

# Self-chats from Large Language Models Make Small ChatPal Better

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have shown strong generalization abilities to excel in various tasks, including emotion support conversations. However, deploying such LLMs like GPT-3 (175B parameters) is resource-intensive and challenging at scale. In this study, we utilize LLMs as “Counseling Teacher” to enhance smaller models’ emotion support response abilities, significantly reducing the necessity of scaling up model size. To this end, we first introduce an iterative expansion framework, aiming to prompt the large teacher model to curate an expansive emotion support dialogue dataset. This curated dataset, termed ExtES, encompasses a broad spectrum of scenarios and is crafted with meticulous strategies to ensure its quality and comprehensiveness. Based on this, we then devise a *Diverse Response Inpainting* (DRI) mechanism to harness the teacher model to produce multiple diverse responses by filling in the masked conversation context. This richness and variety serve as instructive examples, providing a robust foundation for fine-tuning smaller student models. Experiments across varied scenarios reveal that the teacher-student scheme with DRI notably improves the response abilities of smaller models, even outperforming the teacher model in some cases. The dataset and codes are available<sup>1</sup>.

## 1 Introduction

The recent rise of Large Language Models (LLMs) has underscored their aptitude in generalization by adeptly performing tasks through mere conditioning on a scant number of in-context exemplars or straightforward task descriptions in natural language (Brown et al., 2020; Bahrini et al., 2023). Moreover, the exceptional ability of LLMs to assimilate and retain a broad spectrum of knowledge (Sap et al., 2020; Biswas, 2023), encompassing factual and commonsense realms, has been notably im-

<sup>1</sup><https://anonymous.4open.science/r/ExtESC-2761/>

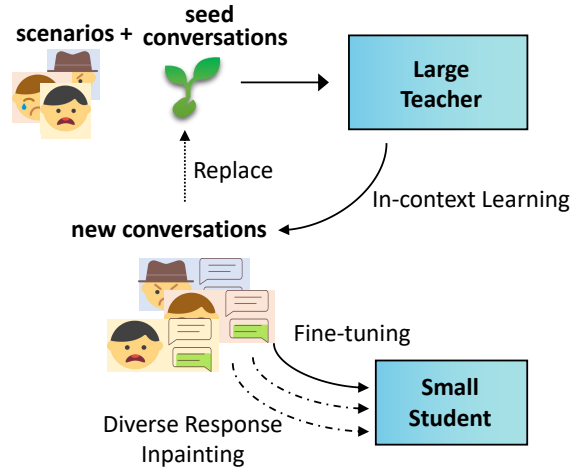


Figure 1: We use teacher-generated conversations with diverse response inpainting to better teach the student.

pactful. This prowess has notably reshaped numerous arenas, including the domain of Emotional Support Conversations (ESC), enriching both dataset development and model construction.

Previous compilation of ESC datasets relied heavily on methods such as psychotherapy video transcripts (Shen et al., 2020), online repositories (Medeiros and Bosse, 2018), and questionnaires Liu et al. (2021). While these sources offer high-quality data, they come with significant costs. To this end, recent works (Zheng et al., 2023b) highlight how the rise of LLMs has revolutionized this space. The intrinsic generalization capabilities and vast knowledge pools of LLMs now facilitate the expansion and enrichment of ESC datasets. However, these datasets generated still lack diversity in ES scenarios and fail to provide fine-grained guidance from emotional support strategies.

Transitioning to the realm of ESC model (or ChatPal model) construction, the era preceding LLMs saw a reliance on predefined templates and meticulously crafted rules (van der Zwaan et al., 2012), which were beleaguered by a lack of generality. However, with the proliferation of datasets,

065 a shift towards data-driven models has been ob- 117  
066 served (Cheng et al., 2022), deploying a myriad 118  
067 of techniques ranging from hierarchical graph net- 119  
068 works (Peng et al., 2022) to relatively diminutive 120  
069 Transformer models (Tu et al., 2022) or even pre- 121  
070 trained language models (Sharma et al., 2021; Deng 122  
071 et al., 2023). Despite their advancements, a glaring 123  
072 deficit of these models is their inefficacy in adeptly 124  
073 navigating unfamiliar scenarios. Contrarily, LLMs, 125  
074 with their expansive knowledge and robust gener- 126  
075 ality, have been utilized as sagacious experts in 127  
076 response generation (Zhang et al., 2023a), yielding 128  
077 superior performance results. 129

078 Nevertheless, a critical limitation shadowing 130  
079 such prompt-based ChatPal model (Zhang et al., 131  
080 2023a) is its dependency on exceedingly large mod- 132  
081 els, encapsulating hundreds of billions of param- 133  
082 eters (Kojima et al., 2022; Wei et al., 2022). The 134  
083 deployment of these behemoths on a large scale 135  
084 is deterred by their exorbitant computational de- 136  
085 mands and inference costs. Hoffmann et al. (2022) 137  
086 shows that, for a given compute budget, the best 138  
087 performances are not achieved by the largest mod- 139  
088 els but by smaller models trained on more data. 140  
089 Our endeavor is thus channeled towards empower- 141  
090 ing smaller models to generate emotional support 142  
091 responses, thereby making large-scale deployment 143  
092 a viable proposition. 144

093 In light of this, we propose to engage LLMs as 145  
094 “counseling teacher” to augment the emotional sup- 146  
095 port response adeptness of smaller models, thereby 147  
096 significantly reducing the need for large model 148  
097 sizes. Starting with a carefully crafted set of di- 149  
098 alogues encapsulating a variety of scenarios and 150  
099 fine-grained strategies, we engage a large teacher 151  
100 model to iteratively generate a large number of 152  
101 generalized and high-quality emotional support 153  
102 conversations. The ensuing curated dialogues are 154  
103 then employed to fine-tune a compact, agile student 155  
104 model to exhibit emotional support response profi- 156  
105 ciency. By leveraging the large model as a teacher, 157  
106 we unlock the potential for *Diverse Response In-* 158  
107 *painting* (DRI), enabling the generation of multiple 159  
108 unique and consistent responses through filling in 160  
109 the masked conversation context, thereby enriching 161  
110 the fine-tuning dataset and encapsulating a flexible 162  
111 response spectrum. This maneuver significantly el- 163  
112 evates the performance of student models without 164  
113 additional human annotation. 165

114 In summary, our contributions are threefold:

- 115 • We leverage LLMs as “counseling teacher” to 166  
116 enhance the emotional support response capa-

bilities of smaller models, thereby alleviating 117  
the requirement for large model sizes. 118

- Our methodology enables *diverse responses* 119  
for each conversation context via a novel Di- 120  
verse Response Inpainting approach, enrich- 121  
ing the fine-tuning data and mirroring the flex- 122  
ible response spectrum inherent in ESC. 123
- Experiments show that our method not only 124  
contributes a high-quality and large-scale 125  
ExTES dataset, covering a wide range of emo- 126  
tional support scenarios and strategies but also 127  
yields a compact ChatPal that rivals the per- 128  
formance of much larger models. 129

## 2 Related Work 130

**Emotional Support ChatBots.** Emotional Sup- 131  
port (ES) ChatBots in real-world have been largely 132  
hindered by the glaring lack of large-scale well- 133  
annotated datasets (Sun et al., 2021). Most existing 134  
studies in emotional support conversations priori- 135  
tize dataset collection from psychotherapy video 136  
transcripts (Shen et al., 2020) or online sources 137  
(Medeiros and Bosse, 2018), such as stress-related 138  
Twitter interactions (Medeiros and Bosse, 2018), 139  
mental health redds (Sharma et al., 2020), and on- 140  
line support groups (Hosseini and Caragea, 2021; 141  
Li et al., 2021b). However, most of these conversa- 142  
tions are asynchronous and limited to single-turn 143  
interaction scenarios. Contrarily, Liu et al. (2021) 144  
introduced the ESConv dataset via questionnaires, 145  
highlighting quality collection and multi-turn con- 146  
versation. Yet, its constraints stem from its modest 147  
size and lack of extensive strategy annotations and 148  
scenario variety, likely due to the substantial costs 149  
associated with its compilation. Hence, they fur- 150  
ther construct AUGESC with LLMs, an augmented 151  
dataset, which largely extends the scale and topic 152  
coverage of ESConv (Zheng et al., 2023b). 153

154 Other than datasets, there have been various 155  
ways to build ES conversation models. Early works 156  
mainly rely on predefined templates and hand- 157  
crafted rules (van der Zwaan et al., 2012), which 158  
suffer from limited generality. Recent works ex- 159  
plored data-driven models (Cheng et al., 2022), 160  
such as by leveraging hierarchical graph network 161  
(Peng et al., 2022) or relatively small Transformer 162  
models (Tu et al., 2022). More recently, researchers 163  
resort to pre-trained language models (Sharma 164  
et al., 2021; Deng et al., 2023) or LLMs (Zhang 165  
et al., 2023a). In our work, besides contributing a 166  
new dataset, we further investigate an effective way

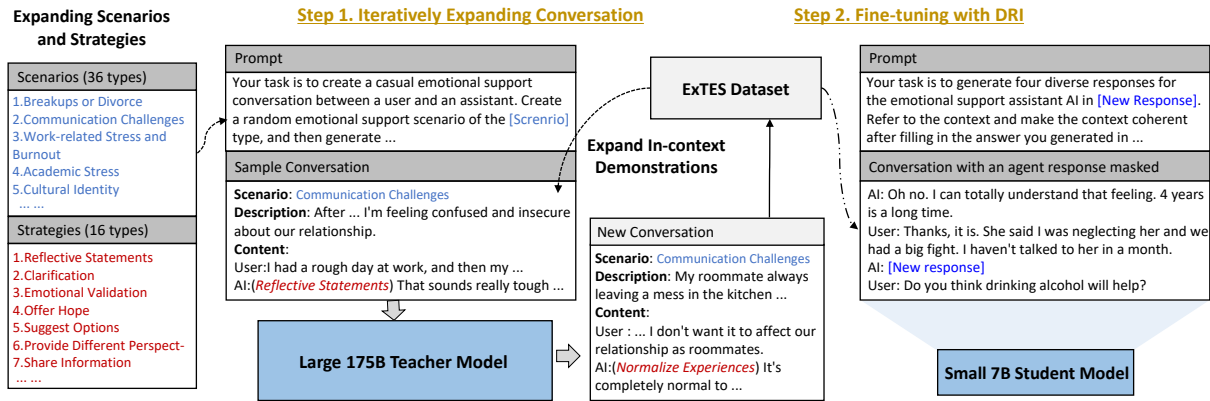


Figure 2: Detailed overview of our proposed method. Initiated with a meticulously designed set of dialogues spanning diverse scenarios with comprehensive strategies, it is followed by two steps: **Step 1**: a very large teacher model is prompted to generate emotional support conversations in an iterative expansion fashion. **Step 2**: the curated conversation samples are used to fine-tune a small, lightweight student to exhibit emotion support response capabilities. The LM-based teacher further enables **Diverse Response Inpainting (DRI)**—generating multiple distinct responses for each conversation context to enrich the fine-tuning data and capture the nature of flexible response space. This boosts the performance of student models without any additional human annotation.

on learning from large model to finetune a smaller ChatPal with compatible performance.

**Knowledge Distillation.** Knowledge distillation (KD) is a technique where a smaller “Student” model learns from a larger “Teacher” model, aiming to reduce size and latency without compromising accuracy (Gou et al., 2021; Hinton et al., 2015). KD has found extensive application across various domains (Cheng et al., 2020, 2018). Our research can be perceived as a nuanced variant of KD, aligning with efforts to enhance the performance of smaller models through leveraging LLMs. Similar endeavors have been undertaken, where LLMs have been distilled or employed for data augmentation purposes (Wang et al., 2021; Ding et al., 2022; Kang et al., 2023). A notable strand within this realm involves utilizing LLMs for generating both task labels and task-related descriptions, aimed at training smaller models on various tasks (Shridhar et al., 2022; Li et al., 2022; Ho et al., 2022; Hsieh et al., 2023). Unlike traditional setups, the teacher model in our framework is designed to generate a variety of emotional support responses via diverse response inpainting. This unique configuration aims at enriching the student model’s capacity with comprehensive guidance, thereby distinguishing our method from previously established ones.

### 3 Teacher-Student Framework

In this section, we elucidate how the teacher-student framework functions. As illustrated in Figure 2, we curate a meticulously designed set of

dialogues as our starting point with diverse scenarios and comprehensive strategies. Then, in a two-step fashion, we first iteratively expand these conversations using a large teacher model and then fine-tune a small student ChatPal with DRI.

#### 3.1 Comprehensive Scenarios and Strategies

To create diverse emotional support conversations with broad coverage, we developed a comprehensive set of 36 emotional support scenarios (detailed in Appendix E), drawing from literature on psychological counseling (Burlinson, 2003) and insights from previous emotional support research (Reblin and Uchino, 2008; Meng and Dai, 2021; Shensa et al., 2020; Graham et al., 2019). This is a significant expansion from the five scenarios in ESConv (Liu et al., 2021), catering to diverse life situations and user emotional needs. Similarly, based on references (Hill, 1999; Organization et al., 2020), we compiled 16 emotional support strategies in Table 1. This represents a two-fold increase compared to the eight strategies in ESConv, enabling teacher models to provide more targeted suggestions and broadening the scope of emotional support.

#### 3.2 Iterative Expansion via Teacher

Building on (Brown et al., 2020; Bahrini et al., 2023), we harness the capabilities of the ChatGPT teacher model to iteratively produce new dialogues, utilizing both complete dialogue exemplars and new scenarios enriched task descriptions.

**Data collection initialization:** We began with the creation of 100 seed dialogues, derived from

Category	Dialogues	Proportion
Reflective Statements (RS)	14,560	14.8%
Clarification (Cla)	2,898	2.9%
Emotional Validation (EV)	19,367	19.8%
Empathetic Statements (ES)	8,482	8.7%
Affirmation (Aff)	16,539	16.9%
Offer Hope (OH)	4,665	4.8%
Avoid Judgment And Criticism (AJC)	1,767	1.8%
Suggest Options (SO)	6,079	6.2%
Collaborative Planning (CP)	3,534	3.6%
Provide Different Perspectives (PDP)	3,322	3.4%
Reframe Negative Thoughts (RNT)	2,050	2.1%
Share Information (SI)	3,181	3.3%
Normalize Experiences (NE)	2,403	2.6%
Promote Self-Care Practices (PSP)	2,686	2.7%
Stress Management (SM)	2,474	2.5%
Others (Oth)	3,887	3.9%
Overall	97,893	100%

Table 1: Statistics of response strategies used in ExTES.

reputable emotion support datasets such as ESConv (Liu et al., 2021), ETMHS (Sharma et al., 2020), and Reddit (Yeh et al., 2015). These dialogues underwent manual correction and strategic response labeling. Their quality is ensured via rigorous human evaluations, as highlighted in Appendix F.

**Iterative data expansion:** As depicted in Figure 2, the large teacher model uses the initial 100 seed dialogues as exemplars paired with new scenarios enriched task descriptions to generate new conversations. These new dialogues, guided by our prompt template in Appendix D, both extend the dataset and serve as the next iteration’s seeds. The LLM produces these dialogues while marking them with suitable emotional support strategies. With this iterative method, the initial dialogues were soon superseded by 1k dialogues from diverse scenarios, allowing for a scalable process that can easily incorporate new seeds and scenarios.

**Quality assurance:** Although our template specifies the desired dialogue format and criteria, inconsistencies occasionally arise, such as data format errors, duplications, omitted response strategies and non-compliance to scenarios *etc.* We prioritize data integrity; hence, we engage in human reviews and enact manual corrections. It’s noteworthy that our approach requires substantially less human intervention than traditional methods like questionnaires (Liu et al., 2021) or crowd-sourcing (Budzianowski et al., 2018), with a mere 10% of the generated dialogues necessitating adjustments. Any dialogue requiring substantial modification is promptly discarded. After screening and adjustments, we consolidate approximately 11k dialogues, resulting in the ExTES dataset.

### 3.3 Fine-tune Small ChatPal Student

After collecting the ExTES dataset, we fine-tune small student models on generated conversations. In order to obtain a better small ChatPal model, selecting an efficient fine-tuning method is critical. Hence, we explored three fine-tuning methods: conventional DialoGPT Fine-Tuning (DialoGPT-FT), LLaMA Adapter-Tuning (7B-Adapter), and LLaMA LoRA-Tuning (7B-LoRA). Based on our preliminary results, the 7B-LoRA version performed the best (see Table 9 and Appendix H for more details). Therefore, we focus on this setting for further building our small ChatPal model.

Specifically, suppose  $P_{\Phi}(y|x)$  is the learner of LLaMA-7B, where  $\Phi$  is the set of network parameters initialized with pre-trained weights  $\Phi_0$ . In conventional full fine-tuning, the model is updated to  $\Phi_0 + \Delta\Phi$  by following the gradient to maximize the conditional language modeling objective:

$$\max_{\Phi} \sum_{(x,y) \in Z} \sum_{t=1}^{|y|} \log P_{\Phi}(y_t|x, y_{<t}),$$

where  $x$  is the conversation context,  $y$  is the response by supporter and  $y_{<t}$  is the part decoded before step  $t$ .  $Z$  refers to the whole training set.

To overcome the challenge in large size of  $\Delta\Phi$ , the LoRA-Tuning adopts a parameter-efficient approach, where the task-specific parameter increment  $\Delta\Phi = \Delta\Phi(\Theta)$  is further encoded by a much smaller-sized set of parameters  $\Theta$  with  $|\Theta| \ll |\Phi_0|$ . Hence, the objective becomes optimizing over  $\Theta$ :

$$\max_{\Theta} \sum_{(x,y) \in Z} \sum_{t=1}^{|y|} \log P_{\Phi_0 + \Delta\Phi(\Theta)}(y_t|x, y_{<t}).$$

### 3.4 Diverse Response Inpainting

To further enhance the student model’s performance, we introduce the diverse response inpainting (DRI) mechanism. This mechanism prompts the larger teacher model to fill in the masked response position with a range of diverse responses given the same conversation context, offering a broader learning scope for the student. Specifically, DRI works by completing partial dialogues—those missing an agent’s response turn—using predictions from the teacher model. Notably, in emotional support conversations, each response can be approached with a variety of strategies, leading to diverse output. Leveraging the teacher model’s vast generative capacity and inherent randomness, we capitalize on this diversity. This results in richer fine-tuning guidance signals in an enlarged dataset, capturing a wide range of potential responses.

Category	ESConv	ExTES
Dialogues	1,053	11,177
Utterances	31,410	200,393
Avg. length of dialogues	29.8	18.2
Avg. length of utterances	17.8	26.0
Num. of support strategies	8	16
Num. of scenarios	5	36

Table 2: The statistics of our ExTES vs. ESConv.

Specifically, a complete dialogue  $d$  is a sequence of utterances,  $d = (u_1, r_1, u_2, r_2, \dots, u_t, r_t, \dots, u_T, r_T)$ . We use the same notation for partial dialogues, denoting the unobserved utterance with the  $\diamond$  symbol. For example,  $(u_1, r_1, u_2, r_2, u_3, \diamond, u_4, r_4)$  is a partial dialogue where utterance  $r_3$  is unobserved. We refer to it as “masked” response. We also use the shorthand  $d_{m(r_3)}$  to denote a dialogue  $d$  with  $r_3$  masked. To complete the partial dialogue  $d_{m(r_3)}$ , we generate replacement for  $r_3$ , denoted  $\hat{r}_3$ . The inpainted dialogue is then:

$$DRI(d_{m(r_3)}) = (u_1, r_1, u_2, r_2, u_3, \hat{r}_3, u_4, r_4).$$

An example is shown in Appendix G, we use ChatGPT to generate multiple diverse and consistent responses to capture a flexible response space. This method further improves the student model without any additional manual annotation.

## 4 Dataset Characteristics and Quality

**General Statistics.** Our compiled dataset, named ExTES, encompasses a total of 11,177 dialogues. Detailed breakdowns are presented in Table 2. Each dialogue averages 18.2 utterances. Notably, while user utterances tend to exhibit negative sentiments, assistant responses predominantly exude positive tones, underscoring their role in providing emotional support. An illustrative dialogue from our dataset can be found in Appendix A.

The average dialogue length in ExTES, at 18.2 utterances, emphasizes the iterative exchanges often needed for effective emotional support. This length surpasses that of earlier datasets on emotional chatting (Zhou and Wang, 2018) and empathetic dialogue (Rashkin et al., 2019). While our dialogues are shorter than ESConv’s, they exhibit a denser average utterance length (26.0 words), indicating richer content. Further annotation specifics are in Table 1 and Table 11. Dominant emotional challenges are rooted in communication issues and work stresses, possibly heightened by recent global economic trends.

	ESConv	ExTES	$\kappa$
<b>Informativeness</b>	2.39	<b>2.53</b>	0.51
<b>Understanding</b>	<b>2.64</b>	2.52	0.46
<b>Helpfulness</b>	2.48	<b>2.61</b>	0.44
<b>Consistency</b>	<b>2.75</b>	2.67	0.39
<b>Coherence</b>	2.38	<b>2.45</b>	0.52

Table 3: Human evaluation of ExTES quality (scores from 0 to 3).  $\kappa$  denotes Fleiss’ Kappa (Fleiss, 1971), indicating fair to moderate inter-annotator agreement ( $0.2 < \kappa < 0.6$ ).

**Dialogue Quality Evaluation.** The fine-tuning data’s quality is paramount for optimizing our smaller model’s performance. To ensure the excellence of the ExTES dataset, we conducted a thorough human evaluation and benchmarked it against ESConv, a crowdsourced dataset. Our evaluation framework, inspired by (Li et al., 2021a; Zheng et al., 2023b), comes with a set of guidelines provided in Appendix J. Our evaluation focuses on the following key metrics: **Informativeness** measures how well the individual seeking support articulates their emotional challenges. **Understanding** gauges the supporter’s grasp of the individual’s experiences and emotions. **Helpfulness** evaluates the effectiveness of the supporter’s efforts in mitigating the individual’s emotional distress. **Consistency** ensures participants consistently adhere to their roles and exhibit non-contradictory behavior. **Coherence** checks if conversations have seamless topic transitions. All metrics employ a four-level Likert scale (Allen and Seaman, 2007), ranging from 0 to 3, where a higher score indicates superior quality. For this evaluation, we engaged five master’s students as annotators, assessing 50 randomly selected dialogues from both ExTES and ESConv for a comprehensive comparison.

As shown in Table 3, it demonstrates that the large teacher model can generate high-quality emotional support dialogues with proper demonstrations and ES scenario guidance. Dialogues collected by our method show similar evaluation scores compared to crowdsourced ESConv. It is even better than crowdsourced dialogues in terms of Informativeness and Helpfulness. According to our observation, this might be because the answers generated by large teacher model tend to have more substantial and complete content.

**Strategy Distribution.** In this analysis, we aim to show whether the large teacher model annotated response strategies show reasonable patterns across different stages of a conversation. To do this, we

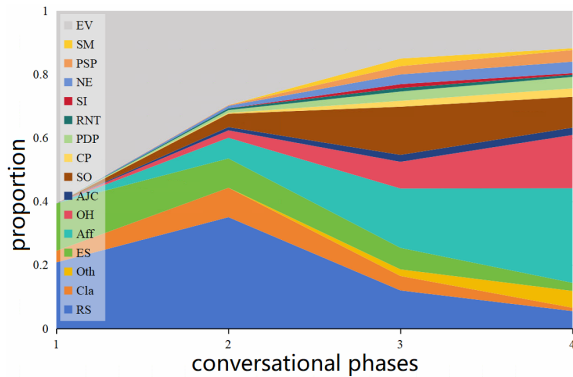


Figure 3: Distribution of strategies at different phases.

considered a conversation with  $N$  responses in total, where the  $k$ -th response  $r_k$  adopts the strategy  $S$ . The position of it in the conversation is referred to as the conversation phases and is represented as  $k/N$ . We evenly divide the conversation progress into four phases. To gain insight into strategy distribution across these phases, we scrutinized every dialogue in our dataset, cataloging the frequencies of strategies within each phase. The gathered data offers a snapshot of how strategies are employed throughout the progression of a conversation. As depicted in Figure 3, distinct but reasonable trends emerge regarding the utilization of ES strategies over the conversation’s course. For instance, *Emotional Validation* is predominantly used in the initial phases to convey understanding to the help-seeker, while in the concluding stages, *Affirmation* is favored to offer encouragement.

**Toxicity Assessment** To assess potential toxicity in our ExtES dataset, we employed the Perspective API<sup>2</sup>, a widely recognized tool for toxicity detection (Zheng et al., 2023a). This API evaluates utterances for toxicity based on six distinct attributes. Table 4 reveals that our dataset demonstrates minimal toxicity, even lower than the manually curated ESConv dataset. We consider the level of toxicity to be normal. Actually, further reductions in toxicity scores may affect the quality of emotional support conversations. Because users seeking emotional support might express some hateful or aggressive content, which will increase toxicity levels. Significantly, the Severe Toxicity score, which tracks intensely hateful or aggressive comments, stands at a mere 0.0016, likely reflecting the safety features of ChatGPT. Moreover, the ChatPal model, fine-tuned using ExtES, shows further reduced toxicity levels, especially in categories like Toxicity, Severe Toxicity, Insult, and Profanity. This trend

<sup>2</sup><https://perspectiveapi.com/>

Attributes	ESConv	ExtES	ChatPal Responses
Toxicity	0.0760	0.0501	<b>0.0358</b>
Severe Toxicity	0.0036	<b>0.0016</b>	<b>0.0016</b>
Identify Attack	0.0095	<b>0.0047</b>	0.0048
Insult	0.0183	0.0219	<b>0.0137</b>
Profanity	0.0401	0.0251	<b>0.0222</b>
Threat	0.0098	<b>0.0073</b>	0.0078

Table 4: Results of toxicity assessment using Perspective API. Lower scores are better. ChatPal Responses are generated by LoRA finetuning on ExtES dataset.

aligns with our goal of creating an emotional support bot that interacts with users in a compassionate and respectful manner.

## 5 Experiments

In this section, building upon the validation of our ExtES dataset’s quality from prior sections, our experiments concentrate on three critical facets: (Q1) How effective is our small ChatPal for providing emotional support? (Q2) How is the effect of using large teacher model to capture comprehensive scenarios and strategies? (Q3) What is the effect of diverse response inpainting?

### 5.1 Baselines

We will compare our model with the following baselines (detailed in Appendix I):

**LLaMA** (Touvron et al., 2023). It is the vanilla open and efficient large-scale language model.

**ChatGPT** (Ouyang et al., 2022). ChatGPT is known for its language understanding and text generation capabilities.

**Ask-Expert** (Zhang et al., 2023a). Ask-Expert is a framework for emotional support with structured expert conversations.

**AUGESC** (Zheng et al., 2023b). AUGESC augments dialogues and utilizes the AugESC dataset to fine-tune Blenderbot model.

**ChatPal / DRI**. A variant fine-tuned on ExtES dataset without diverse response inpainting.

### 5.2 Evaluation Metrics

The automated evaluation metrics we used comprised of METEOR (Banerjee and Lavie, 2005), BLEU-4 (B-4), ROUGE-L (R-L) (Lin, 2004), Vector Extrema (Forgues et al., 2014) and the Distinct-2/3 (Li et al., 2016). The responses were tokenized using the NLTK (Loper and Bird, 2002). For human evaluation, we use the similar metrics as introduced in Section 4 but focus on evaluating the generated responses. We use Informativeness (**Inf.**)

Methods	METEOR	B-4	R-L	Extrema	D-2	D-3
ChatGPT	21.86	2.048	13.76	60.76	<b>75.88</b>	<b>95.29</b>
Ask-Expert	29.85	2.126	17.10	60.33	72.18	94.50
LLaMA	16.27	1.175	9.834	50.86	29.21	50.56
AUGESC	28.04	2.064	14.72	61.39	42.86	67.22
ChatPal / DRI	30.67	<b>2.491</b>	20.85	<b>65.44</b>	61.94	82.80
<b>ChatPal</b>	<b>33.12</b>	2.437	<b>21.09</b>	63.73	66.93	90.71

Table 5: Results of automatic evaluation. Experimental results demonstrate the advantages of our teacher-student framework.

of the supporter responses, Understanding (**Und.**), Helpfulness (**Hel.**), Consistency (**Con.**), Coherence (**Coh.**), and a new **Overall (Ove.)** which evaluates how good the emotion support model is in general.

### 5.3 Overall Evaluation (Q1)

#### 5.3.1 Automatic Evaluation Results

To demonstrate the effectiveness of our teacher-student framework, we compare our ChatPal with other methods and report results in Table 5.

Firstly, regarding the content-based metrics (*incl*, METEOR, B-4, R-L, and Extrema), it is evident that our ChatPal consistently outperforms other baselines. Among them, ChatGPT exhibits a significant superiority over LLaMA. Ask-Expert further improves the performance by excelling in offering more specific advice than the vanilla ChatGPT. Built upon a small language model, AUGESC can achieve competitive performance as Ask-Expert, indicating the advantages of distilling the knowledge from large models. Overall, our method integrates a broader range of emotional support strategies and scenarios that are distilled from the large teacher, allowing for a more generalizable ChatPal model.

Secondly, when assessing diversity-based metrics (namely, *incl*, D-2, and D-3), it’s evident that methods rooted in ChatGPT naturally generate responses that are both lengthier and richer in content compared to others. The Ask-Expert method, with its fixed guiding prompts, somewhat restricts ChatGPT’s response diversity. Yet, extreme diversity isn’t always advantageous. By tailoring ChatGPT to specific emotional support scenarios, our student model not only elicits a range of responses for its own education but also strikes a balance in diversity. This makes it more diverse than the original LLaMA and more measured than Ask-Expert. Overall, our teacher-student framework delivers dual benefits: it produces a sizable, high-quality ESC dataset and refines a smaller ChatPal that rivals the performance of its larger counterparts.

Methods	Inf.	Und.	Hel.	Con.	Coh.	Ove.
ChatGPT	2.47	2.07	2.34	<b>2.41</b>	<b>2.55</b>	2.40
Ask-Expert	2.15	1.34	1.78	1.94	1.84	1.84
LLaMA	1.59	1.21	1.68	1.44	1.58	1.71
AUGESC	2.16	1.83	2.09	1.85	2.40	2.23
ChatPal / DRI	2.31	2.20	2.46	2.36	2.37	2.43
<b>ChatPal</b>	<b>2.49</b>	<b>2.31</b>	<b>2.51</b>	2.39	2.41	<b>2.48</b>
$\kappa$	0.42	0.33	0.37	0.35	0.40	0.41

Table 6: Human evaluation results. The scores (from 0 to 3) are averaged over all the samples rated by three annotators.  $\kappa$  denotes Fleiss’ Kappa (Fleiss, 1971), indicating fair or moderate inter-annotator agreement ( $0.2 < \kappa < 0.6$ ).

#### 5.3.2 Human Evaluation Results

We further conduct human evaluation on the generated responses with five annotators. We randomly sample 50 conversations from ExTES’s test data for comparison. The annotators were asked to rate the performance of different models. The outcomes of comparison (as shown in Table 6) demonstrate the following findings. (1) It reveals that our final ChatPal (student model) trained on our ExTES dataset achieves better performances than the vanilla ChatGPT (teacher model) on most metrics. It also confirms the high quality and practicality of our ExTES dataset in enhancing emotional support capabilities. (2) We find that Ask-Expert, due to its reliance on fixed formats, is only suitable for providing specific actionable advice and cannot offer comprehensive emotional support, hence it received lower scores. On the other hand, AUGESC may provide unhelpful responses to unfamiliar situations, resulting in lower scores on the Understanding and Helpfulness metrics. Based on our expanded wide-ranging scenarios and comprehensive strategies, our ChatPal outperforms other models in almost all metrics. In general, the results show the effectiveness of our teacher-student framework, enhancing the ability of smaller models to provide emotional support.

### 5.4 Advantages of ExTES Dataset (Q2)

#### 5.4.1 Performance on New Scenarios

The collected ExTES dataset covers a wide range of new ES scenarios, which provides a valuable testbed for the analysis of generalization ability. The automatic and human evaluation of various methods in new scenarios are shown in Table 7 and Table 8. For large language models, ChatGPT and Ask-Expert are less sensitive to varying scenarios, thanks to ChatGPT’s generation capabilities while Ask-Expert further instructs ChatGPT to re-

Methods	METEOR	B-4	R-L	Extrema	D-2	D-3
ChatGPT	22.29	2.114	12.52	60.56	<b>74.96</b>	<b>94.13</b>
Ask-Expert	24.61	2.190	17.13	59.85	72.10	93.38
LLaMA	14.46	1.256	10.24	50.11	27.76	48.04
AUGESC	21.96	1.789	15.57	50.09	48.51	76.07
<b>ChatPal</b>	<b>32.56</b>	<b>2.425</b>	<b>20.98</b>	<b>61.63</b>	68.07	92.25

Table 7: Automatic evaluation results in new scenarios. It reveals that our student model outperforms other methods on most metrics.

Methods	Inf.	Und.	Hel.	Con.	Coh.	Ove.
ChatGPT	<b>2.41</b>	2.04	2.36	2.42	2.37	2.39
Ask-Expert	1.80	1.65	1.79	1.52	1.89	1.93
LLaMA	1.24	1.22	1.14	1.86	1.65	1.55
AUGESC	1.68	1.74	1.72	2.03	1.82	1.92
<b>ChatPal</b>	<b>2.37</b>	<b>2.38</b>	<b>2.42</b>	<b>2.47</b>	<b>2.39</b>	<b>2.46</b>
$\kappa$	0.45	0.31	0.35	0.33	0.47	0.42

Table 8: The human evaluation in new scenarios (scores from 0 to 3). The Fleiss’ Kappa is a fair or moderate inter-annotator agreement ( $0.2 < \kappa < 0.6$ ).

spond by using tailored prompts. On the other hand, vanilla LLaMA and AUGESC struggle to provide specific advice in unseen scenarios, due to limited generation capabilities of relatively small models. Our approach ChatPal, which involves venturing into new scenarios and fine-tuning from high-quality datasets, equips it to address a wide range of user emotional issues with greater empathy and provide more detailed guidance.

#### 5.4.2 Effect of Strategy Guidance

To show the effect of fine-grained strategies in ExTES for helping finetuning, we conduct an ablation study on all three fine-tuning schemes. Results are presented in Table 9. We observe that the variants with strategies are generally better than those without strategies in all schemes, except for their performance on D-2/3 metrics. This is because, under the guidance of specific strategies, the response generation space becomes more constrained, reducing the diversity of responses in certain extent. Therefore, we refer our final ChatPal model as the version trained with strategy annotation and enhanced with DRI.

#### 5.5 Diverse Response Inpainting Effect (Q3)

Table 5 and 6 also show the comparison of performance between our student model and its variant w/o DRI. Additionally, Figure 4 demonstrates the impact of generating varying numbers of diverse re-

Method	Stra?	METEOR	B-4	R-L	Extrema	D-2	D-3
DialoGPT-FT	✗	26.03	1.721	13.37	53.27	49.29	62.92
	✓	26.82	1.966	13.23	55.71	53.11	77.47
7B-Adapter	✗	28.48	1.944	16.95	64.47	60.43	82.62
	✓	29.71	1.987	16.39	62.73	60.83	82.24
7B-LoRA (ChatPal / DRI)	✗	30.31	2.333	19.60	65.06	63.64	84.90
	✓	30.67	<b>2.491</b>	20.85	<b>65.44</b>	61.94	82.80
ChatPal	✗	31.05	2.402	20.94	64.51	<b>69.88</b>	<b>91.96</b>
	✓	<b>33.12</b>	2.437	<b>21.09</b>	63.73	66.93	90.71

Table 9: Comparison of fine-tuning methods. We compare the no-strategy (✗) and with-strategy (✓) variants.

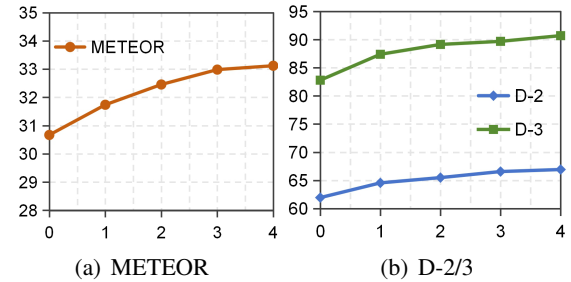


Figure 4: The impact of the number of diverse responses k, ranging from k=0 (w/o DRI) to 4.

sponses during DRI for later finetuning. Compared to the variant w/o DRI, the student model exhibits a significant performance improvement. But our ChatPal scores lower on B-4 and Extrema metrics than ChatPal w/o DRI. This is understandable, under the support of diverse responses, the student model can provide a wider range of emotional support replies. Additionally, generating diverse responses further expands the data scale based on our ExTES dataset, which effectively enhances the quantity of high-quality data. Overall, leveraging the teacher model to generate diverse responses, the performance of small student model can further elevate the performance and help building a more powerful and versatile emotional support chatbot.

## 6 Conclusion

In this paper, we proposed a teacher-student framework and demonstrated the potential of LLMs as “counseling teacher” in enhancing the emotional support response-abilities of smaller models. By leveraging the in-context generalization and extensive knowledge reservoirs of LLMs, we curated a large-scale emotional support conversation dataset (ExTES) and deliberately fine-tuned smaller models with diverse response inpainting mechanism to exhibit proficiency in providing emotional support. Extensive experiments validate the advantages of the ExTES dataset as well as the superiority of the proposed teacher-student framework.



## 603 Limitations

604 Our proposed approach relies heavily on LLMs and  
605 is subject to the same limitations, namely, known  
606 biases in the training data and the ability to hal-  
607 lucinate incorrect information. Since our student  
608 model (ChatPal) is trained on conversations gener-  
609 ated by LLM, it is possible that such characteristics  
610 of the teacher model can get passed along to the  
611 student. Additionally, it is known that for different  
612 cultures, the emotional support strategies can be  
613 very diverse which requires cultural background  
614 knowledge and reasoning processes (Atkins et al.,  
615 2016). And our fine-tuning data is only available  
616 in English and cannot provide support for other  
617 languages at this moment.

## 618 Ethical Considerations

619 Working in the field of emotional support requires  
620 additional ethical considerations. Regarding safety,  
621 we acknowledge the limitations of the current  
622 framework proposed and the potential risks associ-  
623 ated with deploying them directly for emotionally  
624 vulnerable individuals. We do not recommend the  
625 direct deployment of the fine-tuned models from  
626 this work into real-life situations; currently, they  
627 are only suitable for academic research. While we  
628 intend to develop models for the greater good of  
629 society, it is crucial to recognize that the dataset  
630 contains potentially problematic content, includ-  
631 ing toxic or biased material that could be used to  
632 generate negative or offensive content. We openly  
633 provide the dataset collected for this work to assist  
634 in supporting future improvements in ESC.

635 On the other hand, our proposed system relies  
636 heavily on large language models and therefore in-  
637 herits their well-known problems centered around  
638 societal biases learned through pretraining, halluci-  
639 nations, and expensive use of resources (Weidinger  
640 et al., 2021). Various controls are included to con-  
641 strain the LLMs to the emotional support task, but  
642 these are unlikely to fully wash away their inher-  
643 ent issues. Significant further progress needs to  
644 be made in areas like debiasing, grounding in ac-  
645 tuality, and efficient serving before we can safely  
646 deploy this type of system in a production setting.

## 647 References

648 Armen Aghajanyan, Sonal Gupta, and Luke Zettle-  
649 moyer. 2021. Intrinsic dimensionality explains the  
650 effectiveness of language model fine-tuning. In *Pro-*

*ceedings of the 59th Annual Meeting of the Associa-  
tion for Computational Linguistics and the 11th Inter-  
national Joint Conference on Natural Language Pro-  
cessing (Volume 1: Long Papers)*, pages 7319–7328,  
Online. Association for Computational Linguistics.

I Elaine Allen and Christopher A Seaman. 2007. Likert  
scales and data analyses. *Quality progress*, 40(7):64–  
65.

Aram Bahrini, Mohammadsadra Khamoshifar, Hos-  
sein Abbasimehr, Robert J. Riggs, Maryam Esmaeili,  
Rastin Mastali Majdabadkohne, and Morteza Pase-  
hvar. 2023. ChatGPT: Applications, Opportunities,  
and Threats. *arXiv preprint arXiv:2304.09103*.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR:  
An automatic metric for MT evaluation with im-  
proved correlation with human judgments. In *Pro-  
ceedings of the ACL Workshop on Intrinsic and Ex-  
trinsic Evaluation Measures for Machine Translation  
and/or Summarization*, pages 65–72.

Som S Biswas. 2023. Role of chat gpt in public health.  
*Annals of biomedical engineering*, 51(5):868–869.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, et al. 2020. Language models are few-shot  
learners. *Advances in neural information processing  
systems*, 33:1877–1901.

Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang  
Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ra-  
madan, and Milica Gasic. 2018. Multiwoz-a large-  
scale multi-domain wizard-of-oz dataset for task-  
oriented dialogue modelling. In *Proceedings of the  
2018 Conference on Empirical Methods in Natural  
Language Processing*, pages 5016–5026.

Brant Burleson. 2003. Emotional support skills. *Hand-  
book of communication and social interaction skills*,  
pages 180–399.

Jian Cheng, Peisong Wang, Gang Li, Qinghao Hu, and  
Hanqing Lu. 2018. Recent advances in efficient  
computation of deep convolutional neural networks.  
*Frontiers of Information Technology & Electronic  
Engineering*, 19:64–77.

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 Empirical  
Methods in Natural Language Processing, EMNLP  
2022, Abu Dhabi, United Arab Emirates, December  
7-11, 2022*, pages 3014–3026.

Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. 2020.  
*A survey of model compression and acceleration for  
deep neural networks*.

704	Yang Deng, Wenxuan Zhang, Yifei Yuan, and Wai Lam.	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	756
705	2023. Knowledge-enhanced mixed-initiative dia-	Bruna Morrone, Quentin de Laroussilhe, Andrea Ges-	757
706	logue system for emotional support conversations.	mundido, Mona Attariyan, and Sylvain Gelly. 2019.	758
707	In <i>Proceedings of the 61st Annual Meeting of the</i>	Parameter-efficient transfer learning for nlp.	759
708	<i>Association for Computational Linguistics (Volume</i>		
709	<i>1: Long Papers)</i> , <i>ACL 2023, Toronto, Canada, July</i>	Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh,	760
710	<i>9-14, 2023</i> , pages 4079–4095.	Hootan Nakhost, Yasuhisa Fujii, Alexander J. Ratner,	761
		Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister.	762
711	Bosheng Ding, Chengwei Qin, Linlin Liu, Lidong Bing,	2023. Distilling step-by-step! outperforming larger	763
712	Shafiq R. Joty, and Boyang Li. 2022. Is gpt-3 a good	language models with less training data and smaller	764
713	data annotator? <i>ArXiv</i> , abs/2212.10450.	model sizes. <i>ArXiv</i> , abs/2305.02301.	765
714	Joseph L. Fleiss. 1971. Measuring nominal scale agree-	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	766
715	ment among many raters. <i>Psychological Bulletin</i> ,	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,	767
716	76:378–382.	and Weizhu Chen. 2021. Lora: Low-rank adap-	768
		tation of large language models. <i>arXiv preprint</i>	769
717	Gabriel Forgues, Joelle Pineau, Jean-Marie	<i>arXiv:2106.09685</i> .	770
718	Larchevêque, and Réal Tremblay. 2014. Bootstrap-	Junmo Kang, Wei Xu, and Alan Ritter. 2023. Distill	771
719	ping dialog systems with word embeddings. In <i>Nips,</i>	or annotate? cost-efficient fine-tuning of compact	772
720	<i>modern machine learning and natural language</i>	models. <i>ArXiv</i> , abs/2305.01645.	773
721	<i>processing workshop</i> .		
		Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	774
722	Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	775
723	Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui	guage models are zero-shot reasoners. <i>Advances in</i>	776
724	He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. 2023.	<i>neural information processing systems</i> , 35:22199–	777
725	Llama-adapter v2: Parameter-efficient visual instruc-	22213.	778
726	tion model.		
		Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason	779
727	Jianping Gou, Baosheng Yu, Stephen J. Maybank, and	Yosinski. 2018. Measuring the intrinsic dimension of	780
728	Dacheng Tao. 2021. Knowledge distillation: A	objective landscapes.	781
729	survey. <i>International Journal of Computer Vision</i> ,	Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao,	782
730	129(6):1789–1819.	and Bill Dolan. 2016. A diversity-promoting objec-	783
		tive function for neural conversation models.	784
731	Sarah Graham, Colin Depp, Ellen E Lee, Camille	Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen,	785
732	Nebeker, Xin Tu, Ho-Cheol Kim, and Dilip V Jeste.	Xinlu Zhang, Zekun Li, Hong Wang, Jingu Qian,	786
733	2019. Artificial intelligence for mental health and	Baolin Peng, Yi Mao, Wenhui Chen, and Xifeng Yan.	787
734	mental illnesses: an overview. <i>Current psychiatry</i>	2022. Explanations from large language models	788
735	<i>reports</i> , 21:1–18.	make small reasoners better. <i>ArXiv</i> , abs/2210.06726.	789
736	Clara E. Hill. 1999. Helping skills: Facilitating explo-	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	790
737	ration, insight, and action. In <i>American Psychologi-</i>	Optimizing continuous prompts for generation.	791
738	<i>cal Association</i> .		
		Xinmeng Li, Wansen Wu, Long Qin, and Quanjun Yin.	792
739	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.	2021a. How to evaluate your dialogue models: A	793
740	<i>Distilling the knowledge in a neural network</i> .	review of approaches.	794
		Yanran Li, Ke Li, Hongke Ning, Xiaoqiang Xia, Ya-	795
741	Namgyu Ho, Laura Schmid, and Se-Young Yun. 2022.	long Guo, Chen Wei, Jianwei Cui, and Bin Wang.	796
742	Large language models are reasoning teachers. <i>arXiv</i>	2021b. Towards an online empathetic chatbot with	797
743	<i>preprint arXiv:2212.10071</i> .	emotion causes. In <i>Proceedings of the 44th Inter-</i>	798
		<i>national ACM SIGIR Conference on Research and</i>	799
744	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch,	<i>Development in Information Retrieval</i> . ACM.	800
745	Elena Buchatskaya, Trevor Cai, Eliza Rutherford,	Chin-Yew Lin. 2004. ROUGE: A package for automatic	801
746	Diego de Las Casas, Lisa Anne Hendricks, Johannes	evaluation of summaries. In <i>Text Summarization</i>	802
747	Welbl, Aidan Clark, Tom Hennigan, Eric Noland,	<i>Branches Out</i> , pages 74–81.	803
748	Katie Millican, George van den Driessche, Bogdan		
749	Damoc, Aurelia Guy, Simon Osindero, Karen Si-	Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand	804
750	monyuan, Erich Elsen, Jack W. Rae, Oriol Vinyals,	Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie	805
751	and Laurent Sifre. 2022. <i>Training compute-optimal</i>	Huang. 2021. Towards emotional support dialog	806
752	<i>large language models</i> .	systems. In <i>Proceedings of the 59th Annual Meet-</i>	807
		<i>ing of the Association for Computational Linguistics</i>	808
753	Mahshid Hosseini and Cornelia Caragea. 2021. It takes		
754	two to empathize: One to seek and one to provide. In		
755	<i>AAAI Conference on Artificial Intelligence</i> .		

809				
810		<i>and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> ,		
811		pages 3469–3483. Association for Computational Linguistics.		
812				
813	Edward Loper and Steven Bird. 2002. Nltk: The natural			
814	language toolkit.			
815	Lenin Medeiros and Tibor Bosse. 2018. Using crowd-			
816	sourcing for the development of online emotional			
817	support agents. In <i>Practical Applications of Agents</i>			
818	<i>and Multi-Agent Systems</i> .			
819	Jingbo Meng and Yue (Nancy) Dai. 2021. Emotional			
820	Support from AI Chatbots: Should a Supportive Part-			
821	ner Self-Disclose or Not? <i>Journal of Computer-</i>			
822	<i>Mediated Communication</i> , 26(4):207–222.			
823	World Health Organization et al. 2020. Mental health			
824	and psychosocial considerations during the covid-19			
825	outbreak, 18 march 2020. Technical report, World			
826	Health Organization.			
827	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,			
828	Carroll L. Wainwright, Pamela Mishkin, Chong			
829	Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray,			
830	John Schulman, Jacob Hilton, Fraser Kelton, Luke			
831	Miller, Maddie Simens, Amanda Askell, Peter Welin-			
832	der, Paul F. Christiano, Jan Leike, and Ryan Lowe.			
833	2022. Training language models to follow instruc-			
834	tions with human feedback. In <i>NeurIPS</i> .			
835	Wei Peng, Yue Hu, Luxi Xing, Yuqiang Xie, Yajing Sun,			
836	and Yunpeng Li. 2022. Control globally, understand			
837	locally: A global-to-local hierarchical graph network			
838	for emotional support conversation. In <i>Proceedings</i>			
839	<i>of the Thirty-First International Joint Conference on</i>			
840	<i>Artificial Intelligence, IJCAI 2022, Vienna, Austria,</i>			
841	<i>23-29 July 2022</i> , pages 4324–4330.			
842	Hannah Rashkin, Eric Michael Smith, Margaret Li, and			
843	Y-Lan Boureau. 2019. Towards empathetic open-			
844	domain conversation models: A new benchmark and			
845	dataset. In <i>Proceedings of the 57th Annual Meet-</i>			
846	<i>ing of the Association for Computational Linguistics</i> ,			
847	pages 5370–5381. Association for Computational			
848	Linguistics.			
849	Maija Reblin and Bert N Uchino. 2008. Social and emo-			
850	tional support and its implication for health. <i>Current</i>			
851	<i>opinion in psychiatry</i> , 21(2):201.			
852	Maarten Sap, Vered Shwartz, Antoine Bosselut, Yejin			
853	Choi, and Dan Roth. 2020. Commonsense reason-			
854	ing for natural language processing. In <i>Proceedings</i>			
855	<i>of the 58th Annual Meeting of the Association for</i>			
856	<i>Computational Linguistics: Tutorial Abstracts</i> , pages			
857	27–33.			
858	Ashish Sharma, Inna W Lin, Adam S Miner, David C			
859	Atkins, and Tim Althoff. 2021. Towards facilitating			
860	empathic conversations in online mental health sup-			
861	port: A reinforcement learning approach. In <i>Proceed-</i>			
862	<i>ings of the Web Conference 2021</i> , pages 194–205.			
	Ashish Sharma, Adam Miner, David Atkins, and Tim Al-			
	thoff. 2020. A computational approach to understand-			
	ing empathy expressed in text-based mental health			
	support. In <i>Proceedings of the 2020 Conference on</i>			
	<i>Empirical Methods in Natural Language Processing</i>			
	<i>(EMNLP)</i> , pages 5263–5276. Association for Com-			
	putational Linguistics.			
	Siqi Shen, Charles Welch, Rada Mihalcea, and Verónica			
	Pérez-Rosas. 2020. Counseling-style reflection gen-			
	eration using generative pretrained transformers with			
	augmented context. In <i>SIGdial 2020</i> , pages 10–20.			
	Ariel Shensa, Jaime E. Sidani, César G. Escobar-Viera,			
	Galen E. Switzer, Brian A. Primack, and Sophia			
	Choukas-Bradley. 2020. Emotional support from			
	social media and face-to-face relationships: Assoc-			
	iations with depression risk among young adults. <i>Journal of Affective Disorders</i> , 260:38–44.			
	Kumar Shridhar, Alessandro Stolfo, and Mrinmaya			
	Sachan. 2022. Distilling reasoning capabilities into			
	smaller language models. <i>ArXiv</i> .			
	Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju,			
	Eric Michael Smith, Stephen Roller, Megan Ung,			
	Moya Chen, Kushal Arora, Joshua Lane, Morteza			
	Behrooz, William Ngan, Spencer Poff, Naman Goyal,			
	Arthur Szlam, Y-Lan Boureau, Melanie Kambadur,			
	and Jason Weston. 2022. <a href="#">Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage</a> .			
	Hao Sun, Zhenru Lin, Chujie Zheng, Siyang Liu, and			
	Minlie Huang. 2021. Psyqa: A chinese dataset for			
	generating long counseling text for mental health			
	support.			
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier			
	Martinet, Marie-Anne Lachaux, Timothée Lacroix,			
	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal			
	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard			
	Grave, and Guillaume Lample. 2023. Llama: Open			
	and efficient foundation language models.			
	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 sup-			
	port conversation. In <i>Proceedings of the 60th Annual</i>			
	<i>Meeting of the Association for Computational Lin-</i>			
	<i>guistics (Volume 1: Long Papers), ACL 2022, Dublin,</i>			
	<i>Ireland, May 22-27, 2022</i> , pages 308–319.			
	Janneke M van der Zwaan, Virginia Dignum, and			
	Catholijn M Jonker. 2012. A conversation model			
	enabling intelligent agents to give emotional support.			
	In <i>Modern Advances in Intelligent Systems and Tools</i> ,			
	pages 47–52.			
	Shuohang Wang, Yang Liu, Yichong Xu, Chenguang			
	Zhu, and Michael Zeng. 2021. Want to reduce la-			
	beling cost? GPT-3 can help. In <i>Findings of the</i>			
	<i>Association for Computational Linguistics: EMNLP</i>			
	<i>2021</i> , pages 4195–4205, Punta Cana, Dominican Re-			
	public. Association for Computational Linguistics.			

919 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,  
920 Barret Zoph, Sebastian Borgeaud, Dani Yogatama,  
921 Maarten Bosma, Denny Zhou, Donald Metzler, et al.  
922 2022. Emergent abilities of large language models.  
923 *arXiv preprint arXiv:2206.07682*.

924 Laura Weidinger, John Mellor, Maribeth Rauh, Conor  
925 Griffin, Jonathan Uesato, Po-Sen Huang, Myra  
926 Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh,  
927 Zac Kenton, Sasha Brown, Will Hawkins, Tom  
928 Stepleton, Courtney Biles, Abeba Birhane, Julia  
929 Haas, Laura Rimell, Lisa Anne Hendricks, William  
930 Isaac, Sean Legassick, Geoffrey Irving, and Iason  
931 Gabriel. 2021. [Ethical and social risks of harm from  
932 language models](#).

933 Chun-Hung Yeh, Anuradha Welivita, and Pearl Pu Falt-  
934 ings. 2015. A dialogue dataset containing emo-  
935 tional support for people in distress. *arXiv preprint  
936 arXiv:1503.08895*.

937 Qiang Zhang, Jason Naradowsky, and Yusuke Miyao.  
938 2023a. Ask an expert: Leveraging language mod-  
939 els to improve strategic reasoning in goal-oriented  
940 dialogue models. In *Findings of the Association for  
941 Computational Linguistics: ACL 2023*, pages 6665–  
942 6694, Toronto, Canada. Association for Computa-  
943 tional Linguistics.

944 Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Ao-  
945 jun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hong-  
946 sheng Li, and Yu Qiao. 2023b. Llama-adapter: Effi-  
947 cient fine-tuning of language models with zero-init  
948 attention.

949 Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen,  
950 Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing  
951 Liu, and Bill Dolan. 2020. Dialogpt: Large-scale  
952 generative pre-training for conversational response  
953 generation.

954 Chujie Zheng, Pei Ke, Zheng Zhang, and Minlie Huang.  
955 2023a. Click: Controllable text generation with se-  
956 quence likelihood contrastive learning.

957 Chujie Zheng, Sahand Sabour, Jiaxin Wen, Zheng  
958 Zhang, and Minlie Huang. 2023b. Augesc: Dialogue  
959 augmentation with large language models for emo-  
960 tional support conversation. In *Findings of the As-  
961 sociation for Computational Linguistics: ACL 2023*,  
962 pages 1552–1568.

963 Xianda Zhou and William Yang Wang. 2018. MojiTalk:  
964 Generating emotional responses at scale. In *Proceed-  
965 ings of the 56th Annual Meeting of the Association for  
966 Computational Linguistics (Volume 1: Long Papers)*,  
967 pages 1128–1137. Association for Computational  
968 Linguistics.

969 **Appendix**

970 **A Data Example from our dataset**

971 An detailed example can be seen in Figure 5. Each  
972 conversation is labeled with its scene category and  
973 a brief of description of the user. In the context  
974 of each conversation, the strategies used by the  
975 assistant are labeled red.

<p><b>Scene:</b> Communication Challenges</p>
<p><b>Description:</b> With my roommate always leaving a mess in the kitchen, it's becoming increasingly hard to address the issue and maintain a harmonious living environment.</p>
<p><b>User:</b> Hey.  <b>AI:</b> Hi, how are you?  <b>User:</b> I'm alright, thanks. But I've been having a hard time with my roommate lately. They always leave a mess in the kitchen, and it's getting frustrating to deal with.  <b>AI (Reflective Statements):</b> That sounds really irritating. It must be tough to constantly clean up after someone else.  <b>User:</b> Exactly. I've tried talking to them about it, but it doesn't seem to have any effect. I'm not sure how to approach the situation without causing conflict.  <b>AI (Suggest Options):</b> Maybe you could try setting up some agreed-upon rules for kitchen cleanliness. That way, both you and your roommate are on the same page.  <b>User:</b> That's a good idea. I'll give that a try. It's just been really stressful, and I don't want it to affect our relationship as roommates.  <b>AI (Normalize Experiences):</b> It's completely normal to have disagreements or conflicts with roommates. The key is finding a way to address the issues together that works for both of you.  <b>User:</b> You're right. I'll try to approach it from a calm and understanding perspective. Hopefully, we can come to a resolution.  <b>AI (Affirmation):</b> I believe in your ability to handle this situation maturely. Remember to take care of yourself throughout this process.  <b>User:</b> Thank you for the support. It means a lot to me. I'll do my best to handle it positively.  <b>AI (Offer Hope):</b> I'm confident that with open communication and a willingness to find a solution, you'll be able to improve your living situation. Don't hesitate to reach out for support if you need it.  <b>User:</b> I really appreciate your positivity and encouragement. It definitely gives me hope. Thank you again for being there.  <b>AI:</b> You're welcome! I'm here whenever you need to talk. Take care and hang in there!  <b>User:</b> Thank you. I'll keep that in mind. Take care too. Goodbye!</p>

Figure 5: Data example from our ExTES dataset.

976 **B Definitions of Strategies**

977 **Reflective Statements (RS)** Repeat or rephrase  
978 what the User has expressed to show that you're ac-  
979 tively listening. Example: User: "I'm feeling really  
980 overwhelmed with all the work I have to do." As-  
981 sistant: "It sounds like you're feeling overwhelmed  
982 with your workload."

983 **Clarification (Cla)** Seek clarification to ensure  
984 a clear understanding of the User's emotions and  
985 experiences. Example: User: "I just can't shake  
986 off this feeling of sadness." Assistant: "Could you

help me understand what might have triggered this  
feeling of sadness?"

**Emotional Validation (EV)** Acknowledge and  
validate the User's emotions without judgment. Ex-  
ample: User: "I'm so frustrated with myself for  
making the same mistake again." Assistant: "It's  
completely understandable to feel frustrated when  
you make a mistake."

**Empathetic Statements (ES)** Express under-  
standing and empathy towards the User's experi-  
ences. Example: User: "I'm really struggling  
with my self-confidence right now." Assistant: "I  
can imagine how challenging it must be to navigate  
through situations that affect your self-confidence."

**Affirmation (Aff)** Provide positive reinforce-  
ment and encouragement to uplift the User's spirits.  
Example: User: "I feel like I'm not good enough."  
Assistant: "You've accomplished so much already,  
and your abilities speak for themselves. Don't un-  
derestimate your capabilities."

**Offer Hope (OH)** Share optimistic perspectives  
or possibilities to instill hope. Example: User: "I  
don't know if things will ever get better." Assis-  
tant: "Remember that change is constant, and there  
are always opportunities for growth and positive  
change."

**Avoid Judgment and Criticism (AJC)** It's im-  
portant to create a non-judgmental and safe space  
for the User to express their emotions without fear  
of criticism. Refrain from passing judgment or be-  
ing overly critical of their experiences or choices.  
Example: User: "I'm feeling so guilty for taking  
time off work to focus on my mental health." Assis-  
tant: "Taking care of your mental health is crucial,  
and it's not something to feel guilty about. Your  
well-being should always be a priority, and I'm  
glad you recognized that. Is there anything I can  
do to support you during this time?"

**Suggest Options (SO)** Offer practical sugges-  
tions or alternative perspectives for addressing the  
issue at hand. Example: User: "I'm having trou-  
ble managing my stress." Assistant: "Have you  
considered trying relaxation techniques like deep  
breathing or mindfulness exercises?"

**Collaborative Planning (CP)** Work together  
with the User to develop an action plan. Exam-  
ple: User: "I want to improve my time manage-  
ment skills." Assistant: "Let's brainstorm some  
strategies together. How about breaking tasks into  
smaller, more manageable chunks?"

**Provide Different Perspectives (PDP)** Offer al-  
ternative ways of looking at the situation to help

the User gain new insights. Example: User: "I'm devastated that my project didn't succeed." Assistant: "Sometimes setbacks can lead to unexpected opportunities for learning and growth. It's a chance to reassess and try again."

**Reframe Negative Thoughts (RNT)** Help the User reframe negative thoughts into more positive or realistic ones. Example: User: "I'm such a failure." Assistant: "Instead of thinking that way, let's focus on what you've learned from this experience and how you can apply it moving forward."

**Share Information (SI)** Provide educational or factual information about emotions, coping mechanisms, or self-care practices. Example: User: "I'm struggling to manage my anxiety." Assistant: "Did you know that deep breathing exercises and grounding techniques can help reduce anxiety symptoms? Would you like me to explain how to practice them?"

**Normalize Experiences (NE)** Explain that certain emotions or reactions are common and part of the human experience. Example: User: "I feel so guilty for taking time for myself." Assistant: "It's common to feel guilty about self-care, but it's essential for your well-being. Remember, you deserve to prioritize your needs too."

**Promote Self-Care Practices (PSP)** Advocate for engaging in activities that promote well-being and self-care. Example: "Make sure to take some time for yourself and do something that brings you joy and relaxation."

**Stress Management (SM)** Provide suggestions for stress management techniques like exercise, meditation, or spending time in nature. Example: "Engaging in regular physical activity can help reduce stress and improve mood."

**Others (Oth)** Interact with friendly greetings and employ additional supportive techniques that are not covered by the previously mentioned categories.

## C Strategy Transition

We present the top-5 most frequent strategy transitions with 3-5 hops in Table 10. These transitions indicate that supporters usually ask questions and explore the user's situation before comforting the user. Emotional support supporters usually first understand the cause of the user's distress and then say some words of comfort or express sympathy for the user's experience. This is generally as expected. It also might not be wise enough to make actionable

	Strategy Transition	Proportion
<b>3-Hop</b>	EV → RS → EV	17.19 %
	EV → RS → SO	16.23 %
	EV → RS → ES	14.49 %
	RS → EV → SO	11.03 %
	EV → ES → RS	9.75 %
<b>4-Hop</b>	EV → RS → ES → SO	7.08 %
	EV → RS → SO → Aff	6.61 %
	EV → ES → RS → NE	6.04 %
	RS → Aff → ES → RS	5.27 %
	EV → RS → SO → Cla	4.36 %
<b>5-Hop</b>	EV → RS → EV → Aff → SO	1.97 %
	EV → RS → SO → Aff → RS	1.34 %
	RS → EV → SO → OH → SO	0.89 %
	EV → RS → ES → SO → Aff	0.45 %
	EV → ES → RS → NE → Cla	0.27 %

Table 10: Proportions of top-5 strategy transitions in responses. The adjacent same strategies are merged. Abbreviations are consistent with the Appendix B.

suggestions at the beginning of the whole dialogue.

## D Template of Expanding Conversation

The template for ChatGPT to iteratively expand conversations (Figure 2) is as follows:

*Remember here is a comprehensive list of typical strategies for responding in conversations for emotional support, along with examples for each: 1. Reflective Statements: Repeat or rephrase what the person has expressed to show that you're actively listening. 2. Clarification: Seek clarification to ensure a clear understanding of the person's emotions and experiences. 3. Emotional Validation: Acknowledge and validate the person's emotions without judgment. ... 15. Stress Management: Provide suggestions for stress management techniques like exercise, meditation, or spending time in nature. 16. Others: Other strategies. Example:  $\{SEED\}$  EXAMPLE*

*Your task is to create a casual emotional support conversation between a user and an assistant. Create a random emotional support scenario of the ' $\{SCENE\}$ ' type, write it in the Description, and then generate a complete set of dialogue. Make the conversation more like a real-life chat and be specific. Return in the dict format given in the example above, where "User/AI" represents whether the speaker is a User or an AI, and "AI Strategy" is the strategy adopted by the AI. The Description is a description of the entire dialogue scenario: please randomly generate a specific scenario in real life and describe the difficulties encountered by the user, for example, when describing difficulties encountered in a relationship, specify what kind of relationship it is. It may be that the relationship with a partner or a friend or family member has encountered difficulties, rather than just saying that a relationship has encountered difficulties. The return format is a dict ...*

## E Details of Scenarios

Below are 36 emotional support scenarios and examples that we have compiled. And Table 11 is the statistics of all ES scenarios.

Category	Dialogues	Proportion	Category	Dialogues	Proportion
Breakups or Divorce	710	6.3%	Navigating Gender Identity and Transitioning	202	1.8%
Conflicts or Communication Problems	1,109	9.9%	Moving to a New City or Country	202	1.8%
Communication Challenges	1,008	9.0%	Career Transitions	202	1.8%
Coping with the Death of a Loved One	593	5.3%	Parenthood and Parenting Challenges	202	1.8%
Dealing with the Loss of a Pet	601	5.4%	Low Self-Esteem or Lack of Confidence	302	2.7%
Work-related Stress and Burnout	403	3.6%	Body Image Concerns and Eating Disorders	101	0.9%
Financial Worries and Uncertainty	403	3.6%	LGBTQ+ Identity	101	0.9%
Unemployment-related Stress	403	3.6%	Cultural Identity and Belonging	101	0.9%
Academic Stress	403	3.6%	Academic Stress or Pressure	202	1.8%
Spirituality and Faith	202	1.8%	Job Loss or Career Setbacks	202	1.8%
Managing Bipolar Disorder	202	1.8%	Parenting Challenges and Parental Guilt	202	1.8%
Anxiety and Panic	202	1.8%	Sibling Rivalry or Family Conflict	403	3.6%
Depression and Low Mood	403	3.6%	Surviving and Recovering from Physical or Emotional Abuse	101	0.9%
Adjusting to a New Job or Role	302	2.7%	Healing from Sexual Assault or Domestic Violence	101	0.9%
Chronic Illness or Pain Management	302	2.7%	Post-Traumatic Stress Disorder (PTSD)	101	0.9%
Coping with a Diagnosis or Medical Treatment	202	1.8%	Healing from Abuse	202	1.8%
Caregiver Support	202	1.8%	Addiction and Recovery	202	1.8%
Finding Meaning and Purpose in Life	202	1.8%	Support for Loved Ones or Friends	202	1.8%

Table 11: Statistics of all 36 emotional support scenarios covered in our ExTES dataset.

1098	<b>Breakups or Divorce</b>	Example 1: Processing the emotions and grief following the end of a long-term relationship. Example 2: Seeking guidance on how to navigate a recent breakup and move forward.	1134
1099			1135
1100			1136
1101			1137
1102	<b>Conflicts or Communication Problems</b>	Example 1: Dealing with a misunderstanding or disagreement with a close friend or family member. Example 2: Seeking advice on resolving conflicts with a romantic partner and improving communication.	1138
1103			1139
1104			1140
1105			1141
1106			1142
1107	<b>Communication Challenges</b>	Example: Helping a person find effective ways to express their needs and concerns to their partner, fostering open and constructive communication.	1143
1108			1144
1109			1145
1110			1146
1111	<b>Coping with the Death of a Loved One</b>	Example 1: Navigating the stages of grief and finding ways to honor the memory of the deceased. Seeking support in managing the emotional impact of losing a close family member or friend.	1147
1112			1148
1113			1149
1114			1150
1115			1151
1116	<b>Dealing with the Loss of a Pet</b>	Example 1: Processing the deep sadness and emptiness after the death of a beloved pet. Example 2: Seeking understanding and comfort while grieving the loss of a long-time companion animal.	1152
1117			1153
1118			1154
1119			1155
1120			1156
1121	<b>Work-related Stress and Burnout</b>	Example 1: Coping with excessive workload, pressure, and a demanding work environment. Example 2: Seeking strategies to manage stress and achieve a healthier work-life balance.	1157
1122			1158
1123			1159
1124			1160
1125			1161
1126	<b>Financial Worries and Uncertainty</b>	Example 1: Navigating financial challenges such as debt, job loss, or unexpected expenses. Example 2: Seeking emotional support and practical advice to alleviate financial stress and regain stability.	1162
1127			1163
1128			1164
1129			1165
1130			1166
1131	<b>Unemployment-related stress</b>	Example: Encouraging someone who is about to lose their job due to poor company performance, discussing the possibility of changing jobs, prioritizing self-care, and staying positive.	1167
1132			1168
1133			1169

1170	<b>Career Transitions</b> Example: Assisting someone who is considering a career change, helping them explore their passions, and transferable skills and develop a plan for transitioning into a new field.	1222
1171		1223
1172		1224
1173		1225
1174	<b>Parenthood and Parenting Challenges</b> Example: Supporting a new parent who is feeling overwhelmed and sleep-deprived, offering reassurance, and sharing tips for self-care and coping strategies for the demands of parenthood.	1226
1175		1227
1176		1228
1177		1229
1178		1230
1179	<b>Low Self-Esteem or Lack of Confidence</b> Example 1: Addressing negative self-perceptions and building self-worth. Example 2: Seeking techniques for cultivating self-compassion and improving self-esteem.	1231
1180		1232
1181		1233
1182		1234
1183		1235
1184	<b>Body Image Concerns and Eating Disorders</b> Example 1: Dealing with body dissatisfaction and the impact it has on self-image and overall well-being. Example 2: Seeking support in recovering from an eating disorder and developing a healthy relationship with food and body.	1236
1185		1237
1186		1238
1187		1239
1188		1240
1189		1241
1190	<b>LGBTQ+ Identity</b> Example: Assisting someone in the process of coming out as gay, offering support, connecting them with LGBTQ+ community resources, and being a source of understanding.	1242
1191		1243
1192		1244
1193		1245
1194	<b>Cultural Identity and Belonging</b> Example: Engaging in discussions with someone exploring their mixed-race identity and helping them embrace and celebrate their diverse heritage.	1246
1195		1247
1196		1248
1197		1249
1198	<b>Academic Stress or Pressure</b> Example 1: Coping with academic expectations, exam anxiety, or perfectionism. Example 2: Seeking strategies for time management, study techniques, and reducing academic stress.	1250
1199		1251
1200		1252
1201		1253
1202		1254
1203	<b>Job Loss or Career Setbacks</b> Example 1: Navigating the emotions and challenges of losing a job or facing career setbacks. Example 2: Seeking guidance and encouragement for career transitions or exploring new professional opportunities.	1255
1204		1256
1205		1257
1206		1258
1207		1259
1208	<b>Parenting Challenges and Parental Guilt</b> Example 1: Managing parental responsibilities, parenting styles, and dealing with parental guilt. Example 2: Seeking advice on effective communication with children and finding a balance between work and family.	1260
1209		1261
1210		1262
1211		1263
1212		1264
1213		1265
1214	<b>Sibling Rivalry or Family Conflict</b> Example 1: Resolving conflicts and improving relationships with siblings or other family members. Example 2: Seeking guidance on navigating family dynamics, establishing healthy boundaries, and fostering understanding.	1266
1215		1267
1216		1268
1217		1269
1218		1270
1219		1271
1220	<b>Surviving and Recovering from Physical or Emotional Abuse</b> Example 1: Processing the trauma of past abuse and seeking support for healing and recovery. Example 2: Finding resources and coping strategies for managing the emotional impact of abuse.	1272
1221		1273
	<b>Healing from Sexual Assault or Domestic Violence</b> Example 1: Navigating complex emotions, seeking support, and developing coping mechanisms after experiencing sexual assault or domestic violence. Example 2: Accessing information on trauma-informed therapy and support networks for survivors of assault or violence.	
	<b>Post-Traumatic Stress Disorder (PTSD)</b> Example: Creating a safe and non-judgmental space for military veteran with PTSD to share their experiences and providing resources for trauma-focused therapy and support groups.	
	<b>Healing from Abuse</b> Example: Assisting someone who has recently left an abusive relationship, connecting them with local support services, and offering encouragement as they rebuild their life.	
	<b>Navigating Gender Identity and Transitioning</b> Example 1: Seeking support and resources while exploring gender identity and considering transitioning. Example 2: Accessing guidance on navigating social, medical, and legal aspects of transitioning.	
	<b>Chronic Illness or Pain Management</b> Example 1: Coping with the emotional impact of a chronic illness, including pain, limitations, and lifestyle adjustments. Example 2: Seeking support in managing daily challenges, finding self-care strategies, and connecting with others facing similar health issues.	
	<b>Coping with a Diagnosis or Medical Treatment</b> Example 1: Processing the emotions surrounding a new medical diagnosis and navigating treatment options. Example 2: Seeking emotional support and practical guidance to cope with medical procedures, side effects, and lifestyle changes.	
	<b>Caregiver Support</b> Example: Offering guidance and resources to a caregiver of an elderly parent, discussing techniques for managing caregiver stress, and suggesting respite care options.	
	<b>Finding Meaning and Purpose in Life</b> Example 1: Exploring questions related to the meaning of life, personal values, and finding purpose. Example 2: Assisting someone who is questioning their life's purpose and exploring different avenues for finding meaning, discussing their values and interests, and encouraging self-reflection.	
	<b>Spirituality and Faith</b> Example: Offering guidance and resources to someone who is questioning	



1274 their faith or seeking spiritual fulfillment, providing  
1275 support as they explore their beliefs and values.

1276 **Addiction and Recovery** Example: Offering em-  
1277 pathy and understanding to someone battling addic-  
1278 tion, discussing treatment options, and providing  
1279 emotional support during their journey to recovery.

1280 **Support for Loved Ones or Friends** Example:  
1281 Supporting a parent who has a child dealing with  
1282 addiction, offering a listening ear, and connecting  
1283 them with support groups and counseling services.

## 1284 F The quality of Seed Dialogues

1285 Table 12 shows the results of human evaluation on  
seed dialogues and ExTES.

	Seeds	ExTES	$\kappa$
<b>Informativeness</b>	2.39	<b>2.53</b>	0.51
<b>Understanding</b>	<b>2.64</b>	2.52	0.46
<b>Helpfulness</b>	2.48	<b>2.61</b>	0.44
<b>Consistency</b>	<b>2.75</b>	2.67	0.39
<b>Overall</b>	2.38	<b>2.45</b>	0.52

1286 Table 12: Human evaluation of seed dialogues quality  
1287 and ExTES quality. The scores (from 0 to 3) are aver-  
1288 aged over all the samples rated by three annotators.  $\kappa$   
1289 denotes Fleiss’ Kappa (Fleiss, 1971), indicating fair to  
moderate inter-annotator agreement ( $0.2 < \kappa < 0.6$ ).

## 1287 G Diverse Response Inpainting Example

1288 Figure 6 shows the process of diverse response  
1289 inpainting.

## 1290 H Fine-tune Methods

### 1291 H.1 Fine-tune Methods

1292 We explore the following three methods to  
1293 fine-tune our ChatPal (student model):

1294 **DialoGPT Fine-Tuning** DialoGPT (Zhang et al.,  
1295 2020) is a medium-sized GPT2 Model trained  
1296 on 147M conversation-like exchanges extracted  
1297 from Reddit. It was trained with a causal language  
1298 modeling (CLM) objective on conversational data  
1299 and is therefore powerful at response generation  
1300 in open-domain dialogue systems. In order to  
1301 fine-tune DialoGPT, we use CLM training. We  
1302 follow the OpenAI GPT-2<sup>3</sup> to model a multiturn  
1303 dialogue session as a long text and frame the  
1304 generation task as language modeling.

1305 **LLaMA Adapter-Tuning** LLaMA-Adapter  
1306 (Zhang et al., 2023b) is a form of prefix-tuning

<sup>3</sup>[https://huggingface.co/docs/transformers/model\\_doc/gpt2](https://huggingface.co/docs/transformers/model_doc/gpt2)

1307 that prepends a learnable adaption-prompt to  
1308 the inputs of the attention blocks in LLaMA.  
1309 There are only 1.2M parameters to update during  
1310 finetuning, which significantly reduces the memory  
1311 footprint and speeds up training. Recently,  
1312 LLaMA-Adapter v2 (Gao et al., 2023) is developed  
1313 to further include more trainable parameters.  
1314 We use LLaMA-Adapter v2 to demonstrate  
1315 instruction-tuning LLaMA 7B on our dataset.  
1316 Inspired by prefix tuning (Li and Liang, 2021) and  
1317 the original adapter method (Houlsby et al., 2019),  
1318 Adapter-Tuning introduces some new sublayers  
1319 (i.e., adapter layers) acting as low-rank bottlenecks  
1320 within each Transformer layer. Generally, instead  
1321 of tuning all parameters, Adapter-Tuning focuses  
1322 on tuning mainly the adapter layers.

1323 **LLaMA LoRA-Tuning** Low-rank adaption  
1324 (LoRA) (Hu et al., 2021) is a technique to approxi-  
1325 mate the update to the linear layers in a LLM with  
1326 a low-rank matrix factorization. This significantly  
1327 reduces the number of trainable parameters and  
1328 speeds up training with little impact on the final  
1329 performance of the model. We demonstrate this  
1330 method by instruction-tuning LLaMA 7B on our  
1331 dataset. The authors take inspiration from (Li  
1332 et al., 2018; Aghajanyan et al., 2021) which show  
1333 that the learned over-parametrized models in fact  
1334 reside on a low intrinsic dimension. Based on  
1335 the inherent low-rank characteristics of the large  
1336 model, the bypass matrix is added to simulate the  
1337 fine-tuning of the full model parameters. LoRA  
1338 achieves the purpose of lightweight fine-tuning  
1339 through a simple and effective solution. It turns  
1340 various large models into professional models in  
1341 different fields through light fine-tuning.

### 1342 H.2 Experimental Setup

1343 We select LLaMA LoRA-Tuning to build our small  
1344 ChatPal model. During the fine-tuning phase, we  
1345 set the maximum input sequence length to 256 and  
1346 the rank  $K$  in LoRA to 8. We use LLaMA-7B  
1347 model and initialize the checkpoints with the 8-bit  
1348 integer format (int8) parameters released by Tou-  
1349 vron et al. (2023). These parameters remain fixed  
1350 throughout training, reducing GPU memory con-  
1351 sumption and improving training speed. We use  
1352 the Adam optimizer to update LoRA parameters  
1353 with a batch size of 128 and learning rates of  $3e-4$ ,  
1354 respectively. The trainable LoRA parameters are  
1355 fine-tuned on NVIDIA A100-40GB GPUs, and the  
1356 training duration is approximately 15 hours. Fi-  
1357 nally, we yield a small ChatPal with compatible

<p><b>Below is an emotional support dialog. Your task is to generate four diverse responses for the emotional support assistant AI in [New Response]. Refer to the context and make the context coherent after filling in the answer you generated in [New Response].</b></p> <p><b>User:</b> Hi I am feeling blue, things aren't going well at all  <b>AI:</b> Hello. I am sorry to hear that. Can I ask you more about it?  <b>User:</b> My girlfriend and I split after 4 years!  <b>AI:</b> Oh no! I am sorry. How do you feel? Sad? Joyful?  <b>User:</b> Very sad and lonely. I loved her. What should I do?  <b>AI:</b> Oh no. I can totally understand that feeling. 4 years is a long time.  <b>User:</b> Thanks, it is. She said I was neglecting her and we had a big fight. I haven't talked to her in a month.  <b>AI:</b> [New response]  <b>User:</b> Do you think drinking alcohol will help?  <b>AI:</b> Ohhhh. Well, I never ditch my partner, always my partner ditches me, so I do not know why they do that. I think it will worsen the situation.</p> <p><b>Diverse responses:</b>  1.Focus on taking care of yourself physically and emotionally. Make sure you're eating well, getting enough sleep, and engaging in activities that bring you joy.  2.Share your feelings with people you trust. They can provide emotional support and be there for you during this tough time.  3.Take some time to think about what you've learned from this relationship and how you've grown. It can be a valuable opportunity for personal growth.  4.If you find it challenging to cope with your emotions, consider talking to a therapist or counselor who can provide guidance and support.</p>	<p><b>User:</b> You had a similar experience? "I want to drink till I forget about her.  <b>AI:</b> Yes, I quit alcohol 3 years ago. It took me 20 years.  <b>User:</b> How did you do it?  <b>AI:</b> It took me 1 whole year to forget about my ex. My ex told me I am a party girl because I drink too much, so I stopped. I was trying to get him to like me more. Since then I have not been drinking.  <b>User:</b> Did you find someone else?  <b>AI:</b> Nope. I am cool without a man.  <b>User:</b> That's the way to be. But I am too lonely. I have a secret account so I can still see her on Facebook. Do you think that's a good idea?  <b>AI:</b> Haha, smart, yes. No need to force yourself to forget about her. I heard men cannot do without a woman but women don't really need a man.  <b>User:</b> Yeah, that way maybe I can get her back at the right time. I feel a little better - thanks.</p>
--	---

Figure 6: An example of generation diverse responses. The DRI task description and the conversation context are given in ChatGPT to generate multiple diverse responses. The square below is the four different responses generated in [New Response].

1358 performance to much larger models, thereby signif- 1381  
1359 icantly alleviating the requirement for large model 1382  
1360 sizes. 1383

## 1361 I Baselines 1384

1362 We will compare our model with five different base- 1385  
1363 lines: 1386

1364 **LLaMA** (Touvron et al., 2023). LLaMA is an open 1387  
1365 and efficient large-scale base language model that 1388  
1366 sources publicly available datasets. This model is 1389  
1367 trained on a large amount of unlabeled data, mak- 1390  
1368 ing it well suited for fine-tuning a variety of tasks, 1391  
1369 and can be run on a single V100 GPU<sup>4</sup>. 1392

1370 **ChatGPT** (Ouyang et al., 2022). ChatGPT is a 1393  
1371 model for processing sequential data with amazing 1394  
1372 language understanding and text generation capa- 1395  
1373 bilities, and in particular, it trains the model by 1396  
1374 connecting it to a large corpus of real-world con- 1397  
1375 versations. ChatGPT can be used for a wide range 1398  
1376 of domains, including emotional support tasks. 1399

1377 **Ask-Expert** (Zhang et al., 2023a). Ask-Expert is a 1400  
1378 framework in emotional support domain, where the 1401  
1379 structure of expert conversation is outlined by pre- 1402  
1380 specified prompts which reflect a reasoning strat- 1403

egy taught to practitioners in the field. Blenderbot 1381  
model (Shuster et al., 2022) utilizing “Ask-Expert” 1382  
shows quality improvements across all expert sizes. 1383  
**AUGESC** (Zheng et al., 2023b). Zheng et al. 1384  
(2023b) prompt a fine-tuned LLM to complete full 1385  
dialogues from available dialogue posts of vari- 1386  
ous topics, which are then postprocessed based on 1387  
heuristics. They proposed AugESC dataset and 1388  
then fine-tuned Blenderbot model, which is supe- 1389  
rior to strong baselines of dialogue augmentation. 1390  
**Our Chatpal w/o DRI** We only fine-tune LLaMA 1391  
on our ExtES dataset w/o diverse response inpaint- 1392  
ing, which is an original variant of our small Chat- 1393  
Pal and can help us understand the influence of 1394  
diverse responses in Section 5.5. 1395

## 1396 J Guideline of Human Evaluation 1397

We present the guideline of human evaluation in 1397  
Figure 7. Before showing them the final evaluation 1398  
materials, we first train our human evaluators by 1399  
providing them this form, together with detailed 1400  
instructions on how to carefully do the evaluations, 1401  
what these metrics and corresponding scores mean 1402  
*etc.* 1403

<sup>4</sup>We chose the LLaMA-7B version based on the needs of the emotional support task.

<b>Guideline of Human Evaluation</b>	
You need to score the conversation between the help seeker (User) and the emotional support assistant (AI). Read the definitions and examples of evaluation metrics below to rate the results generated by different models. These examples illustrate how each metric can be applied to evaluate an emotional support conversation.	
Scores	3 (Excellent) , 2 (Good) , 1 (Accepted) , 0 (Unsatisfactory)
(1) Informativeness	
Definition	Informativeness measures how well the individual seeking support articulates their emotional challenges.
Examples	1. Low Informativeness: "I'm feeling really bad today." 2. High Informativeness: "I've been feeling overwhelmed because of work. I have tight deadlines, and my boss has been giving me extra tasks. I don't have much time for myself, and it's really stressing me out."
(2) Understanding	
Definition	Understanding gauges the supporter's grasp of the individual's experiences and emotions.
Examples	1. Low Understanding: "That sucks." 2. High Understanding: "I can imagine how stressful it must be to have such a heavy workload and demanding boss. It sounds like you're going through a tough time right now."
(3) Helpfulness	
Definition	Helpfulness evaluates the effectiveness of the supporter's efforts in mitigating the individual's emotional distress.
Examples	1. Low Helpfulness: "I'm sorry to hear that. I hope you feel better soon." 2. High Helpfulness: "It sounds like you could use some time management strategies to handle your workload more effectively. Have you considered talking to your boss about your workload or seeking support from colleagues?"
(4) Consistency	
Definition	Consistency ensures participants consistently adhere to their roles and exhibit non-contradictory behavior.
Examples	1. Inconsistent Behavior: Initially providing empathetic responses and later becoming dismissive or indifferent about the person's feelings. 2. Consistent Behavior: Maintaining a supportive and empathetic tone throughout the conversation, showing genuine care and concern.
(5) Coherence	
Definition	Coherence checks if conversations have seamless topic transitions.
Examples	1. Low Coherence: Frequent topic changes without exploring any of them in depth. For example, discussing work stress, then suddenly switching to talking about hobbies without any connection. 2. High Coherence: A focused conversation that explores a specific issue thoroughly before transitioning to a related topic. For instance, discussing work stress and then gradually shifting the conversation to coping mechanisms or self-care strategies.

Figure 7: Guideline of human evaluation for dialogue quality.