Xia, W.; Lu, C.; and Liu, J. 2020. DFEW: A Large-Scale Database for Recognizing Dynamic Facial Expressions in the Wild. In Chen, C. W.; Cucchiara, R.; Hua, X.; Qi, G.; Ricci, E.; Zhang, Z.; and Zimmermann, R., eds., *MM '20*, 2881–2889.
- Khattab, O.; and Zaharia, M. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In *SIGIR 2020*, 39–48. ACM.
- Kruse, C. S.; Goswamy, R.; Raval, Y.; and Marawi, S. 2016. Challenges and Opportunities of Big Data in Health Care: A Systematic Review. *JMIR Medical Informatics*, 4(4): e38.
- Le, H. M.; Luong, V. T.; and Luong, N. H. 2023. Data Augmentation with Large Language Models for Vietnamese Abstractive Text Summarization. In *MAPR 2023*, 1–6.
- Le, H. M.; Tien, D. N.; Duy, K. L.; Quang, T. N. D.; Toan, N. K.; Nguyen, T.; and Nguyen, B. T. 2025a. Fusionista: Fusion of 3-D Information of Video in Retrieval System. In *MMM 2025*, volume 15524, 278–285.
- Le, H. M.; et al. 2025b. Fustar: Divide and Conquer Query in Video Retrieval System. In *SOICT 2024*, volume 2353.
- Li, H.; Niu, H.; Zhu, Z.; and Zhao, F. 2023. Intensity-Aware Loss for Dynamic Facial Expression Recognition in the Wild. In Williams, B.; Chen, Y.; and Neville, J., eds., *AAAI 2023*, 67–75.
- Li, Z.; Chen, G.; Shao, R.; Xie, Y.; Jiang, D.; and Nie, L. 2024. Enhancing Emotional Generation Capability of Large Language Models via Emotional Chain-of-Thought. *arXiv:2401.06836*.
- Lin, B.; Ye, Y.; Zhu, B.; Cui, J.; Ning, M.; Jin, P.; and Yuan, L. 2024a. Video-LLaVA: Learning United Visual Representation by Alignment Before Projection. In *EMNLP 2024*, 5971–5984.
- Lin, B.; Ye, Y.; Zhu, B.; Cui, J.; Ning, M.; Jin, P.; and Yuan, L. 2024b. Video-LLaVA: Learning United Visual Representation by Alignment Before Projection. In Al-Onaizan, Y.; Bansal, M.; and Chen, Y.-N., eds., *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 5971–5984. Miami, Florida, USA: Association for Computational Linguistics.

- Liu, S.; Zheng, C.; Demasi, O.; Sabour, S.; Li, Y.; Yu, Z.; Jiang, Y.; and Huang, M. 2021. Towards Emotional Support Dialog Systems. In *ACL/IJCNLP 2021*, 3469–3483.
- Malgaroli, M.; Schultebrucks, K.; Myrick, K. J.; Loch, A. A.; Ospina-Pinillos, L.; Choudhury, T.; Kotov, R.; De Choudhury, M.; and Torous, J. 2025. Large language models for the mental health community: framework for translating code to care. *The Lancet Digital Health*.
- Marquez, P. V.; and Saxena, S. 2016. Making Mental Health a Global Priority. *Cerebrum: The Dana Forum on Brain Science*, 2016: cer–10–16.
- MetaAI. 2025. The Llama 4 Herd: The Beginning of a New Era of Natively Multimodal AI Innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>.
- Miner, A. S.; Milstein, A.; Schueller, S.; Hegde, R.; Mangurian, C.; and Linos, E. 2016. Smartphone-Based Conversational Agents and Responses to Questions About Mental Health, Interpersonal Violence, and Physical Health. *JAMA Internal Medicine*, 176(5): 619–625.
- Mitsui, K.; Mitsuda, K.; Wakatsuki, T.; Hono, Y.; and Sawada, K. 2024. Pslm: Parallel generation of text and speech with llms for low-latency spoken dialogue systems. *EMNLP Findings 2024*.
- Mudgal, S. K.; Agarwal, R.; Chaturvedi, J.; Gaur, R.; and Ranjan, N. 2022. Real-world application, challenges and implication of artificial intelligence in healthcare: an essay. *Pan African Medical Journal*, 43: 3.
- NAMI. 2023. Mental Health By the Numbers. <https://www.nami.org/about-mental-illness/mental-health-by-the-numbers/>. Accessed: 2025-11-11.
- Nguyen, D. M.; Diep, N. T.; Nguyen, T. Q.; Le, H.-B.; Nguyen, T.; Nguyen, T.; Nguyen, T.; Ho, N.; Xie, P.; Wattenhofer, R.; et al. 2024. Logra-med: Long context multi-graph alignment for medical vision-language model. *arXiv preprint arXiv:2410.02615*.
- OpenAI. 2023. GPT-4 Technical Report. *CoRR*, abs/2303.08774.
- Poria, S.; Hazarika, D.; Majumder, N.; Naik, G.; Cambria, E.; and Mihalcea, R. 2019. MELD: A Multimodal Multi-Party Dataset for Emotion Recognition in Conversations. In *ACL 2019*, 527–536.
- Radford, A.; Kim, J. W.; Hallacy, C.; and et al., A. R. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *ICML 2021*, volume 139, 8748–8763.
- Rashkin, H.; Smith, E. M.; Li, M.; and Boureau, Y.-L. 2019. Towards Empathetic Open-domain Conversation Models: A New Benchmark and Dataset. In *ACL 2019*, 5370–5381.
- Reddy, S. 2023. Evaluating large language models for use in healthcare: A framework for translational value assessment. *Informatics in Medicine Unlocked*, 41: 101304.
- Richmond, J.; Khodyakov, D.; Barber, C.; Maurer, M.; Scholle, S.; Pillemer, F.; Thakore, S.; Brown, J.; Federman, A.; and Shrank, W. 2022. Development and Validation of the Trust in My Doctor, Trust in Doctors in General, and Trust in the Health Care Team Scales. *Social Science & Medicine*, 298: 114827.
- Rogers, C. R. 1957. The Necessary and Sufficient Conditions of Therapeutic Personality Change. *Journal of Consulting Psychology*, 21(2): 95–103.
- Saffaryazdi, N.; Gunasekaran, T. S.; Laveys, K.; Broadbent, E.; and Billinghamurst, M. 2025. Empathetic Conversational Agents: Utilizing Neural and Physiological Signals for Enhanced Empathetic Interactions. *arXiv preprint arXiv:2501.08393*.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal Policy Optimization Algorithms. *CoRR*, abs/1707.06347.
- Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Zhang, M.; Li, Y. K.; Wu, Y.; and Guo, D. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *CoRR*, abs/2402.03300.
- Sheng, Y.; Yang, J.; Yang, L.; Shi, Y.; Hu, J.; and Jiang, W. 2023. Muffin: A Framework Toward Multi-Dimension AI Fairness by Uniting Off-the-Shelf Models. In *DAC 2023*, 1–6.
- Sim, K. Y. H.; Fortuno, K. T.; and Choo, K. T. W. 2024. Towards Understanding Emotions for Engaged Mental Health Conversations. In *DIS 2024*.
- Tan, Z.; Li, D.; Wang, S.; Beigi, A.; Jiang, B.; Bhattacharjee, A.; Karami, M.; Li, J.; Cheng, L.; and Liu, H. 2024. Large Language Models for Data Annotation and Synthesis: A Survey. In *EMNLP 2024*, 930–957.
- Tong, Z.; Song, Y.; Wang, J.; and Wang, L. 2022. VideoMAE: Masked Autoencoders are Data-Efficient Learners for Self-Supervised Video Pre-Training. In Koyejo, S.; Mohamed, S.; Agarwal, A.; Belgrave, D.; Cho, K.; and Oh, A., eds., *NeurIPS 2022*.
- Wang, Y.; Li, X.; Yan, Z.; He, Y.; Yu, J.; and et al., X. Z. 2025. InternVideo2.5: Empowering Video MLLMs with Long and Rich Context Modeling. *CoRR*, abs/2501.12386.
- Wei, A.; Haghtalab, N.; and Steinhardt, J. 2023. Jailbroken: How Does LLM Safety Training Fail? In *NeurIPS 2023*.
- WHO. 2022. Mental disorders. <https://www.who.int/news-room/fact-sheets/detail/mental-disorders>. Accessed: 2025-11-11.
- xAI. 2025. Grok-2 Beta Release. <https://x.ai/news/grok-2>.
- Yang, A. 2024. Qwen2 Technical Report. *CoRR*, abs/2407.10671.
- Yang, Z.; Ren, Z.; Wang, Y.; Peng, S.; Sun, H.; Zhu, X.; and Liao, X. 2024. Enhancing Empathetic Response Generation by Augmenting LLMs with Small-scale Empathetic Models. *CoRR*, abs/2402.11801.
- Zhai, X.; Mustafa, B.; Kolesnikov, A.; and Beyer, L. 2023. Sigmoid Loss for Language Image Pre-Training. In *ICCV 2023*, 11941–11952.