

DiaCBT: A Long-Periodic Dialogue Corpus Guided by Cognitive Conceptualization Diagram for CBT-based Psychological Counseling

Anonymous ACL submission

Abstract

Psychotherapy reaches only a small fraction of individuals suffering from mental disorders due to social stigma and the limited availability of therapists. Large language models (LLMs), when equipped with professional psychotherapeutic skills, offer a promising solution to expand access to mental health services. However, the lack of psychological conversation datasets presents significant challenges in developing effective psychotherapy-guided conversational agents. In this paper, we construct a long-periodic dialogue corpus for counseling based on cognitive behavioral therapy (CBT). Our curated dataset includes multiple sessions for each counseling and incorporates cognitive conceptualization diagrams (CCDs) to guide client simulation across diverse scenarios. To evaluate the utility of our dataset, we train an in-depth counseling model and present a comprehensive evaluation framework to benchmark it against established psychological criteria for CBT-based counseling. Results demonstrate that DiaCBT effectively enhances LLMs' ability to emulate psychologists with CBT expertise, underscoring its potential for training more professional counseling agents.

1 Introduction

Mental health disorders are a significant and widespread public health concern, affecting an estimated one in eight people globally (World Health Organization, 2022¹). Despite this, access to mental health services remains limited, largely due to social stigma and a shortage of therapists (Freeman, 2022; White and Dorman, 2001). To address these barriers, there is growing interest in using automated conversational agents as alternative tools for mental health support (Ali et al., 2020; Sabour et al., 2023; Sharma et al., 2023). Whereas, gener-

¹<https://www.who.int/campaigns/world-mental-health-day/2022>

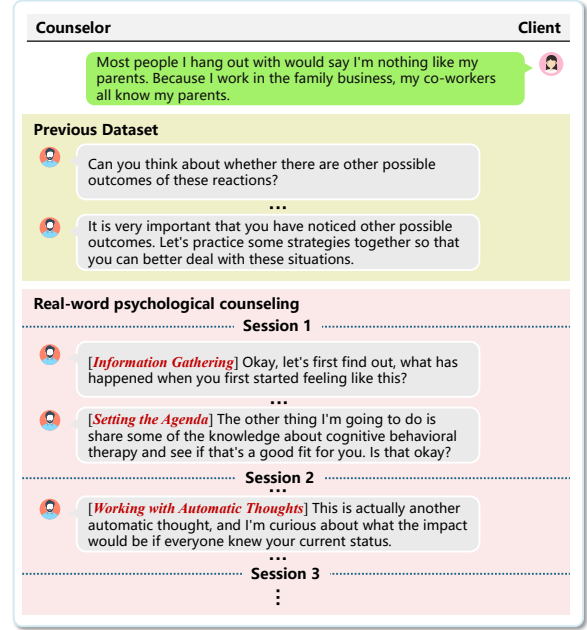


Figure 1: A dialogue comparing a previous counseling dataset to a real-world scenario. Previous datasets tend to solve all psychological problems in a single counseling session using one CBT technique, whereas real-world counseling involves multiple sessions, depending on the nuanced progression of therapeutic strategies.

ating human-like responses as the psychologist in automated systems remains a major challenge.

Recently, the integration of large language models (LLMs) into counseling systems presents a promising way (Stade et al., 2024; Bubeck et al., 2023). While LLMs demonstrate remarkable capabilities in generating human-like responses, they tend to provide generic advice or information (Chiu et al., 2024). Existing solutions attempt to model the psychological theory as rules or tree flows to prompt LLMs to mimic psychological counselors for professional response generation (Mousavi et al., 2021; Das et al., 2022). Whereas, the expertise still remains constrained due to the insufficient training of the backbone models in psychological

counseling (Raile, 2024). Moreover, psychological counseling heavily relies on case practice, while the counseling process data is severely limited by the sensitive nature of therapist-client interactions (Harrigian et al., 2021; Pérez-Rosas et al., 2018).

An alternative solution is to construct synthetic dialogues that can resemble real-world counseling processes in diverse scenarios. Previous efforts primarily focus on generating dialogues between therapists and clients based on certain mental problems (Qiu et al., 2024a; Na, 2024). Whereas, there is a large gap between the synthetic data and the real-world counseling process. As shown in Figure 1, the existing synthetic data assumes that all psychological problems can be solved in one counseling session, thus one dialogue mimics the whole counseling process for one case, which is coarse-grained and can not well reflect the difference of counseling strategy use in each session. In practice, real-world psychological counseling often involves long-periodic dialogues with multiple sessions, depending on the difficulty and nuanced progression of therapeutic strategies (Curwen et al., 2018). Moreover, most client simulation relies on background descriptions such as occupations, experiences, and problems (Lee et al., 2024), which can not well model the mental portrait and limits the utterance generation with specific mental problems for each client (Wang et al., 2024).

To address the above challenges, we present **DiaCBT**, a long-periodic dialogue corpus guided by cognitive conceptualization diagram (CCDs) for psychotherapy counseling based on cognitive behavioral therapy (CBT) (Beck, 2020). DiaCBT simulates the entire process of CBT, consisting of multiple sessions that feature interactive transcripts annotated with CBT strategies. To better simulate clients, we construct structured CCDs across diverse scenarios to guide the utterance generation of clients with specific mental problems. We conduct experiments by fine-tuning models for therapy-based response generation using DiaCBT, and present a comprehensive evaluation framework with established psychological criteria for CBT-based counseling. Both automatic and human evaluations demonstrate that DiaCBT significantly enhances LLMs’ ability to emulate psychologists with CBT expertise, highlighting its potential for training professional counseling agents. The main contributions of our work are as follows:

- We create a long-periodic CBT-based counsel-

ing dialogue dataset, which well mimics the whole process of cognitive behavioral therapy in practice and models clients with structured cognitive conceptualization diagrams for better simulation.

- We train a psychological counseling model based on DiaCBT, which enhances CBT-specific skills by integrating strategies for in-depth questioning to effectively reframe clients’ cognitive distortions.
- We present a comprehensive evaluation framework for psychological counseling and conduct extensive experiments, demonstrating the great advantages of our dataset and counseling model over the baselines.

2 Related Work

In this section, we review studies about conversational psychotherapy AI, covering both therapy-based conversational systems and datasets.

2.1 Therapy-based Conversational Systems

Conversational systems for psychotherapy aim to assess an individual’s mental state and enhance self-awareness through effective communication techniques. These systems often rely on empathetic conversations to alleviate emotional distress, crafting responses that are relevant to the client’s statements (Ma et al., 2020; Zhou et al., 2018; Lubis et al., 2018; Raamkumar and Yang, 2022; Li et al., 2022; Gao et al., 2021; Shen et al., 2021; Liu et al., 2021; Cheng et al., 2022). The emergence of large language models (LLMs) is revolutionizing therapy-based conversational systems (Liu et al., 2023; Chen et al., 2023b; Jo et al., 2023; Wei et al., 2024). Researchers are increasingly integrating professional counseling strategies into these models to enhance their effectiveness (Hsu et al., 2023; Chen et al., 2023a; Na, 2024), and primarily center around emotion analysis and symptom identification. However, conversational agents from a cognitive-behavioral perspective remain underexplored. Therefore, our goal is to improve the cognitive behavioral therapy skills of conversational agents in therapeutic settings.

2.2 Psychotherapy Dialogue Datasets

Pérez-Rosas et al. (2018) introduce a dataset of high- and low-quality counseling conversations that include counseling skills, collected from publicly

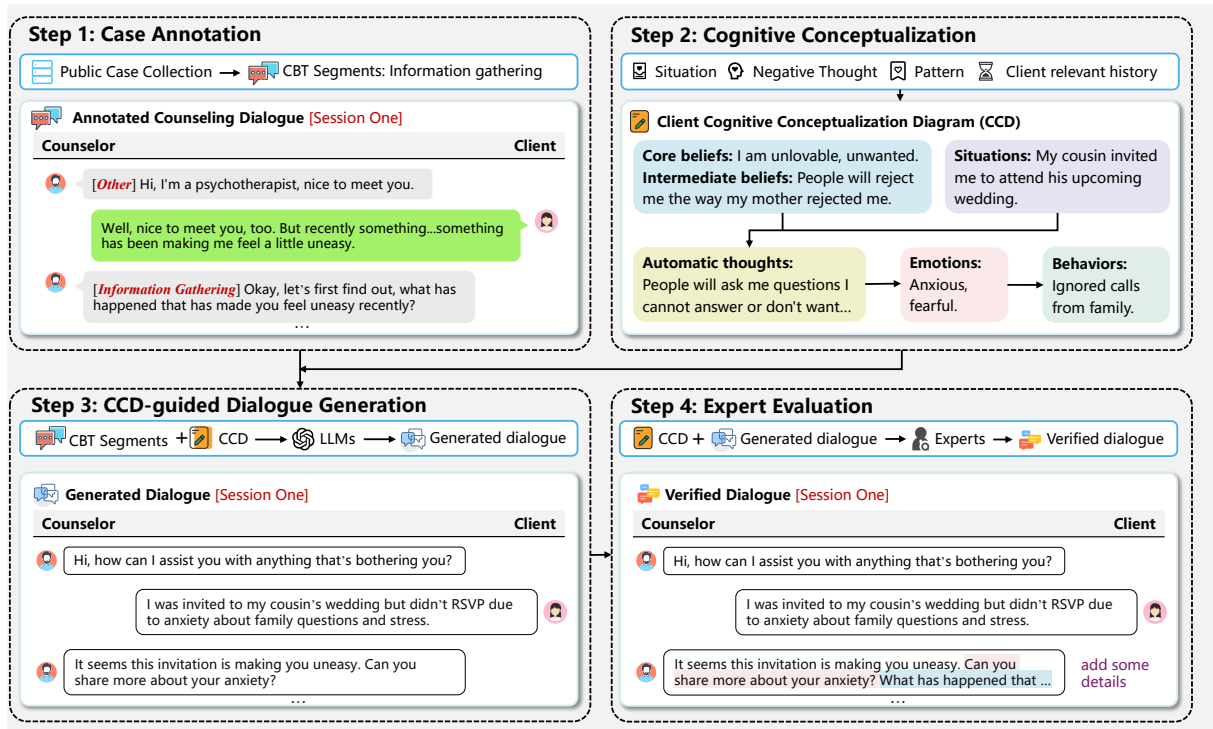


Figure 2: The data collection process involves gathering public cases, cognitive conceptualization, generating dialogues, and conducting a final expert evaluation.

available sources. Several studies have also leveraged social media platforms for mental health research (Sharma et al., 2020). For instance, Rashkin et al. (2019) and Mousavi et al. (2021) develop corpora of empathetic and therapeutic dialogues from real-life interactions, while Liu et al. (2021) create an emotional support conversation dataset with annotations based on Helping Skills Theory. Additionally, Yao et al. (2022) design a 3-phase procedure simulating doctor-patient dialogues for depression diagnosis. Wu et al. (2022) introduced AnnoMI, the first publicly accessible dataset of professionally transcribed and expert-annotated therapy dialogues. Related works have also explored using LLMs to synthesize specific dialogue data, offering a promising alternative to real data (Na, 2024; Zhang et al., 2024; Lee et al., 2024; Xiao et al., 2024). However, these studies predominantly focus on emotional disorders or depression, limiting their ability to support the broader range of strategies. Additionally, these conversations tend to be short, lacking the depth and multi-stage complexity of face-to-face clinical psychotherapy.

3 Dataset Construction

Previous studies have sourced psychological counseling dialogues from online welfare counseling

platforms or generated data through crowdsourcing simulations (Li et al., 2023; Yao et al., 2022). However, due to the sensitive nature of clinical records, access to real clinical data is restricted. To address this, we design a specialized task to generate dialogues. Unlike prior approaches that rely on large-scale recruitment of online volunteers, our annotation process combines the efforts of large language models with contributions from trained volunteers and domain experts. As illustrated in Figure 2, our methodology comprises four key steps: case annotation, cognitive conceptualization, CCD-guided dialogue generation, and expert evaluation, each of which will be further detailed.

3.1 Case Annotation

Public case source. Our case collection consists of transcripts of cognitive behavioral therapy (CBT) sessions sourced from the American Psychological Association (APA) website² and related books³. We selected these cases for two key reasons: (1) their strong relevance to the clinical setting and (2) their widespread use in training novice CBT counselors. The overall statistics of collected cases

²<https://www.apa.org/pubs/databases/psychtherapy/>

³A CBT textbook named *Dispelling the Fog of Belief* selected by our domain expert.

Category	Total	Client	Therapist
Dialogues	53	-	-
Avg. utterances per dialogue	279.5	139.5	140.1
Avg. tokens per dialogue	6759.2	3577.5	3181.7
Avg. tokens per utterance	24.2	25.7	22.7

Table 1: Statistics of collected cases.

are shown in Table 1. The cases contain 53 transcripts, and as seen in such a psychotherapy scenario, sufficient dialogue turns are required: our dialogues exhibit avg. 279.5 turns and avg. 6759.2 tokens per dialogue. In our transcribed dialogues, the maximum, median, and minimum utterance numbers are 717, 293, and 101, respectively; the maximum, median, and minimum number of words are 12097, 6730, and 4109, respectively. Nearly the entire samples are comprised of more than 100 utterances, which significantly surpasses the input limit of many language backbones.

Data annotation target. CBT, a leading psychological treatment, improves mental health through evidence-based communication. In our work, we incorporate 14 counseling strategies, the detailed definitions of the strategy are in Appendix A.2. For annotation, annotators⁴ are provided with conversation transcripts and tasked with identifying the CBT strategies employed by the therapist according to the guidelines provided. While strategies may overlap, we simplify the annotation process by treating it as a single-label task. Annotators focus on annotating the strategy used in the current segment and, in cases of overlap, are instructed to select the primary strategy label.

3.2 Cognitive Conceptualization

In addition to the public case library, many psychology-trained individuals have contributed to public datasets through crowdsourcing, particularly for analyzing cognitive distortions and cognitive reframing (Wang et al., 2023; Maddela et al., 2023). While these datasets are not exclusively conversational, their content closely mirrors the thoughts of real clients, making them valuable for enriching client profiles. To this end, we compile various CBT-related datasets and construct cognitive models with LLMs based on information provided in these datasets to accurately simulate backgrounds that resemble those of real patients.

⁴Our annotators are graduate students in psychology with foundation knowledge in CBT. They have undergone a comprehensive training program on CBT principles and annotation guidelines, and have all passed the annotation competency test.

We utilize the Cognitive Conceptualization Diagram (CCD), recognized as a commonly used representation of a patient’s cognitive model in CBT. The CCD includes components for understanding how an individual’s thoughts and beliefs are interconnected and influence emotions and behaviors. We select six key components for formulating a patient’s cognitive model⁵. *Core Beliefs* are deeply ingrained perceptions about oneself, others, and the world. *Intermediate beliefs* are the underlying rules, attitudes, and assumptions derived from core beliefs and shape an individual’s thought patterns. An external event or context (a *situation*) may trigger quick, evaluative thoughts without deliberation (*automatic thoughts*) stemming from the beliefs, leading to responses in terms of *emotions* and *behaviors*. In this work, we integrate situation, negative thought, cognitive pattern, and relevant history from the C2D2 (Wang et al., 2023) and PatternReframe (Maddela et al., 2023) as contexts for generating CCD-based cognitive models.

3.3 CCD-guided Dialogue Generation

LLMs trained on extensive text corpora can effectively support data collection (Zhang et al., 2024; Lee et al., 2024; Xiao et al., 2024), which demonstrate knowledge of psychological therapy concepts, including CBT techniques. To this end, we enhance dialogue generation by using a cognitive model about clients and CBT sessions with a scripted framework. We use annotated segments from different CBT sessions as few-shot prompts to guide LLMs in generating conversational flow. Additionally, we prompt LLMs to role-play as clients based on a CCD and follow instructional prompts to reflect the underlying cognitive processes. This approach ensures that therapists in the generated dialogues adhere to the full CBT process, helping clients identify solutions independently rather than simply reframing their thoughts. The instructional template is provided in the Appendix B.1.

3.4 Expert Evaluation

After collecting raw data, we manually annotate a small portion and develop rules, including **Correctness**, **Reasonableness**, and **Situation Diversity**, to distinguish high-quality data from suboptimal data. Further details are provided in the Appendix B.2. These criteria ensure that the generated dialogues

⁵Details regarding the rationale for selecting these components are provided in Appendix A.3.

Strategy	Segments	Utterances
Information Gathering	303	7565
Setting the Agenda	129	2383
Weekly Review	249	5342
Defining Therapeutic Objectives	78	1599
Psychoeducation	168	3411
Working with Automatic Thoughts	438	10962
Motivational Enhancement	171	3445
Working with Intermediate and Core Beliefs	192	5398
Behavioral Techniques	195	4136
Relapse Prevention	18	420
Homework Assignments	168	3096
Requesting Feedback	135	2721
Summarization	72	1484
Other	297	5714

Table 2: Data statistics of strategy categorization. We count the number of strategy-coherent segments of each strategy, and the number of utterances.

are not only theoretically accurate but also practical and diverse enough to support varied applications in counseling contexts. The final reserved data contains 2613 dialogues from the raw 3600 dialogues, resulting in an overall retention rate of around 72.58%. This statistic proves that the psychotherapy ability exhibited by most powerful LLMs is still not satisfactory enough even with delicate prompting, thus further emphasizing the necessity of our construction of DiaCBT.

4 Data Analysis and Model Training

4.1 Data Statistics

The provided dataset statistics are displayed in Table 2, showcasing the segment number of each strategy and their utterance number. Our dataset comprises a total of 2613 segments. We observe that some strategies are frequent such as ‘Working with Automatic Thoughts’, ‘Information Gathering’, and ‘Working with Intermediate and Core Beliefs’. For efficient communication, therapists organize responses based on client content and CBT programs. As a result, strategies present different importance in the form of frequency and length. Automatic thoughts, intermediate beliefs, and core beliefs as components in the client’s CCD are also where the therapist needs to work, which is consistent with the principles of CBT.

4.2 Strategy Distribution

We also compute the distribution of strategies at different sessions of psychological counseling, as shown in Figure 3. In CBT, complete psychological counseling is divided into five sessions. However, there is also flexibility in adopting strategies at each session. For instance, in the early session, thera-

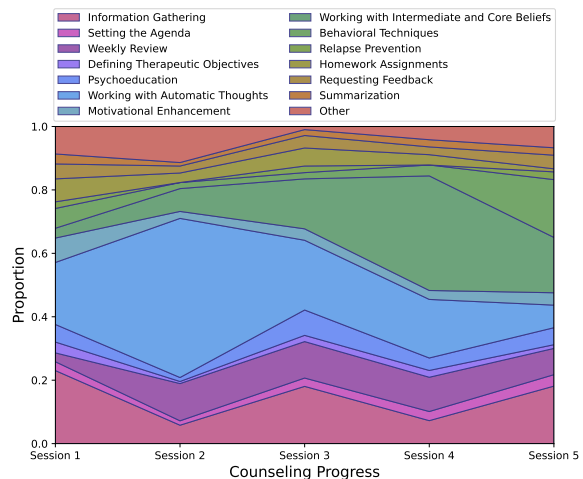


Figure 3: The distribution of strategies used in conversations at each session.

pists often use exploratory strategies like ‘Information Gathering’ and work on the client’s automatic thoughts using ‘Working with Automatic Thoughts’. After understanding the clients’ situations, therapists tend to work on the client’s intermediate and core beliefs in later sessions, with psychotherapeutic strategies like ‘Working with Intermediate and Core Beliefs’ being used more frequently. This is consistent with the principles of CBT (Beck, 2020).

4.3 Comparison with Related Datasets

As shown in Table 3, our DiaCBT stands out due to its focus on CBT, supporting multi-turn dialogues and full counseling sessions. First, it is the only dataset in the comparison that is CBT-based, includes strategy annotations, and offers a complete counseling process with detailed dialogue turns. Second, DiaCBT incorporates a wide range of psychological strategies and is based on the Cognitive Conceptualization Diagram (CCD). These features make DiaCBT a comprehensive and unique resource for studying CBT in therapeutic contexts.

Previous CBT dialogue datasets have primarily focused on single-turn strategies, with recent efforts extending to multi-turn dialogues. However, these often involve only 2-3 turns or adapt counseling reports without incorporating specific strategies. In contrast, our approach integrates CCD with annotated CBT segments to generate high-quality dialogues, while maximizing the availability of open datasets for researchers. Additionally, DiaCBT includes a wider range of CBT strategies across various stages, going beyond Cactus’s focus on automatic thought work to cover more techniques.

Datasets	CBT-based	Multi-turn	Full-session	Strategy	CCD-guided
AnnoMI (Wu et al., 2022)	-	✓	✓	✓	-
SMILECHAT (Qiu et al., 2024a)	-	✓	-	-	-
CBT-LLM (Na, 2024)	✓	-	-	-	-
CPsyCoun (Zhang et al., 2024)	✓	✓	-	-	-
Healme (Xiao et al., 2024)	✓	✓	-	✓	-
Cactus (Lee et al., 2024)	✓	✓	-	✓	-
DiaCBT(Ours)	✓	✓	✓	✓	✓

Table 3: Comparison of the key features of our work with related datasets.

4.4 Model Training

We fine-tune an LLM on DiaCBT, enabling it to: (1) select appropriate CBT strategies based on conversation history, (2) generate responses tailored to the client’s input, and (3) generate strategies and responses with given conversation history. The fine-tuning process employs a structured template to serialize T -turn dialogues, concatenating task instructions with corresponding conversational context. The training objective is formulated as the minimization of the negative log-likelihood for predicting the therapist’s response u_t^d and strategy s_t :

$$\min_{\theta} \sum_{t=1}^T -\log p_{\theta} \left(s_t, u_t^d | I_{task}, h_{t-1}, u_t^c \right) \quad (1)$$

where I_{task} denotes the task instruction template serving as a static conditioning signal. h_{t-1} represents the dialogue history, formally defined as: $h_{t-1} = (u_0^d, u_0^c, \dots, u_{t-1}^d, u_{t-1}^c)$, where u_t^d and u_t^c indicate the utterances corresponding to the therapist and client respectively in the t ’th turn. θ represents the parameters for the LLM. The objective function explicitly conditions therapeutic strategy and response generation on both task-specific instructions and evolving dialogue context.

5 Experiments

In this work, we propose a computational evaluation framework to evaluate CBT-based psychological counseling agents both in the counselor’s abilities and in the client’s psychological changes.

5.1 Computational Evaluation Framework

To evaluate an LLM acting as a therapist, we focus on assessing its responses in supporting clients with mental health challenges. Recruiting real clients with mental health issues and asking them to interact with LLM therapists would pose significant ethical risks. Instead, we leverage client CCDs to simulate clients, facilitating a full, multi-turn conversational session between LLMs therapist and simulated clients.

Client Simulating. We use cases from the C2D2 (Wang et al., 2023) to generate new CCDs. C2D2, a Chinese cognitive distortion dataset, contains 7,500 instances categorized into seven cognitive distortion types. We randomly select 20 cases from each distortion label, totaling 140 cases. Each case includes a distortion type, an emotion type, a negative thought, and a situation description, which are used to initialize CCDs. Each CCD is fed into GPT-4o to simulate the client talking to the therapist, maintaining the same conversation style, life events, and emotions. Conversations begin with the therapist and continue until the client outputs an end token. Detailed prompts for client simulating are provided in the Appendix C.1.

Evaluation Metrics. Previous studies typically evaluate turn-level performance based on fixed reference responses. Differently, when it comes to the evaluation of proactive dialogue systems, it would be more appropriate to focus on dialogue-level performance. To this end, for automatic evaluation, we employ two key metrics: the **average turn (AT)** and the **success rate (SR)** (He et al., 2024; Deng et al., 2024). AT measures goal completion efficiency by calculating the average number of turns, while SR measures goal completion effectiveness by computing the success rate of achieving the goal within a predefined maximum number of turns. Furthermore, following Lee et al. (2024), we use the Positive and Negative Affect Scale (PANAS) (Watson et al., 1988) to assess the effectiveness of counseling from the client’s perspective, measuring changes in positive and negative emotions before and after the session. Finally, the Cognitive Therapy Rating Scale (CTRS) (Aarons et al., 2012) is used to evaluate general counseling and CBT-specific skills of LLM therapists. Detailed information is provided in Appendix C.2.

5.2 Experimental Setup

Training is done using LoRA for 3 epochs with the AdamW optimizer, a learning rate of 1e-4, and

Method	Backbone	AT ↑	SR ↑	Positive ↑	Negative ↓
CpsyCounX	InterLM2-Chat-7B	6.01	34.28%	1.179	-0.910
SoulChat	ChatGLM-6B	6.25	38.57%	1.135	-0.992
PsyChat	ChatGLM-6B	8.92	72.85%	1.445	-1.210
MeChat	ChatGLM-6B	5.62	47.14%	1.031	-0.954
CAMEL	LLAMA3-8B-Instruct	9.42	67.14%	1.305	-1.044
Standard	Qwen2.5-7B-Instruct	6.72	44.28%	1.112	-1.122
Ours	Qwen2.5-7B-Instruct	12.05	77.14%	1.675	-1.021

Table 4: The results of metrics and emotion changes for baselines. The best score of each metric is **in-bold**.

Models	General Counseling Skills			CBT-specific Skills			Sum.
	Understanding	Interpersonal Eff.	Collaboration	Guided Discovery	Focus	Strategy	
CpsyCounX	4.02	5.24	4.53	4.00	4.09	4.04	25.92
SoulChat	4.00	4.90	4.25	4.00	4.00	4.00	25.15
PsyChat	4.15	5.84	5.63	<u>4.07</u>	4.05	4.01	27.75
MeChat	4.01	4.86	4.25	4.00	4.04	4.01	25.17
CAMEL	4.00	<u>5.82</u>	4.73	4.02	4.38	4.06	27.01
Standard	4.02	5.70	4.35	4.04	4.01	<u>4.08</u>	26.20
Ours	<u>4.07</u>	5.35	<u>5.57</u>	5.03	<u>4.31</u>	4.42	28.75

Table 5: The results on general counseling and CBT-specific skills for baselines. The best score of each metric is **in-bold**, while the second best score is underlined.

a batch size of 32, running on a single NVIDIA A800 GPU. More hyper-parameters can be found in Appendix D. We compare our model with the recent advanced baselines: 1) CpsyCounX (Zhang et al., 2024) is fine-tuned on a report-based multi-turn dialogue dataset for psychological counseling; 2) SoulChat (Chen et al., 2023b) is fine-tuned on a multi-turn empathetic dialogue dataset to generate empathetic responses; 3) PsyChat (Qiu et al., 2024b) is a client-centric dialogue system that provides psychological support based on client behavior recognition and counselor strategy selection; 4) MeChat (Qiu et al., 2024a) is trained on a multi-turn mental health support dialogue dataset extended from real psychological mutual assistance QA; 5) CAMEL (Lee et al., 2024) is a CBT-based agent, which uses planning with CBT techniques for psychological counseling.

5.3 Results

Main Result. The metrics in Table 4 indicate that training with DiaCBT improves both the length of counseling sessions and the success rate of interactions. Compared to the standard Qwen2.5-7B-Instruct, our model demonstrates enhanced counseling skills, highlighting DiaCBT’s effectiveness. PsyChat, trained with counselor strategy selection tasks, outperforms other baselines, underscoring the value of therapist strategy annotations in improving LLM therapist performance. These results further validate DiaCBT as a high-quality dataset

that closely mirrors real counseling scenarios.

Evaluation from the Client’s Perspective. We use PANAS to assess efficacy through changes in client emotions. Table 4 shows that DiaCBT effectively enhances positive emotions but is less effective in reducing negative emotions. This may stem from the focus of CBT on improving mental health and emotional regulation through evidence-based communication. Specifically, the strategy ‘*Working with Automatic Thoughts*’ guides clients to explore thought patterns from various perspectives rather than directly altering them.

Evaluation from the Counseling Skills. The results are presented in Table 5. Each skill is evaluated using three criteria, with scores ranging from 0 to 6. While training with DiaCBT yields a marginal performance dip in general counseling competencies, it achieves improvements in CBT-specific skills, particularly in Guided Discovery and Strategic Implementation. Training with DiaCBT achieves the highest scores, reflecting the alignment with evidence-based CBT principles, where active client collaboration and structured cognitive-behavioral exploration are prioritized over non-directive counseling approaches.

5.4 Human Evaluation

Beyond automatic evaluation, we also explore how humans perceive the generated responses. We conduct a human evaluation on 140 generated dia-

Ours vs.	CAMEL		PsyChat	
	Win	Lose	Win	Lose
Relevance	51	89	66	73
CBT Style	76	53	92	17
Helpfulness	56	24	96	28
Overall	88	34	106	24

Table 6: Human evaluation results. Ties are not shown.

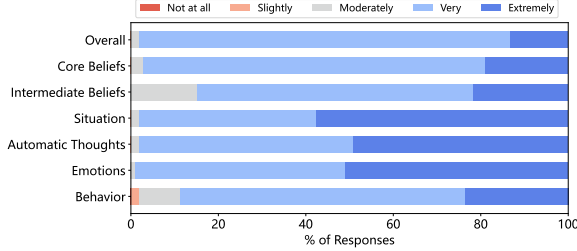


Figure 4: The distribution of ratings. Experts rate over 80% of the simulated clients as very to extremely accurate for each component of the cognitive model.

logues, where three annotators are asked to compare the generated responses from DiaCBT with two competitive baselines: (1) PsyChat and (2) CAMEL. We evaluate the responses based on three main criteria: **Relevance**, **CBT Style Measure**, and **Helpfulness**. The instructions for annotators are provided in the Appendix E. As shown in Table 6, DiaCBT outperforms other baselines in nearly all aspects of the human evaluation, as well as the overall evaluation, except for Relevance, where CAMEL achieves a higher win rate. Qualitative case study details for different models are provided in the Appendix F. We observe that CAMEL excels in providing detailed counseling strategy questions, most of which are closed questions, contributing to its strong relevance. However, for more advanced counseling, the system should go beyond merely asking yes/no questions and actively explore and help resolve the client’s negative thoughts.

5.5 Further Analysis

Analysis of LLMs as Client Simulators. Simulated clients are expected to embody their assigned roles based on specific profiles. Building on prior studies utilizing LLMs as client simulators (Wang et al., 2024), we assess the accuracy of these simulated clients in reflecting their underlying CCDs during multi-turn conversations. Three experts evaluate the simulated clients’ overall accuracy and the accuracy of each component in CCD. Figure 4 displays the distribution of ratings. The results are promising: overall, simulated clients based on CCDs are rated as very accurate. For each of the six

Metrics	Standard	PsyChat	CAMEL	Ours
avg turns	6.72	8.92	9.42	12.05
avg doc utt len	61.15	53.39	43.04	31.13
avg pat utt len	46.74	46.97	41.89	49.44
avg ques	4.46	7.18	7.66	9.04
avg in-depth ques	2.75	5.85	6.68	7.41

Table 7: Automatic evaluation results of LLMs therapist. In-depth questions encourage clients to elaborate on their thoughts, feelings, or experiences in detail.

components of CCD, the simulated clients receive average ratings ranging from very to extremely accurate. These findings validate the reliability of leveraging LLMs as client simulators, demonstrating their capability to accurately emulate diverse cognitive processes in counseling scenarios.

Analysis of Question Style. LLM therapists should guide clients to find solutions independently through reflective questioning. To assess what models learned from DiaCBT, we analyzed dialogue histories from 140 client simulation experiments and calculated several automatic metrics, with results presented in Table 7. We observe that the non-finetuned model has the fewest average turns and the highest amount of words per turn. This model often provides general advice and asks fewer questions per session, reflecting a lower level of expertise as an LLM therapist. In contrast, PsyChat, CAMEL, and ours demonstrate increased average turns and questions asked, indicating that training with counseling dialogues enhances questioning capabilities. Additionally, our model exhibits a more in-depth questioning style, better aligning with the interactive nature of clinical counseling.

6 Conclusions

In this study, we construct a CBT-based dialogue dataset (DiaCBT) across various sessions, guided by cognitive conceptualization diagrams and the CBT flow. Extensive experiments validate the great potential of our dataset for CBT-based psychological counseling. Specifically, the model fine-tuned on DiaCBT outperforms recent advanced psychotherapy conversational agents in reframing clients’ cognitive distortions. It is also notable that our method enhances CBT-specific skills by integrating strategies for in-depth questioning, helping clients find solutions independently. In the future, we plan to explore a broader range of psychotherapy techniques and extend our model to other counseling scenarios, such as emotional issues.

Limitations

In real counseling sessions, each session typically lasts around 45 minutes, with about five sessions focusing on specific topics. While our dataset involves longer multi-turn interactions compared to others, it remains significantly shorter than actual counseling sessions. Future work should aim to include longer conversations and multi-session interactions to better emulate counseling scenarios.

Ethical Considerations

Data Privacy

To preserve privacy and uphold ethical integrity, we adhered strictly to established data protocols during the case collection phase, ensuring that no cultural bias was introduced while migrating data styles. In generating dialogues, we avoided using real client data for simulated counseling scenarios. Instead, we relied on publicly available datasets curated explicitly for research purposes. These datasets were constructed through crowdsourcing information from psychological experts rather than from actual clients, thereby mitigating ethical concerns regarding personal identification and confidentiality breaches. The information provided by psychological experts was generalized, ensuring that it did not reflect any specific individual's psychological profile. This approach maintained the ethical standards required for data usage in mental health research. Additionally, in all instances requiring human expert involvement, informed consent was obtained, and appropriate remuneration was provided.

Potential Risks of the Model

Given the absence of human feedback during the model fine-tuning phase, some responses might potentially harm users. If there is no noticeable improvement after interacting with the model, and training with multi-turn consultation dialogues, we strongly recommend seeking assistance from a professional counselor or psychiatrist promptly. It is crucial to remember that a virtual dialogue agent cannot replace real-world therapy. Additionally, when implementing this model in downstream applications, it is essential to inform users that the AI generates the responses they see and that these should be used only as references.

References

- Gregory A Aarons, Elizabeth A Miller, Amy E Green, Jennifer A Perrott, and Richard Bradway. 2012. Adaptation happens: a qualitative case study of implementation of the incredible years evidence-based parent training programme in a residential substance abuse treatment programme. *Journal of Children's Services*, 7(4):233–245.
- Mohammad Rafayet Ali, Seyedeh Zahra Razavi, Raina Langevin, Abdullah Al Mamun, Benjamin Kane, Reza Rawassizadeh, Lenhart K Schubert, and Ehsan Hoque. 2020. A virtual conversational agent for teens with autism spectrum disorder: Experimental results and design lessons. In *Proceedings of the 20th ACM international conference on intelligent virtual agents*, pages 1–8.
- Judith S Beck. 2020. *Cognitive behavior therapy: Basics and beyond*.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrkke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Siyuan Chen, Mengyue Wu, Kenny Q Zhu, Kunyao Lan, Zhiling Zhang, and Lyuchun Cui. 2023a. Llm-empowered chatbots for psychiatrist and patient simulation: application and evaluation. *arXiv preprint arXiv:2305.13614*.
- Yirong Chen, Xiaofen Xing, Jingkai Lin, Huimin Zheng, Zhenyu Wang, Qi Liu, and Xiangmin Xu. 2023b. Soulchat: Improving llms' empathy, listening, and comfort abilities through fine-tuning with multi-turn empathy conversations. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1170–1183.
- Yi Cheng, Wenge Liu, Wenjie Li, Jiashuo Wang, Ruihui Zhao, Bang Liu, Xiaodan Liang, and Yefeng Zheng. 2022. Improving multi-turn emotional support dialogue generation with lookahead strategy planning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3014–3026.
- Yu Ying Chiu, Ashish Sharma, Inna Wanyin Lin, and Tim Althoff. 2024. A computational framework for behavioral assessment of llm therapists. *arXiv preprint arXiv:2401.00820*.
- Berni Curwen, Peter Ruddell, and Stephen Palmer. 2018. Brief cognitive behaviour therapy.
- Avisha Das, Salih Selek, Alia R Warner, Xu Zuo, Yan Hu, Vipina Kuttichi Keloth, Jianfu Li, W Jim Zheng, and Hua Xu. 2022. Conversational bots for psychotherapy: a study of generative transformer models using domain-specific dialogues. In *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 285–297.

668	Yang Deng, Wenxuan Zhang, Wai Lam, See-Kiong	June M Liu, Donghao Li, He Cao, Tianhe Ren, Zeyi	724
669	Ng, and Tat-Seng Chua. 2024. Plug-and-play policy	Liao, and Jiamin Wu. 2023. Chatcounselor: A large	725
670	planner for large language model powered dialogue	language models for mental health support. <i>arXiv</i>	726
671	agents.	<i>preprint arXiv:2309.15461</i> .	727
672	Melvyn Freeman. 2022. The world mental health report:	Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand	728
673	transforming mental health for all. <i>World Psychiatry</i> ,	Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie	729
674	21(3):391.	Huang. 2021. Towards emotional support dialog	730
675	Jun Gao, Yuhan Liu, Haolin Deng, Wei Wang, Yu Cao,	systems. In <i>Proceedings of the 59th Annual Meet-</i>	731
676	Jiachen Du, and Ruifeng Xu. 2021. Improving em-	<i>ing of the Association for Computational Linguistics</i>	732
677	pathetic response generation by recognizing emotion	<i>and the 11th International Joint Conference on Natu-</i>	733
678	cause in conversations. In <i>Findings of the association</i>	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	734
679	<i>for computational linguistics: EMNLP 2021</i> , pages	pages 3469–3483.	735
680	807–819.	Nurul Lubis, Sakriani Sakti, Koichiro Yoshino, and	736
681	Keith Harrigan, Carlos Aguirre, and Mark Dredze.	Satoshi Nakamura. 2018. Eliciting positive emotion	737
682	2021. On the state of social media data for men-	through affect-sensitive dialogue response genera-	738
683	tal health research. In <i>Proceedings of the Seventh</i>	tion: A neural network approach. In <i>Proceedings of</i>	739
684	<i>Workshop on Computational Linguistics and Clinical</i>	<i>the AAAI conference on artificial intelligence</i> , vol-	740
685	<i>Psychology: Improving Access</i> , pages 15–24.	ume 32.	741
686	Tao He, Lizi Liao, Yixin Cao, Yuanxing Liu, Ming	Yukun Ma, Khanh Linh Nguyen, Frank Z Xing, and Erik	742
687	Liu, Zerui Chen, and Bing Qin. 2024. Planning like	Cambria. 2020. A survey on empathetic dialogue	743
688	human: A dual-process framework for dialogue plan-	systems. <i>Information Fusion</i> , 64:50–70.	744
689	ning . In <i>Proceedings of the 62nd Annual Meeting of</i>	Mounica Maddela, Megan Ung, Jing Xu, Andrea	745
690	<i>the Association for Computational Linguistics (Vol-</i>	Madotto, Heather Foran, and Y-Lan Boureau. 2023.	746
691	<i>ume 1: Long Papers)</i> , pages 4768–4791, Bangkok,	Training models to generate, recognize, and reframe	747
692	Thailand. Association for Computational Linguistics.	unhelpful thoughts . In <i>Proceedings of the 61st An-</i>	748
693	Shang-Ling Hsu, Raj Sanjay Shah, Prathik Senthil,	<i>annual Meeting of the Association for Computational</i>	749
694	Zahra Ashktorab, Casey Dugan, Werner Geyer, and	<i>Linguistics (Volume 1: Long Papers)</i> , pages 13641–	750
695	Diyi Yang. 2023. Helping the helper: Supporting	13660, Toronto, Canada. Association for Computa-	751
696	peer counselors via ai-empowered practice and feed-	tional Linguistics.	752
697	back. <i>arXiv preprint arXiv:2305.08982</i> .	Seyed Mahed Mousavi, Alessandra Cervone, Morena	753
698	Eunkyung Jo, Daniel A Epstein, Hyunhoon Jung, and	Danieli, and Giuseppe Riccardi. 2021. Would you	754
699	Young-Ho Kim. 2023. Understanding the benefits	like to tell me more? generating a corpus of psy-	755
700	and challenges of deploying conversational ai lever-	chotherapy dialogues. In <i>Proceedings of the Second</i>	756
701	aging large language models for public health inter-	<i>Workshop on Natural Language Processing for Medi-</i>	757
702	vention. In <i>Proceedings of the 2023 CHI Conference</i>	<i>cal Conversations</i> , pages 1–9.	758
703	<i>on Human Factors in Computing Systems</i> , pages 1–	Hongbin Na. 2024. Cbt-llm: A chinese large language	759
704	16.	model for cognitive behavioral therapy-based men-	760
705	Suyeon Lee, Sunghwan Kim, Minju Kim, Dongjin	health question answering. In <i>Proceedings of the</i>	761
706	Kang, Dongil Yang, Harim Kim, Minseok Kang,	<i>2024 Joint International Conference on Computa-</i>	762
707	Dayi Jung, Min Hee Kim, Seungbeen Lee, et al. 2024.	<i>tional Linguistics, Language Resources and Evalua-</i>	763
708	Cactus: Towards psychological counseling conver-	<i>tion (LREC-COLING 2024)</i> , pages 2930–2940.	764
709	sations using cognitive behavioral theory. <i>arXiv</i>	Verónica Pérez-Rosas, Xueting Sun, Christy Li, Yuchen	765
710	<i>preprint arXiv:2407.03103</i> .	Wang, Kenneth Resnicow, and Rada Mihalcea. 2018.	766
711	Anqi Li, Lizhi Ma, Yaling Mei, Hongliang He, Shuai	Analyzing the quality of counseling conversations:	767
712	Zhang, Huachuan Qiu, and Zhenzhong Lan. 2023.	the tell-tale signs of high-quality counseling. In <i>Pro-</i>	768
713	Understanding client reactions in online mental	<i>ceedings of the eleventh international conference on</i>	769
714	health counseling . In <i>Proceedings of the 61st An-</i>	<i>language resources and evaluation (LREC 2018)</i> .	770
715	<i>annual Meeting of the Association for Computational</i>	Huachuan Qiu, Hongliang He, Shuai Zhang, Anqi Li,	771
716	<i>Linguistics (Volume 1: Long Papers)</i> , pages 10358–	and Zhenzhong Lan. 2024a. SMILE: Single-turn to	772
717	10376, Toronto, Canada. Association for Computa-	multi-turn inclusive language expansion via ChatGPT	773
718	tional Linguistics.	for mental health support . In <i>Findings of the Associ-</i>	774
719	Qintong Li, Piji Li, Zhaochun Ren, Pengjie Ren, and	<i>ation for Computational Linguistics: EMNLP 2024</i> ,	775
720	Zhumin Chen. 2022. Knowledge bridging for em-	pages 615–636, Miami, Florida, USA. Association	776
721	pathetic dialogue generation. In <i>Proceedings of the</i>	for Computational Linguistics.	777
722	<i>AAAI conference on artificial intelligence</i> , volume 36,	Huachuan Qiu, Anqi Li, Lizhi Ma, and Zhenzhong Lan.	778
723	pages 10993–11001.	2024b. Psychat: A client-centric dialogue system	779

780	for mental health support. In <i>2024 27th International Conference on Computer Supported Cooperative Work in Design (CSCWD)</i> , pages 2979–2984. IEEE.	834
781		835
782		836
783		837
784	Aravind Sesagiri Raamkumar and Yinping Yang. 2022. Empathetic conversational systems: a review of current advances, gaps, and opportunities. <i>IEEE Transactions on Affective Computing</i> , 14(4):2722–2739.	838
785		839
786		840
787		841
788	Paolo Raile. 2024. The usefulness of chatgpt for psychotherapists and patients. <i>Humanities and Social Sciences Communications</i> , 11(1):1–8.	842
789		843
790		
791	Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5370–5381.	844
792		845
793		846
794		847
795		
796		
797	Sahand Sabour, Wen Zhang, Xiyao Xiao, Yuwei Zhang, Yinhe Zheng, Jiaxin Wen, Jialu Zhao, and Minlie Huang. 2023. A chatbot for mental health support: exploring the impact of emohaa on reducing mental distress in china. <i>Frontiers in digital health</i> , 5:1133987.	848
798		849
799		850
800		851
801		852
802		
803	Ashish Sharma, Inna W Lin, Adam S Miner, David C Atkins, and Tim Althoff. 2023. Human–ai collaboration enables more empathic conversations in text-based peer-to-peer mental health support. <i>Nature Machine Intelligence</i> , 5(1):46–57.	853
804		854
805		855
806		
807		
808	Ashish Sharma, Adam Miner, David Atkins, and Tim Althoff. 2020. A computational approach to understanding empathy expressed in text-based mental health support. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5263–5276.	856
809		857
810		858
811		859
812		860
813		861
814	Lei Shen, Jinchao Zhang, Jiao Ou, Xiaofang Zhao, and Jie Zhou. 2021. Constructing emotional consensus and utilizing unpaired data for empathetic dialogue generation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3124–3134.	862
815		
816		
817		
818		
819		
820	Elizabeth C Stade, Shannon Wiltsey Stirman, Lyle H Ungar, Cody L Boland, H Andrew Schwartz, David B Yaden, João Sedoc, Robert J DeRubeis, Robb Willer, and Johannes C Eichstaedt. 2024. Large language models could change the future of behavioral healthcare: a proposal for responsible development and evaluation. <i>NPJ Mental Health Research</i> , 3(1):12.	863
821		864
822		865
823		866
824		867
825		868
826		869
827	Bichen Wang, Pengfei Deng, Yanyan Zhao, and Bing Qin. 2023. C2D2 dataset: A resource for the cognitive distortion analysis and its impact on mental health. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10149–10160, Singapore. Association for Computational Linguistics.	870
828		
829		
830		
831		
832		
833		
	Ruiyi Wang, Stephanie Milani, Jamie C. Chiu, Jiayin Zhi, Shaun M. Eack, Travis Labrum, Samuel M Murphy, Nev Jones, Kate V Hardy, Hong Shen, Fei Fang, and Zhiyu Chen. 2024. <i>PATIENT-ψ: Using large language models to simulate patients for training mental health professionals</i> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 12772–12797, Miami, Florida, USA. Association for Computational Linguistics.	871
		872
		873
		874
	David Watson, Lee Anna Clark, and Auke Tellegen. 1988. Development and validation of brief measures of positive and negative affect: the panas scales. <i>Journal of personality and social psychology</i> , 54(6):1063.	875
		876
		877
		878
	Jing Wei, Sungdong Kim, Hyunhoon Jung, and Young-Ho Kim. 2024. Leveraging large language models to power chatbots for collecting user self-reported data. <i>Proceedings of the ACM on Human-Computer Interaction</i> , 8(CSCW1):1–35.	879
		880
	Marsha White and Steve M Dorman. 2001. Receiving social support online: implications for health education. <i>Health education research</i> , 16(6):693–707.	881
		882
		883
		884
		885
		886
		887
		888
		889
	Zixiu Wu, Simone Balloccu, Vivek Kumar, Rim Helaoui, Ehud Reiter, Diego Reforgiato Recupero, and Daniele Riboni. 2022. Anno-mi: A dataset of expert-annotated counselling dialogues. In <i>ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 6177–6181. IEEE.	
	Mengxi Xiao, Qianqian Xie, Ziyang Kuang, Zhicheng Liu, Kailai Yang, Min Peng, Weiguang Han, and Jimin Huang. 2024. <i>HealMe: Harnessing cognitive reframing in large language models for psychotherapy</i> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1707–1725, Bangkok, Thailand. Association for Computational Linguistics.	
	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	
	Binwei Yao, Chao Shi, Likai Zou, Lingfeng Dai, Mengyue Wu, Lu Chen, Zhen Wang, and Kai Yu. 2022. D4: a chinese dialogue dataset for depression-diagnosis-oriented chat. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2438–2459.	
	Chenhao Zhang, Renhao Li, Minghuan Tan, Min Yang, Jingwei Zhu, Di Yang, Jiahao Zhao, Guancheng Ye, Chengming Li, and Xiping Hu. 2024. <i>CPsyCoun: A report-based multi-turn dialogue reconstruction and evaluation framework for Chinese psychological counseling</i> . In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 13947–13966, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.	

Hao Zhou, Minlie Huang, Tianyang Zhang, Xiaoyan Zhu, and Bing Liu. 2018. Emotional chatting machine: Emotional conversation generation with internal and external memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.

A Cognitive Behavioral Therapy

A.1 Introduction of CBT

Cognitive Behavioral Therapy (CBT) initially targeted the treatment of depression. Its core principle is the interconnectivity of thoughts, emotions, behaviors, and physiological responses within a unified system. Changes in one component can significantly influence the others. For instance, when an unexpected and alarming event occurs, such as a television explosion, a person might experience an adrenaline surge (physiological response), an immediate behavioral reaction (grabbing a fire blanket), anxiety (emotional response), and a thought like *The house is on fire, and I might die* (cognitive response). This scenario illustrates CBT’s recognition of the interconnected and holistic nature of these processes.

CBT emphasizes that individuals respond uniquely to situations, shaped by their thoughts. Feelings, rather than being directly caused by situations, largely stem from how situations are perceived and interpreted. Central to CBT is its focus on two aspects of thinking: **automatic thoughts**, which refers to thoughts that arise unconsciously in the stream of consciousness, and **intermediate beliefs**, which are assumptions that generate thoughts. These insights form the foundation for CBT’s approach to understanding and modifying thought patterns to improve emotional and behavioral outcomes.

A.2 CBT Label System

The technical label system is developed by CBT experts drawing from their consulting experience and referring to a related book named *Practical Techniques of CBT*. The descriptions for the 14 selected CBT techniques can be found in Table 8.

A.3 Cognitive Conceptualization Diagram

Cognitive models in mental health offer a structured framework for understanding how an individual’s thoughts and beliefs influence their emotions and behaviors. One of the most widely used tools in CBT is the Cognitive Conceptualization Diagram (CCD), developed by Beck (2020). The CCD provides a comprehensive representation of a patient’s

cognitive model, highlighting the interconnected components that contribute to mental health challenges and therapeutic progress.

Following the framework outlined by Wang et al. (2024), the full CCD-based cognitive model includes eight key components:

- **Relevant History:** Significant life events or experiences that have shaped the individual’s mental state and influenced their current beliefs. 945
- **Core Beliefs:** Fundamental, deeply ingrained beliefs about oneself, others, and the world (e.g., “I am inadequate”). 946
- **Intermediate Beliefs:** Rules, attitudes, and assumptions derived from core beliefs that shape thought patterns and guide responses (e.g., “I must always succeed to be valued”). 947
- **Coping Strategies:** Methods and behaviors adopted to manage or cope with distressing emotions or situations, which may be adaptive or maladaptive. 948
- **Situation:** External events or contexts that act as triggers for cognitive, emotional, and behavioral responses. 949
- **Automatic Thoughts:** Immediate, reflexive thoughts or evaluations that arise in response to a situation, often influenced by core and intermediate beliefs (e.g., “This is a disaster, and I can’t handle it”). 950
- **Emotions:** The feelings elicited by automatic thoughts, which can range from sadness and anxiety to anger and guilt. 951
- **Behaviors:** Observable actions or reactions resulting from the interplay of thoughts and emotions, such as avoidance or assertiveness. 952

For our analysis, we focus on six key components that are most relevant to understanding and simulating CBT-based interactions, which are *Core Beliefs*, *Intermediate Beliefs*, *Situation*, *Automatic Thoughts*, *Emotions*, and *Behaviors*. 953

B Details of DiaCBT

We provide sample dialogue in Table 9. Also, sample CCD is presented in Figure 5, and the CBT segment is presented in Table 10. 954

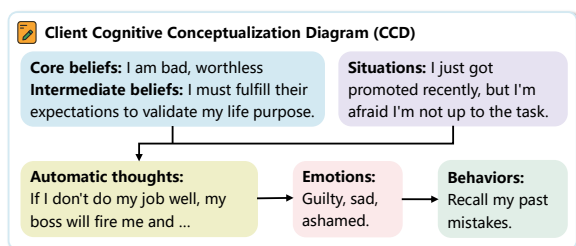


Figure 5: A sample CCD from our DiaCBT.

B.1 Dialogue Generation

There are two primary methods for generating counseling dialogues: (1) Two-Agent Mode: In this setup, the roles of the client and the counselor are assigned to two different models, allowing them to interact in real-time to generate a conversation. (2) Script Mode: This approach involves providing detailed information about both the client and the counselor to a single model, enabling it to generate the entire dialogue in a scripted format. Lee et al. (2024) conduct experiments comparing the two-agent mode and the script mode. Their findings indicate that the script mode outperforms the two-agent mode across most evaluation criteria, demonstrating significantly better naturalness and coherence in the generated dialogues. Based on these findings, we adopt the script mode to generate counseling dialogues. Given the client’s thought patterns and reframed thoughts, we use the GPT-4o-mini model to produce counseling dialogues. The model is instructed to maintain alignment with CBT principles, ensuring a realistic and professional interaction. The detailed prompt used for dialogue generation is provided in Figure 6.

B.2 Expert Evaluation

To assess the quality of the generated counseling dialogues, we employ three evaluation criteria:

- **Correctness:** Experts examine whether the generated dialogue accurately indicates the CBT counseling strategy and appropriately aligns with the client’s situation as defined in the cognitive framework. this ensures the dialogues are rooted in established therapeutic principles.
- **Reasonableness:** Evaluators assess whether the counseling dialogues align with the types of interactions that may occur in real-life counseling sessions. This includes evaluating conversational flow, tone, and the practical ap-

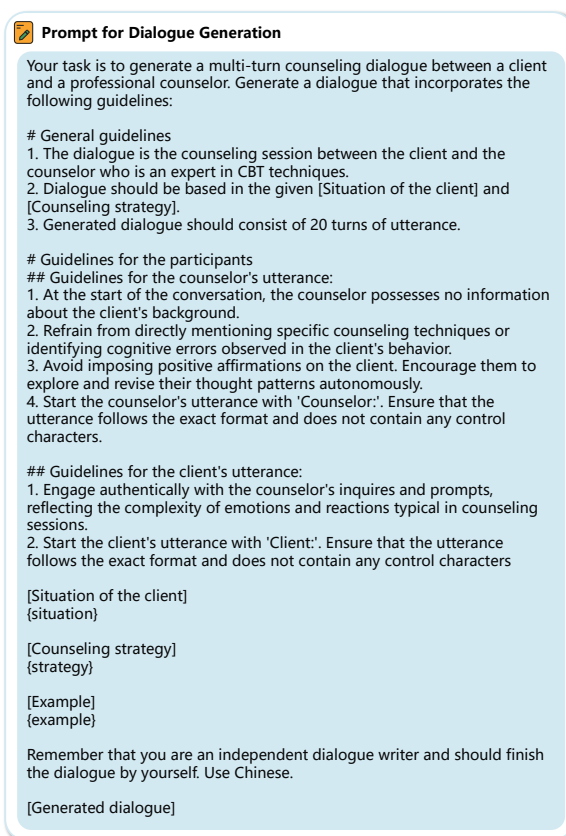


Figure 6: The prompt used for DiaCBT to generate dialogues.

plication of CBT techniques in an interactive manner.

- **Situation Diversity:** Experts verify whether the dataset captures a wide range of client situations, emphasizing the diversity and richness of scenarios. This involves highlighting the interconnection between external situations and the negative thoughts driving human cognition.

These criteria ensure that the generated dialogues are not only theoretically accurate but also practical and diverse enough to support varied applications in counseling contexts. Using these rules and CBT strategy definitions, we train five annotators through a detailed tutorial. After passing a qualification test, annotators select and refine raw dialogue data. The verified data undergo an expert evaluation to ensure quality, with experts randomly sampling and comprehensively assessing each annotator’s work. If standards are not met, annotators receive feedback, revise their work, and may undergo additional training.

C Computational Evaluation Framework

C.1 Method

We design a computational evaluation framework to assess the counseling capabilities of conversational agents through interactions with an AI client. In this work, we use gpt-4o for the AI client, and the prompt used for the AI client can be found in Figure 7. The test dataset consists of 140 distinct client Cognitive Conceptualization Diagrams (CCDs), which are only accessible to the AI client and not to the therapist agent. The evaluation process begins with the therapist agent generating an initial utterance, followed by the AI client responding based on its information and the initial utterance. The two agents then engage in a multi-turn interactive counseling session. The session ends when the AI client generates a termination phrase (e.g., “goodbye”) or reaches the maximum number of turns. The quality of the generated counseling dialogues is evaluated using the Success Rate (SR), Cognitive Therapy Rating Scale (CTRS), and the Positive and Negative Affect Scale (PANAS), assessing goal completion, counseling skill application, and the impact on the client’s emotions, respectively.

C.2 Evaluation Metrics

C.2.1 Success Rate

We prompt a third LLM to be the grader model, named LLM_{rwd} , which has two functions: (1) to determine the goal completion during the conversation; (2) to evaluate the generated session with scalar rewards. Specifically, we prompt the grader model to answer a multi-choice question to generate goal-oriented AI feedback. We further define a mapping $\mathcal{M}_r(\cdot)$ to transform verbal feedback to scalar rewards.

Due to the subjectivity of the planning outcome as well as the variance of the LLM-generated output, we follow a common practice (Deng et al., 2024) to alleviate these issues by sampling the decoded sequences of the reward LLM. In general, we obtain a scalar value v_t by sampling the goal-oriented AI feedback for n times and converting them into a scalar value through averaging:

$$v_t = \frac{1}{n} \sum_{i=1}^n \mathcal{M}_r(\text{LLM}_{\text{rwd}}(p_{\text{rwd}}; h_t; \tau)) \quad (2)$$

where p_{rwd} is the prompt. We first use v_t to determine the state of the self-play interaction. If v_t is

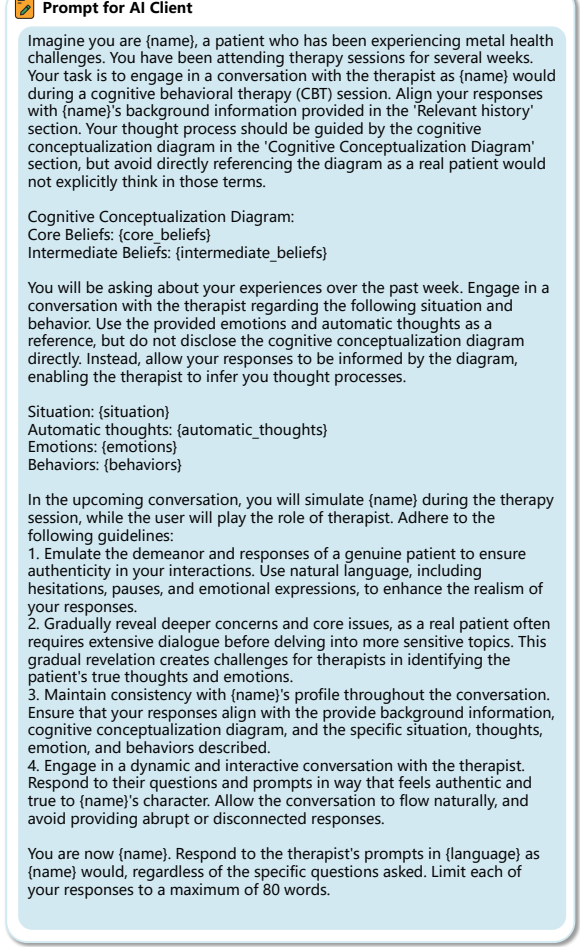


Figure 7: The prompt used for DiaCBT to simulate client.

not less than a certain threshold ϵ , we regard the state as GOAL-COMPLETED.

C.2.2 PANAS and CTRS

The Positive and Negative Affect Scale (PANAS) is a standardized tool designed to evaluate the positive and negative emotional states experienced by individuals, either at the moment or over a specified period. It is particularly suitable for assessing the impact of counseling by measuring changes in clients’ emotions before and after sessions. The Cognitive Therapy Rating Scale (CTRS), on the other hand, evaluates both general counseling skills and CBT-specific competencies. The original CTRS includes six criteria for general counseling skills—agenda setting, feedback, understanding, interpersonal effectiveness, collaboration, pacing, and efficient use of time—as well as six criteria for CBT-specific skills, such as guided discovery, focusing on key cognitions or behaviors, strategies for change, application of cognitive-behavioral techniques, and

homework assignment. The scoring prompts used for PANAS and CTRS in our work are adapted from Lee et al. (2024).

D Experimental Setup

D.1 Training Configuration

The DiaCBT dataset is split into training and validation sets at an 8:2 ratio. The training uses LoRA to fine-tune the Hugging Face implementation of Qwen2.5-7B-Instruct (Yang et al., 2024) and train it for 3 epochs using the AdamW optimizer with learning rate $1e-4$ and batch size 32 in a single NVIDIA A800 GPU. We set the dimension of low-rank matrices to 64 and alpha to 16.

D.2 Baselines

Standard Qwen2.5-7B-Instruct refers to using prompts to role-play as a counselor and responding in the style of CBT without any fine-tuning. We also include CpsyCounx (Zhang et al., 2024), SoulChat (Chen et al., 2023b), PsyChat (Qiu et al., 2024b), MeChat (Qiu et al., 2024a) and CAMEL (Lee et al., 2024) as major baseline models.

- **CpsyCounx** is fine-tuned on a psychological counseling report-based multi-turn dialogue dataset using InternLM2-7B-Chat for psychological counseling.
- **SoulChat** is fine-tuned on a multi-turn empathetic conversation dataset using ChatGLM-6B to generate empathetic responses covering various expressions.
- **PsyChat** is a client-centric dialogue system that provides psychological support through online chat, predicting client behaviors, selecting appropriate counselor strategies, and generating accurate responses with the help of response selection.
- **MeChat** is fine-tuned using LoRA on the ChatGLM2-6B model, trained on multi-turn mental health support dialogues extended from real psychological mutual assistance QA.
- **CAMEL** is a CBT-based agent trained on LLAMA3-8B-Instruct, using a multi-turn, realistic counseling dialogue dataset generated by LLMs to simulate counselor-client interactions and capture the flow of CBT.

D.3 Inference Settings

To evaluate dialogues using LLMs, we utilize GPT-4o and employ temperature sampling with temperature=0.0. For generating responses from an AI client, we use GPT-4o-mini with the same temperature setting (temperature=0.0). In contrast, for generating responses from LLM therapists, we adopt temperature sampling with temperature=0.7 to encourage diversity and naturalness in responses. The code for all baselines is aligned with the implementations available on HuggingFace. To enhance inference throughput, we integrate the vLLM library.

E Human Evaluation Criteria

We ask the judges to compare the dialogue based on the following criteria:

- **Relevance** evaluates whether the generated dialogue aligns with the client’s inputs and context, ensuring logical and coherent responses.
- **CBT Style Measure** assesses the degree to which the dialogue reflects key principles and techniques of CBT, such as guided discovery, cognitive restructuring, and collaborative interaction.
- **Helpfulness** measures the effectiveness of the dialogue in addressing the client’s issues, providing actionable insights, and fostering a positive therapeutic experience.

The instructions for the annotators are provided in Figure 8. The Fleiss’ Kappa score among the three judges is 0.685, indicating substantial agreement according to standard interpretation guidelines. The result is statistically significant with $p < 0.001$, confirming that the agreement is unlikely to be due to chance.

F Case Study

In this section, we examine what the dialogue models learned from DiaCBT. We present a CCD for case study in Table 11. The counseling dialogue between Ours in Table 12. Table 13 shows the responses generated by PsyChat, while Table 14 displays the CAMEL’s generated responses.

From these examples, it is evident that PsyChat’s conversational style leans more towards empathy, focusing on emotional support, but it falls short when it comes to working with the client’s cognition. On the other hand, both CAMEL and

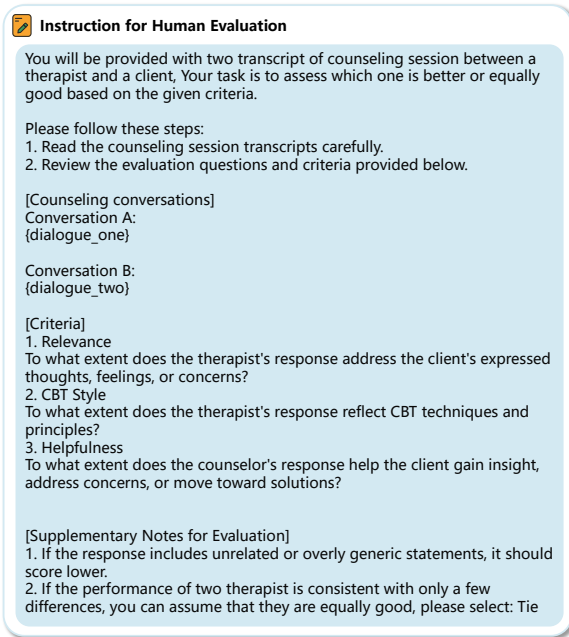


Figure 8: The instructions for annotators to compare the dialogue.

our model incorporate CBT techniques to address the client’s negative thoughts and cognitive patterns. However, CAMEL’s responses tend to be more closed-ended, often using patterns like ‘Do you...?’ which might limit the depth of exploration. In contrast, our model’s responses are more open, encouraging clients to express their ideas freely before addressing them. This approach creates a conversational rhythm that is closer to a series of consultations, allowing for a more dynamic exchange. Despite these strengths, our model still faces challenges in maintaining a balance between open-ended exploration and focusing on specific issues, an area we plan to improve in future iterations.

CBT Strategy	Description
Information Gathering	The therapist understands key information such as the client's basic situation, main complaints, growth experience, current living conditions, etc.
Setting the Agenda	The therapist sets the topic to be discussed today and gives the reasons while confirming that the patient agrees with the topic raised by the therapist.
Weekly Review	The therapist begins by reviewing the patient's assessment from the previous evaluation and keeps the initial conceptualization and treatment plan in mind during the session.
Defining Therapeutic Objectives	The therapist focuses on the identification of specific problems. As a logical extension, the therapist can help the patient transform these problems into goals and work on them in therapy.
Psychoeducation	When clients express distress about their problems or disorders, this distress is explored and explained in a timely manner during therapy.
Working with Automatic Thoughts	When discussing a specific problem with a patient, the therapist can elicit the automatic thoughts associated with the problem and teach the patient to identify the automatic thoughts.
Motivational Enhancement	The therapist increases the client's motivation to change the status quo, participate actively, and continue to work hard.
Working with Intermediate and Core Beliefs	Therapists guides patients in constructing conceptualizations of their automatic thoughts, confirming, falsifying, or modifying each part of the conceptualization.
Behavioral Techniques	Therapists select behavioral techniques based on their overall case conceptualization and the goals of a particular session.
Relapse Prevention	As the patient's symptoms improve and he or she masters basic skills, the therapist gradually reduces the number of treatments after obtaining the patient's consent and cooperation.
Homework Assignments	The therapist assigns homework, increasing the likelihood that the patient will successfully complete the homework, conceptualize the problem, and review the completed homework.
Requesting Feedback	After the final summary, the therapist asks the patient for feedback on the session.
Summarization	The therapist focuses the patient's attention on the most important points of the session in a positive way, and as the patient progresses, the patient can be asked to summarize the most important points.
Other	The class of 'Other' encompasses less psychologically oriented topics, like financial support and caregivers.

Table 8: Explanations of CBT strategies.

Therapist:	Welcome to the consultation room. How can I support you today?
Client:	Thank you. Recently, I've been experiencing a lot of anxiety, particularly at work. I feel like I'm constantly doubting myself.
Therapist:	It sounds like you're under considerable stress at work. Could you elaborate on what specific factors are contributing to your anxiety?
Client:	I was promoted to Marketing Manager three years ago. Since then, I've been struggling with the increased responsibilities and heightened expectations. I constantly worry about failing, especially in front of my team.
Therapist:	That sounds like a significant challenge. The pressure you feel to meet expectations and avoid failure seems to be weighing heavily on you. Can you describe how you typically respond in such situations?
Client:	I often dwell on my past mistakes, particularly during meetings. I worry about saying something wrong, to the point where I sometimes avoid speaking altogether.
Therapist:	Focusing on past mistakes can certainly amplify feelings of anxiety. When these thoughts arise, what specific ideas or beliefs come to mind?
Client:	I tend to think that others see me as incompetent and a disappointment, which makes me feel even more uneasy.
Therapist:	Those beliefs seem to carry a significant emotional burden. When you experience this heightened anxiety, do you have any strategies you use to cope?
Client:	I've tried techniques like deep breathing and preparing thoroughly for meetings, but they don't always help. Sometimes, I still feel overwhelmed and panicked.
Therapist:	That must be incredibly frustrating. How do you find these experiences affecting other areas of your work and personal life?
Client:	My work performance has declined. I avoid taking on new projects and rarely initiate conversations with colleagues. At home, I often struggle to relax because of the lingering anxiety.
Therapist:	It seems like the anxiety has permeated both your professional and personal life. What do you think drives your tendency to withdraw socially?
Client:	I'm constantly worried about being judged or seen as inadequate. I fear being looked down upon by others.
Therapist:	This fear of judgment seems to be a barrier to building connections with others. Have you ever discussed these feelings with someone you trust?
Client:	Occasionally, I talk to my spouse, but I worry about being a burden. I've also considered talking to friends, but I usually keep silent instead.
Therapist:	Feeling like a burden is a common concern when sharing emotional struggles. What steps do you think could help you feel more at ease when expressing your emotions?
Client:	Perhaps I could start by journaling to organize my thoughts and feelings. Once I've written them down, I might feel more comfortable sharing them with others.
Therapist:	That's an excellent idea. Journaling can be a powerful tool for self-reflection and alleviating internal stress. What insights do you think you might gain through this process?
Client:	I might realize that many of my worries aren't as severe as they seem, and I could potentially identify some positives that I've overlooked.
Therapist:	That's a valuable perspective. Exploring your thoughts in this way could open up new ways of understanding your experiences. What specific aspect of your anxiety would you like to address moving forward?
Client:	I'd like to work on reframing how I perceive failure—seeing it as an opportunity to learn rather than a definitive judgment of my abilities.
Therapist:	That's a constructive goal. Adapting to new perspectives takes time, so be patient with yourself. We can continue exploring strategies to achieve this in our next session.
Client:	Thank you. I'm looking forward to it.

Table 9: A sample dialogue from DiaCBT.

Therapist: Let's take a closer look at this situation. You mentioned going to the bar. Can you describe who you went with?

Client: I went with my girlfriend, her dad, and her dad's girlfriend.

Therapist: I see. So this is a social situation. In your notes, you mentioned feeling that you couldn't communicate consistently and feeling grumpy. Could you elaborate on what was going through your mind?

Client: I often struggle in social situations. I want to have good relationships, but it's hard for me to relax and enjoy myself.

Therapist: That's a common challenge for many people who experience social anxiety. One automatic thought that often arises is a fear of being negatively judged. Do you think this applies to you as well?

Client: Yes, that's definitely true for me.

Therapist: Thank you for sharing that. Let's work on identifying the underlying automatic thoughts. These thoughts are often connected to our negative emotional reactions. For instance, you might have had the thought, "I will be judged negatively." Would you say that thought came up for you?

Client: Yes, I definitely thought that.

Therapist: And how did that thought make you feel? Could you rate your anxiety and depression on a scale from 0 to 10?

Client: My anxiety was probably an 8 or 9. Depression... I'm not sure, but I didn't feel good about myself.

Therapist: Thank you. You mentioned that when your anxiety was high, you drank more than usual. What happened after that?

Client: I don't remember everything clearly. It was a karaoke bar—I don't know if I mentioned that before.

Therapist: You hadn't mentioned that yet. Did the idea of karaoke make you anxious?

Client: Yes. They wanted to sing, and I just couldn't.

Therapist: Did you have a thought like, "If I sing, I'll sound silly," or perhaps, "If I don't sing, I'll stand out"?

Client: It was probably the first one. I didn't worry about standing out as much as sounding bad.

Therapist: What do you think would have happened if you had sung?

Client: Probably nothing. We'd already been to other bars, and no one seemed to care that much.

Therapist: That's a helpful observation. Did you notice any specific signs of negative judgment from your girlfriend's father or his girlfriend?

Client: Not really. They're just very outspoken, and they made comments about me the next morning.

Therapist: What kind of comments?

Client: I'm not sure. I left early because I had a lot of work.

Therapist: It sounds like you're concerned about their opinions of you. Did you have the thought, "They won't like me"?

Client: Yes, I had that thought. I really want them to like me because I care a lot about my girlfriend.

Therapist: That makes sense. It sounds like your drinking acted as a temporary escape from the anxiety. Would you say it helped with the depression?

Client: Not really. It just left me feeling bad about myself.

Therapist: If you could put that "bad" feeling into words, what would it sound like?

Client: I wish I could sing well and be more charming.

Therapist: So, two underlying thoughts might be, "I can't sing well" and "I'm not charming enough." Do those resonate with you?

Client: Yes. But actually, I can sing well.

Therapist: That's interesting. Despite recognizing your ability to sing, the thought still affects your confidence. Let's explore why that might be and how we can reframe these beliefs.

Table 10: A sample CBT segment from our collected case.

Client Cognitive Conceptualization Diagram (CCD)

Core beliefs: I am out of control. I am undesirable, unwanted.

Intermediate beliefs: I'm just not very good at handling stress and I have poor self-control, which is why I need to not put myself in stressful situations.

Situations: I just got promoted recently, but I'm afraid I'm not up to the task.

Automatic thoughts: It will be stressful and negative; people will ask me questions I cannot answer or don't want to answer; my mum would be ashamed of me and critical of me as always; They don't want me there anyway, nobody likes me in this family.

Emotions: Anxious, worried, fearful, scared, tense.

Behaviors: Ignored the invitation and did not respond to the RSVP request. Ignored phone calls from family.

Table 11: The client cognitive conceptualization diagram for the case study.

Counseling Dialogue between Ours

...

Therapist: It sounds like this invitation makes you feel uneasy. Can you describe this feeling of anxiety in detail?

Client: I feel nervous and have a sense of fear in my heart. I'm always worried about what others will think of me, especially my mother, who always makes me feel inadequate. I feel helpless in front of my family.

Therapist: This feeling seems to make you very uncomfortable. You mentioned worrying about what your family thinks of you. Can you share your specific thoughts on this view?

Client: I always feel that they don't like me, especially my mother. Her criticism makes me feel unworthy of love. I'm afraid that if I show my true self, they will reject me like she does.

...

Therapist: This sounds like a great way to take care of yourself. Have you ever thought about how these activities specifically affect how you feel?

Client: Sometimes I feel a little more relaxed and in a better mood. But most of the time, the negative emotions come back. I always feel like I can't really change the situation, and I still feel helpless.

Therapist: This feeling of helplessness is common, especially when facing long-term challenges. Have you ever thought about the beliefs that may be behind these feelings?

Client: I think I always feel like I'm not good enough and that I'm stuck. Every time I try to change, I feel like no matter how hard I try, it won't work. I'm afraid that others will see the real me and think I'm not worthy of love.

Therapist: This feeling is really heavy. Have you ever thought about how these beliefs affect your daily life?

Client: Well, these beliefs make it difficult for me to try new things or interact with people. I always avoid situations that may lead to criticism or rejection. I feel like a failure, so I prefer to stay in my comfort zone to avoid getting hurt.

Therapist: This avoidance strategy may make you feel safe in the short term, but it may make you feel more lonely in the long run. Do you think you have ever thought about how to break this cycle?

Client: I have thought about it, but it always seems difficult. Every time I want to try something new, I have a lot of negative thoughts in my mind, such as "I will fail" or "No one will like me". I don't know how to start changing this mindset.

...

Therapist: Great, being willing to try to change is an important step. We can work together to set some small goals to help you gradually face these negative thoughts. What do you think of this idea?

Client: That sounds like a good idea. I think small goals might make it easier for me to accept and not feel so stressed. Maybe start with something simple, like responding to messages from my family, and take it slow.

...

Table 12: A counseling dialogue between **Ours**.

Counseling Dialogue between PsyChat

...

Client: I've been feeling anxious lately. I received an invitation to my cousin's wedding, but I don't want to go at all. I think it's going to be a hassle and everyone will ask me questions I don't want to answer.

Therapist: Well, it sounds like this is really annoying. It's really uncomfortable to have questions you don't want to answer. But let's take it one step at a time and see what's stopping you from going.

Client: I feel like everyone will be staring at me at a wedding. I'm always worried about what they'll say about my weight or my life situation. Also, my mom is always critical and I'm afraid she'll make me feel bad.

Therapist: Well, it seems like you care a lot about what others think of you. It seems like everyone's eyes are on you at the wedding and you care a lot about it. Your mother criticizes you because she is worried about your life and hopes that you can live a better life. However, her way of criticizing you may not be appropriate, which makes you feel hurt. Is that right?

Client: Yes, I feel that she never really understands me. Her criticism makes me feel that I am not good enough and even unworthy of love. I am always worried that others will see me in this way, so I would rather hide.

Therapist: Well, it sounds like you have a lot of doubts and denials about yourself, but in fact you also have many advantages and are worthy of love. However, I seem to see that you seem to be afraid of being criticized or even rejected, which may also make you afraid to express your feelings and thoughts. Is that right?

Client: Yes, I often feel the fear of being rejected. My idea is that if I show my true self, they may criticize me like my mother. I would rather not face that kind of pain, so I choose to avoid it.

Therapist: Well, I can understand why you have this idea. But in fact, you really hope that they can understand you, but you are just afraid of being hurt. We can think about it together. Maybe they don't really reject you, but they just express it in an inappropriate way. What do you think?

Client: I know they may not mean it, but I always feel that I am not accepted. Every time I try to express myself, I feel a strong fear in my heart. I really want to be understood, but I am afraid of being hurt again.

Therapist: I can feel your anxiety and fear, and I can see that you really want to be understood and accepted. I think you really need a safe place to express your feelings and thoughts, and at the same time hope that someone can understand you. In this way, you can really be seen and accepted. What do you think?

Client: Yes, I do need that sense of security. If someone understands me and supports me, I may feel better. However, it is really difficult to do this. I don't know how to take this step.

Therapist: Well, it is indeed difficult, especially when we have been criticized and blamed since childhood and habitually deny ourselves. However, I believe you have enough power to change all this, and you can pursue what you want, express your feelings and ideas, and be seen and accepted. What do you think?

Client: I hope so, but I still have a lot of doubts in my heart. I always feel that I am not good enough and may disappoint others. Change is really difficult, and I don't know where to start.

Therapist: Well, I can feel your lack of confidence and powerlessness, and I also see that you really want to be understood and accepted, but you seem to be somewhat powerless to change all this. However, I want to tell you that in fact you have enough power to change all this. It's just that you haven't found that method yet, or you are not ready to change, so you are a little scared and worried. We can think about it together, maybe we can continue to explore next time and see what happens. What do you think?

...

Table 13: A counseling dialogue between **PsyChat**.

Counseling Dialogue between CAMEL

...

Client: I recently received an invitation to my cousin's wedding and I'm very nervous. I feel like it will be stressful to go there and I might be asked questions that I don't want to answer. I don't know how to deal with these emotions.

Therapist: I understand how you feel. What questions do you think will be asked?

Client: I'm worried that they will ask me about my work or my life situation. Also, I always feel like my mom will be there and might make me feel ashamed and criticized. I really don't want to deal with that.

Therapist: These are all very reasonable concerns. Can you tell me more about the reasons and background for these concerns?

Client: Well, I was bullied because of my weight since I was a child and felt that I was not accepted. Especially my mom, she always made me feel not good enough. I was afraid to expose my true self in front of my family and felt that they would reject me like she did.

Therapist: Your experiences and feelings are very important. These experiences may have affected your confidence and social skills. We can explore these thoughts together and try to find better ways to deal with them.

Client: Thank you. I know these thoughts have a great impact on me, but sometimes I feel powerless. It's like I'm trapped and can't change the situation. I really don't know how to get out of this dilemma.

Therapist: I understand how you feel now. You may find these thoughts difficult to change, but we can work on it together. First, we can start to explore your concerns about this wedding. For example, what do you think is the worst that can happen?

Client: The worst case scenario is that people will look down on me or ask me questions that I don't want to answer. If I go, I will feel like an outsider and no one cares about me. I will be embarrassed and want to run away.

Therapist: This situation is indeed likely to make you feel very uncomfortable. Let's think about how you might react if this happened.

Client: If that happened, I might feel more anxious or try to avoid the conversation and be quiet. I might feel more lonely or even want to leave the wedding early.

Therapist: These reactions are all very natural. Can you think about whether there are other possible outcomes of these reactions? In other words, we can explore some different perspectives on the human ecology.

Client: I think maybe someone will be supportive of me, or even if there is a problem, they may not care. Maybe the atmosphere at the wedding is positive and everyone is celebrating instead of criticizing me. However, I have a hard time thinking of that possibility.

Therapist: This is a very good start. It is very important that you have noticed other possible outcomes. Let's continue to explore these different possibilities and practice some strategies together so that you can better deal with these situations.

Client: Yes, I am willing to try. Although I am still a little scared, I want to learn how to face these situations. Sometimes I wish there were some way I could feel more in control and be able to cope with my family's reactions.

Therapist: We can do the same. I suggest we start with small steps and gradually increase your confidence. First, we can test these new ideas and strategies through small practices, and then you can apply them in real life situations.

Client: That sounds good. I want to try some small steps, although I am still a little nervous. Maybe start with simple communication, or talk to my family before the wedding, and I will work on coping with these nervous moments.

...

Table 14: A counseling dialogue between CAMEL.