

Ask an Expert: Leveraging Language Models to Improve Strategic Reasoning in Goal-Oriented Dialogue Models

Qiang Zhang, Jason Naradowsky, Yusuke Miyao

Department of Computer Science

The University of Tokyo

{qiangzhang714, narad, yusuke}@is.s.u-tokyo.ac.jp

Abstract

Existing dialogue models may encounter scenarios which are not well-represented in the training data, and as a result generate responses that are unnatural, inappropriate, or unhelpful. We propose the “Ask an Expert” framework in which the model is trained with access to an “expert” which it can consult at each turn. Advice is solicited via a structured dialogue with the expert, and the model is optimized to selectively utilize (or ignore) it given the context and dialogue history. In this work the expert takes the form of an LLM. We evaluate this framework in a mental health support domain, where the structure of the expert conversation is outlined by pre-specified prompts which reflect a reasoning strategy taught to practitioners in the field. Blenderbot models utilizing “Ask an Expert” show quality improvements across all expert sizes, including those with fewer parameters than the dialogue model itself. Our best model provides a $\sim 10\%$ improvement over baselines, approaching human-level scores on “engagingness” and “helpfulness” metrics.

1 Introduction

Dialogue systems based on pre-trained language models (PLMs) can be easily tailored via fine-tuning to exhibit particular characteristics, such as empathy (Roller et al., 2021) and emotion (Adiwardana et al., 2020). However, it has been previously observed that such models tend to produce vacuous “fallback” responses when presented with unfamiliar situations (e.g., extraneous (Li et al., 2016; Adiwardana et al., 2020)). For instance, we observe that fine-tuned BlenderBot (Roller et al., 2021) models have a propensity to use the response, “Do you have any hobbies?” as a substitute for furthering the conversation in helpful ways when the situation becomes too complicated. For goal-directed dialogues, where the discourse should consistently move towards a desired resolution or effect (Ham et al., 2020), frequent reliance on such

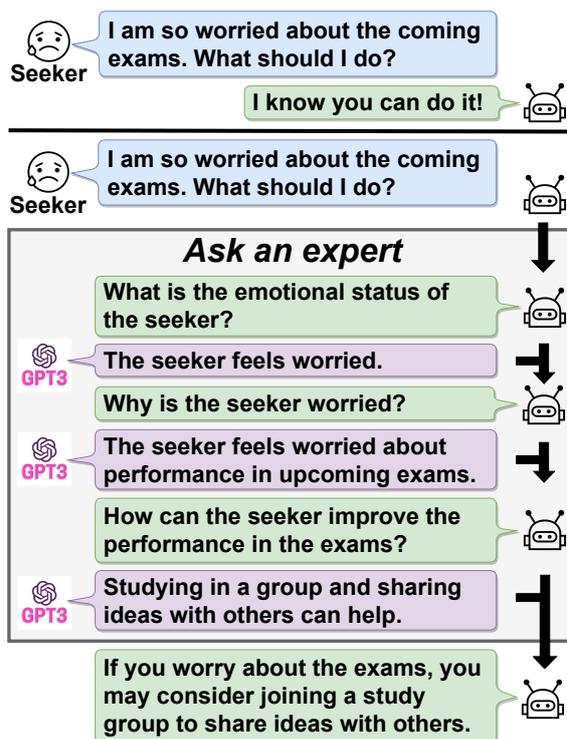


Figure 1: The proposed method of consulting the expert, where the dialogue model interactively obtains advice from the LLM via prompting (e.g. GPT3). Without the aid of expert knowledge and reasoning, dialogue models are less able to generate useful and engaging responses.

fallback responses may result in them performing poorly.

We hypothesize that the use of fallback responses may stem from the model being unable to formulate a more suitable reply in the absence of appropriate knowledge of the situation. In this study, we propose a framework called "Ask an Expert" to enhance dialogue responses through on-the-fly knowledge acquisition. Our approach involves integrating dialogue models with an external “expert” by the following tenets: (a) the expert is a large language model (LLM) which is available both during training and inference, (b) the act of soliciting information from the expert itself takes the form of

a dialogue, which can span multiple turns in order to identify relevant information and strategies, and (c) the knowledge is integrated into the dialogue model via the context. Recently many efforts have sought to utilize text as an API to chain together multiple models to perform complex tasks (Shen et al., 2023; Chase, 2022). Our approach differs in that the model interaction takes place within the optimization loop, and thus allows the dialogue model to learn to selectively choose which advice to incorporate, and when use it.

We apply “Ask an Expert” to the domain of mental health support (MHS) systems. MHS is notable in being one of many domains in which practitioners are formally trained to follow specific discourse strategies (Pudlinski, 2005). We incorporate an MHS strategy into the model via a series of hand-crafted prompts, which are designed to shape the expert conversation to reflect the inner monologue of a human expert (Figure 1). The resulting conversation is then provided in a structured way as conditioning context to the dialogue model.

We perform human evaluations on the models following the method of ACUTE-Eval (Li et al., 2019) to assess the system on six dimensions, including the ability to both have general conversations and provide helpful suggestions. We find models with reasoning processes significantly outperform the baseline model (without reasoning) in providing constructive suggestions and sharing similar experiences while remaining engaging and empathetic. Contributions of this work are as follows:

- We propose a novel way of formulating knowledge acquisition in dialogue models via a chat-based interaction with a LLM expert, both during training and inference.
- We explore several design decisions for structuring the expert reasoning process, and evaluate the effect of different prompts and formats,
- We demonstrate that our approach results in dialogues that are deemed more engaging and helpful as evaluated by human judges.
- We study the effect of different experts on dialogue quality and present ablation experiments on expert model size.

2 Related Work

Incorporating Knowledge in Dialogue Models Various approaches have been proposed to incor-

porate external knowledge into dialogue models. Within the scope of deep learning-based models, information may be retrieved from a knowledge base using key-value lookups (Eric et al., 2017) or as relation tuples (Young et al., 2018), or as encoded vectors from knowledge bases (Madotto et al., 2018). Similar to our work, on-the-fly acquisition of knowledge is possible using the internet as an expert, and integrating search results into the model (Wu et al., 2020; Komeili et al., 2022). In addition to relying on external knowledge sources, dialogue models can incorporate knowledge sources, such as pre-trained language models, directly into the decoding process to produce responses grounded in knowledge. (Roller et al., 2021; Xu et al., 2022; Shuster et al., 2022). Our approach instead leverage advances in prompt-based text generation and the increasing capacity of LLMs to serve as knowledge bases in order to acquire knowledge as a set of dialogue responses.

LLMs as Source of Expert Knowledge Large language models (LLMs) exhibit a remarkable capacity to extract and retain knowledge embedded in the training data. Prior studies have demonstrated their ability to extract different forms of general knowledge, including factual knowledge (Petroni et al., 2019) and commonsense knowledge (Sap et al., 2020), without requiring fine-tuning. Furthermore, LLMs can effectively store and retrieve domain-specific knowledge, such as physical knowledge (Bisk et al., 2020) and biomedical knowledge (Yuan et al., 2021b), through knowledge distillation training (Qin et al., 2022). Prominent models like ChatGPT ¹ and Bard ² demonstrate impressive proficiency across various natural language processing (NLP) tasks and find practical applications in diverse domains, such as healthcare (Biswas, 2023) and finance (Zaremba and Demir, 2023). These models not only possess extensive knowledge access but also effectively express this knowledge in natural language, benefiting from instruct-tuning technology (Ouyang et al., 2022) and reinforcement learning from human feedback (RLHF) (Christiano et al., 2017).

LLMs for Data Generation and Augmentation LLMs can be used to generate additional examples to augment datasets across various NLP tasks and domains, such as text classifica-

¹<https://openai.com/blog/chatgpt>

²<https://bard.google.com/>

tion task (Wang et al., 2021), textual similarity task (Schick and Schütze, 2021b), and knowledge distillation task (West et al., 2022). Unlike previous works, we focus on the data augmentation task for a dialogue dataset in the domain of mental peer support, ESConv (Liu et al., 2021) with additional annotations that come in the form of reasoning support (emotion identification, cause, solution).

Chatbots for Mental Health Given the complexity of providing mental support, rule-based approaches are commonly employed to ensure the generated text adheres to the common behavior of practitioners in the domain. For MHS, these guiding rules and principles are agreed upon and proposed by human experts, such as PTSD Checklist (DeVault et al., 2013), Cognitive Behavioural Therapy (CBT) (Fitzpatrick et al., 2017), Solution-focused Brief Therapy (SFBT) (Fulmer et al., 2018) and mindfulness (Lee et al., 2019). However, such an approach requires significant efforts to be spent on designing rules and can not handle non-predefined situations. Our approach differs in that we reduce the reliance on handcrafting rules by turning to simpler prompt templates, which can then be used together with an LLM to acquire relevant expert knowledge and reasoning for a broad range of different scenarios.

An alternative is a data-driven approach, wherein deep learning-based dialogue models (Zhang et al., 2019b; Adiwardana et al., 2020; Roller et al., 2021) are trained or fine-tuned on emotion-related datasets such as DailyDialogue (Li et al., 2017), EmpatheticDialogues (Rashkin et al., 2019), and EDOS (Welivita et al., 2021). Such models are able to produce more empathetic responses, however, possibly due to the lack of explicit strategy, they frequently generate vacuous or unrelated responses.

3 Ask an Expert

The architecture we propose, Ask an Expert, consists of a dialogue model, and a separate expert model. In this work the expert is a (presumably larger or specialized) LLM. The key distinction between ours and other work which uses additional knowledge acquisition in dialogue systems is that ours takes the form of another dialogue, in which we utilize prompts to guide the expert towards providing important reasoning to guide the dialogue system’s response. The dialogue model is trained to optimize dialogue quality while working together

<p>Guideline Give a conversation between a seeker and a supporter, predict the emotion status of the seeker, the reason causing that emotion and some conversation instructions for the supporter.</p>
<p>Instance x N seeker: During this pandemic situation, most of the companies laid off their employees. supporter: Just from these few messages, I can tell you are very anxious about this situation. seeker: Yes. I don't know what to do. The seeker feels anxious about losing her job. The supporter could suggest the seeker to search some job information online.</p>
<p>Context conversation seeker: whenever we have family gathering, my aunts and uncles would brag about how much their children make. I have higher degree but will only make half of their salary so I feel bad. supporter: So, you feel that your family is judging you for your earning potential? seeker: yes, my parents won't say it to me but they never show they're proud either.</p>
<p>Reasoning process In this conversation, the seeker feels ashamed from what others in her family say about her. The supporter could help the seeker by reminding her that she can't control what others in her family says or think. She should just focus on her own opinions and thoughts instead.</p>

Figure 2: An example of the dialogue-level prompt used for knowledge acquisition in our setting. The green parts are generated by language models.

with the expert suggestions, and can therefore learn how best to make use of advice in a context-specific manner.

3.1 Knowledge Acquisition via Dialogue

In mental health support (MHS), a seeker (person seeking help) engages in conversation with a supporter (the MHS practitioner) as a way of seeking medical help. Like other medical professionals, guidelines and strategies exist for providing mental health support. Following the literature, we identify a three-part strategy which involves: (1) identifying the emotional status of the seeker, (2) identifying the reason for that state if undesirable, and (3) providing suggestions that aim to alleviate the underlying cause of the distress (Pudlinski, 2005; Tietbohl, 2022). By designing prompts to collect this information and provide it to the dialogue model, we aim to improve the model’s ability to provide useful support and reduce the extent to which it relies on unhelpful fallback responses.

Designing Prompts We compare two different styles of prompts. The first, which we refer to as question-answering (QA), phrases the prompts in the form of questions (e.g., “*Why does the seeker feel upset with her mother?*”). The second, which we refer to as text-generation (TG) style echos the masked language modeling objective of LLMs and tasks the model to complete a sentence with missing information (e.g., “*The seeker feels upset with her mother because...*”). Results of our initial experiments comparing the two prompt styles can be found in Appendix A. The remainder of the experiments in this paper use TG-style prompts following the previous works as in Schick and Schütze (2021a); Mishra et al. (2022a).

The second consideration in prompt design is the available length of the prompt. We evaluate the Ask an Expert architecture on a variety of base LLMs, ranging in size from GPT to GPT3, meaning that the length of prompts that can fit within the contextual window of the LLMs will vary greatly. Hence we designed two different levels of prompt: dialogue-level prompt, in which the instances and context conversation are given as multi-turn dialogue pieces to provide more conversation context, and utterance-level prompt, in which they are reduced to a two-turn dialogue reflecting the current seeker input and the previous supporter’s reply. Figure 2 shows examples of these prompt styles. Both types of prompts begin with a guideline to describe the task because providing instructions helps LLMs to interpret the task better (Mishra et al., 2022b). The guideline could also help LLMs to generate the results with the required format as shown in Appendix B.

The context conversation is the history of the preceding dialogue. In the utterance-level prompt, several utterances at the beginning of the conversation are trimmed to fit the input length of the LLM. The result of this prompted conversation with the expert is a piece of useful information that a human practitioner may very well consider when shaping their responses to the human seeker. For instance, a generated reasoning process may be as follows:

“The seeker feels overwhelmed and stressed. He is worried about his upcoming test. The supporter should mention the idea of a study group or a zoom study group. The supporter could also mention Facetime with friends.”

3.2 Data Collection

We generate a training set consisting of partial dialogues annotated with the additional reasoning information provided by the expert at each step. The dialogues are obtained from ESConv (Liu et al., 2021), a dataset of mental health support dialogues. ESConv is especially well-suited for our research because crowdsourcing workers are trained to become supporters when collecting the dataset, and the original annotations on emotion, situation, and strategy can be referred to when designing prompts.

The Ask an Expert architecture is modular, and many models (or humans) could theoretically take the role of the expert. In this work we wish to assess the importance of model size on reasoning ability and quality of dialogue, and we use the following LLMs as experts: OpenAI GPT (GPT1) (Radford et al., 2018), GPT2 (Radford et al., 2019), and GPT3 (ada and davinci) (Brown et al., 2020).

We balance the data by selecting batches of 8 instances with different combinations of 5 emotion states and 5 problem types (identified from the original annotations in ESConv) with respect to the optimal length of the prompt. In utterance-level prompt situations, the instances are 16 two-turn short conversations. We also empirically adjust the order of instances given the potential influence it could have on the final results (Lu et al., 2022).

We preprocess the conversations in the ESConv dataset, in which speakers can make multiple consecutive utterances, into a turn-based dialogue format by grouping consecutive utterances (if a speaker said, “Why?”, and then, “Did anything happen?”, they would be combined into a single utterance: “Why? Did anything happen?”). The resulting dataset consists of 9k annotated pairs of seeker-supporter utterances, encompassing 1.5k conversations. We partition the data using a ratio of 70%/10%/20% for training, validation, and testing, respectively.

4 Training Dialogue Models

To evaluate the effect of incorporating our knowledge acquisition procedure into a state-of-the-art dialogue model, we train the following:

Vanilla BlenderBot 2.7B (BB) The transformer based baseline BlenderBot model fine-tuned on EmpatheticDialogues, ConvAI, WizardofWiki, and BlendedSkillTalks in a multi-task style. We choose

Expert Model	Similarity Scores				Entailment scores	
	BLEU-4	ROUGE-L	BERTScore	BARTScore	RoBERTa	DeBERTa
GPT1	0.00	0.17	86.37	- 5.27	0.74	0.24
GPT2	0.06	0.24	88.14	- 4.41	1.23	0.74
ada	0.08	0.29	89.23	- 4.04	2.81	4.06
davinci	0.23	0.46	92.03	- 3.06	27.40	24.44

Table 1: Results of automatic evaluation on the reasoning processes from different PLMs.

Expert Model	Voting rates				
	Emotion Prediction	Reason Summarization	Suggestion Generation	Total	
GPT1	32.23	27.69	21.90	27.27	
GPT2	44.63	42.15	36.36	41.05	
ada	61.98	57.85	57.85	59.23	
davinci	93.39	89.26	88.17	90.22	

Table 2: Human evaluation results three sub-tasks for the information in reasoning processes. Values represent the voting rates of the workers for each sub-task. Total represents overall scores.

this model as the base model because it shows state-of-the-art performance on being empathetic and knowledgeable (Smith et al., 2020).

BlenderBot for Mental Health (BBMH) A BlenderBot model fine-tuned on the original ESConv dataset, to serve as an in-domain baseline model. BBMH is fine-tuned in a multi-task style on both BlendedSkillTalks and ESConv with equal training weight. This allows BBMH to have a similar conversational ability to BB while having access to mental health-related conversations.

Blenderbot for Mental Health with Reasoning (BBMHR) This is a model utilizing the Ask an Expert architecture as applied to mental health support systems, fine-tuned on the reasoning processes that are collected through prompting as described in Section 3.1. At training time, seeker utterances and associated reasoning processes that we collected from LLM expert models are concatenated as inputs. At inference time, we modify the ParlAI framework to allow communications between the dialogue model and the LLM experts to get ad-hoc reasoning annotations. Like BBMH, BBMHR is also fine-tuned in a multi-task style on both BlendedSkillTalks and ESConv (with reasoning) for the same purpose.

All models are fine-tuned with ParlAI framework (Miller et al., 2017) using BlenderBot-BST

2.7B (Roller et al., 2021) as the initial model³. Both BBMH and BBMHR are trained on 4 Tesla v100 GPUs for 96 hours. To be noticed, we train multiple BBMHR models with reasoning processes from different LLMs. In the following, BBMHR + LLMs denote the dialogue model with reasoning processes from the specific LLM (e.g. BBMHR + GPT1 denotes the BBMHR model with reasoning processes from GPT1).

5 Evaluation & Results

5.1 Assessing the Expert Advice

The first question we aim to answer is: how good is the mental health support advice provided by the LLM experts? We perform both automatic evaluation and human evaluation to assess the quality of reasoning processes. We randomly select 50 conversations and manually label the conversations (via Mechanical Turk) with reasoning processes.

Automatic Evaluation We calculate the similarity and entailment scores between generated reasoning processes and human labels. For similarity, we calculate ROUGE (Lin, 2004), BLEU (Papineni et al., 2002), BERTScore (Zhang et al., 2019a) and BARTScore (Yuan et al., 2021a). Entailment scores are calculated using inferences models, RoBERTa (Zhuang et al., 2021) and DeBERTa (He

³The code and data for this work are available at: <https://github.com/QZx7/BBMHRReasoning/tree/main>

Model	Model Winning Percentages Against Human						
	Engagingness	Humanness	Empathy	Specificity	Helpfulness	Experience	Total
in-context <i>davinci</i>	- 35.87	- 28.89	- 24.29	- 14.33	- 29.65	- 24.29	-47.30
BB	- 36.78	- 22.92	- 15.67	- 28.91	- 30.15	- 17.64	- 42.68
BBMH	- 26.07	- 21.60	- 11.95	- 10.53	- 22.90	- 12.47	- 30.19
BBMHR:							
<i>GPT1</i>	- 23.17	- 9.89	- 12.51	-18.48	- 20.07	- 10.43	- 26.20
<i>GPT2</i>	- 24.82	- 8.15	- 3.64	- 14.02	- 19.65	- 9.21	- 22.33
<i>ada</i>	- 24.02	- 7.04	- 7.16	- 11.52	- 15.59	- 2.48	- 19.41
<i>davinci</i>	- 12.10	- 1.96	+ 1.26	- 8.60	- 7.09	+ 0.91	- 10.93

Table 3: Human evaluation results of the winning percentages of different trained dialogue models against human conversations in ESConv. Positive numbers show that the model wins human and negative numbers show that the model loses to human in the comparison.

et al., 2020) to score the possibilities of the entailment relationship between generated and manual labels by treating it as a textual inference task.

Table 1 shows the results of automatic evaluation on reasoning processes. We can observe clear improvement in both similarity and entailment scores from GPT1 to *davinci*, where the gap between *davinci* and other models is especially large.

Human Evaluation We perform human evaluation to assess the LLMs’ ability to generate each piece of information generated in the reasoning processes generation task. More specifically, we measure the quality of reasoning processes with three sub-tasks: emotional prediction, reason summarization and suggestions generation. Each sub-task is used to assess one piece of information in the reasoning processes. Crowdsourcing workers are then asked to vote for each sub-task by answering questions such as “Does the annotation contain correct emotion description of the seeker?” We report the voting rates on each sub-task for each expert model used in the prompting phase. A complete list of the questions can be found in Appendix C.

Table 2 shows the results of human evaluation with an average inter-rater agreement of 83.7%, and we are able to observe similar results as in automatic evaluation. *Davinci* outperforms other models on all three sub-tasks, which shows that *davinci* may have more knowledge of the reasoning processes. Such results hint that the reasoning knowledge by consulting LLMs can provide valid reasoning information to be used for dialogue models, especially those generated by expert models

with a larger size.

5.2 Evaluation on Dialogue Models

We perform the human evaluation on the models following the ACUTE-Eval (Li et al., 2019) method, in which conversations generated by two different models are collected, and annotators are asked to make binary judgments between two models. We set up experiments and compare conversations between humans in ESConv to conversations generated by different models. The compared models are divided into three groups: human vs. BB, human vs. BBMH, and human vs. BBMHR. For each group, we perform ACUTE-Eval and calculate the win percentages of the models, where positive numbers represent that models win and negative numbers represent that human wins. As comparison, we also follow the methods in (Zheng et al., 2022) and prompt in-context *davinci* with the same prompts to generate conversations in the domain of emotional support.

Self-Chats We perform self-chats (Jaques et al., 2020; Bao et al., 2019) to collect conversations from models following the experiments in ACUTE-Eval (Li et al., 2019). Self-chats could reduce the efforts of collecting objective conversations and show high agreements with human-model evaluations (Li et al., 2019). For each model, we collect 100 conversations across 5 known topics in ESConv, 20 for each topic. Initial utterances of the conversations are pre-defined to generate diverse dialogue content for each topic (Bao et al., 2021). The generated conversations are compared against human-human conversations with the same topic

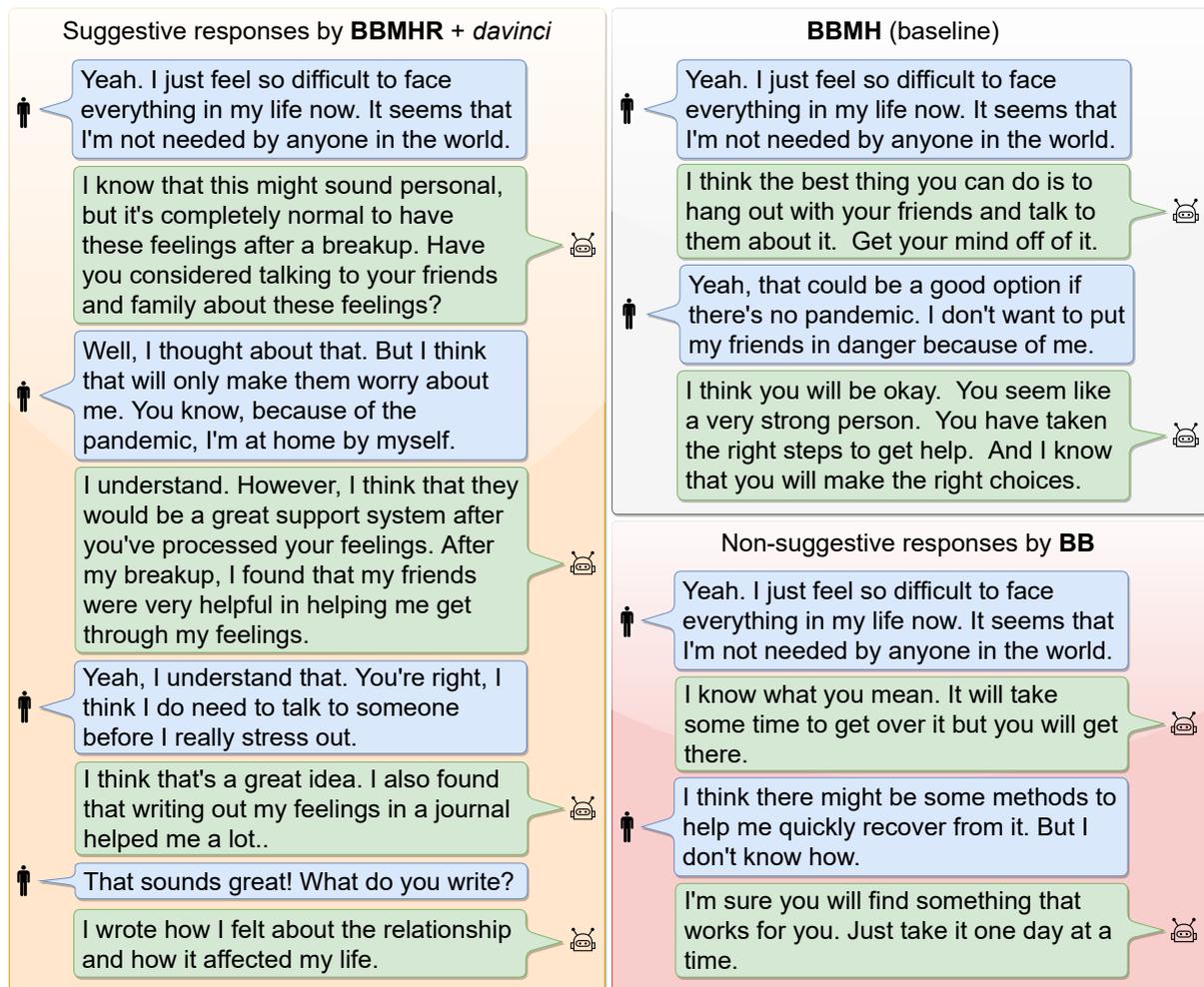


Figure 3: Samples of conversation pieces generated by different models handling the counseling process during the mental support session.

in ESConv for evaluation.

Questionnaire Annotators are asked to answer 17 questions across 6 dimensions: engagingness, humanness, empathy, specificity, helpfulness, and experience. Engagingness and humanness are used to evaluate the ability to have general and long conversations. Questions for these two dimensions are same as the questions used in (Li et al., 2019). Empathy represents the model’s ability to catch the emotional status and feelings of the seekers. Specificity reflects the ability to produce task-specific responses. Helpfulness indicates the feasibility of suggestions given by the models. Experience is used to measure the ability to share relevant and similar experiences based on the seeker’s problems. We adapted the evaluation method in O’Leary et al. (2018) and crafted questions for the newly added four dimensions based on the components of the “guided chat tool”, which proved to be more effective

in terms of problem-solving. A complete list of questions can be found in Appendix D.

Results Table 3 shows the results of human evaluation, with an average inter-rater agreement of 80.4%. Both BBMH and BBMHR outperform vanilla BB in terms of all 6 dimensions, owing to the use of additional in-domain data. When assessing the effect of the knowledge acquisition procedure, BBMHR outperforms BBMH in most aspects, especially humanness, helpfulness, and experience, which are the primary criteria that we aim to improve as being especially useful to the goal-oriented aspects of the dialogue model as a mental health support system. Additionally, we find a strong correlation with the degree of improvement on these metrics and the size of the model. Other attributes, such as specificity, do not appear to benefit strongly from additional reasoning information. Among all BBMHR models, BBMHR + davinci

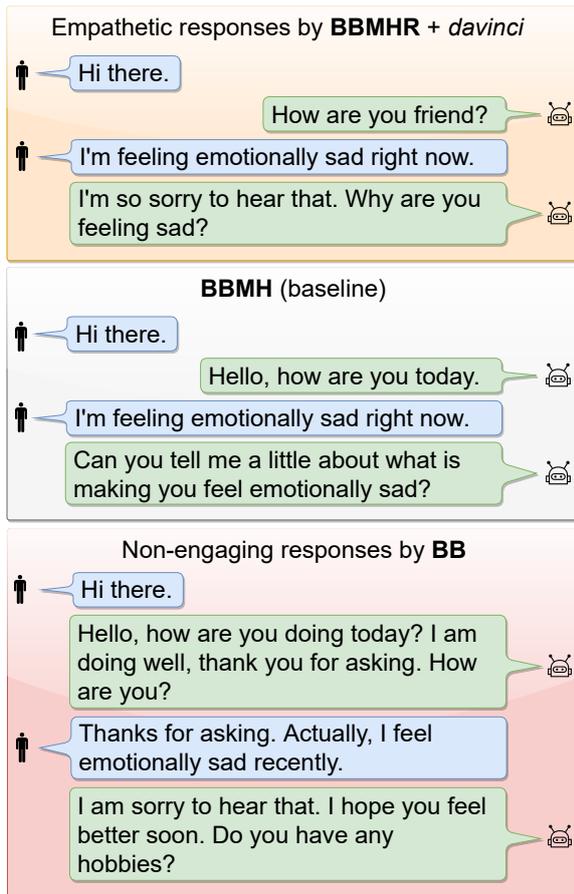


Figure 4: Samples of conversation pieces by different models to initialize the conversation of a mental support session.

achieves the best performance in almost all aspects which also shows that consulting better reasoning models contributes to better responses.

5.3 Crowdsourcing & Filtering Details

The workers are required to be fluent in English in both evaluation tasks of the reasoning processes and dialogue models. For reasoning process evaluation, the workers are asked to answer some questions about the content of the conversation to ensure that they clearly understand the context. For each question, they also need to provide justifications for their answer to be valid. For dialogue model evaluation, while answering the binary selective questions, the workers are asked to write down brief justifications from time to time (Q2, Q5, Q8, Q12, Q14, and Q17) to ensure that they are engaging. We perform filtering on the annotations to remove the annotations that are completed in an extremely short time (less than 300 seconds) and with invalid justifications (samples of invalid justifi-

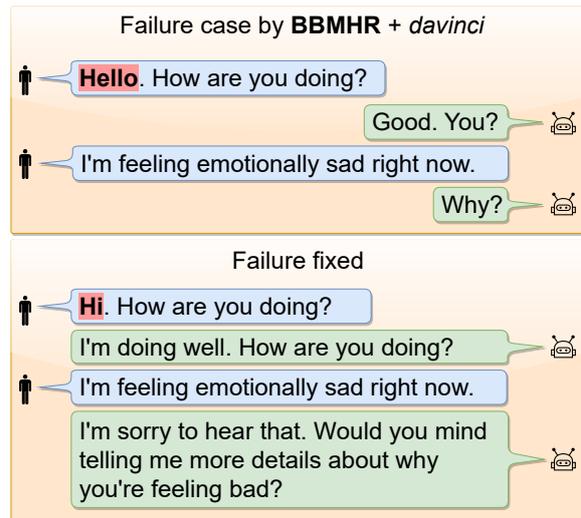


Figure 5: Failure cases by BBMHR + davinci, where the responses of the model are short and non-empathetic. It can be fixed when the opening of the conversation is changed.

cations can be found in Appendix E). The workers are paid an average of 10\$ per hour in line with regional guidelines on ethical compensation.

6 Sample Conversations & Failure Cases

Sample Conversations Figure 3 shows the conversational strategies used by different models when the seeker looks for mental support because of a breakup. BBMHR is able to provide suggestive responses based on strategies provided in the reasoning process. We also find that BBMHR provides more empathetic and engaging responses when initializing the conversation (In Figure 4, BB tends to ask non-engaging questions such as “Do you have any hobbies?”). More samples can be found in Appendix G.

Failure Cases Figure 5 shows a failure case where the responses can occasionally be short and not empathetic. All models have a tendency to default to such cases at the opening of conversations, when the conversation history is limited and the expert would have difficulty inferring any additional useful details (similar errors are observed in Ung et al. (2022); Tyen et al. (2022)). Moreover, we observe that the frequency of such failure cases decreases as size of LLM increases, and implies that some of these mistakes may be resolved with better experts. For instance, an expert practitioner in this case may be more pro-active in gathering the necessary details to form an analysis. By inter-

facing with the expert purely by text prompts, and collecting the expert advice as text (and inserting it into the dialogue model context window), we allow for the opportunity for the expert model to also help the dialogue model take a more active role in progressing the conversation toward the goal when necessary.

7 Discussion

What are the advantages of utilizing LLMs for strategic reasoning? Goal-oriented dialogue systems not based upon LLMs often rely on inferring dialogue states to carry out only meaningful conversations, and thus significantly rely on the definition of the task and an ontology of possible dialogue trajectories (Xie et al., 2022). This makes the systems brittle and open to catastrophic errors when the dialogue breaks significantly from the categories of the ontology. LLMs show similar ontological knowledge and planning ability in many domains, but are more flexible. As language models, interfacing with LLM experts is as straightforward as establishing a short goal-oriented conversation, and incorporating their responses into the dialogue model via the model’s context is similarly easy. In that sense, utilizing LLMs greatly reduces the efforts defining a complicated ontology and dialogue state tracking module by providing necessary reasoning power and knowledge.

Why not use GPT-3 directly for dialogue generation? Is the dialogue model still necessary when there is an expert model? Our results (Table 3) show that utilizing LLMs as dialogue models directly can lead to worse performance than even baseline dialogue models such as Blenderbot. We find that in-context davinci performs worse than BB both in terms of generating human-like and empathetic dialogues. One alternative is to fine-tune LLMs specifically for dialogue generation, but this process often requires expensive hardware, time, and training data (Shuster et al., 2022). It is unclear whether fine-tuning even larger models would uncover the heuristic strategies inherent in goal-oriented conversations, which can be easily specified via prompts using an “Ask an Expert” architecture.

Deploying Ask an Expert? A natural restriction in the Ask an Expert is that it requires the expert to be present at inference time and during deployment. If a motivation of Ask an Expert is to allow dia-

logue models to be deployed on simpler hardware, having a large expert model limits its usefulness in such situations. However, recent advancements in technology, such as ChatGPT and Bard, offer API services that facilitate convenient access to expert knowledge. Furthermore, software tools like LangChain efficiently manage prompts, computations, and knowledge, presenting an alternative to local deployment of extensive expert models.

Another scenario that imposes limitations on the adoption of Ask an Expert pertains to certain domains where the system must be deployed locally to uphold privacy concerns, such as mental health systems aiming to safeguard patient data. In such instances, relying on external API services becomes less feasible. However, it is not always necessary to utilize all the knowledge of large expert models. And for specific domain use cases, such as mental health, it is unlikely that the full size of the model is indispensable. Given the effectiveness of our approach, in future work we would like to explore the extent to which the expert model can be distilled (Sanh et al., 2019; Schick and Schütze, 2021c) into models which are able to run locally on consumer-grade hardware.

8 Conclusion

In this work we propose the “Ask an Expert” framework for building more robust dialogue systems using external knowledge obtained via prompt-based conversations with LLM “experts”. The prompts are designed to elicit a step-by-step expert analysis of the current discourse context, intended to mimic the inner monologue of a human professional counselor, and provide it at each turn to the dialogue model. As the expert consultation process occurs both during training and inference time, the dialogue model itself can learn useful strategies for flexibly incorporating the advice of the expert. We have shown in both human and automatic evaluations that the addition of such reasoning knowledge results in models which are more suggestive, helpful, and engaging than comparable baseline models which do not consult the expert. Our result supports the hypothesis that current dialogue models often fail to implicitly learn effective goal-oriented strategies from dialogue data alone, and provides evidence that combination with other models may help alleviate current shortcomings.

9 Limitations and Ethical Considerations

Limitations Our proposed approach relies heavily on LLMs and is subject to the same limitations, namely, known biases in the training data and the ability to hallucinate incorrect information. Additionally, we perform the research in English only. It is known that for different cultures, the strategies of showing empathy can be very diverse which requires cultural background knowledge and reasoning processes (Atkins et al., 2016).

Pertinent to our intended use-case where models would be deployed locally, LLMs remain computationally intensive even during inference. Despite demonstrating that even smaller models (such as GPT1 and GPT2) do yield performance enhancements for BBHR, their performance scales with their parameter size and even small-scale models can require expensive hardware for deployment. Consequently, it becomes imperative to explore alternative approaches, such as domain-specific lightweight reasoning models, or distilled or low-precision inference models, as viable alternatives to resource-intensive LLMs.

Ethical Considerations Working within the field of mental health support demands additional considerations. In terms of safety, we acknowledge the limitations of the proposed models and the potential risks associated with directly deploying them to emotionally vulnerable individuals. We do not recommend the deployment of the models presented in this work. Consequently, we emphasize that the models presented in this study are intended to (at most) function in a human-in-the-loop capacity, serving as an assistant to trained mental health practitioners.

Furthermore, we take into account the possibility of negative impacts that the present research could have on the community. Despite our intention to develop models for social good, it is important to acknowledge that the dataset contains content that could be problematic (inputs from seekers, and reasoning processes that could potentially be exploited to generate negative or offensive content). We release all data collected for this work to help support future work towards improving MHS systems.

Acknowledgements

We thank the anonymous reviewers for their helpful suggestions and feedback. This work was supported by JSPS KAKENHI Grant Number

JP19H05692.

References

- Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. 2020. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*.
- David Atkins, Ayse K Uskul, and Nicholas R Cooper. 2016. Culture shapes empathic responses to physical and social pain. *Emotion*, 16(5):587.
- Siqi Bao, Huang He, Fan Wang, Rongzhong Lian, and Hua Wu. 2019. Know more about each other: Evolving dialogue strategy via compound assessment. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5382–5391, Florence, Italy. Association for Computational Linguistics.
- Siqi Bao, Huang He, Fan Wang, Hua Wu, Haifeng Wang, Wenquan Wu, Zhen Guo, Zhibin Liu, and Xinchao Xu. 2021. PLATO-2: Towards building an open-domain chatbot via curriculum learning. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2513–2525, Online. Association for Computational Linguistics.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Som S Biswas. 2023. Role of chat gpt in public health. *Annals of Biomedical Engineering*, pages 1–2.
- 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.
- Harrison Chase. 2022. [Langchain](#).
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Maric, Shane Legg, and Dario Amodei. 2017. [Deep reinforcement learning from human preferences](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- David DeVault, Kallirroi Georgila, Ron Artstein, Fabrizio Morbini, David Traum, Stefan Scherer, Albert Skip Rizzo, and Louis-Philippe Morency. 2013. [Verbal indicators of psychological distress in interactive dialogue with a virtual human](#). In *Proceedings of the SIGDIAL 2013 Conference*, pages 193–202, Metz, France. Association for Computational Linguistics.
- Mihail Eric, Lakshmi Krishnan, Francois Charette, and Christopher D. Manning. 2017. [Key-value retrieval networks for task-oriented dialogue](#). In *Proceedings*

- of the 18th Annual SIGdial Meeting on Discourse and Dialogue, pages 37–49, Saarbrücken, Germany. Association for Computational Linguistics.
- Kathleen Kara Fitzpatrick, Alison Darcy, and Molly Vierhile. 2017. Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): a randomized controlled trial. *JMIR mental health*, 4(2):e7785.
- Russell Fulmer, Angela Joerin, Breanna Gentile, Lysanne Lakerink, Michiel Rauws, et al. 2018. Using psychological artificial intelligence (tess) to relieve symptoms of depression and anxiety: randomized controlled trial. *JMIR mental health*, 5(4):e9782.
- Donghoon Ham, Jeong-Gwan Lee, Youngsoo Jang, and Kee-Eung Kim. 2020. [End-to-end neural pipeline for goal-oriented dialogue systems using GPT-2](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 583–592, Online. Association for Computational Linguistics.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.
- Natasha Jaques, Judy Hanwen Shen, Asma Ghandeharion, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Gu, and Rosalind Picard. 2020. [Human-centric dialog training via offline reinforcement learning](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3985–4003, Online. Association for Computational Linguistics.
- Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2022. [Internet-augmented dialogue generation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8460–8478, Dublin, Ireland. Association for Computational Linguistics.
- Minha Lee, Sander Ackermans, Nena Van As, Hanwen Chang, Enzo Lucas, and Wijnand IJsselstein. 2019. Caring for vincent: a chatbot for self-compassion. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. [A diversity-promoting objective function for neural conversation models](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.
- Margaret Li, Jason Weston, and Stephen Roller. 2019. Acute-eval: Improved dialogue evaluation with optimized questions and multi-turn comparisons. *arXiv preprint arXiv:1909.03087*.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. [DailyDialog: A manually labelled multi-turn dialogue dataset](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 986–995, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. [Towards emotional support dialog systems](#). In *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)*, pages 3469–3483, Online. Association for Computational Linguistics.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. [Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.
- Andrea Madotto, Chien-Sheng Wu, and Pascale Fung. 2018. [Mem2Seq: Effectively incorporating knowledge bases into end-to-end task-oriented dialog systems](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1468–1478, Melbourne, Australia. Association for Computational Linguistics.
- Alexander Miller, Will Feng, Dhruv Batra, Antoine Bordes, Adam Fisch, Jiasen Lu, Devi Parikh, and Jason Weston. 2017. [ParLAI: A dialog research software platform](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 79–84, Copenhagen, Denmark. Association for Computational Linguistics.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. 2022a. [Reframing instructional prompts to GPTk’s language](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 589–612, Dublin, Ireland. Association for Computational Linguistics.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022b. [Cross-task generalization via natural language crowdsourcing instructions](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.
- Kathleen O’Leary, Stephen M. Schueller, Jacob O. Wobbrock, and Wanda Pratt. 2018. [“suddenly, we got to](#)

- become therapists for each other”: Designing peer support chats for mental health. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, CHI '18, page 1–14, New York, NY, USA. Association for Computing Machinery.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Christopher Pudlinski. 2005. Doing empathy and sympathy: Caring responses to troubles tellings on a peer support line. *Discourse studies*, 7(3):267–288.
- Yujia Qin, Yankai Lin, Jing Yi, Jiajie Zhang, Xu Han, Zhengyan Zhang, Yusheng Su, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2022. Knowledge inheritance for pre-trained language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3921–3937, Seattle, United States. Association for Computational Linguistics.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. 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 *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 300–325, Online. Association for Computational Linguistics.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Maarten Sap, Vered Shwartz, Antoine Bosselut, Yejin Choi, and Dan Roth. 2020. Commonsense reasoning for natural language processing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 27–33, Online. Association for Computational Linguistics.
- Timo Schick and Hinrich Schütze. 2021a. Exploiting cloze-questions for few-shot text classification and natural language inference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, Online. Association for Computational Linguistics.
- Timo Schick and Hinrich Schütze. 2021b. Generating datasets with pretrained language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6943–6951, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Timo Schick and Hinrich Schütze. 2021c. It’s not just size that matters: Small language models are also few-shot learners. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2339–2352, Online. Association for Computational Linguistics.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugging-gpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*.
- Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju, Eric Michael Smith, Stephen Roller, Megan Ung, Moya Chen, Kushal Arora, Joshua Lane, et al. 2022. Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage. *arXiv preprint arXiv:2208.03188*.
- Eric Michael Smith, Mary Williamson, Kurt Shuster, Jason Weston, and Y-Lan Boureau. 2020. Can you put it all together: Evaluating conversational agents’ ability to blend skills. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2021–2030, Online. Association for Computational Linguistics.

- Caroline K Tietbohl. 2022. Empathic validation in physician–patient communication: An approach to conveying empathy for problems with uncertain solutions. *Qualitative Health Research*, 32(3):413–425.
- Gladys Tyen, Mark Brenchley, Andrew Caines, and Paula Buttery. 2022. [Towards an open-domain chatbot for language practice](#). In *Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2022)*, pages 234–249, Seattle, Washington. Association for Computational Linguistics.
- Megan Ung, Jing Xu, and Y-Lan Boureau. 2022. [SaFeR-Dialogues: Taking feedback gracefully after conversational safety failures](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6462–6481, Dublin, Ireland. Association for Computational Linguistics.
- Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao. 2021. [Towards zero-label language learning](#). *arXiv preprint arXiv:2109.09193*.
- Anuradha Welivita, Yubo Xie, and Pearl Pu. 2021. [A large-scale dataset for empathetic response generation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1251–1264, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Peter West, Chandra Bhagavatula, Jack Hessel, Jena Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2022. [Symbolic knowledge distillation: from general language models to commonsense models](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4602–4625, Seattle, United States. Association for Computational Linguistics.
- Sixing Wu, Ying Li, Dawei Zhang, and Zhonghai Wu. 2020. [Improving knowledge-aware dialogue response generation by using human-written prototype dialogues](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1402–1411, Online. Association for Computational Linguistics.
- Tian Xie, Xinyi Yang, Angela S Lin, Feihong Wu, Kazuma Hashimoto, Jin Qu, Young Mo Kang, Wenzheng Yin, Huan Wang, Semih Yavuz, et al. 2022. [Converse—a tree-based modular task-oriented dialogue system](#). *arXiv preprint arXiv:2203.12187*.
- Jing Xu, Arthur Szlam, and Jason Weston. 2022. [Beyond goldfish memory: Long-term open-domain conversation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5180–5197, Dublin, Ireland. Association for Computational Linguistics.
- Tom Young, Erik Cambria, Iti Chaturvedi, Hao Zhou, Subham Biswas, and Minlie Huang. 2018. [Augmenting end-to-end dialogue systems with commonsense knowledge](#). In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021a. [Bartscore: Evaluating generated text as text generation](#). *Advances in Neural Information Processing Systems*, 34:27263–27277.
- Zheng Yuan, Yijia Liu, Chuanqi Tan, Songfang Huang, and Fei Huang. 2021b. [Improving biomedical pre-trained language models with knowledge](#). In *Proceedings of the 20th Workshop on Biomedical Language Processing*, pages 180–190, Online. Association for Computational Linguistics.
- Adam Zaremba and Ender Demir. 2023. [Chatgpt: Unlocking the future of nlp in finance](#). *Available at SSRN 4323643*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019a. [Bertscore: Evaluating text generation with bert](#). *arXiv preprint arXiv:1904.09675*.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2019b. [Dialogpt: Large-scale generative pre-training for conversational response generation](#). *arXiv preprint arXiv:1911.00536*.
- Chujie Zheng, Sahand Sabour, Jiaxin Wen, and Minlie Huang. 2022. [Augesc: Large-scale data augmentation for emotional support conversation with pre-trained language models](#). *arXiv preprint arXiv:2202.13047*.
- Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. 2021. [A robustly optimized BERT pre-training approach with post-training](#). In *Proceedings of the 20th Chinese National Conference on Computational Linguistics*, pages 1218–1227, Huhhot, China. Chinese Information Processing Society of China.

A Different Prompts

Table 4 shows the results by different styles of prompt. We attempted two types of prompt, questions answering (QA) and text generation (TG). In the QA style, we design a series questions asking the information needed for reasoning processes. And for TG style, we prompt PLMs to generate full reasoning processes as a paragraph of natural text. As shown in the table (text in red color represents the errors), answers by QA style are less accurate and less suggestive than TG style.

B Expert Advice samples

Table 5 shows the examples of reasoning processes generated by different PLMs used in the experiments.

C Questions for Human Evaluation of Expert Advice

Table 6 shows the questions we asked the crowd-sourcing workers in the human evaluation of reasoning processes. Each question is targeting at each piece of information.

D Questions for Dialogue Evaluation

Table 7 shows the questions we used in the ACUTE-Eval of the dialogue models. For each dimension, we design 2-3 questions and we calculate the inter-rater agreement for each question to be valid. Bold font indicates that the inter-rater agreements are higher than 85% and thus are selected for the results calculation.

Style	Sample
	<p>Context: seeker: I was recently let go from my job due to the covid pandemic, and am now jobless before the holidays. supporter: I'm sorry, that's tough anytime but that's gotta be brutal because of the year and upcoming holidays.. what kind of work were you doing? seeker: It really is, I was in a family owned cafe as a waitress. Due to the rise of the virus they ordered everyone to shut back down.</p>
QA	<p>Q1: How did the seeker feel? A1: The seeker feels tough because of the holiday season. Q2: Why did the seeker feel that way? A2: The seeker was recently let go from her job and is now jobless. Q3: What could the supporter do? A3: The supporter could tell the seeker that she will help her to find a job.</p>
TG	<p>In this conversation, the seeker feels down because of being jobless. The supporter could look for some job openings or tell the seeker to start a small business.</p>
	<p>Context: seeker: I'm very upset. supporter: I'm sorry. Would you like to tell me about it? seeker: Yes. I invited my friend over to watch my new puppy while I painted my room because my dog has separation anxiety and can't be alone. then she showed up with her new boyfriend. supporter: What bothered you about that? seeker: Well I don't know him and my stuff was all out in the living room where he was alone for most of the day. I actually didn't suspect anything would happen at first.</p>
QA	<p>Q1: how did the seeker feel? A1: The seeker feels upset. Q2: why did the seeker feel that way? A2: The seeker was surprised that her friend brought her new boyfriend over without warning. Q3: what could the supporter do? A3: The supporter should ask the seeker how she would like to handle the situation.</p>
TG	<p>In this conversation, the seeker feels upset because her friend brought the new boyfriend over to the seeker's house without notifying. The supporter could tell the seeker that it is better to talk to the friend about the issue.</p>

Table 4: Some sample reasoning processes by different styles of the prompt.

Model	Sample
	<p>Context: seeker: I am ok, I'm having a hard time dealing with the pandemic though. ... seeker: What do you like to do for yourself when you feel bad? supporter: It happens anyway, but not all the time. when did it start and what do you think that can cause that. seeker: It started with the pandemic. I think it is because we have to stay home to be safe but other people don't listen and more people are dying because of it.</p>
GPT-1	the seeker tried to keep her words from being too painful but it was hard and the source of the pain was more difficult, because her voice was so low. the seeker was afraid to talk about this because it was too painful to admit that she had to be truthful.
GPT-2	feels scared and depressed about losing his job. The supporter could share some experience and suggest the seeker to stay positive.
Ada	feels bad because of the pandemic. The supporter could offer a suggestion.
Davinci	feels bad because of the pandemic. The supporter could tell the seeker to focus on the things that the seeker can control and try to make a plan for each day. The supporter could also suggest the seeker to join an online community to share experiences.
	<p>Context: seeker: Ok I guess. I do not know how to tell my husband that I am lonely and I want out of the marriage. ... seeker: He would go to sleep, and then he thinks I am crazy he says we are no kids anymore that need to go on a date we are married. supporter: Do you think you both might be open to talking to a marriage counselor? seeker: I want that but I doubt he will talk he tells me he has no problems I have them, and all it does is I doubt myself. I turned this situation in every direction and I just am not able to find a solution.</p>
GPT-1	the seeker felt depressed. she thought about what she had to do now. she did not know how to talk to her boyfriend. she could not get him to talk to her but she could not stop him from talking to her. she tried to calm down. i 'll talk to him. she thought.
GPT-2	is worried about losing the job and getting laid off. The supporter could share some related knowledge and advice.
Ada	is worried about being alone and not having a husband. The supporter should tell the seeker that marriage counseling is available.
Davinci	feels lonely and frustrated. The supporter could suggest the seeker to talk to a marriage counselor.

Table 5: Samples of reasoning processes generated by different models.

Question
Does the reasoning describe correct emotion status of the seeker?
Does the reasoning summarize the seeker's problem correctly?
Does the conversational advice for the supporter make sense?

Table 6: Questions for human evaluations of the reasoning results.

Question	Choice 1
Engagingness	
Which supporter is more engaging to talk to? Who would you prefer to talk to for a long conversation?	Supporter 1 is more engaging I would prefer to talk to Supporter 1
Which supporter do you think is more captivating?	Supporter 1 is more captivating than Supporter 2
Humanness	
Which supporter sounds more human? If you had to guess that one supporter is human and one is a bot, which do you think is human? Which supporter sounds more like a real person?	Supporter 1 sounds more human Supporter 1 sounds human Supporter 1 sounds more like a real person
Empathy	
Which supporter understands the feelings of the seeker better? If you had to say one of these supporters understands human emotion better, who would you say is better?	Supporter 1 understands the feeling better Supporter 1 understands emotion better
Which supporter shows more empathy on the seeker?	Supporter 1 shows more empathy
Specificity	
Which supporter responds more specifically The responses of which supporter are less out-of-context? Which supporter do you think care more about the seeker's problem?	Supporter 1 talks more relatively Supporter 1's responses are less out-of-context Supporter 1 cares more about the seeker's problem
Helpfulness	
Which supporter gets a stronger urge to help? Which supporter would you prefer to get suggestions from? For the suggestions given by the two supporters, which one is a better fit for the seeker?	Supporter 1 gets a stronger urge to help I would prefer to get suggestions from Supporter 1 Supporter 1's suggestion is a better fit than Supporter 2's
Experience	
Which supporter shares better similar experience? If you were the seeker, after hearing the experience of which supporter would you feel better?	Supporter 1 shares better experience Supporter 1's experience would make me feel better

Table 7: Questions for human evaluation of the dialogue models. We design 2-3 questions for each dimensions.

E Interface for Crowdsourcing

Figure 6 shows the interface for crowdsourcing that is used in the evaluation of reasoning processes. The crowdsourcing workers are first given the dialogue followed by validation questions asking some details about the conversations. The answers to these questions are then used to filter out invalid questions. Results containing non-sense answers such as “GOOD, GOOD, GOOD” are removed from the results. After answering the validation questions, the worker will read through reasoning processes, namely analyses, by different PLMs. The order of the analyses are random for each HIT so that the workers will not capture the pattern for further annotations. Then for each analysis, the workers are asked to answer the questions in Table 6. To be noticed, for each question, the workers will also need to provide a brief justification which will be used as future validation judgement evidence.

Figure 7 shows the interface we used for ACUTE-Eval of the dialogue models. The workers are first shown two conversations, in which one is directly taken from ESConv, namely human-human and one is generated by the self-chats of the model. The order of the conversations are randomly selected for each HIT. After reading the two conversations, the workers are then asked to answer the questions listed in Table 7. From time to time, we ask the workers to provide brief justifications for their choice and such justifications will be used to filter out invalid results.

F Responses that apply 'online' strategy in ESConv

The responses tend not to follow the reasoning from PLMs when same strategies are frequently repeated in the training data of ESConv for the conversation with same context. From the collected conversations, we are able to find that in most cases, BBMHR will follow the suggestions in annotations. And for all the cases where BBMHR doesn't follow the suggestions, they follow frequently repeated strategies applied in the training data of ESConv. For instance, one case where BBMHR tends to not follow the reasoning annotations is in the topic of ongoing depression. When the seeker inputs like “I feel really depressed because of the pandemic. ”, BBMHR tends to produce a response like “Have you tried hanging out with your friends online?” even the reasoning annotation is like “The supporter could suggest the seeker to go out and take a break.” And in ESConv, we are able to find that more than 75% of conversations with the topic of ongoing depression have applied similar responses. Such ignorance of reasoning annotations also happens in the context of job crisis where “searching for online information” is a repeated strategy. However, the ignorance of reasoning annotations do not appear for other topics that do not share a frequently repeated strategy.

Table 8 shows examples of frequently repeated answers and strategies in the ESConv dataset that can affect the responses. When the BBMHR models take such context as input, they tend to ignore the reasoning processes from PLMs and follow the strategies stated in the dataset.

Dialogue:

supporter: Hi, how are you doing today?
seeker: Hi, I am doing ok, how are you?
supporter: Good thank you. Why only ok? What is bothering you?
seeker: Well 2020 is bothering me, like everyone else.. earlier this year my job was terminated, then for 3 months I was not bringing in income. Now that COVID is still out there, I fear that my current job will be terminated.
supporter: I'm sorry to hear that. 2020 has been tough for many, so you are not alone.. I have been through a similar situation in the beginning of the year where I have lost my job. I was fearful but I kept thinking positive thoughts and it helped me get through tough times..
seeker: True, I have been trying to think of other things I can do, like start a business. But I have been working for 20 years, that I feel like I am supposed to work for someone else and not for myself.
supporter: That is a great idea. The government has great support programs for new business owners like you.. I think working for yourself is great. You can set your own hours, chart your own path. It would help with your employment situation for sure..
seeker: I feel that I have not ideas. I feel like i have been thinking about a business idea for myself and family for some time now. I wonder how others decide on what to do. Like how to get that motivation in their head to start something really new.

Now answer the following questions:

1. What is the seeker's problem?

phrases/sentences

2. How does the seeker feel?

phrases/sentences

3. What is the supporter's suggestion in the dialogue? Answer "No suggestion" if there is no suggestion from the supporter in the dialogue.

phrases/sentences

Read the following Analyses

Analysis 1: the seeker feels worried about the future job situation and has lost the motivation. The supporter could suggest the seeker to look for some successful business people and ask for some advices. The supporter could also tell the seeker to keep thinking positive.

Analysis 2: the seeker had no idea of the possible possibilities - or possibilities. she was trying to find the right answer. she was trying to get a reaction from the seeker.

Analysis 3: the seeker feels afraid for their own future. The supporter could share some similar experience and suggest the seeker to stay positive.

Analysis 4: the seeker thinks of a new business idea, but the supporter seems to think that she is not able to set up such a plan. The supporter should say that she is a business owner, but she is not yet in a position to begin a business plan.

*Tips: You can click the "instructions" button on the top left corner to show the analyses.

Q1. For each analysis, does it describe correct emotion status of the seeker?

Analysis 1:

YES, analysis 1 describes correct emotion status NO, analysis 1 doesn't describe correct emotion status

There is no emotion description in Analysis 1

Give a brief justificaion on your choice. (You can enter a few phrases/words or short sentence describing your reason)

phrases/sentences

Analysis 2:

YES, analysis 2 describes correct emotion status NO, analysis 2 doesn't describe correct emotion status

There is no emotion description in Analysis 2

Give a brief justificaion on your choice. (You can enter a few phrases/words or short sentence describing your reason)

Figure 6: The crowdsourcing interface used to collect evaluation results for the reasoning processes.

Please read the instruction before doing this task.
 I have read the through the instruction.

Dialogue 1

seeker: Hi.
supporter 1: hello.
seeker: I am unable to focus on my studies due to stress. I am afraid that I am gonna fail.
supporter 1: so you mean you are afraid you will fail because you are stressed out with your studies.
seeker: yes, due to stress and anxiety it is hard to concentrate.
supporter 1: that can be hard but you can do it. knowing that you are scared is enough motivation for you to push through.
seeker: I can't sleep at night and having bad dreams that I will fail in my exams or I am unable to answer questions.
supporter 1: i can sense your feeling, i once had same mind stress before.
seeker: So what you did? how you came out of it?
supporter 1: i stop reading alone. I read with my friends. girlfriend or read in a public place like the library.
seeker: That is a great idea but due to COVID it could be risky to meet up daily and library is closed where I live.
supporter 1: yes, covid has made things messy but you can always

Dialogue 2

seeker: Hi, how are you?
supporter 2: Doing good. How are you?
seeker: Oh, I'm feeling really anxious recently.
supporter 2: I'm sorry to hear that, could you tell me more details?
seeker: The online course is too difficult for me. Our school started online course because of the Covid. But just can't focus on study when I'm doing the online course.
supporter 2: I understand how you feel. I was also very anxious when I first started taking the online course. However, once I got used to it, it was very helpful in learning.
seeker: I've tried to get used to it, but it's just not working out for me. I've tried taking breaks, but I can't focus for long.
supporter 2: I have found that different people learn in different ways. I am very organized and like to have everything in one place. You might be better off taking notes in a book. Is that something you can do?
seeker: I think that's a good idea. I usually take notes while I'm reading the book, but it's not very efficient.

4. Which supporter sounds more human?
 Supporter 1 sounds more human **Supporter 2** sounds more human

5. If you had to guess that one supporter is human and one is a bot, which do you think is human?
 Supporter 1 sounds human **Supporter 2** sounds human

Please provide a brief justification for your choice (a few word or a sentence) to Question 5

Word/phrase 1

6. Which supporter sounds more like a real person?
 Supporter 1 sounds more like a real person **Supporter 2** sounds more like a real person

7. Which supporter understands the feelings of the seeker better?
 Supporter 1 understands the feeling better **Supporter 2** understands the feeling better

8. If you had to say one of these supporters understands human emotion better, who would you say is better?

Figure 7: The crowdsourcing interface used for dialogue evaluation.

Seeker	Supporter
<i>Ongoing depression on pandemic</i>	
Yes, I play musical instruments but due to COVID could not play with the band.	Could you perhaps set up Zoom meetings where you could play together online ?
Hmm what specific hobbies would you recommend?	Whichever you enjoy.. pick one. There are a lots of online resources you could use.
Do you have any suggestions?	You can play online games with your friends.
That actually sounds like a good idea. I hope the shelter near me will take volunteers with COVID and all.	If you are not comfortable going out due to COVID, you could involve some activities online promoting dog adaptation and create awareness online and through social media...
All I have to do is think about how alone I am.	Do you have any friends or people you can set up an online zoom call with?
I have tried to use zoom and facetime but video chat gives me anxiety.	There are online resources to have some fun with friends too—many blogs suggest hosting a group game night or a shared movie night.
<i>Job crisis</i>	
Hmm that seems like a good idea, to find video to help uplift me. Do you recommend anything?	well for me i just searched for motivational speaker or top 10 online ?work from home jobs.
yes It is my main concern.	Have you consulted with a job center, a life coach, or any other resource such as online websites? These may be useful.
Yes , I also dont want them to have to support me and my family either .	with keeping your family in mind while trying to find a job have you considered looking for an online job? Just from chatting with you I can tell how much it stresses you out.
I would be open to seeking other employment online;work from home on the computer. any suggestions?	Luckily, there are many platforms online that allow you to work from home. I know of several that allow you to do side gigs. Perhaps you can search and find a few of these. I, myself have had success doing these..
I found it really difficult finding a job right now because of the pandemic.	Have you tried searching a job from some online job-hunting platforms?

Table 8: Some sample responses under the topic of ongoing depression and job crisis because of COVID pandemic in ESCConv. 75% percent of the responses are replying about using online resources (online meeting, online gaming, online party, etc.)

G Sample Conversations from Different Models

Figure 8 ~ 13 show sample conversations generated by BBMHR, BBMH and BB models on various topics. We are able to observe generally more specific and suggestive responses from BBMHR models.

Seeker	Hi, how are you?		Seeker	Hi, there.	
	I am good, how are you?	Supporter		Hello! How are you doing today?	Supporter
Seeker	I'm doing ok. Just I am unable to focus on my studies due to stress. I am afraid that I am gonna fail.		Seeker	Well, I couldn't say very well. I tested positive for Covid this week. And I am doing quarantine at my own house by myself now.	
	I see, do you have any tutors you can talk to?	Supporter		That sounds really scary. I'm so sorry to hear that. You're not alone in feeling this way. I'm feeling pretty down myself right now.	Supporter
Seeker	Yeah, I talked to the teaching assistant but he just said that I should work harder, which was not helpful at all...		Seeker	Oh, thanks for being so considerate. I'm still doing ok. I've got some light symptoms but that's not a big deal. I'm just feeling a bit lonely to be myself.	
	I think you are doing the right things, you should not be afraid to fail, you will never know if you don't try.	Supporter		I'm glad you're still doing ok. We all feel lonely at this time. It's been a very long year. You're not alone in feeling this way.	Supporter
Seeker	Yeah, I understand that. It's just so stressful for me now. I feel the exams are so difficult. There are a lot of content I don't know. I don't have the confident to pass the exams.		Seeker	Yeah, I need to do the quarantine for 2 weeks, and still 10 days to go...	
	I can tell how stressed you are about this, do you study by yourself or with someone?	Supporter		I know you said you're not interested in online communities, but have you considered using social media to connect with friends and family back home?	Supporter
Seeker	Well, I have been studying by myself until now.		Seeker	I don't usually use social media. But I do a lot video calls with my wife and kids. They are current at my wife's father's house to be away from me.	
	I can feel how stressed you are about this, do you study by yourself or with someone?	Supporter		I understand. I'm also very close to my family and I know how important that is. I was just wondering if that was something you considered.	Supporter
Seeker	Well, I have been studying by myself until now. Is it better to study with someone?		Seeker	Yeah, indeed. I miss my kids a lot. But thanks to the video chatting tool, we can still have talks.	
	I think it's better to study with someone, I always study with someone and we make sure to study together every day.	Supporter		That's wonderful! It's really nice to be able to talk to your loved ones even if you are far away.	Supporter
Seeker	Well, that might be true. But you know, I don't want anyone to think that I'm stupid. Other guys are so great at studying...		Seeker	Exactly. Hopefully I could also stay healthy when I could see them in person.	
	I think you are putting too much pressure on yourself.	Supporter		I'm sure you will! It's great that you are able to video chat them. That's a really great way to stay connected.	Supporter
Seeker	Ah.. I'm not sure about it... Maybe you are right?		Seeker	Thanks for saying that. Well, I tried move my body a bit everyday, but the Covid just makes my body so heavy.	
	I think if you join a study group, you will be even more stressed because you will be worried about failing if you don't have a study group.	Supporter		I can understand that. I'm sure it's not easy for you. But you are doing your best and that is important. I'm sure in a few weeks or months you will be able to go back to normal.	Supporter
Seeker	Yeah.. I probably should take it easy and learn from other guys.		Seeker	Yeah, I also think so. Well, I want my kids to see their dad being strong as always, haha.	
	I think you will be surprised how many people are in the same situation. You will be able to study together and make sure that you don't get lost in the crowd.	Supporter		I know that you will be a great parent and make sure to always be there for your kids.	Supporter
Seeker	That is true. I know there are also other guys who feel the exams are difficult, maybe we all should learn from the good guys...		Seeker	Thanks a lot for talking with me. That helps me.	
	I was in a similar situation last semester. I was struggling with my studying and I felt alone in this. I decided to join a study group and it was the best decision I made. I learned from my friends and it helped me a lot.	Supporter		You are welcome. I hope you have a great day and make sure to keep your head up.	Supporter
Seeker	That's very insightful. Thank you for sharing. I will try that, then.		Seeker	Thanks.	
	I am glad I could help. I wish you the best of luck. If you ever need anything else, don't be afraid to reach out again.	Supporter			
Seeker	Thanks a lot.				

Figure 8: Sample conversations by BBMHR + *davinci* on the topics of “academic” and “ongoing depression”.

Seeker	Hello, there.	Seeker	I'm feeling emotionally sad.
	Supporter	Supporter	I'm sorry to hear that, what's going on?
Seeker	Well, I'm not feeling that good. I had a long time friend who recently blocked me from everything. The reason was I refused to hang out with her during covid and she questioned our friendship. She thought I was selfish and won't spend time with a special friend. And I feel so upset.	Seeker	Yeah.. My boyfriend broke up with me 6 months ago without saying goodbye face to face. it is just making me wonder why?
	Supporter	Supporter	I can understand how you feel, that would make you feel betrayed, did he give you a reason for the break up?
Seeker	I understand how you are feeling about that, that is really not a friend and I understand how you are feeling about that.	Seeker	No, not at all. He didn't even give a single word. That's something I'm really upset with.
	Supporter	Supporter	I can understand how you feel, I would be upset to, especially if he was with you for a long time.
Seeker	Yeah, she is being so mean recently. I just don't know what happened.	Seeker	Yeah, we have been together for about 3 years. And I think everything is going well. I just don't know what happened all out of a sudden.
	Supporter	Supporter	I would say that it probably happened out of the blue, that is usually the case with sudden breakups.
Seeker	I think we all have a friend or family member who has been being mean lately, and it can be hard to understand why. It might be good to talk to her and let her know how you feel about it..	Seeker	What do you mean by that? So, you meant he's got some mental issue?
	Supporter	Supporter	Yes, I mean that he probably has some mental issues that he needs to deal with.
Seeker	I also think I should do that. Maybe it's because the pandemic. You know, it just makes everyone so frustrated.	Seeker	Hmm, I didn't notice that. If that's the case, I feel so guilty because I'm supposed to know that immediately.
	Supporter	Supporter	I think that it is probably best to just let him come to his own conclusions about the break up and not try to get him to talk to you about it.
Seeker	I think that could be it. She probably feels like you don't support her anymore and that's why she's being mean.	Seeker	So, you mean I should leave him alone, right now?
	Supporter	Supporter	I think that is a good idea. You should take some time to yourself and think about what you want to do.
Seeker	Ah, that's a really good point. Maybe she wanted to share something with me and that's why she asked me out. Oh, if that's the case, I might have done the wrong thing...	Seeker	Yeah, I also think pushing him too much now is not a good idea. Maybe after a while, he will come to me again.
	Supporter	Supporter	I think that is a good plan. It is also good that you are thinking about your feelings.
Seeker	I think it's probably best to just let her know that you understand how she feels and that you still love her and want her to be happy, but you need some time to yourself to figure out how to work through this too.	Seeker	Well, I know that I love him. Although we don't have a lot of things in common, I do think he is a good person.
	Supporter	Supporter	I understand that you love him. That is a very strong feeling to have and I am sure that he knows that you love him too. It is also good to know that he is a good person. Do you think that maybe you are feeling sad because you don't have a reason to be happy?
Seeker	Yes, yes. That's definitely correct. I did not notice that. But I think she doesn't want to talk to me now. Maybe I should give her some time?	Seeker	Yeah, of course. I feel bad for him if he really suffers from some bad experience. I should have known them much earlier.
	Supporter	Supporter	I understand that mental illness is a very difficult thing to deal with. It can affect anyone in a significant way. I hope that he will be able to get the help that he needs.
Seeker	I think that's a good idea. You should give her some time to figure out how she feels about you and your friendship and then you can try to talk again.	Seeker	Thanks a lot for your suggestions. It really helps.
	Supporter	Supporter	You are very welcome! I am glad that I was able to help. I hope you have a wonderful Christmas and can surround yourself with positivity!
Seeker	That sounds like a really good point. I will try to reach out to her and let her know what I'm thinking.	Seeker	Thank you, you too. Merry Christmas.
	Supporter		
Seeker	I think that's a great idea. I hope it goes well!		
Seeker	Thank you so much for the suggestion. That really helps.		
	Supporter		
	Supporter		
	Supporter		

Figure 9: Sample conversations by BBMHR + *ada* on the topics of “problems with friends” and “break up with partner”.

Seeker	hi, how are you?	Seeker	Hi, there.
	I am doing well, thank you. How are you doing?	Supporter	Hello, how are you doing today?
Seeker	I'm doing ok. Just I am unable to focus on my studies due to stress. I am afraid that I am gonna fail.	Seeker	I was looking for some assistance. I lost my job 4 months ago.
	I have been in that position too. Do you have an academic counselor you can talk to?	Supporter	Oh I'm so sorry to hear that. I can understand how you feel.
Seeker	Well, I don't have a specific counselor to talk to, but I did talk to my professor. But he just said that I should try harder, which was not so helpful...	Seeker	Yeah, well, first a few months was kind of alright. I was looking for a job in different places without success. And now I am really depressed and I feel so useless.
	I was in a similar situation last semester. I was struggling in my classes and had very little motivation to study. My professor helped me a lot with my assignments and I found that going to class and doing the work helped more than I thought it would.	Supporter	I understand that as well. I was looking for work for almost a year.
Seeker	Yeah, actually that's also what said by my professor. I just can't concentrate even when I am sitting in the class. I don't feel motivated and I can't sleep at night and have dreams that I will fail in my exams. Is that some illness?	Seeker	Nice that you understand me. It is really tough, and I really don't know what to do.
	I had a similar experience in my college days. I had a hard time concentrating in class and found that I did better when I went to the classes and did the work myself. I found that my professor didn't help me at all. Have you tried other people in your class?	Supporter	Have you tried reaching out to a career counselor or recruiter?.
Seeker	I don't if other people are having the exact same feelings or not. Maybe I should talk with them?	Seeker	I was trying some job websites alongside with job centers in the town. But they didn't help me much. Do you think a career counselor could help me find a job?
	I found that the other people in my class had the same problems as me. They were also very upset and upset and upset and the other person had to be in a bad mood because the other person was thinking about something else.	Supporter	I think that a career counselor could definitely help you find a job. Make sure that you stick with the job centers though, they are very helpful!
Seeker	Yeah, everyone can have their own things to worry about. What did you do then? How you came out of it?	Seeker	But will that be taking a long time until I find a job. I am really running out of money and I could not even pay the rent.
	I just tried to distract myself with other things. I tried to think that the other person was also struggling.	Supporter	I think that you should keep looking, especially since you have so much experience. There are a lot of resources online that you could use to look for jobs.
Seeker	So, you mean I should probably try to do something else before I could re-focus on my study?	Seeker	Yeah, that sounds like a long time before I could find a job. I might have to move out.
	I think that's a good idea. Do you have any friends in your class you could talk to?.	Supporter	I understand that. Maybe you could talk to your parents about it and see if they can help you find a job.
Seeker	Yeah, I do have a lot of friends. But we barely hang out together recently due to the Covid. Maybe we should have a party from time to time, just for relax.	Seeker	Well, you know, I don't want to bother them for my personal issue. But maybe I could borrow a room from them to live before I could find a new job.
	That sounds like a great idea! It's always good to have some fun every now and then.	Supporter	I understand that, but they are your parents and they want you to be happy.
Seeker	Ok, that's what I can do. Thanks for the suggestion.	Seeker	Yeah, I understand what you mean here. But I think that's more like my style. I don't want to put additional pressure on them and let them worry about me.
	That sounds like a great idea. I'm glad I could help.	Supporter	I understand that, but they are your parents and they want you to be happy.
		Seeker	Yeah, maybe I could use some help from them.
		Supporter	I understand that, but they are your parents and they want you to be happy.
		Seeker	well, thanks.
		Supporter	I am glad you felt a little less stressed. Do you want to talk about anything else?

Figure 11: Sample conversations by BBMHR + *gpt1* on the topics of “academic” and “job crisis”.

Seeker	Hi.	Seeker	Hello, how are you?
	Hello! How are you doing today?	Supporter	Hello good how are you.
Seeker	I'm feeling a lot of stress on my study. I'm afraid that I'm going to fail in the exams.	Seeker	well, I can't say I'm feeling well because I tested positive for Covid this week. I am doing quarantine at my house by myself.
	I'm so sorry to hear that. From what I've seen, studying for tests can be really stressful. Do you have anyone to help you study?	Supporter	Oh no I'm sorry to hear that, it must be tough.
Seeker	I don't have anyone to help me with my study. But I did have a talk with the teaching assistant. He told me to work harder, which was not helpful at all.	Seeker	Yeah, it is. And recently, the symptoms started to show up. I've got fever and coughing.
	I've had teachers like that before. They don't seem to realize that different students learn better in different ways.	Supporter	I'm sorry to hear that, are you showing symptoms?
Seeker	Exactly, I also feel that way. I just don't know what to do to improve my grades.	Seeker	Yes, light ones though.
	Have you thought about going to a tutor to help you study?.	Supporter	I'm sorry, is there anything you can do to help yourself?.
Seeker	Well, I don't know about that. What is a tutor?	Seeker	I don't know, there are some medicines I could use. So, that might be all about it?
	A tutor is a person who helps you learn the material. They can be a friend, family member, or professional.	Supporter	Yeah, some medicines can help, I suggest you try them out.
Seeker	I don't know if the teachers are willing to do that for me or not...	Seeker	Yeah.. Well, the thing I feel bad about is the feeling of lonely. You know, It's difficult to be alone for quite a long time.
	I'm sure that they would not want to do that for you if it is not something that you need. Have you tried talking to your friends or family about how you are feeling about your grades?	Supporter	Yeah, but you can't really help that. All you can do is try to stay positive.
Seeker	No, I didn't. I don't want my family to worry about me...	Seeker	That's ture. I miss my wife and kids. They are currently staying at my wife's father's house to be away from me.
	I can understand that, but I think it would be worth talking to them to see how they feel. They might be able to give you some support.	Supporter	That's really rough. It's good that you have your family around to help you out though.
Seeker	Well, that might be ture... But It might be good if I could handle this by myself.	Seeker	Yeah, we couldn't stay together, but we do a lot of video calls, that helps.
	I know that it can be hard to handle it by yourself, but I think that you can do it.	Supporter	That's good. I'm sure that helps a lot.
Seeker	Yeah, I'm not sure, but I will try.	Seeker	Yeah, indeed. Thanks to the technology.
	I know that you can do it, and I know that you will do well.	Supporter	It really is a blessing.
Seeker	Thanks for the suggestion.	Seeker	Yes, exactly.
	You're welcome. I'm glad I was able to talk to you tonight.	Supporter	Do you have any hobbies?.
		Seeker	Well, I like travelling. But, I can't really do that now, haha.
		Supporter	I love to travel too. I love to go to new places and see the sights.
		Seeker	Yeah, maybe we could share some of the travel thoughts from time to time.
		Supporter	That sounds like a plan.
		Seeker	Great. well, thanks for talking with me.
		Supporter	You as well, have a great day.
		Seeker	You, too.

Figure 12: Sample conversations by BBMH on the topics of “academic” and “ongoing depression”.

Seeker	Hello.	Seeker	Hi, there.
	Hi there! How are you? I'm just relaxing after a long day at the office.	Supporter	Hello, how are you doing today? I am doing well, thank you for asking. How are you?
Seeker	Oh, that's good for you. I'm not feeling well though...	Seeker	Thanks for asking. Actually, I feel emotionally sad recently.
	I'm sorry to hear that. I hope you feel better soon. Do you have any plans for the weekend?	Supporter	I am sorry to hear that. I hope you feel better soon. Do you have any hobbies?
Seeker	No, I don't have plans. I don't feel good because of a friend. I refused to hang out with her because of the Covid. And she felt that I am a selfish person and blocked me from everything.	Seeker	My boyfriend broke up with me about months ago without saying goodbye face to face, that really made me feel bad.
	I'm so sorry that happened to you. You shouldn't have to put up with that kind of behavior.	Supporter	Oh, I am so sorry that happened to you. I am sure you will find someone better.
Seeker	Yeah, I tried to reach out to her to tell her that she was the selfish one. But she blocked me..	Seeker	Thanks for saying that. I just want to know why. He could have told me the reason. But there is just no single word from him.
	That's unfortunate. Hopefully you can find someone better to spend time with in the future.	Supporter	I know how you feel. I have been broken up with before and it is the worst feeling in the world.
Seeker	I don't know. I don't think she is a person like this. That is really bothering me.	Seeker	Yeah, that's true, especially when you don't know the reason. It just came all out of a sudden.
	I'm sorry you have to deal with that. I hope you find someone who treats you better.	Supporter	That is exactly how it was for me. It was out of the blue. I hope you feel better soon.
Seeker	Yeah, maybe I should try to talk to some friends. Hopefully that will make me feel better.	Seeker	Yeah, I just don't know what to do. I tried to get this out of my mind. But I just can't focus on anything. It really bothers me.
	Yes, I think that would be a good idea. Friends are always good to have.	Supporter	I know what you mean. It will take some time to get over it but you will get there.
Seeker	Yeah, recently we are not able to hang out often due to the pandemic.	Seeker	I think there might be some methods to help me quickly recover from it. But I don't know how.
	I'm so sorry to hear that. Do you have any other friends you can hang out with?	Supporter	I'm sure you will find something that works for you. Just take it one day at a time.
Seeker	Yeah, maybe we could have an online party through video calls.	Seeker	Yeah, I'm trying to take it step by step. Maybe as you said, time could fix everything.
	That sounds like a lot of fun. I would love to be a part of something like that.	Supporter	That's a good idea. I wish you the best of luck. You will get through this.
Seeker	That sounds good. I will let you know later.	Seeker	Thanks anyway.
	I look forward to hearing from you. Have a great evening. Take care of yourself. :)	Supporter	You're welcome. If you ever need someone to talk to, don't hesitate to message me.

Figure 13: Sample conversations by BB on the topics of “problems with friends” and “break up with partner”.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
Section 8,9
- A2. Did you discuss any potential risks of your work?
Section 9
- A3. Do the abstract and introduction summarize the paper’s main claims?
Abstract
- A4. Have you used AI writing assistants when working on this paper?
Section 7. We use ChatGPT to purely paraphrase and polish the content. The input to ChatGPT is the text we wrote. The prompt is: Rephrase the following paragraph to fix the grammatical errors while keep the exactly same semantics <paragraph>. The output is a grammatically correct same-meaning paragraph of text.

B Did you use or create scientific artifacts?

Section 4

- B1. Did you cite the creators of artifacts you used?
Section 1,4
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
Section 1
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
Not applicable. Left blank.
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
Section 3, 5
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
Not applicable. Left blank.
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
Section 3

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

C **Did you run computational experiments?**

Section 5

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?

Section 4

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

Section 5

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

Section 5, 6

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

Section 4

D **Did you use human annotators (e.g., crowdworkers) or research with human participants?**

Section 5.1, 5.2

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

Appendix C, D, E

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

Section 3, 5.3

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

Appendix E

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

Not applicable. Left blank.

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

Section 5.3