

Context-Emotion Aware Therapeutic Dialogue Generation: A Multi-component Reinforcement Learning Approach to Large Language Models for Mental Health Support

Eric Hua Qing Zhang and Julia Ive

University College London, UCL, United Kingdom
eric.zhang.24@ucl.ac.uk

Abstract. Mental health illness represents a substantial global socioeconomic burden, with COVID-19 further exacerbating accessibility challenges and driving increased demand for telehealth mental health support. While large language models (LLMs) offer promising solutions through 24/7 availability and non-judgmental interactions, pre-trained models often lack the contextual and emotional awareness necessary for appropriate therapeutic responses. This paper investigated the application of supervised fine-tuning (SFT) and reinforcement learning (RL) techniques to enhance GPT-2's capacity for therapeutic dialogue generation. The methodology restructured input formats to enable simultaneous processing of contextual information and emotional states alongside user input, employing a multi-component reward function that aligned model outputs with professional therapist responses and annotated emotions. Results demonstrated improvements through reinforcement learning over baseline GPT-2 across multiple evaluation metrics: BLEU (0.0111), ROUGE-1 (0.1397), ROUGE-2 (0.0213), ROUGE-L (0.1317), and METEOR (0.0581). LLM evaluation confirmed high contextual relevance and professionalism, while reinforcement learning achieved 99.34% emotion accuracy compared to 66.96% for baseline GPT-2. These findings demonstrate reinforcement learning's effectiveness in developing therapeutic dialogue systems that can serve as valuable assistive tools for therapists while maintaining essential human clinical oversight. The code and appendices are publicly available at: <https://github.com/ez-anthro-tech-design/RLforGPT2MentalHealth>

Keywords: Reinforcement Learning, Supervised Fine-tuning, Large Language Models

1 Introduction

Mental illness imposes a significant socioeconomic burden, with one in five adults in the United States experiencing a mental health condition each year [1, p.]. Patients are further deterred from seeking care due to societal stigma [2], attitudinal barriers [3], and high costs and lack of knowledge about accessing these services [4].

Direct patient-facing AI systems without clinical supervision face significant obstacles, as patients demonstrate lack of trust and lower adherence towards autonomous non-human consultations [5]. Concerns exist regarding patient safety and accuracy when AI-assisted dialogue tools operate without clinical oversight [6]. Limited laws governing ethical use, data privacy, and transparent communication in autonomous AI systems compound these issues [7]. Additionally, LLMs exhibit hallucination and may be pretrained on datasets of unknown validation status, creating liability uncertainties on AI-generated clinical guidance as standard-of-care evidence despite absence of expert review protocols [8]. Given these constraints, LLMs may be most effectively deployed as supervised dialogue generation tools that support healthcare professionals across various therapeutic contexts, including telehealth and in-person sessions, by providing real-time conversational suggestions rather than replacing clinical judgment.

LLMs arise from advancements in natural language processing (NLP) using statistical methods to model word sequences [9]. Studies have demonstrated LLMs' potential to support healthcare professionals in clinical settings by providing consistent, non-judgmental dialogue suggestions during patient interactions [10]. These systems offer conversational guidance when practitioners need alternative phrasing options [10]. Evidence indicates AI-assisted dialogue generation tools reduce communication barriers, with patients showing increased comfort discussing sensitive topics when practitioners use supportive frameworks [11] [12].

Beyond these deployment challenges, additional technical limitations exist in LLMs for domain-specific tasks, particularly in capturing user emotions and comprehending context-specific cues [13]. Two methods stand out as fine-tuning strategies to enhance LLMs' capabilities in domain-specific settings: supervised fine-tuning and reinforcement learning. SFT operates on a straightforward supervised learning paradigm, where pre-trained models learn from input-output pairs through standard cross-entropy loss minimisation by updating model parameters to minimise prediction errors on labelled datasets [14] [15]. Yet, Stiennon et al. [16] identified three fundamental SFT limitations: (1) equally penalising critical errors and minor lexical variations, (2) reinforcing biases in low-quality training data, and (3) producing repetitive outputs. These SFT limitations highlight a broader challenge: whilst SFT improves task-specific performance [17], it often struggles with aligning to human preferences [18]. Luo et al. [19] show that whilst LLM responses are rated as more empathetic than human responses in 72.85% of cases, patients generally prefer human responses, indicating that genuine empathy is deeply valued in clinical contexts.

In contrast, RL offers a robust approach to align model outputs with human preferences through human preference collection, reward model learning, and policy optimisation [20]. RL learns through interactions with an environment, receiving reward signals that guide policy improvement through techniques like proximal policy optimisation (PPO) [18]. RL has shown success in state-of-the-art LLMs, demonstrating that reward-driven optimisation can enhance reasoning capabilities to meet user preferences for transparent explanations.

In this paper, we present three key contributions:

- 1) a novel multi-component reward function that evaluates emotional appropriateness and contextual relevance in therapeutic dialogue

- 2) a reinforcement learning optimisation framework that enhances SFT models for affective empathy, and
- 3) a practical implementation using GPT-2 that enables local deployment in clinical settings.

GPT-2, as an open-source model, offers low resource requirements without API access, and with relatively low parameters (124M - 1.5B), can be accomplished practically with commercially available computers for local deployment and affordable fine-tuning [34]. This integrated approach enables the model to generate transparent examples of contextual and emotionally aligned responses with explicit emotion states, providing therapists with clear demonstrations of ideal therapeutic communication for use across supervised clinical settings.

Results demonstrated substantial improvements through reinforcement learning over baseline GPT-2 across multiple evaluation metrics: BLEU improved by 428% (0.0021 to 0.0111), ROUGE-1 by 192% (0.0478 to 0.1397), ROUGE-2 by 233% (0.0064 to 0.0213), ROUGE-L by 223% (0.0408 to 0.1317). LLM-as-a-judge confirmed high contextual relevance and professionalism, whilst reinforcement learning achieved 99.34% emotion accuracy, representing a 48% improvement over baseline GPT-2's 66.96%.

2 Related Work

2.1 AI-Assisted Therapeutic Dialogue Tools in Clinical Practice

The integration of AI-assisted dialogue tools in mental health services has emerged as a promising approach to enhance therapeutic communication while maintaining necessary clinical supervision. These tools are designed to support healthcare professionals across various therapeutic modalities, including telehealth and in-person sessions, by providing real-time dialogue suggestions and conversational guidance.

Telehealth services have demonstrated effectiveness as scalable and cost-effective alternatives to traditional face-to-face mental health services, helping to overcome barriers such as social stigma, limited accessibility, and high costs [21] [22]. During the pandemic, the American Psychological Association [23] found that 76% of respondents were exclusively delivering services via telephone, digital platforms, or video conferencing. Similarly, text-based therapeutic interventions have also shown promise, with email-based exercises incorporating positive psychology interventions demonstrating improvements in stress levels and psychological well-being [24], indicating that digital therapeutic communication can effectively support psychological health.

However, telehealth services face several challenges including unstable internet connections, confidentiality and privacy concerns over online transmission, and video conferencing quality issues [25]. AI-assisted dialogue generation tools, when integrated into supervised therapeutic settings, offer potential solutions to enhance communication quality while addressing these limitations. These tools can support therapists by providing contextually appropriate dialogue suggestions in real-time, whether during telehealth sessions or in-person consultations.

Studies have demonstrated the effectiveness of LM-backed therapeutic support tools, which are valued for their non-judgmental nature and ability to provide structured conversational frameworks that enhance patient comfort [10]. When used under clinical supervision, these tools can offer consistent dialogue suggestions and emotional support frameworks during periods when therapists need alternative phrasing options or guidance on addressing sensitive topics [10]. The key distinction lies in positioning these tools as augmentative rather than replacement technologies, ensuring that clinical judgment and human empathy remain central to therapeutic practice while leveraging AI capabilities to enhance communication effectiveness.

2.2 Fine-tuning LLMs for Mental Health Support

Much of the literature exploring the use of pre-trained LLMs for mental health support (MHS) relies heavily on prompt engineering approaches. For example, Gabriel et al. [26] found that GPT-4 can generate responses with nearly 45% greater overall empathy than human peer-to-peer responses.

Beyond prompt engineering, supervised fine-tuning approaches have shown significant promise [27]. Jin et al. [28] demonstrated that fine-tuned models outperform pre-trained models on emotional support tasks, yet argued that fine-tuned models struggle with human-like emotional support and risk producing empty, repetitive outputs. These findings highlight the need for more sophisticated training approaches that go beyond traditional supervised learning.

Reinforcement learning has emerged as a promising complementary approach to address these limitations. Sharma et al. [29] applied reinforcement learning to DialoGPT, a large dialogue generation model based on GPT-2, for empathic rewriting, aiming to transform low-empathy conversational posts to higher empathy responses. While Sharma et al.'s [29] work addressed different objectives, it demonstrated the feasibility of reinforcement learning in improving response generation, achieving a BLEU score of 0.1391 in their Partner model. Similarly, Saha et al. [30] employed RL to fine-tune GPT-2 [31] for empathetic text generation using a combined reward function including BLEU, ROUGE-L, and sentiment scores. However, these studies primarily focused on general empathy metrics rather than clinical-specific emotional appropriateness and contextual relevance, indicating a gap that our work aims to address.

3 Dataset

3.1 Training Dataset

Table 1. Dialogue characteristics of the MESC dataset

	Total Records	Unique Problem Types	Min. Dialogue Turns	Max. Dialogue Turns	Avg. Dialogue Turns
Train	815	15	8	99	28.4
Val	102	15	11	47	26.6
Test	102	15	9	57	28.6

The Multimodal Emotional Support Conversation (MESC) dataset [15] draws from realistic therapy sessions adapted from "In Treatment". Professional actors portray therapists and patients in conversations capturing spoken dialogue, facial expressions, and body language. Mental health experts at Shanghai Mental Health Centre validated these performances as accurate representations of genuine therapeutic interactions [15]. Each dialogue contains multiple exchanges between psychotherapist and patient. This multimodal dataset addresses limitations of previous text-based emotional support datasets by capturing facial expressions and emotional cues alongside spoken dialogue (ibid.).

Chu et al. [15] performed standard train/test/validation split on the dataset encompassing 15 scenarios and 7 emotion categories. Emotional labels were initially classified using GPT-3.5, then manually calibrated by three trained graduate students specialising in emotional support research.

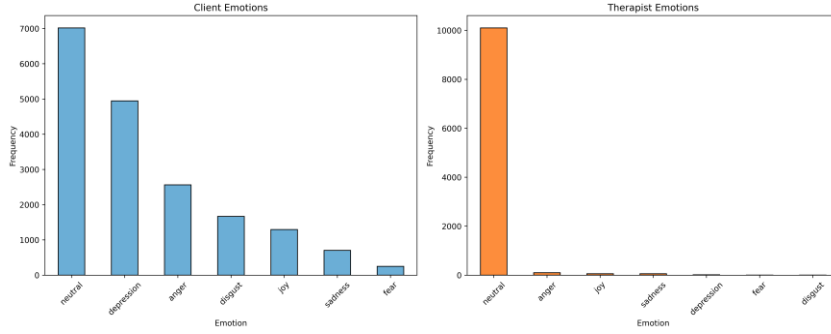


Fig. 1. Distribution of emotion categories by speaker (MESC dataset)

3.2 Evaluation Dataset

This paper will evaluate the fine-tuned models on the ESConv dataset, an emotional support conversation dataset containing annotated dialogues between help-seekers and supporters, with rich metadata including emotion types, problem categories, and support strategies [32].

Table 2. Dialogue characteristics of the ESConv dataset

Metric	Value
Total conversations	1300
Minimum dialogue turns	16
Maximum dialogue turns	120
Mean dialogue turns	29.51
Median dialogue turns	27.0
Mean words per turn	16.4

3.3 Data Preprocessing

During preprocessing, text extraction isolated authentic speech from metadata descriptions using regular expressions to remove patterns beginning with "The speaker" or "The emotion state". Utterances were then concatenated before change in speakers.

Sequence filtering

Sequence filtering preserves computational efficiency during SFT and RL by preventing excessive padding tokens that consume memory and slow training. Analysis revealed a right-skewed distribution with mean length 52.3 tokens exceeding median of 41 tokens. A 128-token threshold achieves ~95% data coverage whilst maintaining efficiency.

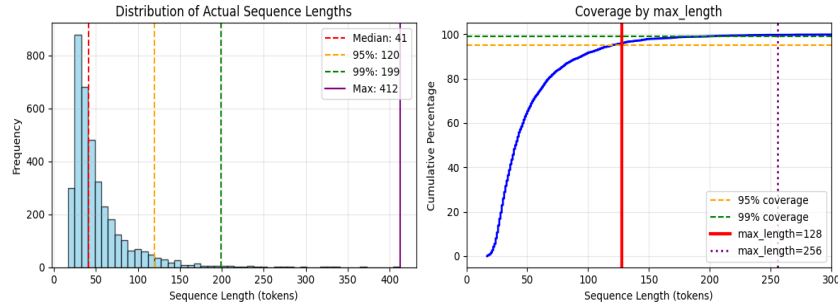


Fig. 2. Figure 2: Sequence Length Distribution and Coverage Analysis for 128-Token Threshold

Tokenizer Setup & Vocabulary Extension

This paper extended the original GPT-2 tokenizer (~50,257 tokens) with custom special tokens, including structural markers (<bos>, <eos>, <pad>, <problem>, <user>, <user_emotion>, <therapist>, <therapist_emotion>) and seven emotion tokens (["anger", "sadness", "depression", "disgust", "fear", "joy", "neutral"]). This extension serves critical functions: structural tokens enable delineation of dialogue components, allowing the model to distinguish between patient context, utterances, and emotional states. Emotion tokens are treated as atomic units to preserve semantic integrity - otherwise, subword tokenisation might fragment "sadness" into ["sad", "ness"], degrading coherent emotional representations. Crucially, it allows explicit generation of emotion tokens.

Table 3. Sequence Construction

<bos><problem>problem_text<user>user_text<user_emotion>emotion <therapist>therapist_text<therapist_emotion>emotion<eos>
--

Label Creation

To train the model for response prediction, a selective labelling strategy masked context tokens (<problem>, <user>, <user_emotion>) and content with label value -100, excluding them from loss calculations whilst preserving input context. This methodology employed HuggingFace's Transformers library, where mask token -100 creates an 'ignore index' without contributing to gradient loss. The response section employed differentiated labelling: whilst <therapist> markers were masked (-100), therapist text and <therapist_emotion> tokens received standard labels for prediction. This enabled response generation including emotional state prediction, allowing generation of content and emotional classification during inference. This contrasted with Wang et al's work [33], who used label masking causing excessive reliance on label embedding rather than contextual information. The Transformers implementation reverses this by predicting only unmasked response content, ensuring learning from contextual understanding rather than superficial patterns.

4 Methodology

4.1 Experiment Setup

All experiments were conducted on Google Colab using NVIDIA A100 GPU (40GB RAM). The baseline GPT-2 model (124M parameters) was imported via HuggingFace Transformers. Reinforcement learning used TRL 0.11.3 with Proximal Policy Optimization (PPO). Training was monitored using TensorBoard for SFT and Weights & Biases for RL experiments

4.2 Baseline GPT-2

A retokeniser was implemented on the baseline GPT-2 model to adapt the input structure to accommodate Table 3's sequence construction.

4.3 Supervised Fine-tuning (With Therapist Emotions)

SFT training employed batches of 32 sequences, each standardised to 128 tokens. During forward propagation, the model processed input sequences through GPT-2's transformer architecture to generate probability distributions over the vocabulary. Attention masks distinguished real tokens from padding, ensuring self-attention focused on meaningful content whilst ignoring padded positions.

Table 4. Supervised fine-tuning hyperparameters

Parameter	Value	Description
Optimizer	AdamW	Adaptive optimizer with weight decay
Learning rate	2e-5	Base learning rate for parameter updates
Batch size	32	Number of training samples processed simultaneously per iteration

Warmup ratio	0.1	Proportion of training steps for learning rate warmup
Training epochs	20	Total number of training epochs
Loss function	Cross-entropy	Computed selectively where labels \neq -100
Early stopping observation	Epoch 10	Diminishing returns observed beyond this point
Validation frequency	Every epoch	Model evaluation after each training epoch

4.4 Supervised Fine-tuning (No therapist emotions)

A separate SFT round was conducted without therapist emotion token generation as an ablation study baseline. This variant used identical training procedures but excluded *< therapist_emotion > emotion_word sequences* from target generation, producing only therapist text without emotional classification. This model helped quantify the joint text-emotion generation task's contribution to dialogue quality. Same hyperparameters were used in Table 4.

4.5 Supervised Fine-tuning + Reinforcement Learning

This paper adapted the standard RL framework consisting of contexts C , response generation actions G , a policy π , and multi-dimensional rewards R . In this framework, given a context $c \in C$, the model generated a response $g \in G$ according to the policy $\pi : C \rightarrow P(G)$, where $P(G)$ represented the probability distribution over possible responses. RL was implemented on top of SFT according to Stiennon et al. 's [16] classical RL approach that utilised SFT as a stable starting policy and established basic task actions.

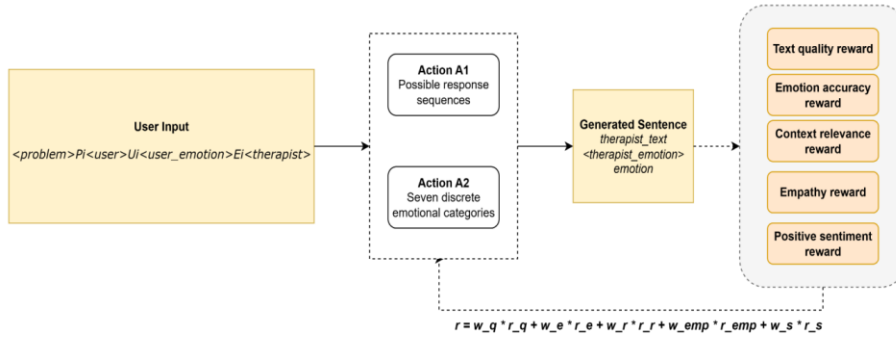


Fig. 3. The overall architectural representation of reinforcement learning framework

By employing a KL divergence penalty in RL, this approach prevented the RL policy from drifting too far from the SFT policy. Additionally, pre-trained token distributions may skew policy updates towards frequent words regardless of rewards. To mitigate this, this paper integrated reward scaling to a range of -10 to 10 into the reward function [34]. A learning rate of 1×10^{-6} , training epochs of 8, batch size of 16, and data shuffling disabled were employed.

Proximal Policy Optimization (PPO)

PPO is a policy gradient method addressing training instability in reinforcement learning by constraining policy updates through a clipped objective function. The algorithm enables multiple epochs of minibatch updates whilst preventing excessive parameter changes that could destabilise training [35]. This stability was crucial for this study, where catastrophic forgetting posed risks to response patterns acquired during SFT.

The PPO objective optimised the clipped surrogate loss:

$$L^{CLIP(\theta)} = \hat{E}_t [\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t)] \quad (1)$$

where

$$r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \quad (2)$$

represented the probability ratio between current and previous policies, \hat{A}_t denoted the advantage estimate that measures how much better an action is compared to the expected value at that state, and ϵ was the clipping parameter that constrained policy updates within a trust region.

Context

The context enabled response generation specific to the patient's presenting issues, utterances, and emotional states at each dialogue turn. The context space consisted of structured user input processed through retokenisation. At each generation step, the context for each user utterance was represented by C_i , containing: problem type P_i , user utterance U_i , and user emotional state E_i , denoted by:

$$C_i = \langle \text{problem} \rangle P_i \langle \text{user} \rangle U_i \langle \text{user}_{emotion} \rangle E_i \quad (3)$$

The C_i is followed by the response prompt $\langle \text{therapist} \rangle$, creating the complete input representation of context as $C_i \langle \text{therapist} \rangle$

Action

The model executed actions through structured response generation, producing both therapist textual responses and corresponding emotional classifications. The action aimed to generate contextually appropriate responses recognising context C_i whilst aligning with appropriate emotional classifications.

The model executed two generation components: $a1$ produced response text, and $a2$ predicted appropriate therapist emotional stance. Action space $A1$ consisted of all possible response sequences within vocabulary constraints. Action space $A2$ encompassed seven discrete emotional categories existing within the dataset: {anger, sadness, depression, disgust, fear, joy, neutral}. The complete action was denoted as $a = (a1, a2) \in A1 * A2$. Following C_i , structured action a generated responses in the sequence $\langle therapist_{emotion} \rangle emotion_{word} \langle eos \rangle$.

Policy

The policy operated by encoding structured input context C_i , then generating text by computing probability distributions over vocabulary tokens through decoder layers that maximised $P(response|context)$. The model continued generation beyond text to produce the emotional alignment sequence. Generation employed sampling parameters (top-p=1.0, top-k=0.0) recommended by TRL documentation to ensure response stability, with sequence completion triggered by the end-of-sequence token.

Rewards

This paper designed a novel composite reward function enhancing response generation through complementary metrics: text quality, emotional alignment, contextual relevance, empathetic content, and sentiment appropriateness. The ensemble approach addressed reward hacking, e.g. pure symbols/punctuation marks for high neutral emotion, where models exploit individual reward functions by generating responses that achieve high scores through rule-following whilst producing clinically meaningless content [36]. This ensemble ensured high-reward responses demonstrated competence across all dimensions simultaneously.

Text Fluency

The text quality reward function $f_{quality}(\cdot)$ measured response coherence and appropriateness.

$$r_q = f_{quality}(therapist_{response}) \quad (4)$$

Implementation incorporated repetition detection and meaningless pattern identification, penalising responses with excessive n-gram repetition and incoherent content.

Emotional Alignment

Emotional alignment rewarded suitable emotional classifications and penalised misalignment. The reward r_e evaluated emotion prediction accuracy and emotion token presence:

$$r_e = f_{emotion}(predicted_{emotion}, target_{emotion}, has_emotion_token) \quad (5)$$

This rewarded exact emotional matches and appropriate categories (positive, negative, neutral), whilst penalising conflicts between incompatible categories.

Contextual Relevance

Contextual relevance ensured responses addressed specific patient concerns rather than generic interventions. Semantic similarity measured contextual appropriateness:

$$r_r = \cos_sim(embed(therapist_response), embed(user_input)) \quad (6)$$

Semantic embeddings were computed using Sentence-BERT (all-MiniLM-L6-v2) (Appendix 1) to measure similarity between responses and user inputs. This model was chosen for its ability to vectorise text for semantic comparison.

Empathetic Content

A pre-trained empathy classification model (paragon-analytics/bert_empathy) (Appendix 1) evaluated empathetic quality of generated responses:

$$r_{emp} = f_{emp}(therapist_response) \quad (7)$$

The empathy model was chosen for its ability to predict empathic language content. This assessed empathetic quality, rewarding responses demonstrating understanding and emotional support.

Sentiment Appropriateness

Sentiment ensured responses conveyed appropriately positive sentiment, distinct from the discrete emotional categories used elsewhere. While emotion classification captured specific emotional states, sentiment analysis measured the overall positive/negative nature of the response

$$r_s = f_{sent}(therapist_response) \quad (8)$$

where $f_{sent}(\cdot)$ employed DistilBERT sentiment classification (distilbert-base-uncased-finetuned-sst-2-english) (Appendix 1).

Total Reward

The final composite reward function combined all components with optimised weighting:

$$r_{total} = w_q \times r_q + w_e \times r_e + w_r \times r_r + w_{emp} \times r_{emp} + w_s \times r_s \quad (9)$$

where text quality ($w_q = 1.1$) and emotional alignment ($w_e = 1.2$) received higher weighting, contextual understanding ($w_r = 1.1$) was emphasised, whilst supportive content components ($w_{emp} = 0.7$, $w_s = 0.7$) received moderate weighting to balance clinical appropriateness with empathetic responsiveness.

Weight selection was based on preliminary experiments evaluating response quality across different configurations. The weights do not sum to 1.0 (total = 4.8) as this approach prioritised absolute reward magnitudes rather than relative proportions, where stronger gradients were needed to overcome the conservative learning rate (1×10^{-6}) used for stability.

4.6 Hyperparameter Tuning

Hyperparameter tuning was conducted for top_p, top_k, and temperature via grid search across SFT (No emotion), SFT (With emotion), and RL-SFT (With emotion) models to optimise combined BLEU and ROUGE-1-2-L scores (see Appendix 2 for grid search results).

4.7 Evaluation

Evaluation was conducted using automated metrics, LLM-as-a-judge evaluation and human evaluation on two datasets: the MESC held-out testing dataset and the ESConv dataset [32].

ESConv contains richly annotated, high-quality conversations providing examples of effective emotional support. It is formatted in JSON and annotated with problem types and situation descriptions. While emotion labels were not originally annotated, this paper employed the roberta-base-go_emotions model (Appendix 1) to classify utterance emotions and map these into the seven emotion categories present in MESC for emotion classification evaluation.

- 1) BLEU: precision against reference responses [37]
- 2) ROUGE-1-2-L: N-gram overlap between generated and reference responses [38]
- 3) METEOR: Automatic metric for machine translation evaluation based on generalised unigram matching [39]

LLM-as-a-judge assessed baseline GPT-2, SFT (With emotion), and RL-SFT (With emotion) models using GPT-4o-mini via OpenAI API with prompt engineering (Appendix 3) to evaluate 100 samples from the testing dataset on a 1-5 scale across six emotional support criteria:

- 1) Therapeutic Rapport - Emotional connection and trust-building
- 2) Active Understanding - Comprehension and reflection of user's state
- 3) Relevance Focus - Addressing specific problems without going off-topic
- 4) Practical Helpfulness - Providing actionable guidance and direction
- 5) Professional Appropriateness - Maintaining boundaries
- 6) Emotional Validation - Acknowledging and validating user emotions

These criteria were adapted from Yuan et al.'s [40] human evaluation metrics to create a comprehensive framework tailored for therapeutic communication assessment. Human validation was conducted through manual inspection of a subset of GPT-4o-mini evaluations to ensure scoring consistency and appropriateness of the automated judgments.

5 Results

5.1 Automated Evaluation

Automated evaluation results (Table 5) from baseline GPT-2 to the reinforcement learning fine-tuned model demonstrate substantial improvements across BLEU, ROUGE-1, ROUGE-2, ROUGE-L and METEOR measures. Baseline GPT-2 achieved minimal scores across BLEU (0.0021) and ROUGE metrics (0.0461, 0.0046, 0.0383), indicating poor alignment with reference responses. Interestingly, baseline GPT-2 achieved the highest METEOR score (0.0670) across all models, likely due to METEOR's emphasis on synonym and stem matching, which may favour more diverse vocabulary generation over precise alignment with references [39].

Table 5. Automated evaluation quality metrics

Model	BLEU	ROUGE-1	ROUGE-2	ROUGE-L	METEOR
GPT-2	0.0021	0.0478	0.0064	0.0408	0.0691
SFT (No therapist emotions)	0.0106	0.1160	0.0159	0.1058	0.0433
SFT (Therapist emotions)	0.0108	0.1164	0.0175	0.1076	0.0485
Reinforcement Learning - SFT	0.0111	0.1397	0.0213	0.1317	0.0581

Both SFT variants demonstrated improvements over baseline across BLEU and ROUGE metrics, though with reduced METEOR scores. The emotion-inclusive SFT model outperformed the no-emotion variant across most metrics, with notable improvements in ROUGE-2 (+10.1%) and ROUGE-L (+1.7%).

The RL-aligned model achieved the strongest overall performance, demonstrating substantial gains in ROUGE-1 (+20.0% over emotion SFT) and ROUGE-L (+22.4% over emotion SFT), suggesting enhanced lexical overlap with reference responses.

Table 6. Automated evaluation quality metrics (ESConv dataset)

Model	BLEU	ROUGE-1	ROUGE-2	ROUGE-L	METEOR
SFT (Therapist emotions)	0.0050	0.1046	0.0089	0.0925	0.0354
Reinforcement Learning - SFT	0.0042	0.1110	0.0092	0.1002	0.0290

Automated evaluation on the ESConv dataset assessed model generalisation beyond the training domain. The SFT emotion-aware model achieved moderate cross-dataset performance. The RL-aligned model demonstrated mixed results, with improvements in ROUGE metrics (+6.1% ROUGE-1, +3.4% ROUGE-2, +8.3% ROUGE-L)

compared to emotion-inclusive SFT, but decreased performance in BLEU (-16%) and METEOR (-18.1%).

Emotion classification accuracy revealed progressive improvements across model variants. Baseline GPT-2 achieved 66.96% accuracy, whilst the no-emotion SFT model reached 80.84%. The emotion-aware SFT model achieved 95.37%, demonstrating the effectiveness of joint text-emotion training. The RL model attained 99.34% accuracy, representing a 4.16% improvement over emotion-aware SFT. This high performance demonstrates that the composite reward function effectively enhanced emotional alignment capabilities.

Table 7. Emotion token generation accuracy performance

Emotion Accuracy	Model
0.6696	GPT-2
0.8084	SFT (No therapist emotions)
0.9229	SFT (Therapist emotions)
0.9934	Reinforcement Learning

5.2 LLM-as-a-judge Evaluation

The LLM-as-a-judge evaluation (GPT-4o-mini) revealed a clear performance hierarchy across all models tested on 100 dialogue samples. The RL model achieved the highest overall average score of 1.718, followed by SFT at 1.655, and baseline GPT-2 at 1.062.

The RL model demonstrated notable superiority in relevance focus (1.830 vs SFT's 1.740), representing 5.2% improvement. This aligns with the contextual relevance component in the composite reward function, suggesting successful optimisation of context-aware response generation. In active understanding, the RL model achieved 1.670 compared to SFT's 1.570 (6.4% improvement), indicating enhanced comprehension of patient concerns. In emotional validation, RL achieved 1.640 versus SFT's 1.530 (7.2% improvement), validating the emotion-focused reward components. However, both advanced models showed similar performance in practical helpfulness (RL: 1.490, SFT: 1.480), with only marginal improvement over baseline GPT-2 (1.020).

Table 8. LLM-as-a-judge evaluation results

Metric	Baseline GPT-2	SFT Model	RL Model
Therapeutic Rapport	1.090	1.700	1.730
Active Understanding	1.030	1.570	1.670

Relevance Focus	1.070	1.740	1.830
Practical Helpfulness	1.020	1.480	1.490
Professional Appropriateness	1.100	1.910	1.950
Emotional Validation	1.060	1.530	1.640

5.3 Human Evaluation

Manual inspection of 50 samples (Appendix 6) revealed that LLM-as-a-judge had underscored around 20 samples. Human evaluation showed that the RL model frequently generated question-oriented responses rather than providing direct therapeutic advice. For instance, when users express intense emotions such as depression regarding bereavement, the generated responses employ open-ended questioning techniques that encourage further exploration without emotional reactivity. However, this likely explains the relatively low ROUGE scores, as reference responses were derived from conversations characterised by natural turn-taking, where not every utterance tends to be a question. Consequently, there was lower lexical overlap with reference responses that contained more declarative statements.

Table 9. Samples of generated responses from the SFT-RL model

Problem Type	User Text	User emotion	Generated Therapist Response	Generated Therapist Emotion
Therapeutic Relationship	Did I offend you?	neutral	No, you didn't offend me. You didn't say that you were upset with me.	neutral
Bereavement	I can't stop thinking about my grandmother who passed away	depression	And you wanna tell me what you wanna do for her?	neutral
Post-Traumatic Stress Disorder	I keep having flashbacks from the accident	fear	Do you want to talk about it?	neutral

Nevertheless, generated responses demonstrated strong contextual relatedness, with every response directly addressing specific keywords in patient utterances rather than providing generic statements. The neutral emotions generated benefit therapeutic

rapport. According to Gelso and Kanninen [41], effective therapists maintain neutrality to remain close to the patient's "zone of proximal development," avoiding manipulation of patient emotional expression through excessive warmth. The RL model's tendency towards neutral emotional stances and seemingly low Emotional Validation scores reflects this professional therapeutic approach.

6 Discussion

This paper presents a proof-of-concept for developing cost-effective, locally deployable therapeutic dialogue models using structured multimodal input formatting in GPT-2 combined with multi-component reward optimisation. The approach addresses critical requirements for mental health applications: data sensitivity through local deployment, computational efficiency for resource-constrained environments, and enhanced safety through smaller model parameters that enable better oversight and control. The developed model serves as a prototype dialogue suggestion tool designed to complement existing telehealth services whilst maintaining patient data privacy and clinical supervision requirements [10].

Automated evaluation demonstrated SFT effectiveness in structuring input format to capture context and emotion, improving over baseline GPT-2 (BLEU: 0.0108, ROUGE-2: 0.0175, ROUGE-L: 0.1076) with enhanced emotion accuracy (GPT-2: 66.96% vs SFT-emotion: 95.37%). The composite reward function in RL further improved generated responses across lexical measures (ROUGE-1: 0.1397, ROUGE-2: 0.0213, ROUGE-L: 0.1317). Notably, these responses are generated upon GPT-2 (124M parameters), a decoder-only model with inherent architectural limitations for complex conversational tasks. Researchers from Anthropic found that 'alignment tax' is inherent in RL training on small LLMs (<10B parameters) where NLP metrics performance may degrade after RL training [42].

Comparing against existing empirical evidence, Sharma et al. [29] performed RL on DialoGPT for empathetic rewriting, transforming low-empathy conversational posts to higher empathy. Whilst serving different purposes, it demonstrated RL's feasibility in improving response generation, achieving a BLEU score of 0.1391. Although this score is significantly higher than our model's, this does not diminish our results as their task differs fundamentally from dialogue generation. Empathetic rewriting resembles machine translation – improving existing text rather than generating new responses - creating an inherent advantage when using BLEU, a metric designed to evaluate machine-translated text against reference translations.

Evaluating Saha et al.'s [30] work, which performed RL on GPT-2 using composite rewards of BLEU, ROUGE-L, and sentiment scores for empathetic text generation, this paper recognises that such metrics may prioritise lexical similarity over contextual appropriateness and often yield sparse, biased signals [34]. Whilst direct performance comparison is not meaningful due to different datasets and tasks, our approach demonstrates that composite rewards incorporating contextual and emotional appropriateness can achieve competitive lexical scores (ROUGE-L: 0.1317) without directly optimising for them.

An important aspect of this model is improving current therapist generation approaches by enhancing affective empathy, where SFT-emotion outperformed SFT-

no-emotion across all metrics (BLEU: 0.0119 vs 0.0106; ROUGE-2: 0.0179 vs 0.0159; ROUGE-L: 0.1084 vs 0.1058; METEOR: 0.0533 vs 0.0433). This suggests that generating emotion tokens before end-of-sentence tokens provides dual benefits: improving lexical generation whilst demonstrating affective empathy, which Rudra et al. [43] emphasised as important for therapists to display genuine emotional engagement rather than surface-level recognition.

This work's primary contribution lies in demonstrating a proof-of-concept for developing smaller, safer therapeutic dialogue models suitable for sensitive healthcare data. By restructuring GPT-2 for contextual-emotional input processing, we present a cost-effective multimodal alternative that addresses critical deployment constraints in mental health applications: data privacy through local deployment, computational efficiency for resource-limited environments, and enhanced oversight capabilities inherent in smaller model architectures. The SFT model's consistent improvements across automated evaluation metrics, validated by LLM-as-a-judge evaluation showing 'Relevance Focus' as the second-highest metric, demonstrates that structured input effectively captures context-emotion relationships without requiring large-scale, resource-intensive models. This approach offers healthcare institutions a viable pathway for implementing AI-assisted therapeutic tools whilst maintaining stringent data security and clinical supervision requirements.

6.1 Limitations

This paper acknowledges several important limitations. Firstly, significant class imbalance in emotion labels within the training dataset resulted in overfitting towards neutral therapist emotion generation. This creates bias where the model defaults to emotionally restrained responses, potentially limiting communication flexibility. Whilst this reflects realistic professional therapist approaches maintaining clinical boundaries, it may reduce empathy and emotional responsiveness in contexts requiring warmer interactions. This limitation appears in LLM-as-a-judge evaluation results, where the model achieved the lowest score for Practical Helpfulness whilst rating highest in Professional Appropriateness, suggesting focus on professional behaviour over empathetic engagement.

Another significant limitation is the absence of empirical hyperparameter optimisation for reward weights prior to training. Given limited computational resources and time, systematic optimisation techniques were omitted. Human evaluation reveals the relevance component may be disproportionately weighted, leading to suboptimal conversational patterns where the model occasionally uses simplistic paraphrasing rather than generating meaningful responses. For instance, when presented with "My ex keeps texting me and I don't know how to respond," the model sometimes produces repetitive responses like "So you're not sure how to respond?" This suggests excessive relevance weighting may encourage surface-level lexical matching rather than deeper engagement.

Significant limitations exist in using LLMs for mental health support. First, while these models simulate empathy effectively, they frequently exaggerate it through lengthy responses, creating misleading emotional connections and potentially fostering inappropriate user reliance [43] [44]. Second, generic prompt engineering fails to capture contextual awareness, resulting in inconsistent responses that undermine

reliability [17] [45]. Third, LLMs lack clinical judgement, making them prone to generating harmful recommendations in legally sensitive psychological contexts [46]. Fourth, pre-trained models exhibit racial and demographic biases, producing stereotyped responses and hallucinated inaccuracies that misalign with users' values [26] [18]. Finally, commercially available LLMs raise privacy concerns regarding API vulnerabilities and data leakage when handling sensitive information [47].

6.2 Future Research

Whilst this paper employed a relatively small LLM with only 124M parameters, qualitative examples demonstrated GPT-2's capability to produce responses that effectively capture context and remain highly relevant to user input. Fine-tuning on therapy session data shows its effectiveness in enabling LLMs to replicate professionalism and appropriate speech tones, potentially addressing limitations of generalised LLMs that produce superficial empathetic responses through lengthy, generic outputs. This paper highlights the critical need for higher quality datasets to better align LLMs with human needs and serve vulnerable groups requiring specialised support.

The structured input approach demonstrated success and reveals potential for improving emotional support CAs to capture multimodal information. Whilst multimodality here involves preprocessing video and audio signals into text, future research could employ RL on graph-convolutional-network-based multimodal classifiers [48], restructuring input formats for multimodal response generation.

7 Ethics

All datasets are open source from published articles. Recognising potential harms of LLMs in mental health support [49], including retraumatisation, unreliable advice, overreliance, and safety bypass risks, this paper positions the model as a clinical assistance tool rather than direct patient interface. Implementation follows a tiered ethical framework [50]: Level 1 Safeguards (expert-use disclaimers, local deployment), Level 2 Protections (human oversight, transcript review), and Level 3 Responsibility (transparent documentation of limitations and capabilities).

Declaration on Generative AI: GenAI has been used for supporting drafting and language refinement

References

- [1] ‘Mental Illness - National Institute of Mental Health (NIMH)’. Accessed: Sep. 13, 2025. [Online]. Available: <https://www.nimh.nih.gov/health/statistics/mental-illness>
- [2] P. Corrigan, ‘How stigma interferes with mental health care.’, *Am. Psychol.*, vol. 59, no. 7, pp. 614–625, Oct. 2004, doi: 10.1037/0003-066X.59.7.614.
- [3] B. M. Coêlho, G. L. Santana, M. C. Viana, Y.-P. Wang, and L. H. Andrade, “‘I don’t need any treatment’ – barriers to mental health treatment in the general population of a megacity”, *Braz. J. Psychiatry*, vol. 43, no. 6, pp. 590–598, Dec. 2021, doi: 10.1590/1516-4446-2020-1448.
- [4] J. Goodwin, E. Savage, and A. O’Donovan, “‘I Personally Wouldn’t Know Where to Go’: Adolescents’ Perceptions of Mental Health Services”, *J. Adolesc. Res.*, vol. 39, no. 5, pp. 1384–1412, Sep. 2024, doi: 10.1177/07435584221076056.
- [5] S. Li, M. Chen, P. L. Liu, and J. Xu, ‘Following Medical Advice of an AI or a Human Doctor? Experimental Evidence Based on Clinician-Patient Communication Pathway Model’, *Health Commun.*, vol. 40, no. 9, pp. 1810–1822, Jul. 2025, doi: 10.1080/10410236.2024.2423114.
- [6] E. Jo *et al.*, ‘Assessing GPT-4’s Performance in Delivering Medical Advice: Comparative Analysis With Human Experts’, *JMIR Med. Educ.*, vol. 10, no. 1, p. e51282, Jul. 2024, doi: 10.2196/51282.
- [7] A. Følstad *et al.*, ‘Future directions for chatbot research: an interdisciplinary research agenda’, *Computing*, vol. 103, no. 12, pp. 2915–2942, Dec. 2021, doi: 10.1007/s00607-021-01016-7.
- [8] D. O. Shumway and H. J. Hartman, ‘Medical malpractice liability in large language model artificial intelligence: legal review and policy recommendations’, *J. Osteopath. Med.*, vol. 124, no. 7, pp. 287–290, Jul. 2024, doi: 10.1515/jom-2023-0229.
- [9] D. Beeferman, A. Berger, and J. Lafferty, ‘Statistical Models for Text Segmentation’.
- [10] Z. Ma, Y. Mei, and Z. Su, ‘Understanding the Benefits and Challenges of Using Large Language Model-based Conversational Agents for Mental Well-being Support’, *AMIA. Annu. Symp. Proc.*, vol. 2023, pp. 1105–1114, Jan. 2024.
- [11] K. Witkowski, R. B. Dougherty, and S. R. Neely, ‘Public perceptions of artificial intelligence in healthcare: ethical concerns and opportunities for patient-centered care’, *BMC Med. Ethics*, vol. 25, no. 1, p. 74, Jun. 2024, doi: 10.1186/s12910-024-01066-4.
- [12] M. E. te Pas, W. G. M. M. Rutten, R. A. Bouwman, and M. P. Buise, ‘User Experience of a Chatbot Questionnaire Versus a Regular Computer Questionnaire: Prospective Comparative Study’, *JMIR Med. Inform.*, vol. 8, no. 12, p. e21982, Dec. 2020, doi: 10.2196/21982.
- [13] E. M. Boucher *et al.*, ‘Artificially intelligent chatbots in digital mental health interventions: a review’, *Expert Rev. Med. Devices*, vol. 18, no. sup1, pp. 37–49, Dec. 2021, doi: 10.1080/17434440.2021.2013200.
- [14] C. Raffel *et al.*, ‘Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer’.
- [15] Y. Chu, L. Liao, Z. Zhou, C.-W. Ngo, and R. Hong, ‘Towards Multimodal Emotional Support Conversation Systems’, Oct. 19, 2024, *arXiv*: arXiv:2408.03650. doi: 10.48550/arXiv.2408.03650.
- [16] N. Stiennon *et al.*, ‘Learning to summarize with human feedback’, in *Advances in Neural Information Processing Systems*, Curran Associates, Inc., 2020, pp. 3008–3021. Accessed: Apr. 08, 2025. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2020/hash/1f89885d556929e98d3ef9b86448f951-Abstract.html
- [17] Z. Guo, A. Lai, J. H. Thygesen, J. Farrington, T. Keen, and K. Li, ‘Large Language Model for Mental Health: A Systematic Review’, Feb. 18, 2024. doi: 10.2196/preprints.57400.
- [18] L. Ouyang *et al.*, ‘Training language models to follow instructions with human feedback’.

- [19] M. Luo, C. J. Warren, L. Cheng, H. M. Abdul-Muhsin, and I. Banerjee, ‘Assessing Empathy in Large Language Models with Real-World Physician-Patient Interactions’, in *2024 IEEE International Conference on Big Data (BigData)*, Dec. 2024, pp. 6510–6519. doi: 10.1109/BigData62323.2024.10825307.
- [20] H. Lang, F. Huang, and Y. Li, ‘Fine-Tuning Language Models with Reward Learning on Policy’, Mar. 28, 2024, *arXiv*: arXiv:2403.19279. doi: 10.48550/arXiv.2403.19279.
- [21] R. L. Bashshur, G. W. Shannon, N. Bashshur, and P. M. Yellowlees, ‘The Empirical Evidence for Telemedicine Interventions in Mental Disorders’, *Telemed. E-Health*, vol. 22, no. 2, pp. 87–113, Feb. 2016, doi: 10.1089/tmj.2015.0206.
- [22] E. J. Kraus, B. Nicosia, and D. I. Shalowitz, ‘A qualitative study of patients’ attitudes towards telemedicine for gynecologic cancer care’, *Gynecol. Oncol.*, vol. 165, no. 1, pp. 155–159, Apr. 2022, doi: 10.1016/j.ygyno.2022.01.035.
- [23] ‘Patients with Depression and Anxiety Surge as Psychologists Respond to the Coronavirus Pandemic’.
- [24] M. Torniainen-Holm *et al.*, ‘The effectiveness of email-based exercises in promoting psychological wellbeing and healthy lifestyle: a two-year follow-up study’, *BMC Psychol.*, vol. 4, no. 1, p. 21, Dec. 2016, doi: 10.1186/s40359-016-0125-4.
- [25] J. Crawford, S. Haffar, S. Fernando, H. Stephens, S. B. Harvey, and M. Black, ‘Client perspectives: Telehealth for mental health services’, *Australas. Psychiatry*, vol. 33, no. 1, pp. 96–102, Feb. 2025, doi: 10.1177/10398562241270986.
- [26] S. Gabriel, I. Puri, X. Xu, M. Malgaroli, and M. Ghassemi, ‘Can AI Relate: Testing Large Language Model Response for Mental Health Support’, Oct. 07, 2024, *arXiv*: arXiv:2405.12021. doi: 10.48550/arXiv.2405.12021.
- [27] H. W. Chung *et al.*, ‘Scaling Instruction-Finetuned Language Models’.
- [28] H. Jin, S. Chen, D. Dilixiati, Y. Jiang, M. Wu, and K. Q. Zhu, ‘PsyEval: A Suite of Mental Health Related Tasks for Evaluating Large Language Models’, Jun. 03, 2024, *arXiv*: arXiv:2311.09189. doi: 10.48550/arXiv.2311.09189.
- [29] A. Sharma, I. W. Lin, A. S. Miner, D. C. Atkins, and T. Althoff, ‘Towards Facilitating Empathic Conversations in Online Mental Health Support: A Reinforcement Learning Approach’, in *Proceedings of the Web Conference 2021*, Ljubljana Slovenia: ACM, Apr. 2021, pp. 194–205. doi: 10.1145/3442381.3450097.
- [30] T. Saha, V. Gakhreja, A. S. Das, S. Chakraborty, and S. Saha, ‘Towards Motivational and Empathetic Response Generation in Online Mental Health Support’, in *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, Madrid Spain: ACM, Jul. 2022, pp. 2650–2656. doi: 10.1145/3477495.3531912.
- [31] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, ‘Language Models are Unsupervised Multitask Learners’.
- [32] S. Liu *et al.*, ‘Towards Emotional Support Dialog Systems’, Jun. 02, 2021, *arXiv*: arXiv:2106.01144. doi: 10.48550/arXiv.2106.01144.
- [33] Q. Wang *et al.*, ‘Improving sequence labeling with labeled clue sentences’, *Knowl.-Based Syst.*, vol. 257, p. 109828, Dec. 2022, doi: 10.1016/j.knosys.2022.109828.
- [34] A. Anuchitanukul and J. Ive, ‘SURF: Semantic-level Unsupervised Reward Function for Machine Translation’, in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, M. Carpuat, M.-C. de Marneffe, and I. V. Meza Ruiz, Eds., Seattle, United States: Association for Computational Linguistics, Jul. 2022, pp. 4508–4522. doi: 10.18653/v1/2022.naacl-main.334.
- [35] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, ‘Proximal Policy Optimization Algorithms’, Aug. 28, 2017, *arXiv*: arXiv:1707.06347. doi: 10.48550/arXiv.1707.06347.
- [36] J. Eisenstein *et al.*, ‘Helping or Herding? Reward Model Ensembles Mitigate but do not Eliminate Reward Hacking’, Aug. 16, 2024, *arXiv*: arXiv:2312.09244. doi: 10.48550/arXiv.2312.09244.

- [37] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, 'BLEU: a method for automatic evaluation of machine translation', in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL '02*, Philadelphia, Pennsylvania: Association for Computational Linguistics, 2001, p. 311. doi: 10.3115/1073083.1073135.
- [38] C.-Y. Lin, 'ROUGE: A Package for Automatic Evaluation of Summaries', in *Text Summarization Branches Out*, Barcelona, Spain: Association for Computational Linguistics, Jul. 2004, pp. 74–81. Accessed: Sep. 13, 2025. [Online]. Available: <https://aclanthology.org/W04-1013/>
- [39] S. Banerjee and A. Lavie, 'METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments', in *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, J. Goldstein, A. Lavie, C.-Y. Lin, and C. Voss, Eds., Ann Arbor, Michigan: Association for Computational Linguistics, Jun. 2005, pp. 65–72. Accessed: Sep. 13, 2025. [Online]. Available: <https://aclanthology.org/W05-0909/>
- [40] A. Yuan, E. Garcia Colato, B. Pescosolido, H. Song, and S. Samtani, 'Improving Workplace Well-being in Modern Organizations: A Review of Large Language Model-based Mental Health Chatbots', *ACM Trans. Manag. Inf. Syst.*, vol. 16, no. 1, pp. 1–26, Mar. 2025, doi: 10.1145/3701041.
- [41] C. Gelso and K. Kanninen, 'Neutrality Revisited: On the Value of Being Neutral Within an Empathic Atmosphere', *J. Psychother. Integr.*, vol. 27, pp. 330–341, Feb. 2017, doi: 10.1037/int0000072.
- [42] Y. Bai *et al.*, 'Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback', Apr. 12, 2022, *arXiv*: arXiv:2204.05862. doi: 10.48550/arXiv.2204.05862.
- [43] P. Rudra, W.-T. Balke, T. Kacprowski, F. Ursin, and S. Salloch, 'Large language models for surgical informed consent: an ethical perspective on simulated empathy', *J. Med. Ethics*, p. jme-2024-110652, Mar. 2025, doi: 10.1136/jme-2024-110652.
- [44] V. Sorin *et al.*, 'Large Language Models and Empathy: Systematic Review', *J. Med. Internet Res.*, vol. 26, no. 1, p. e52597, Dec. 2024, doi: 10.2196/52597.
- [45] F. Farhat, 'ChatGPT as a Complementary Mental Health Resource: A Boon or a Bane', *Ann. Biomed. Eng.*, vol. 52, no. 5, pp. 1111–1114, May 2024, doi: 10.1007/s10439-023-03326-7.
- [46] D. Grabb, 'The impact of prompt engineering in large language model performance: a psychiatric example', *J. Med. Artif. Intell.*, vol. 6, no. 0, Oct. 2023, doi: 10.21037/jmai-23-71.
- [47] H. R. Lawrence, R. A. Schneider, S. B. Rubin, M. J. Matarić, D. J. McDuff, and M. J. Bell, 'The Opportunities and Risks of Large Language Models in Mental Health', *JMIR Ment. Health*, vol. 11, no. 1, p. e59479, Jul. 2024, doi: 10.2196/59479.
- [48] Q. Yang, M. Ye, and B. Du, 'EmoLLM: Multimodal Emotional Understanding Meets Large Language Models', Jun. 29, 2024, *arXiv*: arXiv:2406.16442. doi: 10.48550/arXiv.2406.16442.
- [49] I. Song, S. R. Pendse, N. Kumar, and M. D. Choudhury, 'The Typing Cure: Experiences with Large Language Model Chatbots for Mental Health Support', Mar. 06, 2024, *arXiv*: arXiv:2401.14362. doi: 10.48550/arXiv.2401.14362.
- [50] N. Neveditsin, P. Lingras, and V. Mago, 'Clinical Insights: A Comprehensive Review of Language Models in Medicine', Jan. 07, 2025, *arXiv*: arXiv:2408.11735. doi: 10.48550/arXiv.2408.11735.