

Muffin: Mitigating Unhelpfulness in Emotional Support Conversations with Multifaceted AI Feedback

Anonymous ACL submission

Abstract

Emotional support conversation systems are designed to alleviate users’ emotional distress and assist them in overcoming their challenges. While previous studies have made progress, their models occasionally generate unhelpful responses, which are intended to be supportive but instead have counterproductive effects. Since unhelpful responses can hinder the effectiveness of emotional support, it is crucial to mitigate them within conversations. Our solution is motivated by two principal considerations: (1) multiple facets of emotional support are expected to be considered when developing emotional support conversation models, and (2) directly reducing the probability of generating unhelpful responses can effectively mitigate their occurrence. Accordingly, we introduce a novel **model-agnostic** framework named *Mitigating unhelpfulness with multifaceted AI feedback for emotional support (Muffin)*. It first employs a multifaceted AI feedback module designed to assess the helpfulness model responses across various facets of emotional support. Leveraging contrastive learning, Muffin then reduces the unhelpful responses’ likelihoods. To validate the effectiveness of our proposed framework, we apply Muffin to various previous emotional support generation models, including the state-of-the-art. Experimental results demonstrate that Muffin can significantly mitigate unhelpful response generation while enhancing response fluency and relevance.

1 Introduction

Emotional support conversation systems (supporters) are designed to generate responses that can buffer the emotional distress experienced by users (help-seekers) and help users to work through the challenges they are confronting (Liu et al., 2021). Recently, many studies have contributed to this field (Deng et al., 2023; Zhao et al., 2023; Cheng et al., 2022; Tu et al., 2022). Despite great success,

A Conversation between a Help-Seeker and a Supporter

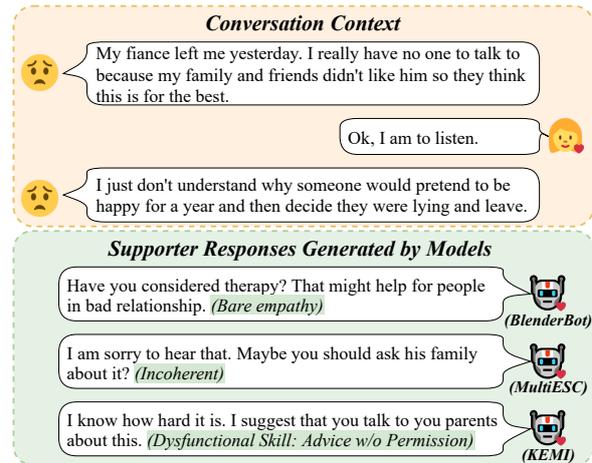


Figure 1: Examples of unhelpful responses generated by recent emotional support conversation models, including BlenderBot (Roller et al., 2021), MultiESC (Cheng et al., 2022), and KEMI (Deng et al., 2023).

their models occasionally generate well-intended responses that produce a counterproductive support effect, i.e., exacerbating the negative emotional states of users or inhibiting effective problem-solving, as shown in Figure 1. In the psychology and communication theories, these failed support attempts are termed “*unhelpful messages*” (Greene and Burleson, 2003; Burleson, 1985).

The frequency of unhelpful responses generated by some of the previous models is not extremely high, e.g, approximately 30% of responses generated by MultiESC (Cheng et al., 2022) on ESConv benchmark (Liu et al., 2021) are identified as unhelpful under strict evaluation criteria. However, their occurrence can significantly undermine earlier supportive efforts and damage the trust between the help-seeker and the supporter (Llewelyn et al., 1988). Therefore, mitigating models’ generation of unhelpful responses is critical. We aim to address this problem with the following two deliberations. **D1 - Consideration of Multiple Facets:** Many pre-

064 various studies generate responses that primarily em- 115
065 phasize a single facet of emotional support, e.g., 116
066 one of empathetic expression (Li et al., 2024), com- 117
067 munication skill efficiency (Cheng et al., 2022; Liu 118
068 et al., 2021), or response coherence (Deng et al., 119
069 2023), in each of their models. However, such a 120
070 singular emphasis on one facet often leads to the 121
071 oversight of the others, potentially resulting in un- 122
072 helpfulness (Greene and Burleson, 2003), as exem-
073 plified in Figure 1. **D2 - Direct Minimization of Un-**
074 **helpful Response Probability:** Previous models are
075 typically optimized by minimizing the negative log-
076 likelihood of the golden responses. Moving beyond
077 this optimization objective, we aim to specifically
078 mitigate unhelpful responses by directly targeting
079 and reducing the probability of their generation.

080 Accordingly, this paper introduces a novel
081 **model-agnostic** framework called Mitigating
082 unhelpfulness with multifaceted AI feedback for
083 emotional support (Muffin). For **D1**, we design a
084 multifaceted AI feedback module within the Muffin
085 framework. This module assesses whether a
086 specific response is unhelpful from multiple facets
087 of emotional support. Leveraging the advanced ca-
088 pabilities of recent large language models (LLMs),
089 we implement this module by instruction-tuning
090 LLaMA, avoiding inefficient and expensive human
091 feedback collection. Then, we continue optimiz-
092 ing an emotional support conversation model with
093 its previous training objective and a new one to
094 implement **D2**. The additional objective is to mini-
095 mize the likelihood of unhelpful responses, which
096 is implemented by contrasting unhelpful responses,
097 identified by the feedback module, and the other
098 (non-unhelpful) ones. Through these two steps,
099 we aim to mitigate the unhelpful responses gener-
100 ated by a given emotional support conversation
101 model. Experimental results highlight the effective-
102 ness of our framework, demonstrating that Muffin
103 can enhance the helpfulness of previous emotional
104 support conversation models, including those rec-
105 ognized as state-of-the-art. The main contributions
106 of this work are:

- 107 1. We recognize and address a crucial problem in 115
108 recent emotional support conversation models, 116
109 i.e., the generation of unhelpful responses, a 117
110 key concern in effective emotional support.
- 111 2. We propose Muffin, a novel model-agnostic 118
112 framework designed to mitigate unhelpful re- 119
113 sponse generation. It incorporates a multi- 120
114 faceted AI feedback module to distinguish 121

115 unhelpful generated responses and mitigates 116
116 responses identified as unhelpful by leverag- 117
117 ing contrastive learning.

- 118 3. We undertake experiments with the latest emo- 119
119 tional support conversation models, including 120
120 state-of-the-art ones, to demonstrate Muffin’s 121
121 effectiveness in mitigating the models’ ten- 122
122 dency to produce unhelpful responses.

123 2 Related Work 123

124 2.1 Emotional Support Conversation 124

125 In the domain of emotional support conversation 125
126 generation, prior studies have achieved some suc- 126
127 cess. Specifically, they have each emphasized 127
128 and incorporated different facets of emotional sup- 128
129 port, such as empathetic expression, communica- 129
130 tion skills, and response coherence, into their re- 130
131 spective models. Some of them consider a sin- 131
132 gle facet. For example, MultiESC (Cheng et al., 132
133 2022) only considers the communication skill ef- 133
134 ficacy of responses by planning response strate- 134
135 gies. KEMI (Deng et al., 2023) incorporates re- 135
136 lated external knowledge for response generation. 136
137 Although the response coherence is enhanced, the 137
138 efficacy of communication skills is ignored. Li 138
139 et al. (2024) enhance the empathetic expression 139
140 of LLMs through chain-of-thought. Beyond these, 140
141 MISC (Tu et al., 2022) generates supportive re- 141
142 sponses considering both commonsense and com- 142
143 munication skills, thereby enhancing two facets of 143
144 emotional support. TranESC (Zhao et al., 2023) 144
145 incorporates commonsense, communication skills, 145
146 and emotional elements into response generation. 146
147 Most all these models are optimized by minimiz- 147
148 ing the negative log-likelihood of golden responses, 148
149 instead of the unhelpful response likelihood. 149

150 2.2 Contrastive Learning for Text Generation 150

151 Contrastive learning was initially employed to learn 151
152 meaningful representations by contrasting positive 152
153 and negative samples in the field of natural lan- 153
154 guage processing (Logeswaran and Lee, 2018; Gut- 154
155 mann and Hyvärinen, 2012). Recently, it has been 155
156 applied to text generation (Jiashuo et al., 2023; 156
157 Zheng et al., 2023; Liu et al., 2022), achieving im- 157
158 pressive performance across various settings. Dur- 158
159 ing training, the model is exposed to a range of 159
160 “hard” negative and positive samples through con- 160
161 trastive learning, enabling the model to distinguish 161
162 preferred outputs from less desirable ones. Con- 162
163 sequently, selecting positive and negative samples 163

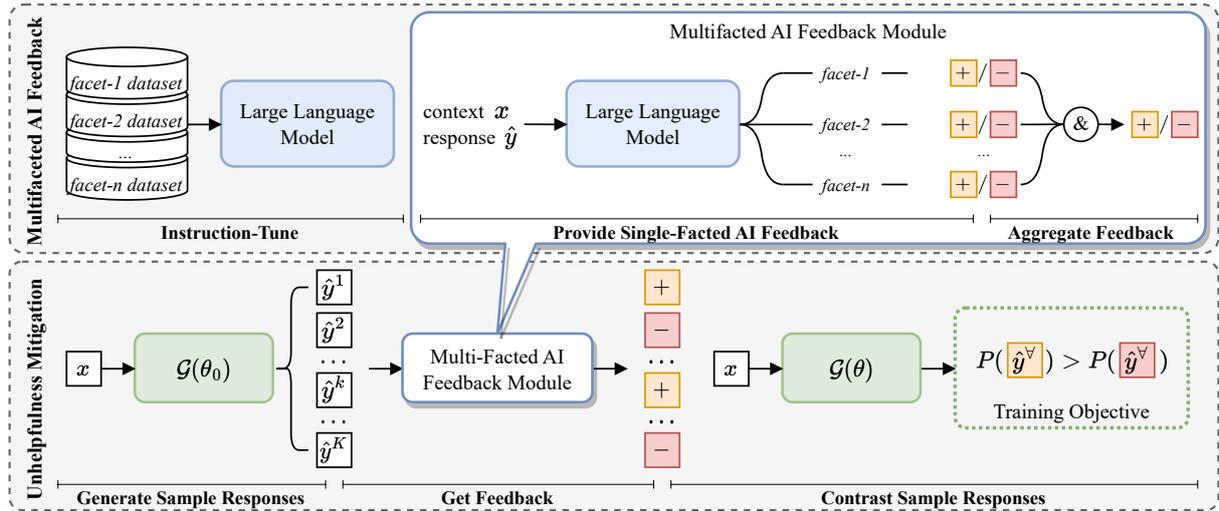


Figure 2: The overview of our proposed model-agnostic framework—Muffin. $+$ and $-$ indicate helpful (non-unhelpful) and unhelpful labels, respectively.

is crucial in this process. In this paper, responses generated by the model that are deemed unhelpful are negative samples, while all other responses are considered positive samples.

3 Preliminaries

3.1 Unhelpfulness of Emotional Support

We draw upon theories in psychology and communication (Greene and Burleson, 2003; Burleson, 1985) and adopt the term “unhelpful” to characterize responses that consistently produce negative outcomes in emotional support conversations. Conversely, responses that yield positive outcomes, or at the very least, do not cause negative effects, are termed “helpful.” These theories also suggest that an unhelpful response often stems from a flaw in merely one specific facet of emotional support. Often, the flaw directly exacerbates the user’s negative feelings or hinders effective problem-solving. For instance, a response can be deemed unhelpful if it either neglects the individual’s feelings and needs (lacking empathy) or portrays the individual’s behavior as problematic (exhibiting a dysfunctional communication skill: confront). In our work, we use this feature to identify unhelpful responses.

3.2 Problem Definition

Our primary goal is to mitigate the generation of unhelpful responses. Rather than training a new model from scratch, we aim to refine a pre-trained emotional support conversation model with the dataset it was originally trained on. This process unfolds as follows. Let $\mathcal{G}(\theta_0)$ represent the model

trained on a dataset \mathcal{D} , where θ_0 denotes the model parameters. Each instance in \mathcal{D} is denoted as (x, y) , with x as the input and y as the expected output. Usually, x is the conversation context, but it contains additional related information in some models. Assume that there are K samples $\{\hat{y}^1, \dots, \hat{y}^K\}$ with labels $\{\hat{l}^1, \dots, \hat{l}^K\}$. These samples are the diverse beam search generation results of $\mathcal{G}(x; \theta_0)$. As for the label $\hat{l}^k \in \{0, 1\}$, it represents feedback to indicate whether the sample \hat{y}^k is unhelpful ($\hat{l}^k = 0$) or not unhelpful ($\hat{l}^k = 1$). Our objective is to refine the model’s parameters θ such that the likelihood of generating unhelpful samples is reduced relative to helpful ones. In this process, we only modify the training process, ensuring that the model’s architecture and the inference mechanism remain untouched. Moreover, our approach is model-agnostic. This implies that $\mathcal{G}(\theta_0)$ can be any deep learning model designed and trained for emotional support conversations.

4 Method

The overall framework of Muffin is outlined in Figure 2. It is composed of two principal components, each specifically designed for deliberations in Section 1: **D1: Consideration of Multiple Facets** and **D2: Direct Minimization of Unhelpful Response Probability**, respectively. The multifaceted AI feedback module aims to identify whether a response from multiple facets of emotional support is unhelpful. The unhelpfulness mitigation module mitigates the likelihood of unhelpful responses by contrasting helpful and unhelpful responses.

4.1 Multifaceted AI Feedback

We distinguish whether a response is unhelpful from multiple facets. However, collecting feedback from humans is inefficient and costly. In addition, recent large language models (LLMs), such as the GPT series (Ouyang et al., 2022) and LLaMA (Touvron et al., 2023), demonstrate remarkable natural language understanding capabilities. Therefore, we decide to obtain feedback from AI.

Instruction-tuning Prompt engineering provides a simple and straightforward approach to obtaining feedback from LLMs. However, our experiments suggest that it is challenging to manifest the full potential of LLMs for emotional support without investing significant effort in prompt design, which will be detailed later. As an alternative, we elicit the desired capabilities of the LLM via instruction tuning (Wei et al., 2022). Specifically, we design task descriptions and instructions tailored to classification tasks related to different emotional support facets. We employ corresponding datasets for these tasks. The prompt employed is illustrated in Figure 3. Notably, the response class can indicate whether the response is unhelpful regarding this facet. During the training phase, all *texts in italics enclosed within curly braces* are provided. During inference, the model is expected to generate the *response class*, based on the other *italicized inputs within the curly braces*.

```

### Instruction:
{task description and instruction}

### Input:
Conversation Context: {context}
The last supporter statement: {response}
{all possible classes}

### Output:
{response class}

```

Figure 3: The prompt used by the Multifaceted AI Feedback for classifying the supporter’s response.

Multifaceted AI feedback module The final feedback for an emotional support response is derived from an aggregation of AI feedback across multiple facets. For the given response and its context, we use the instruction-tuned LLM to provide feedback on all these facets respectively. If feedback from any of these facets suggests the response is unhelpful, the response is accordingly labeled as 0; otherwise, the response is deemed helpful

(non-unhelpful) and labeled as 1, as mentioned in Section 3.1.

4.2 Unhelpfulness Mitigation

We mitigate \mathcal{G}_0 generating unhelpful responses by contrasting helpful and unhelpful responses generated by \mathcal{G}_0 itself, which can be implemented by the following three steps:

Generating sample responses We utilize \mathcal{G}_0 to generate responses on its own training dataset \mathcal{D} using diverse beam search (Vijayakumar et al., 2016). Thus, for each instance $(x, y) \in \mathcal{D}$, there are K sample responses $\{\hat{y}^1, \dots, \hat{y}^K\}$.

Getting feedback These responses can be generated because they have relatively high generation probabilities. However, some of them can be unhelpful responses. Therefore, we adopt the multifaceted AI feedback module to identify whether these responses are unhelpful. Thus, we obtain K labels $\{\hat{l}^1, \dots, \hat{l}^K\}$, where $\hat{l}^k \in \{0, 1\}$.

Contrasting sample responses We expect that the model \mathcal{G} can sign a higher likelihood to the helpful responses than the unhelpful ones. Therefore, we contrast them using the following loss:

$$\mathcal{L}_{cl} = \frac{1}{2K} \sum_i \sum_{j \neq i} \max(0, -(\hat{l}^i - \hat{l}^j) \times (P(\hat{y}^i|x) - P(\hat{y}^j|x) + \lambda)), \quad (1)$$

where λ is the margin hyperparameter. Moreover, $P(\hat{y}^i|x)$ is the length-normalized log-probability of the response \hat{y}^i , and it is computed by:

$$P(\hat{y}^i|x) = \sum_{t=1}^{|\hat{y}^i|} \frac{\log \mathcal{G}(\hat{y}_t^i|x, \hat{y}_{<t}^i; \theta)}{|\hat{y}^i|^\alpha}, \quad (2)$$

where α is the length penalty hyperparameter. In addition to the above loss, we also consider the negative log-likelihood loss to prevent the model’s generation from deviating too much from the ground truth. The loss can be formulated as:

$$\mathcal{L}_{gen} = -\frac{1}{|y|} \sum_{t=1}^{|y|} \log \mathcal{G}(y_t|x, y_{<t}; \theta). \quad (3)$$

The final loss is the combination of the above two losses:

$$\mathcal{L} = \beta_{cl} \mathcal{L}_{cl} + \beta_{gen} \mathcal{L}_{gen}, \quad (4)$$

where β_{cl} and β_{gen} are weight hyperparameters.

5 Experiments

5.1 Experimental Setups

Base models ($\mathcal{G}(\theta_0)$) Our proposed method, i.e., Muffin, is a model-agnostic approach designed to mitigate the unhelpfulness of an existing emotional support conversation model. To examine its effectiveness, we experiment with five recent models: BlenderBot (Vanilla) (Roller et al., 2021), BlenderBot-Joint (Joint) (Liu et al., 2021), Multi-ESC (Cheng et al., 2022), TransESC (Zhao et al., 2023), and KEMI (Deng et al., 2023). We obtain each model’s parameters θ using its official implementation and the default hyperparameters.

When training Muffin $_{\mathcal{G}(\theta_0)}$, we try to use the same training hyperparameters as the base model $\mathcal{G}(\theta_0)$, including batch size and random seed. However, we use a smaller learning rate, i.e., 3×10^{-5} , to help the model converge more efficiently. We set the epoch number as 1 since the training loss converged within one epoch. The margin hyperparameter λ is 0.01, the length penalty hyperparameter α is 1, and the weight hyperparameters β_{cl} and β_{gen} are both 1. In the response generation phase, as described in Section 4.2, we set the number of sample responses K to 10. Consequently, both the beam size and the number of beam groups are configured to be 10, while all other generation hyperparameters are the same as its base model.

ESConv dataset (\mathcal{D}) The ESConv dataset (Liu et al., 2021) is used to train the aforementioned base models. ESConv is a benchmark for emotional support conversations, comprising approximately 1K conversations with 13K utterances. All base models, with the exception of TransESC, follow the original division of ESConv for training, validation, and testing, using an 8:1:1 ratio. TransESC employs a random split while maintaining the same ratio. Notably, each model adopts different data preprocessing methods. We adhere to each base model’s specific data division and pre-processing.

5.2 Facets of Emotional Support

We consider three essential facets of emotional support: empathetic expression, skill efficiency, and response coherence, which the base models incorporate into their models. Here, we would like to describe the unhelpfulness of each facet and detail the corresponding classification dataset.

Empathetic expression Empathetic expressions signify the supporter’s interest in and comprehen-

sion of the help-seeker’s perspective. Conversely, their absence can impede conversation engagement and obstruct establishing a trust-based relationship between the supporter and the help-seeker (Morse et al., 1992). While empathy encompasses various aspects (Paiva et al., 2017; Batson, 2009; Bohart and Greenberg, 1997), we adopt the comprehensive framework proposed by (Sharma et al., 2020). It identifies three empathy communication mechanisms, i.e., emotional reaction, interpretations, and explorations, each assessed across three levels: no communication, weak communication, and strong communication. Our work considers responses exhibiting no empathy across all mechanisms unhelpful. For example, “sleeplessness can result in upsetness,” which inappropriately offers mere information, is considered unhelpful in empathetic expression, especially when responding to a statement like “I am upset”.

For training and testing LLMs in classifying unhelpful responses regarding empathetic expression, we also utilize the dataset compiled by Sharma et al. (2020). It consists of 3K context-response pairs. Each response within this dataset is assessed based on three previously mentioned empathy communication mechanisms. Responses consistently labeled as “no communication” across all these mechanisms are identified as unhelpful.

Skill efficiency The applications of effective strategies can help supporters convey appropriate and impactful support messages (Greene and Burleson, 2003), deepening the understanding of the help-seeker’s state and facilitating problem solution (Hill, 2009). However, some dysfunctional skills lead to opposite effects (Barsky and Coleman, 2001; Burleson and Samter, 1985). For example, while *advice* is typically seen as beneficial, *advice without permission* can be less effective than employing no specific skills at all. This study assesses skill efficiency using motivational interviewing skill codes (Moyers et al., 2003), which include three general categories: MI adherent, MI non-adherent, and others. Responses classified as MI non-adherent category are deemed unhelpful.

We utilize the Motivational Interviewing (MI) dataset proposed by Welivita and Pu (2022) to classify unhelpful responses in the context of skill efficiency. This dataset comprises 17K context-response pairs, where each response is annotated with one of three MI codes: MI adherent, MI non-adherent, or others. Responses labeled as “MI non-

adherent” are considered unhelpful.

Response coherence While response coherence is a fundamental expectation in almost all conversational systems (Huang et al., 2020), it holds particular significance in emotional support conversation systems. Incoherent responses can confuse and impede effective communication. We categorize responses as coherent or incoherent, with the latter labeled unhelpful.

We have synthesized a dataset, derived from the base model’s training set \mathcal{D} , specifically for the purpose of response coherence detection. Specifically, we randomly selected 4K context-response pairs from \mathcal{D} . For each pair, the original response is categorized as coherent. To introduce incoherence, we employ two methods: firstly, by selecting a response from a different conversation, and secondly, by modifying keywords or important information in the original response to create a subtly incoherent variant. These methods result in a total of approximately 12K context-response pairs, which are then utilized for classifying unhelpful responses in terms of their response coherence.

5.3 Multifaceted AI Feedback

Instruction-tuning settings We instruction-tune a 7B LLaMA (Touvron et al., 2023) for the multifaceted AI feedback module and use a low-rank adaptor (LoRA) (Hu et al., 2022) for efficiency. Specifically, we freeze LLaMA’s weight and inject trainable rank decomposition matrices into query, key, value, and output layers. The learning rate is 3×10^{-4} , and the training epoch is 12. We merge the datasets described in Section 5.2, formatting each context-response pair according to the instruction template presented in Figure 3. Consequently, this merged dataset encompasses a total of 22K instances, which are utilized for instruction-tuning. The dataset is partitioned into training, validation, and test sets following an 8:1:1 split ratio.

Module performance The instruction-tuned LLaMA model is subsequently employed within the Multifaceted AI Feedback Module. The module’s effectiveness is evaluated using the instruction-tuning test set, with the results depicted in Figure 4. We report both accuracy and F1 scores for data instances throughout the entire testing set and across distinct facets. Additionally, we extend our evaluation to include the module’s performance when utilizing three other LLMs: GPT-3.5, GPT-4, and the original (vanilla) LLaMA model.

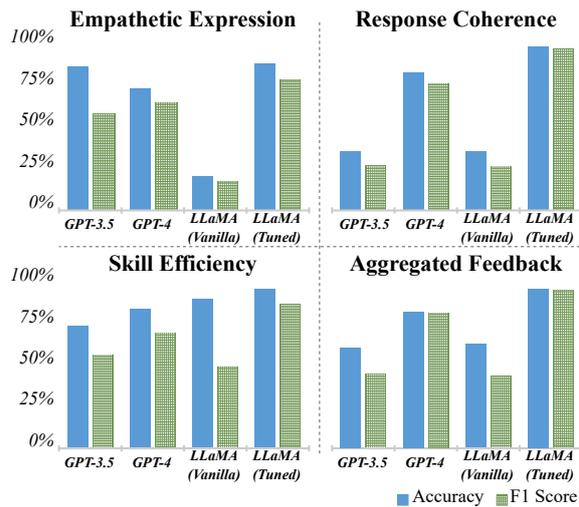


Figure 4: Comparison of performance among various Language Model Models (LLMs) including GPT-3.5, GPT-4, LLaMA (Vanilla), and LLaMA (Tuned) in classification tasks related to different facets of emotional support, as well as the aggregated feedback.

Without fine-tuning, GPT-4 outperforms the other two models in terms of the classification effectiveness. Upon closer examination, GPT-3.5 is prone to classifying responses as non-unhelpful. The vanilla LLaMA model frequently classifies responses into a singular category across different facets, leading to the worst performance. However, after instruction-tuning, LLaMA exhibits a significant enhancement in performance, with an accuracy of 90.72% and an F1 score of 89.86% on the aggregated feedback, thereby exceeding the capabilities of its counterparts. This remarkable improvement provides strong justification for our decision to apply the 7B instruction-tuned LLaMA within the Multifaceted AI Feedback Module.

5.4 Emotional Support Response Generation

Automatic evaluation We evaluate the quality of emotional support responses by a range of automatic evaluation metrics, including BLEU (Papineni et al., 2002) ($B-1/2/3/4$), ROUGE ($R-L$) (Lin, 2004), *METEOR* (Banerjee and Lavie, 2005), *CIDEr* (Vedantam et al., 2015), and BOW Embedding-based matching score (*Extreme*) (Liu et al., 2016). These metrics are good at evaluating the similarity between the generated response and the ground truth.

Table 1 showcases the performance of Muffin with different base models in all automatic evaluation metrics. In general, Muffin demonstrates significant enhancements across nearly all evalua-

Model	B-1	B-2	B-3	B-4	R-L	METEOR	CIDEr	Extreme
Vanilla	18.23	7.02	3.49	1.99	16.09	7.31	14.95	50.48
Muffin _{Vanilla}	19.43*	7.58*	3.66	2.02	16.26	7.72*	13.90*	51.00*
Joint	18.77	7.54	3.79	2.15	17.72	7.59	17.38	50.96
Muffin _{Joint}	20.59*	8.38*	4.26*	2.54*	18.35*	8.18*	19.12*	51.46*
TransESC	17.32	7.10	3.63	2.18	17.47	7.53	22.07	51.33
Muffin _{TransESC}	17.19	7.17*	3.73*	2.25*	17.54*	7.58*	22.72*	51.57*
KEMI	19.85	8.15	4.24	2.52	17.17	7.92	15.09	50.85
Muffin _{KEMI}	20.01*	8.31*	4.36*	2.60*	17.30*	7.99*	15.45*	51.11*
MultiESC	21.79	9.19	4.98	3.05	20.92	8.93	28.84	52.59
Muffin _{MultiESC}	21.83*	9.28*	5.12*	3.21*	21.26*	8.92	31.26*	52.83

Table 1: Automatic evaluation results and AI feedback from the multifaceted AI feedback module. For all metric scores and feedback, a higher value indicates better performance. The values marked with * indicate the results are statistically significant with $p < 0.05$.

tion metrics. Moreover, it can be observed that the performance of Muffin _{$\mathcal{G}(\theta_0)$} is predominantly influenced by its base model $\mathcal{G}(\theta_0)$, assessed through automatic evaluations.

Human evaluation Following previous work (Deng et al., 2023; Liu et al., 2021), we compare the base model and its corresponding Muffin model on five aspects, which are (1) *Fluency*: which model’s response is more fluent? (2) *Identification*: which model’s response is more skillful in identifying the user’s problem? (3) *Comforting*: which model’s response is better at comforting the user? (4) *Suggestion*: which model can give more helpful and informative suggestions? (5) *Overall*: which model’s response is generally more helpful from the aspect of the help-seeker? Specifically, for each $\mathcal{G}(\theta_0)$ -Muffin _{$\mathcal{G}(\theta_0)$} pair, we randomly select 100 instances for comparison. Then, we ask four unique human evaluators to vote which response is better. They can select “tie” if both responses are considered equal. We average their results as the final result.

Figure 5 summarizes the A/B test results on BlenderBot-Joint (Liu et al., 2021), KEMI (Deng et al., 2023), and MultiESC (Cheng et al., 2022), along with their corresponding Muffin models. These three settings are selected for their significant performance in automatic evaluation. The inter-rater agreement, i.e., Fleiss’ Kappa (Fleiss, 1971), is 0.39, implying fair agreement. Our Muffin models are regarded as more helpful in general. Responses generated by Muffin models are slightly more fluent than those generated by base models. Compared with the corresponding base model $\mathcal{G}(\theta_0)$, the Muffin _{$\mathcal{G}(\theta_0)$} model shows more powerful capability in identifying the help-seeker’s problem. Moreover, Muffin _{$\mathcal{G}(\theta_0)$} models can generate responses that have better effects to comfort the users than $\mathcal{G}(\theta_0)$. Annotators also prefer responses

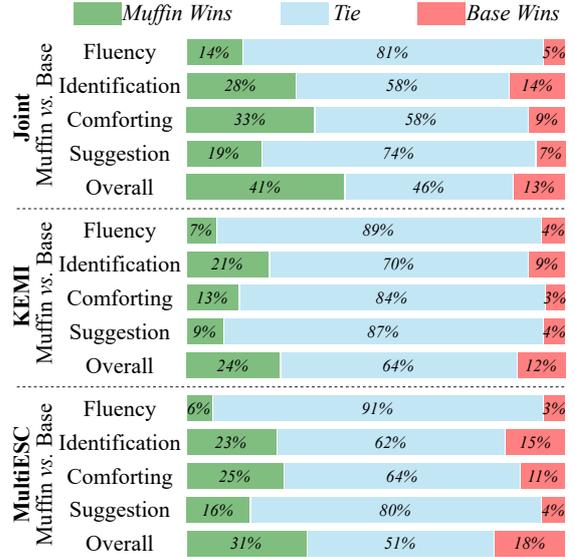


Figure 5: Human A/B test results. Displayed within each bar, from left to right, are the ratios for “Muffin Wins”, “Tie”, and “Base Wins”.

generated by Muffin _{$\mathcal{G}(\theta_0)$} because of their more helpful and informative suggestions. These results prove that Muffin indeed mitigates the unhelpfulness of emotional support conversation models.

Ablation study To assess the impact of different facets of feedback, we undertake an ablation study. Specifically, we employed AI feedback from each individual facet to distinguish helpful from unhelpful responses, subsequently contrasting them to mitigate unhelpfulness. This study is conducted using BlenderBot-Joint as the base model, given its relatively straightforward architecture and the notable improvements exhibited by Muffin_{Joint}.

The findings, as presented in Table 2, reveal that relying solely on AI feedback from a single facet for unhelpfulness mitigation results in diminished performance in automatic evaluations when

Model	B-1	B-2	B-3	B-4	R-L	METEOR	CIDeR	Extreme
Joint	18.77	7.54	3.79	2.15	17.72	7.59	17.38	50.96
Muffin _{Joint}	20.59	8.38	4.26	2.54	18.35	8.18	19.12	51.46
Muffin _{Joint} (<i>emp.</i>)	19.58	8.06	4.04	2.33	18.46	7.87	19.61	51.43
Muffin _{Joint} (<i>skill</i>)	18.68	7.51	3.81	2.19	17.98	7.67	17.84	50.98
Muffin _{Joint} (<i>cohr.</i>)	20.04	8.10	4.04	2.26	18.24	7.95	18.07	51.31

Table 2: Ablation study results. Muffin_{Joint} (X) indicates the mitigation process only uses AI feedback in terms of the facet X . All values are statistically significant with $p < 0.05$.

Conversation History	
Help-seeker	: I'm an alcoholic. I get that thought in my head and its like I become obsessed.
Supporter	: Why not try to put attention on some other things instead of a beer?
Help-seeker	: Hmmm like what?
Model Generated Responses	
Joint	You can try to be normal. That is a great idea. <i>(Lacking empathy)</i>
Muffin _{Joint}	Yes, there are some things that you can do, but you have to be aware of your own mental health, that's why it's important to know you are doing the right thing.
MultiESC	It sounds like you have a lot on your plate and you want to be normal. <i>(Inconsistent)</i> I can understand how that would make you feel.
Muffin _{MultiESC}	I think you should focus on yourself instead of a beer. Do you have any hobbies that you enjoy?
KEMI	I think you should try to find some new things to do to relax.
Muffin _{TransESC}	Have you tried meditation? It can help you relax and relax.

Table 3: Comparison of responses generated by BlenderBot-Joint, MultiESC, and KEMI with their respective Muffin-enhanced counterparts.

compared to the comprehensive Muffin_{Joint} model. This underscores that considering multiple facets of emotional support when building emotional support conversation models (**D1**) is necessary. Another insight from Table 2 is that all three ablated models outperform the base model. It suggests that the proposed solution, directly mitigating unhelpful responses (**D2**), is reasonable and effective. Moreover, the results also indicate that the quality of helpful and unhelpful responses will influence the effects of unhelpfulness mitigation.

Case study To intuitively illustrate the superiority of Muffin over its base model, Table 3 presents a comparative case study, comparing responses generated by three prominent base models and their corresponding Muffin versions. From the comparison of BlenderBot-Joint and Muffin_{Joint}, we can observe that the BlenderBot-Joint implies that the

help-seeker can be abnormal now. Such a statement ignores the help-seeker's feelings, barely expressing empathy. For the comparison of responses generated by BlenderBot-Joint and Muffin_{Joint}, the former tends to state facts more directly, subtly implying that the help-seeker might be experiencing an abnormal state. Such a statement ignores the help-seeker's feelings, barely expressing empathy. In contrast, Muffin_{Joint} conveys concern for the help-seeker's well-being and attempts to solve the problem by shifting the help-seeker's perspective, amplifying the empathetic undertone. In the case of MultiESC, Muffin_{MultiESC} crafts a response that aligns more closely with the context, addressing the inconsistency of the response generated by MultiESC. Lastly, comparing KEMI with Muffin_{KEMI}, even though KEMI's response does not exhibit glaring issues, the Muffin version stands out as more beneficial. This distinction arises because Muffin_{KEMI}, in contrast to KEMI's general advice, offers a more specific and actionable recommendation, aligning closely with the help-seeker's request for precise advice.

6 Conclusion

In this work, we focus on mitigating unhelpful responses generated by recent emotional support conversation models. Such unhelpful responses, despite their well-intentioned nature, can inadvertently counteract prior supportive efforts. Analyzing the potential causes for unhelpful responses, we introduce a novel model-agnostic framework Muffin. In specific, it contains a multifaceted AI feedback module, which can discern helpful and unhelpful responses generated by a specific emotional support conversation model. Then Muffin contrasts helpful and unhelpful responses generated by this model, in order to reduce the likelihood of the unhelpful responses. Experimental results underscore Muffin's efficacy, showcasing enhancements in both automatic evaluations and human ratings. This suggests that Muffin can mitigate helpfulness in emotional support conversations.

598 Limitations

599 Although Muffin’s effectiveness is apparent, oppor-
600 tunities for further refinement exist. Our present ef-
601 forts address general cases of unhelpfulness, specif-
602 ically targeting responses generally perceived as
603 unhelpful by most help-seekers. Nonetheless, it
604 is crucial to acknowledge that certain responses
605 may adversely affect particular individuals under
606 specific circumstances. Consequently, this under-
607 scores the need for personalizing the emotional sup-
608 port conversation system to meet individual user re-
609 quirements. However, this personalization presents
610 substantial challenges. Collecting and processing
611 personal data raises serious privacy and security
612 concerns. Striking a balance between effective per-
613 sonalization and the potential for privacy intrusion
614 remains a delicate issue. Furthermore, reconciling
615 the goals of personalizing individual user experi-
616 ences with the need to generalize models across a
617 broader user base poses a fundamental conflict in
618 system design.

619 Ethical Considerations

620 In our experiments, we adopted open-sourced
621 datasets, including ESConv (Liu et al., 2021), em-
622 pathetic classification dataset (Sharma et al., 2020),
623 and MI dataset (Welivita and Pu, 2022). All per-
624 sonally identifiable information was removed from
625 these datasets. For the human ratings, we empha-
626 sized the comfort and well-being of our annota-
627 tors. Moreover, our research explores the develop-
628 ment of emotional support conversation systems.
629 Compared with existing methods, our proposed
630 approach represents a significant leap towards es-
631 tablishing a more secure emotional support conver-
632 sation framework. Consequently, we confidently
633 assert that our research is conducted in strict ad-
634 herence to the ethical guidelines prescribed by the
635 Association for Computational Linguistics (ACL).

636 References

637 Satanjeev Banerjee and Alon Lavie. 2005. **METEOR:**
638 **An automatic metric for MT evaluation with im-**
639 **proved correlation with human judgments.** In *Pro-*
640 *ceedings of the ACL Workshop on Intrinsic and Ex-*
641 *trinsic Evaluation Measures for Machine Transla-*
642 *tion and/or Summarization*, pages 65–72, Ann Arbor,
643 Michigan. Association for Computational Linguis-
644 tics.

645 Allan Barsky and Heather Coleman. 2001. Evaluating
646 skill acquisition in motivational interviewing: The

development of an instrument to measure practice
skills. *Journal of Drug Education*, 31(1):69–82. 647 648

C Daniel Batson. 2009. These things called empathy:
eight related but distinct phenomena. 649 650

Arthur C Bohart and Leslie S Greenberg. 1997. *Empa-*
thy reconsidered: New directions in psychotherapy.
American Psychological Association. 651 652 653

Brant R Burleson. 1985. The production of comforting
messages: Social-cognitive foundations. *Journal of*
language and social psychology, 4(3-4):253–273. 654 655 656

Brant R Burleson and Wendy Samter. 1985. Consis-
tencies in theoretical and naive evaluations of com-
forting messages. *Communications Monographs*,
52(2):103–123. 657 658 659 660

Yi Cheng, Wenge Liu, Wenjie Li, Jiashuo Wang, Ruihui
Zhao, Bang Liu, Xiaodan Liang, and Yefeng Zheng.
2022. **Improving multi-turn emotional support dia-**
logue generation with lookahead strategy planning.
In *Proceedings of the 2022 Conference on Empiri-*
cal Methods in Natural Language Processing, pages
3014–3026, Abu Dhabi, United Arab Emirates. As-
sociation for Computational Linguistics. 661 662 663 664 665 666 667 668

Alan S Cowen and Dacher Keltner. 2017. Self-report
captures 27 distinct categories of emotion bridged by
continuous gradients. *Proceedings of the national*
academy of sciences, 114(38):E7900–E7909. 669 670 671 672

Yang Deng, Wenxuan Zhang, Yifei Yuan, and Wai Lam.
2023. **Knowledge-enhanced mixed-initiative dia-**
logue system for emotional support conversations.
ArXiv preprint, abs/2305.10172. 673 674 675 676

Joseph L Fleiss. 1971. Measuring nominal scale agree-
ment among many raters. *Psychological bulletin*,
76(5):378. 677 678 679

John O Greene and Brant Raney Burleson. 2003. *Hand-*
book of communication and social interaction skills.
Psychology Press. 680 681 682

Michael U Gutmann and Aapo Hyvärinen. 2012. Noise-
contrastive estimation of unnormalized statistical
models, with applications to natural image statistics.
Journal of machine learning research, 13(2). 683 684 685 686

Clara E Hill. 2009. *Helping skills: Facilitating, explo-*
ration, insight, and action. American Psychological
Association. 687 688 689

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
Weizhu Chen. 2022. **Lora: Low-rank adaptation of**
large language models. In *The Tenth International*
Conference on Learning Representations, ICLR 2022,
Virtual Event, April 25-29, 2022. OpenReview.net. 690 691 692 693 694 695

Minlie Huang, Xiaoyan Zhu, and Jianfeng Gao. 2020.
Challenges in building intelligent open-domain di-
alog systems. *ACM Transactions on Information*
Systems (TOIS), 38(3):1–32. 696 697 698 699

700	WANG Jiashuo, Haozhao Wang, Shichao Sun, and Wenjie Li. 2023. Aligning language models with human preferences via a bayesian approach. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	758
701		759
702		760
703		761
704		762
705	Zaijing Li, Gongwei Chen, Rui Shao, Dongmei Jiang, and Liqiang Nie. 2024. Enhancing the emotional generation capability of large language models via emotional chain-of-thought. <i>ArXiv preprint</i> , abs/2401.06836.	763
706		764
707		765
708		766
709		767
710	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	768
711		769
712		770
713		771
714	Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2122–2132, Austin, Texas. Association for Computational Linguistics.	772
715		773
716		774
717		775
718		776
719		777
720		778
721		779
722	Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. Towards emotional support dialog systems. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3469–3483, Online. Association for Computational Linguistics.	780
723		781
724		782
725		783
726		784
727		785
728		786
729		787
730		788
731	Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022. BRIO: Bringing order to abstractive summarization. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2890–2903, Dublin, Ireland. Association for Computational Linguistics.	789
732		790
733		791
734		792
735		793
736		794
737		795
738	Susan P Llewelyn, Robert Elliott, David A Shapiro, Gillian Hardy, and Jenny Firth-Cozens. 1988. Client perceptions of significant events in prescriptive and exploratory periods of individual therapy. <i>British Journal of Clinical Psychology</i> , 27(2):105–114.	796
739		797
740		798
741		799
742		800
743	Lajanugen Logeswaran and Honglak Lee. 2018. An efficient framework for learning sentence representations. In <i>6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings</i> . OpenReview.net.	801
744		802
745		803
746		804
747		805
748		806
749	Janice M Morse, Joan Bottorff, Gwen Anderson, Beverley O’Brien, and Shirley Solberg. 1992. Beyond empathy: expanding expressions of caring. <i>Journal of advanced nursing</i> , 17(7):809–821.	807
750		808
751		809
752		810
753	Theresa B Moyers, Tim Martin, Jennifer K Manuel, William R Miller, and D Ernst. 2003. The motivational interviewing treatment integrity (miti) code: Version 2.0. Retrieved from <i>Verfügbar unter: www.casaa.unm.edu [01.03. 2005]</i> .	811
754		812
755		813
756		814
757		
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	
	Ana Paiva, Iolanda Leite, Hana Boukricha, and Ipke Wachsmuth. 2017. Empathy in virtual agents and robots: A survey. <i>ACM Transactions on Interactive Intelligent Systems (TiiS)</i> .	
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	
	Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. Recipes for building an open-domain chatbot. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 300–325, Online. Association for Computational Linguistics.	
	Ashish Sharma, Adam Miner, David Atkins, and Tim Althoff. 2020. A computational approach to understanding empathy expressed in text-based mental health support. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5263–5276, Online. Association for Computational Linguistics.	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>ArXiv preprint</i> , abs/2302.13971.	
	Quan Tu, Yanran Li, Jianwei Cui, Bin Wang, Ji-Rong Wen, and Rui Yan. 2022. MISC: A mixed strategy-aware model integrating COMET for emotional support conversation. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 308–319, Dublin, Ireland. Association for Computational Linguistics.	
	Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015</i> , pages 4566–4575. IEEE Computer Society.	
	Ashwin K Vijayakumar, Michael Cogswell, Ramprasaath R Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2016. Diverse beam search: Decoding diverse solutions from neural sequence models.	

815	Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin	demonstrates empathy through emotional re-	867
816	Guu, Adams Wei Yu, Brian Lester, Nan Du, An-	action.	868
817	drew M. Dai, and Quoc V. Le. 2022. Finetuned		
818	language models are zero-shot learners . In <i>The Tenth</i>		
819	<i>International Conference on Learning Representations</i> ,	• Interpretations: This approach reflects an un-	869
820	<i>ICLR 2022, Virtual Event, April 25-29, 2022</i> .	derstanding of feelings and experiences as in-	870
821	OpenReview.net.	ferred from the help-seeker’s statements. For	871
		instance, “I understand that you feel sad” ex-	872
822	Anuradha Welivita and Pearl Pu. 2022. Curating a large-	emplifies the supporter’s interpretations.	873
823	scale motivational interviewing dataset using peer		
824	support forums . In <i>Proceedings of the 29th Inter-</i>	• Explorations: Responses in this category ex-	874
825	<i>national Conference on Computational Linguistics</i> ,	plore feelings and experiences not explicitly	875
826	pages 3315–3330, Gyeongju, Republic of Korea. In-	stated by the help-seeker. An example is the	876
827	ternational Committee on Computational Linguistics.	close question “Why are you feeling alone	877
		right now?”	878
828	Ka Wong, Praveen Paritosh, and Lora Aroyo. 2021.	In the framework and associated dataset outlined	879
829	Cross-replication reliability - an empirical approach	by Sharma et al. (2020) , responses that solely offer	880
830	to interpreting inter-rater reliability . In <i>Proceedings</i>	advice (e.g., “ask your school counselor what re-	881
831	<i>of the 59th Annual Meeting of the Association for</i>	sources they can provide”), merely present factual	882
832	<i>Computational Linguistics and the 11th International</i>	information (e.g., “mindful meditation helps over-	883
833	<i>Joint Conference on Natural Language Processing</i>	come anxiety”), or are offensive or abusive (e.g.,	884
834	<i>(Volume 1: Long Papers)</i> , pages 7053–7065, Online.	“I don’t know what to do; I am also feeling suici-	885
835	Association for Computational Linguistics.	dal right now”) are considered non-empathetic and	886
		characterized as “no communication.” However, as	887
836	Weixiang Zhao, Yanyan Zhao, Shilong Wang, and Bing	noted by Greene and Burluson (2003) , information	888
837	Qin. 2023. TransESC: Smoothing emotional support	and suggestions can often be helpful, particularly	889
838	conversation via turn-level state transition . In <i>Find-</i>	when the help-seeker requires. Consequently, in	890
839	<i>ings of the Association for Computational Linguis-</i>	our work, responses containing pertinent informa-	891
840	<i>tics: ACL 2023</i> , pages 6725–6739, Toronto, Canada.	tion or suggestions are classified as “weak/strong	892
841	Association for Computational Linguistics.	communication” if they align with the help-seeker’s	893
		request in the context. To identify such responses,	894
842	Chujie Zheng, Pei Ke, Zheng Zhang, and Minlie Huang.	we utilize GPT-4 to detect contexts where the help-	895
843	2023. Click: Controllable text generation with se-	seeker explicitly (e.g., “any suggestions”) or im-	896
844	quence likelihood contrastive learning . In <i>Findings</i>	plicitly (e.g., “I don’t know what to do”) requests	897
845	<i>of ACL</i> .	information or suggestions. Subsequently, these	898
		responses are manually evaluated to determine if	899
846	A Multifaceted AI Feedback Module	they provide reasonable information or suggestions.	900
		For the classification task, the specific prompt is	901
847	This section provides more details about the multi-	depicted in Figure 6.	902
848	faceted AI feedback module.		
		Communication Skill Efficiency We adopt the	903
849	A.1 Multiple Facets of Emotional Support	Motivational Interviewing Treatment Integrity	904
		(MITI) code Moyers et al. (2003) , a well-	905
850	For further clarification and better understanding,	established behavioral coding system that differen-	906
851	we have provided additional examples of unhelpful	tiates between favorable and unfavorable responses,	907
852	response instances related to the emotional support	to categorize responses into three classes based on	908
853	facets we emphasize in our work. These examples	their communication skills:	909
854	can be found in Table 4. Furthermore, we have		
855	expanded on the explanations of helpfulness and	• MI Adherent: Responses in this category	910
856	unhelpfulness in relation to each facet and the cor-	support help-seekers with empathetic and	911
857	responding classification dataset.	compassionate statements, fostering a sense	912
		of being heard, respected, and understood. For	913
858	Empathetic Expression The empathetic expres-	example, the response “Well, there is really a	914
859	sion framework introduced by Sharma et al. (2020)		
860	encompasses three empathy communication mech-		
861	anisms, each illustrating a different approach to		
862	conveying empathy:		
863	• Emotional Reaction: This mechanism in-		
864	volves expressing warmth, compassion, and		
865	concern for the help-seeker. An example is		
866	the statement “I am here to help you,” which		

Facet	Context	Unhelpful Response
<i>Empathetic Expression</i>	<i>Help-seeker</i> : Just wanna sleep. Can't motivate at all. I don't know what to do.	<i>Supporter</i> : As Nas says, sleep is the cousin of death. It's the kindest form of oblivion.
	<i>Help-seeker</i> : What do you want to be when you grow up? I don't know what to do with my life or what I truly love when all I really want is not to exist.	<i>Supporter</i> : But what does fully grown up mean? I have my own house and a job etc. Am I fully grown up now?
<i>Communication Skill Efficiency</i>	<i>Help-Seeker</i> : I do nothing but think about how I'm worthless and how I shouldn't be here. I've never tried or contemplated suicide. I've always wanted to fix my issues, but I never get around to it. How can I change my feeling of being worthless to everyone?	<i>Supporter</i> : Still, you should watch out for this because it can be a sign of worsening depression.
	<i>Help-Seeker</i> : I try to have a "heart-to-heart" with my mom, but it usually winds up being messy and frustrating afterward. Thank you for your advice and words of encouragement. Sorry if that sounds silly.	<i>Supporter</i> : Stop being so afraid...you really have nothing to lose :) Well, you are right.
<i>Response Coherence</i>	<i>Help-seeker</i> : Is it afternoon where you are? I'm still waking up; I should have more coffee. And I am just always tired of everything these days.	<i>Supporter</i> : Could you please explain why you are enjoying yourself?
	<i>Supporter</i> : Yes, what happened? Why are you tired? Please can you explain? Everything going ok?	
	<i>Help-seeker</i> : I think it's my depression. I have a very hard time getting motivated to do anything. Nothing seems worth it.	
	<i>Help-seeker</i> : I am upset.	<i>Supporter</i> : Everyone can be happy. Can you tell me more?

Table 4: Unhelpful responses in terms of each facet of emotional support.

lot going on for you right now” qualifies as MI Adherent.

- **Others:** This class encompasses responses that might not immediately elicit a positive support effect. These include responses like a closed question, such as “Did you eat five fruits and vegetables this week?” In our work, we regard both these types of responses and those classified as MI Adherent as *helpful*. Our rationale is that such responses, though seemingly less impactful initially, can be beneficial in the long term. They contribute to the progression of the conversation and strengthen the relationship between the help-seeker and the supporter.
- **MI Non-Adherent:** This category includes responses that involve arguing, confronting, or offering unsolicited advice, which may lead to resistance and impede problem-solving for help-seekers. An example is the response “Yes, you are an alcoholic. You might not think so, but you are,” which is classified as

MI Non-Adherent.

To classify supporter statements based on communication efficiency, we utilize the MI dataset as proposed by [Welivita and Pu \(2022\)](#). Additionally, [Figure 7](#) illustrates the specific prompt employed for this task.

Response Coherence To develop the dataset for response classification in the aspect of response coherence, we initially selected approximately 4K context-response pairs from the ESConv training dataset as the foundation for constructing the coherence data. It is important to note that this training dataset represents the intersection of the training sets used by all base models. As mentioned in [Section 5.2](#), two types of incoherent variants were created: one by replacing the original response with a response from a different dialogue, and the other by altering keywords in the original response. For the latter approach, we utilized GPT-4 to modify keywords or crucial information, ensuring that the altered response either conveyed a contrary meaning or addressed a different topic. [Table 5](#) presents

```

### Instruction:
In the context of empathy, there are three key aspects to consider: (1) Emotional Reactions - expressing emotions like warmth, compassion, and concern that the peer supporter feels after reading the seeker's post; (2) Interpretations - conveying an understanding of the feelings and experiences inferred from the seeker's post; (3) Explorations - seeking a deeper understanding of the seeker by delving into feelings and experiences not explicitly stated in the post. Each aspect can exhibit varying degrees of communication—none, weak, or strong—based on the manner in which related content is expressed. The overall level of empathy is determined by the highest level achieved across these three aspects.
Your task is to identify the level of empathy in the Supporter's response within the provided conversation.

### Input:
Conversation Context: {context}
The last supporter statement: {response}
Identify the empathy level of the Supporter's response. Choose one of the following options: No Communication, Weak Communication, and Strong Communication.

```

Figure 6: The specific prompt utilized by the Multifaceted AI Feedback for classifying supporter statements considering empathetic expression.

an illustrative example of this process. The specific prompt employed for this task is illustrated in Figure 8.

Context	
<i>Help-seeker:</i> Is it afternoon where you are? I'm still waking up; I should have more coffee. And, I am just always tired, tired of everything these days.	
<i>Original Response</i>	Yes, What happened? Why are you tired? Please can you explain?
<i>Incoherent Responses</i>	I am glad to hear you are feeling a little better! Yes, bad management is so toxic. Even with great coworkers in a job you love, horrible management can ruin it quickly. (From another dialogue)
	Why are you asking me to explain when you are feeling exhausted? Can you share your thoughts on this matter? (Keywords changed)

Table 5: The instance of a coherent response and its two incoherent variants.

A.2 Model Tuning and Module Performance

We instruction-tuned LLaMA¹ to equip the model with the capability required for unhelpful response classification tasks.² Furthermore, we initialized the LoRA weights using a low-rank adapter that was fine-tuned on the Stanford Alpaca dataset.³ We evaluated the capabilities of GPT-3.5 and GPT-4 by invoking the gpt-3.5-turbo and gpt-4 models,

¹<https://huggingface.co/decapoda-research/llama-7b-hf>

²The implementation referred to is available at <https://huggingface.co/decapoda-research/llama-7b-hf>

³<https://huggingface.co/tloen/alpaca-lora-7b>

respectively, through their API⁴ during the period from December 2023 to January 2024. We set the temperature parameter to 0 to ensure deterministic output generation. The detailed results of this evaluation are detailed in Table 6.

GPT-3.5 refers to considering the supporter's response not unhelpful. Prior to instruction-tuning, LLaMA exhibited notably poor performance across all facets of the classification tasks. The model consistently favored specific classes, such as "Strong Empathy," "Other," and "Yes," corresponding to empathetic expression, communication efficiency in skill, and response coherence, respectively.

We also conduct human evaluation to examine the module. In particular, we randomly sample 200 responses from the ESConv dataset and model-generated responses. For each response, an annotator was asked to assess whether it is unhelpful. We compared the annotations with the multifaceted AI feedback, and the consistency rate was 88%. We found that the multifaceted AI module is stricter than human annotators. It is evident by the fact that for instances in which the multifaceted AI feedback and human annotations are different, the multifaceted AI feedback tends to consider the response unhelpful.

⁴<https://platform.openai.com/docs/api-reference/chat/create>

```

### Instruction:
Motivational Interviewing involves three distinct strategies. Each strategy can be described
as follows:

1. MI Adherent Strategies:
* Advising: Providing advice when directly requested by the Help-seeker. This may include
indirect forms of permission, such as when the Supporter says to disregard the advice as
appropriate.
* Encouraging: Offering positive remarks or compliments to the Help-seeker.
* Emphasizing Autonomy: Highlighting the Supporter's control, freedom of choice, and ability
to make decisions.
* Compassion Statements: Expressing sympathy or understanding.

2. MI Non-Adherent Strategies:
* Unsolicited Suggestions: Offering solutions or actions without the Supporter's prior
consent.
* Direct Disagreement: Explicitly disagreeing, arguing, blaming, criticizing, or questioning
the Supporter's honesty.
* Commands: Issuing orders or imperatives.
* Cautionary Statements: Warning of potential consequences or serving as a caution.

3. Other Strategies:
* Open Questions.
* Personal Disclosure: The supporter shares their own information or experiences.
* Close-ended Questions: Inquiries answerable with a simple 'yes' or 'no' or a limited set of
responses.
* Open-ended Questions: Questions that allow for a broad range of answers.
* Repetition/Rephrasing: Echoing, rewording, or paraphrasing the seeker's statements.
* Enhanced Repetition: Repeating or rephrasing the Supporter's statement with added emphasis
or meaning.
* Educational Feedback: Providing information, feedback, or opinions without giving direct
advice. Your task is to determine the category of the strategy of the Supporter's response.

### Input:
Conversation Context: {context}
The last supporter statement: {response}
Identify the strategy of the Supporter's response. Choose one of the following options: MI
Adherent, MI Non-Adherent, and Others.

```

Figure 7: The specific prompt used by the Multifaceted AI Feed- back for classifying the supporter’s statement regarding communication skill efficiency.

B Experimental Setup

B.1 Devices and Environment

We employed Pytorch⁵ for the implementation of all models, and conducted our experiments using Nvidia GeForce RTX 3090 GPUs.

B.2 ESConv Dataset

The ESConv dataset comprises a collection of emotional support conversations, each facilitated between a help-seeker and a supporter. While previous studies have partitioned and processed this dataset in various ways, our implementations of Muffin strictly follow each of their methodologies to maintain fairness. However, it is important to note that direct comparisons of performance across these models may not be entirely fair due to the differences in their data preprocessing approaches.

⁵<https://pytorch.org/>

B.3 Base Models

Vanilla is a vanilla BlenderBot (Roller et al., 2021) trained on the dataset ESConv. We used the small version⁶ of BlenderBot in experiments following previous studies.

Joint is built upon the backbone of BlenderBot and is specially trained to generate responses along with an expected communication skill at the beginning of each response (Liu et al., 2021)⁷. The focused facet of this model is skill efficiency.

MultiESC (Cheng et al., 2022)⁸ is an emotional support conversation model that mainly focuses on the communication skill efficiency. It predicts

⁶https://huggingface.co/facebook/blenderbot_small-90M

⁷<https://github.com/thu-coai/Emotional-Support-Conversation>

⁸<https://github.com/lwgkzl/multiesc>

Model		GPT-3.5	GPT-4	LLaMA (Vanilla)	LLaMA (Tuned)
<i>Empathetic Expression</i>	<i>accuracy</i>	81.88	69.58	19.09	83.50
	<i>F1 Score</i>	54.89	61.58	16.03	74.12
<i>Skill Efficiency</i>	<i>accuracy</i>	69.56	79.01	84.96	90.43
	<i>F1 Score</i>	53.22	66.14	45.93	81.83
<i>Response Coherence</i>	<i>accuracy</i>	33.47	78.76	33.31	92.98
	<i>F1 Score</i>	25.26	72.62	24.98	92.06
<i>Aggregated Feedback</i>	<i>accuracy</i>	57.11	78.01	59.08	90.72
	<i>F1 Score</i>	41.99	77.19	40.96	89.86

Table 6: Performance of various Language Model Models (LLMs) on unhelpful response classification in terms of different facets and the aggregated decision. All values are expressed as percentages (%).

```

### Instruction:
Determine if the supporter's response aligns coherently with the seeker's post. A coherent response should maintain a logical flow of ideas in correspondence with the post, often including supporting arguments or evidence directly related to the post's content.

### Input:
Conversation Context: {context}
The last supporter statement: {response}
Identify whether the supporter's last statement is coherent with the help-seeker's post. Answer 'Yes' if the supporter's response is coherent with the help-seeker's post, otherwise answer 'No'.

```

Figure 8: The specific prompt utilized by the Multifaceted AI Feedback for classifying the supporter’s statement in terms of response coherence.

the strategies (communication skills) in the next several turns considering the the user’s state.

KEMI incorporates various knowledge for a mixed-initiative conversation model, which can provide emotional support (Deng et al., 2023)⁹. In this process, the response coherence is also enhanced.

TransESC (Zhao et al., 2023)¹⁰ predicts the transitions of the user’s emotion, the communication skills, and the conversation keywords. Then, such information is used for response generation. This model takes into account more than one facet of emotional support; however, it is not optimized to reduce the likelihood of unhelpful responses like other models.

B.4 Human Evaluation

For each pair of $\mathcal{G}(\theta_0)$ and $\text{Muffin}_{\mathcal{G}(\theta_0)}$, we randomly selected 100 response pairs generated by both $\mathcal{G}(\theta_0)$ and $\text{Muffin}_{\mathcal{G}(\theta_0)}$ under identical conversational contexts. Selecting 100 responses aligns

with the standard practice for human evaluation in dialogue tasks, such as the experiments of TransESC (Zhao et al., 2023) and KEMI (Deng et al., 2023). To prevent annotators from identifying the generation model based on the order of sentences, the sequence in which these two responses are presented is randomized for each evaluation.

Initially, the annotators were briefed about the nature of emotional support conversations to ensure a comprehensive understanding of the task’s objectives. During the rating process, they were suggested to imagine themselves as the help-seekers within the conversations. After being provided with the conversation context, the annotators then proceeded to compare two generated responses, shown as Figure 9. Moreover, we prioritized the comfort and well-being of our annotators, advising them to pause or cease the annotation process if they encountered any content that made them feel uncomfortable. Annotators were paid at a rate of 1.5~2 times their local hourly minimum wage. Based on annotators’ feedback, it was estimated that approximately 40 seconds were spent on evaluating each response pair.

⁹<https://github.com/dengyang17/KEMI>

¹⁰<https://github.com/circle-hit/TransESC>

Here is a conversation between a help-seeker and a support. Imagine that you are the help-seeker, and compare the following two responses. If you find both responses to be equally effective or unsatisfactory, please indicate your assessment as a "Tie".

Conversation History:

Help-seeker: Hi i am okay, a little bit sad though.

Support: Okay. I am very sorry to hear that! Do you want to tell me more about that?

Help-seeker: Well with the holidays coming up I have been very stressed and nervous about what i am going to do.

Two supporter responses:

A. Of course! I am sorry to hear you are stressed about the holidays. Can I ask what are you worried about?

B. Of course! I am sorry to hear you are feeling stressed and anxious about the holidays.

1. Which response is more fluent (grammar errors and inappropriate repetition can decrease the fluency)? A, B or Tie?

>

2. Which response is more skillful in identifying the help-seeker's problem? A, B or Tie?

>

3. Which response is better at comforting the help-seeker? A, B or Tie?

>

4. Which response can give more helpful and information suggestions? A, B or Tie?

>

5. Which response is generally more helpful? A, B or Tie?

>

Figure 9: An example of an instance presented to annotators for evaluation.

The inter-rater agreement among annotators, as measured by Fleiss’s Kappa, is 0.39. This value is relatively high, especially when compared to inter-rater reliability values in most subjective tasks, which typically fall within the range of 0.2 to 0.6 (Wong et al., 2021; Cowen and Keltner, 2017). In the process of aggregating the annotations, we determine the winning response based on the consensus of annotators. If two annotators prefer response A, and another two annotators prefer a tie, we categorize response A as the winner. When an equal number of annotators favor two responses, we label the result as a tie.

C Response Generation Model Performance

C.1 Multifaceted AI Feedback

In addition to presenting the primary results and conducting an ablation study, we provide the multifaceted AI feedback results as a reference, as outlined in Table 7. For each model, we calculate the frequency of generating unhelpful responses on the ESConv testing set. Specifically, we utilize the multifaceted AI feedback module to identify the helpful (non-unhelpful) responses and compute their percentage, displayed in the left subtable. Furthermore, we analyze the helpful response percentage when each model generates ten responses using

diverse beam search, reported in the right subtable. It’s worth noting that the assertion “approximately 30% of responses generated by MultiESC on the ESConv benchmark are identified as unhelpful” is derived from the findings presented in this table.

Overall, Muffin demonstrates enhancements in AI feedback across multiple facets. However, three intriguing phenomena emerge. **(1)**. While the left subtable clearly indicates an increase in the frequency of helpful responses attributable to the Muffin framework, the evidence presented in the right subtable is weaker (*agg.*), particularly when the base model is TransESC, KEMI, or MultiESC. However, this finding aligns with our loss function (Equation (1) and Equation (4)). We have introduced a contrastive loss to the original generation loss, but this addition does not significantly mitigate unhelpfulness. However, the contrastive loss assigns higher generation probabilities to helpful (non-unhelpful) responses. Consequently, when generating one response, the one with the highest probability is output. Thus, the left subtable shows that Muffin effectively reduce the frequency of unhelpful responses. **(2)**. Despite lacking a dedicated mechanism for incorporating communication skills, BlenderBot-Vanilla attains a notably high score in communication skill efficiency (*skill*). Upon closer examination, we observe that this model frequently produces responses such as “I can understand that...”

A. Generating One Response					B. Generating Ten Responses				
Model	<i>emp.</i>	<i>skill</i>	<i>cohr.</i>	<i>agg.</i>	Model	<i>emp.</i>	<i>skill</i>	<i>cohr.</i>	<i>agg.</i>
Vanilla	81.22	90.43	80.69	64.83	Vanilla	81.90	92.03	78.89	64.51
Muffin _{Vanilla}	85.83	92.82	84.33	71.26	Muffin _{Vanilla}	89.84	96.05	82.60	73.78
Joint	80.61	88.93	80.65	61.48	Joint	79.31	90.56	79.75	63.30
Muffin _{Joint}	82.33	90.47	83.04	64.76	Muffin _{Joint}	83.43	92.26	80.02	66.54
TransESC	81.06	91.44	74.76	63.24	TransESC	74.16	92.78	43.93	41.21
Muffin _{TransESC}	81.28	91.49	78.67	66.02	Muffin _{TransESC}	74.20	92.92	44.15	41.43
KEMI	83.01	88.47	85.76	70.90	KEMI	81.16	87.49	81.57	66.23
Muffin _{KEMI}	83.40	88.68	87.15	72.33	Muffin _{KEMI}	81.39	87.64	81.61	66.47
MultiESC	83.39	90.66	85.38	70.00	MultiESC	78.61	89.67	77.33	60.83
Muffin _{MultiESC}	85.24	92.92	86.06	72.38	Muffin _{MultiESC}	78.77	89.73	77.26	60.87

Table 7: The AI feedback provided is sourced from the comprehensive AI feedback module. The left subtable displays the percentage of helpful (non-unhelpful) responses produced by each model employing the decoding strategy of the base model. Conversely, the right table showcases the percentage of helpful responses generated by each model utilizing diverse beam search with a beam size set at 10. *emp.*, *skill*, *cohr.* and *agg.* represent empathetic expression, skill efficiency, response coherence and aggregated feedback, respectively. All values are expressed in percentages (%), where higher percentages signify superior performance. The values in the left subtable are statistically significant with $p < 0.05$.

1125 or “I’ve experienced something similar...”. These
1126 responses are categorized as self-disclosure. While
1127 this strategy doesn’t undermine previous emotional
1128 support efforts and remains undetected as “unhelp-
1129 ful” by the AI feedback module, a conversation
1130 dominated by self-disclosure responses may not
1131 fully align with the objective of providing support.
1132 (3). Another noteworthy observation is that Trans-
1133 ESC receives exceptionally low AI feedback in
1134 response coherence (*cohr.*), despite its incorpora-
1135 tion of external knowledge. However, insights from
1136 TransESC’s ablation experiment (Zhao et al., 2023)
1137 suggest that the generated knowledge may contain
1138 noise. This potential factor could explain the low
1139 feedback score observed in response coherence.

1140 C.2 Cases

1141 For a closer and more detailed examination of our
1142 method, we present additional comparisons be-
1143 tween the base model and the Muffin model in
1144 Table 8. Additionally, we include AI feedback for
1145 each facet as a reference.

Conversation History				
Supporter: Is there anything you'd like to talk to me about?				
Help-seeker: I have a friend who knew me since I was little. She is a little older than me. She still treat me like I'm still little. I tried to explain to her that I'm a grown woman and I grew out of lot of things when I was a kid.				
Model	Response	<i>emp.</i>	<i>skill</i>	<i>cohr.</i>
Joint	I have had a friend like that before. They are a good friend, and you should feel as though you are a better person.	Weak Empathy	Other	No
Muffin _{Joint}	I have been in a similar situation. I had a friend who was a little older than me, and I had to deal with the same thing.	Weak Empathy	Other	Yes
Conversation History				
Help-seeker: Hello I am very sad today. My partner is cheating on me. I was told by a customer at work. I do not know this person well, yet they knew the person he was cheating on me with. He will not admit it himself.				
Supporter: I am sorry if your partner is cheating. So, your partner denies the charge of cheating?				
Help-seeker: Yes. He says the person who told me is lying. We are not married but we live together.				
Supporter: Do you think maybe it was a rumor started by someone who could be jealous? Ash your partner if he loves you enough to propose marriage. If he says no, then ask him why.				
Help-seeker: I am not sure if I am okay with giving him an ultimatum. We are still young and although I love being in a relationship, I am unsure about marriage. But I appreciate your suggestions. And you make a great point about the workplace being hostile.				
Model	Response	<i>emp.</i>	<i>skill</i>	<i>cohr.</i>
MultiESC	I think you should give him an ultimatum and see if he accepts it.	No Empathy	MI Non-Adherent	Yes
Muffin _{MultiESC}	Have you spoken to your friends or family about the situation?	Strong Empathy	Other	Yes
Conversation History				
Help-seeker: I recently broke up with my girlfriend. We had been together for more than 3 years. She said that she simply wasn't in love anymore, and that she didn't have the same feelings she did at the beginning.				
Supporter: You still have your whole life ahead of you. Some people don't find a partner until their 40's or even 50's... I think that's what your girl would want you to do as well.				
Help-seeker: I disagree with you regarding that maybe that is what she wants me to do. I don't think she cares about me at the slightest anymore. I was a burden for her and will probably keep being it for any woman I meet.				
Model	Response	<i>emp.</i>	<i>skill</i>	<i>cohr.</i>
KEMI	That is a great attitude to have. You are a strong person, and I think you will find someone better.	Weak Empathy	MI Non-Adherent	Yes
Muffin _{KEMI}	I think you are a very good person and I think that you will find someone much better.	Weak Empathy	MI Adherent	Yes

Table 8: Comparison of responses generated by BlenderBot-Joint, MultiESC, and KEMI with their respective Muffin-enhanced counterparts.