

Automatic Evaluation for Mental Health Counseling using LLMs

Anonymous ACL submission

Abstract

High-quality psychological counseling is crucial for mental health worldwide, and timely evaluation is vital for ensuring its effectiveness. However, obtaining professional evaluation for each counseling session is expensive and challenging. Existing methods that rely on self or third-party manual reports to assess the quality of counseling suffer from subjective biases and limitations of time-consuming.

To address above challenges, this paper proposes an innovative and efficient automatic approach using large language models (LLMs) to evaluate the working alliance in counseling conversations. We collected a comprehensive counseling dataset and conducted multiple third-party evaluations based on therapeutic relationship theory. Our LLM-based evaluation, combined with our guidelines, shows high agreement with human evaluations and provides valuable insights into counseling scripts. This highlights the potential of LLMs as supervisory tools for psychotherapists. By integrating LLMs into the evaluation process, our approach offers a cost-effective and dependable means of assessing counseling quality, enhancing overall effectiveness.

1 Introduction

Globally, approximately one in five individuals experience mental health problems each year, with many seeking psychological counseling for support (Eysenbach et al., 2004; Steel et al., 2014; Holmes et al., 2018). Timely feedback during counseling sessions can greatly enhance the quality of counseling (Lambert, 2013a). However, obtaining this level of supervision is expensive and challenging. Currently, counseling quality is often assessed through retrospective self-reports from counselors and clients (Goldberg et al., 2020), but these self-assessments are prone to biases, compromising their reliability. Our analysis demonstrates that counselors and clients do not always agree in their

feedback, with counselors exhibiting an optimistic bias and clients influenced by social biases. Therefore, there is a need for an affordable, impartial, and professional third-party method of counseling evaluation and feedback.

Existing research in the field has aimed to provide automatic counseling quality evaluation, but their methods are often not as effective as human evaluation for three reasons. First, they are limited to predicting client-reported scores due to the lack of ratings from third-party experts (Atkins et al., 2014; Wu et al., 2023; Li et al., 2022, 2023). Second, they typically analyze individual turns in the conversation and fail to consider the multi-turn interaction within the entire counseling session (Martinez et al., 2019; Goldberg et al., 2020; Lin et al., 2023, 2022). Third, their evaluations are based on specific linguistic features from counselors' and clients' utterances, limiting interpretability (Martinez et al., 2019; Goldberg et al., 2020).

In this paper, we address these limitations by introducing an effective automatic approach to third-party assessment of the working alliance, achieved through analyzing entire counseling conversations using large language models (LLMs) (Wei et al., 2022a). We collect a large-scale text-based counseling dataset from an online counseling platform, including self-reported working alliance scores from both counselors and clients. Using an observer version of the working alliance scale based on Bordin's theory of therapeutic relationship (Bordin, 1979), we design specific guidelines for each question of the scale. Experts carefully annotate a subset of counseling sessions, highlighting the differences in working alliance assessment from diverse perspectives and emphasizing the need for multiple experts.

We then utilize these guidelines to enhance the proficiency of advanced LLMs, such as ChatGPT (OpenAI, 2023a) and GPT-4 (OpenAI, 2023b), in evaluating the annotated sessions. Experimental

findings demonstrate that the precise guidelines significantly improve GPT-4’s ability to assess the working alliance, ensuring internal consistency and alignment with human evaluations. Additionally, we validate the integration of Chain-of-Thought (CoT) (Wei et al., 2022b) into GPT-4, enabling it to identify supportive evidence for scoring within the conversation, further enhancing its capabilities. Importantly, the evidence extracted by GPT-4 proves instrumental in improving agreement among human annotators, highlighting the value of automatic evaluation as a tool for enhancing human understanding of counseling scripts. This also suggests the potential for LLMs to serve as supervisors for psychotherapists.

2 Related Work

Evaluating Counseling Quality Using NLP.

Many researchers have endeavored to leverage machine learning and NLP techniques for the automatic evaluation of conversations in mental health counseling, including assessing counselors’ therapeutic skills (Cao et al., 2019; Gibson et al., 2016) and treatment fidelity (Atkins et al., 2014), as well as clients’ intervention responses (Tanana et al., 2015; Li et al., 2023). These work mostly focuses on studying individual participants’ behaviors and language features, rather than relational dynamics between them. However, in psychotherapy research, the relationship between counselors and clients are widely investigated (Ribeiro et al., 2013; Norcross, 2010; Falkenström et al., 2014). The working alliance, defined as the collaboration and attachment between counselors and clients, is a crucial researched variable (Bordin, 1979; Norcross, 2010; Falkenström et al., 2014). There are methods that attempt to evaluate the therapeutic relationship between counselors and clients but only limited to clients’ self-reported general alliance, due to the lack of observers’ assessments (Goldberg et al., 2020; Martinez et al., 2019). But the scores of alliance obtained from counselors and their clients appear to be independent and unreliable, influenced by their own biases (Horvath and Greenberg, 1994). Our research is designed to align with the fine-grained observer-rated alliance based on theoretical framework, which demonstrates significant predictive capabilities for objective counseling outcomes in psychology (Fenton et al., 2001).

LLMs for Conversation Evaluation. As the emergence of Large Language Models (LLMs)

showcasing advanced text understanding and reasoning capabilities, recent research has explored to leverage LLMs to evaluate the quality of conversations (Wang et al., 2023; Lin and Chen, 2023; Gong and Mao, 2023). Most works employ LLMs to assess conversation responses, focusing on criteria such as naturalness, coherence and fluency. Wu et al. (2023) specifically investigates the capability of ChatGPT to assess the context-consistency and coherence of the reflection strategy employed in counseling sessions. However, these studies only confine evaluation to the utterance level of conversations. There are some research attempts to assess the overall quality or individual performances exhibited in entire conversations, with the conversation scenarios limited to task-oriented or chit-chat dialogues (Mehri and Eskenazi, 2020; Lin and Chen, 2023). Different from the existing studies, our work aims to empower LLMs to better understand interlocutors’ relationships in the context of mental health counseling. This task presents a greater challenge for LLMs, as it demands a heightened capability to comprehend natural language and human interactions.

3 Framework and Guidelines for Measuring Working Alliance

To understand the working alliance between counselors and clients in psychological counseling, we adapt the existing working alliance definitions and observer-rated scales to the context of online text-only counseling. To facilitate accurate third-party evaluations based on this framework, we carefully design specific guidelines in collaboration with counseling psychology experts.

3.1 Framework

In psychology research, the preeminent definition of alliance is introduced in Bordin (1979)’s theory of therapeutic working alliance. According to this theory, the working alliance represents interactive and collaborative elements of the relationship between counselors and clients in the context of an affective bond of positive attachment (Constantino et al., 2002). This concept consists of three core components – counselors and clients’ mutual agreement on the targets of counseling (*Goal*), abilities to engage in the tasks of counseling (*Approach*), as well as the cultivation of emotional connections (*Affective Bond*) (Bordin, 1979).

In order to measure the working alliance within

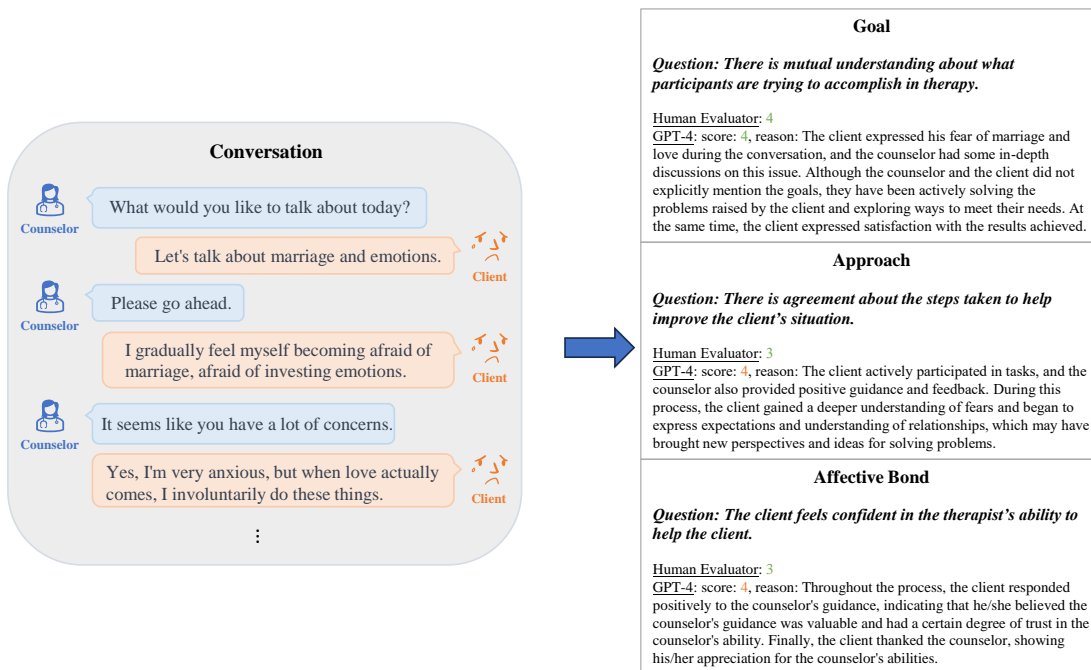


Figure 1: Our framework of working alliance contains three integral components - agreement on goal-setting and the approaches towards goals, alongside the establishment of affective bonds. For each of these components, we design four detailed questions and differentiate the evidence spectrum spanning from considerable evidence against, some evidence against, no evidence against, some evidence for, and considerable evidence for these aspects. All examples shown in this paper have been modified to ensure the absence of any real personal information about the speakers.

182 the aforementioned theoretical framework from the
 183 perspective of observers, we adopt the Observer-
 184 rated Short version of Working Alliance Inventory
 185 (WAI-O-S) (Tichenor and Hill, 1989). This inven-
 186 tory comprises 12 designed questions, where each
 187 dimension of the working alliance is measured by
 188 four questions. Each question is rated ranging from
 189 1 to 5 points. Its reliability and validity has under-
 190 gone thorough and comprehensive verification in
 191 various psychotherapy types (Santirso et al., 2018;
 192 Ribeiro et al., 2021). Table 1 presents the dimen-
 193 sions along with questions that shape the working
 194 alliance.

195 **Goal.** In counseling, goals are important for facil-
 196 itating changes in clients' thoughts, feelings, and
 197 actions. They provide direction for both counselors
 198 and clients during their sessions. Clear agreement
 199 on goals increases adherence and leads to better
 200 outcomes. However, at the beginning of counsel-
 201 ing, there can be a lack of clarity about clients'
 202 issues and differences in goals between clients and
 203 counselors. To address this, counselors should en-
 204 gage in deeper discussions with clients to establish
 205 mutually endorsed and valued objectives.

206 **Approach.** In addition to the agreement on goals,
 207 the strength of the working alliance also depends
 208 on the participants' clear and mutual understand-
 209 ing as well as acceptance on the tasks that their
 210 shared goals impose upon them (Bordin, 1983).
 211 Tasks are usually assigned by counselors based on
 212 their counseling styles, personal experiences and
 213 predispositions. However, clients may not fully un-
 214 derstand the interconnections between the assigned
 215 tasks and the overarching goals. Moreover, clients
 216 may perceive that the demands of tasks exceed their
 217 abilities. In such cases, counselors need to skill-
 218 fully adapt to their clients by offering alternative
 219 or modified tasks, thereby empowering clients to
 220 actively and effectively engage.

221 **Affective Bond.** Apart from cognitive collabora-
 222 tion, emotional connections play a crucial role in
 223 shaping the therapeutic alliance. The concept of
 224 affective bonds embraces the complex network of
 225 positive personal attachments between counselors
 226 and clients, including issues such as mutual trust,
 227 liking, acceptance, and confidence (Horvath and
 228 Marx, 1990). As clients perceive that counselors
 229 genuinely care about and appreciate them, a sense

Dimension	Question	No.
Goal	There is mutual understanding about what participants are trying to accomplish in therapy.	Q1
	The client and counselor are working on mutually agreed upon goals.	Q2
	The client and counselor have same ideas about what the client's real problems are.	Q3
	The client and counselor have established a good understanding of the changes that would be good for the client.	Q4
Approach	There is agreement about the steps taken to help improve the client's situation.	Q5
	There is agreement about the usefulness of the current activity in therapy (i.e., the client is seeing new ways to look at his/her problem).	Q6
	There is agreement on what is important for the client to work on.	Q7
	The client believes that the way they are working with his/her problem is correct.	Q8
Affective Bond	There is a mutual liking between the client and counselor.	Q9
	The client feels confident in the counselor's ability to help the client.	Q10
	The client feels that the counselor appreciates him/her as a person.	Q11
	There is mutual trust between the client and counselor.	Q12

Table 1: The framework of working alliance contains three core components: *Goal*, *Approach*, and *Affective Bond*. Each dimension is assessed through a set of four questions.

of security is established, fostering a greater willingness to delve into deeper self-disclosure during counseling, particularly in discussing their negative behaviors and thoughts. Moreover, clients' confidence in counselors' capabilities to facilitate positive changes make them more inclined to accept counselors' guidance and actively participate in the tasks assigned by the counselors.

3.2 Guidelines

To facilitate the understanding of questions and the differentiation of scores by observers, we have four developers to carefully design specific guidelines for each score associated with each question.

Firstly, we randomly select 15 conversations and ask all the developers to annotate them independently based on general guidelines. After the annotation, the developers discuss the differences and confusions among their annotations in several conversations until reaching a consensus. During this process, they may refine the guidelines by compiling the behavioral indicators of counselors and clients relevant to each question, with the associated degree and frequency at each score level. The developers repeat annotating these conversations based on modified guidelines. After iterating the above step 3 times, the final version of the guidelines is obtained. The intra-class agreement (Koo and Li, 2016) among the four developers in the three iterations are as follows: 0.5267, 0.6084, and 0.6603. The monotonically increasing agreement proves that the iterative process effectively resolves differences among developers. And the moderate agreement ensures the reliability of our guidelines.

The specific guidelines and more details of the development process are presented in Appendix A.

4 Data Collection

To validate the feasibility of our proposed framework, we collect counseling conversations between professional counselors and actual clients, and carefully annotate these conversations according to the framework.

4.1 Data Source

We developed an online text-based counseling platform and enlisted 9 qualified professional counselors (7 females, *Mean age* = 34.67 years old, *SD* = 7.45). We recruited 82 adults (55 females, *Mean age* = 27.62 years old, *SD* = 5.94) as clients who were interested in and eligible for online psycho-counseling. These clients were assessed using the self-report symptom inventory (SCL-90)(Wang et al., 1999) to ensure they did not exhibit severe depressive, anxious, or psychiatric symptoms. Each counseling session lasted 50 minutes, which is a widely accepted standard duration for psychological counseling. Clients were encouraged to attend a minimum of 7 counseling sessions, scheduled weekly or bi-weekly. After each session, counselors and clients completed the therapists' and clients' versions of the Working Alliance Inventory, respectively(Horvath and Greenberg, 1989; Hatcher and Gillaspay, 2006), to assess the therapeutic relationship. These inventories are based on Bordin's theory of alliance, similar to the observer-rated scale used in this study.

We collected total 859 counseling sessions and

728 out of them received the self-reported scales from both counselors and clients. The statistics of the overall conversations are detailed in Table 2. The length of counseling conversations are significantly longer than the existing conversations obtained through crowdsourcing or generated by language models (avg. 76.07 utterances compared to 29.8 utterances in ESConv (Liu et al., 2021) and 6.36 utterances in SMILE (Qiu et al., 2023)). Moreover, each counselor-client pair engages in multiple consecutive counseling sessions (avg. 10.48 sessions compared to 4 sessions in Multi-Session Chat (Xu et al., 2022)), suggesting, in real-world scenarios, an effective resolution of clients’ concerns often requires extended multi-turn interactions and multiple sessions.

Category	Total	Counselor	Client
# Dialogues	859	-	-
# Speakers	91	9	82
# Avg. sessions per speaker	-	95.44	10.48
# Utterances	65,347	32,860	32,487
Avg. utterances per dialogue	76.07	38.25	37.82
Avg. length per utterance	26.84	24.01	29.70

Table 2: Statistics of the overall conversations.

4.2 Annotation Process

To ensure the quality of the annotations, we engaged three experienced developers of the guidelines to annotate a subset of collected conversations. Their extensive knowledge of the working alliance framework and guidelines allowed for a thorough evaluation. Before the annotation process, we took measures to protect the privacy of the counselors and clients by anonymizing their personal information, including names, organizations, addresses, and more.

For the annotation phase, we randomly selected 79 sessions involving 4 counselors and 8 clients. Each conversation was annotated by all three annotators. To determine the final score for each question, we calculated the average of all scores assigned by the annotators.

After obtaining the annotated data, we calculated the intraclass correlation coefficient (ICC)(Koo and Li, 2016) among the annotators for each question. The inter-rater agreement for the dimensions of *Goal*, *Approach*, and *Affective Bond* were found to be 0.7581, 0.6587, and 0.6498, respectively. These values indicate a reliable level of agreement among

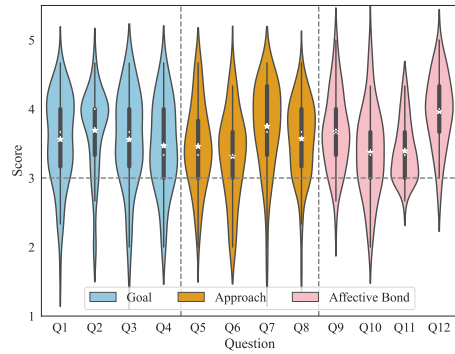


Figure 2: The violin plot of the distribution of scores annotated for each question, with a boxplot inside. The white pentagons within the violins represent the mean values.

the annotators¹. Further details regarding the inter-rater agreement for each question can be found in Appendix B.1.

4.3 Data Characteristics

Figure 2 illustrates the distribution of annotated scores for all the questions. For further insights into the average scores per dimension and question, as well as their corresponding standard deviations, please refer to Appendix B.2.

On average, the scores for each dimension range between 3.5 and 4, indicating a generally high quality in the overall counseling conversations, although there is room for improvement. It is evident that counselors and clients have developed positive emotional connections and are making progress towards shared counseling objectives. Among the three dimensions of alliance, the *Affective Bond* stands out with the highest average score, particularly in the question regarding mutual trust between counselors and clients (Q12), where the score almost reaches 4. However, the *Approach* dimension has the lowest average score, specifically in the question concerning agreement on the usefulness of the current therapy activity (Q6, avg. = 3.32). This signifies that there are still some differences between counselors and clients regarding the steps and tasks to be taken in addressing the client’s psychological issues throughout the counseling process.

¹An ICC value between 0.5 and 0.75 indicates moderate reliability, while a value between 0.75 and 0.9 indicates good reliability.

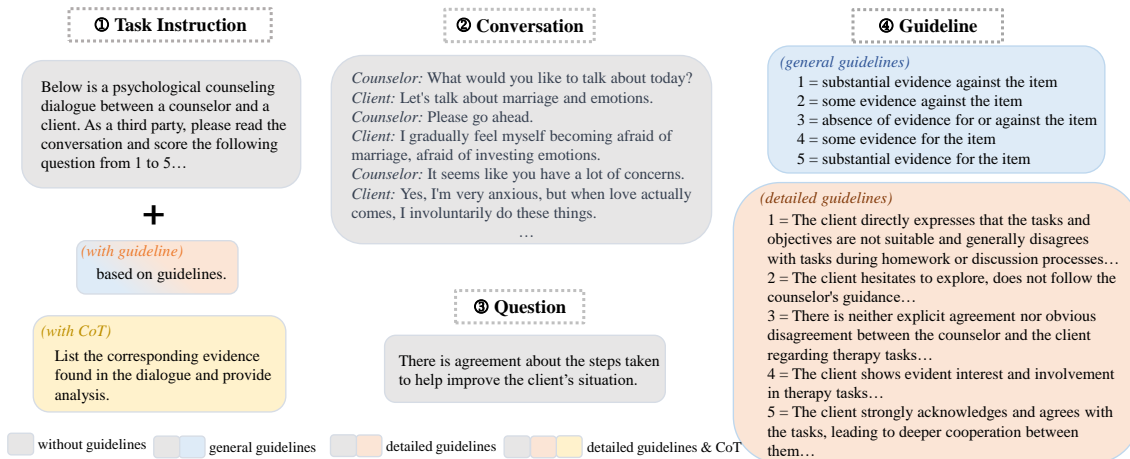


Figure 3: Example prompts for evaluating a giving conversation across different experimental setups (i.e. with different prompt types and with/without CoT) addressing question *There is agreement about the steps taken to help improve the client's situation*. General guidelines remain consistent across different questions, whereas detailed guidelines are intricately tailored to each specific question.

5 LLM Evaluation

To investigate whether LLMs can evaluate the working alliance between counselors and clients, we conduct zero-shot experiments to prompt advanced LLMs including ChatGPT and GPT-4 to assess text-only conversations based on our proposed guidelines.

5.1 Setup

The prompt comprises four key components: the definition of the evaluation task, the entire counseling conversation to be evaluated, the evaluation question and corresponding guidelines. To further investigate the impact of guidelines on the evaluation capabilities of LLMs, we conduct three experimental settings — prompting LLMs without guidelines, with general guidelines, and with our proposed detailed guidelines. We also explore the impact of the Chain-of-Thoughts process on the scoring of LLMs after providing detailed evaluation criteria. In the CoT setting, we require models to provide corresponding evidence for ratings within the dialogue text additionally. We carefully design specific prompts for each experiment setting accordingly. Figure 3 illustrates example prompts designed for LLMs to score a given counseling session across four experimental conditions.

5.2 Models

We select two accessible top-performing GPT-series models – ChatGPT (*gpt-35-turbo-16k* model) and GPT-4 (*gpt-4* model). These models have been

enhanced to follow human instructions through instruction tuning and align with human preferences via reinforcement learning from human feedback (RLHF, (Ouyang et al., 2022)). Compared to ChatGPT, GPT-4 is widely acknowledged as a more robust model that can solve more complex problems with broader knowledge and stronger reasoning capabilities. Our interactions with these models are facilitated using the official API provided by OpenAI. The temperature and nuclear sampling parameter are set as 1.0 for both models. Each model is tasked with rating the same conversation three times independently for thorough evaluation.

6 Results and Analysis

In this section, we answer a set of research questions through the above data collected and the conducted experiments. We first demonstrate the difference among working alliance rated by counselors, clients and observers. We then analyze the potential of LLMs to serve as tools for assessing consultation quality, and effective approaches to enhance their evaluation capabilities. Furthermore, we discuss the feasibility of utilizing evidence generated by GPT-4 to further enhance human annotators' agreement.

6.1 RQ1: Is Retrospective Self-reports Reliable?

To examine differences in the assessment of working alliance from various perspectives, we calculated pairwise Pearsonr correlation among coun-

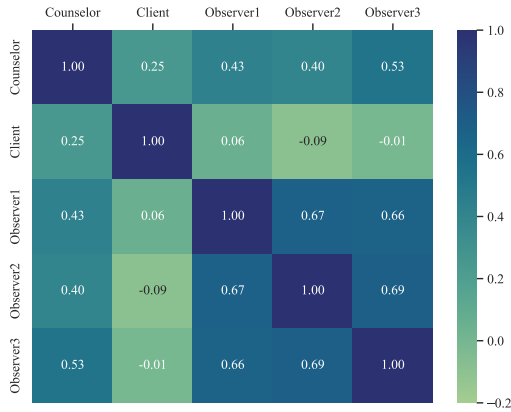


Figure 4: The heatmap of the pairwise correlation among counselors, clients, and three annotators in evaluating the working alliance.

424 selsors, clients, and three annotators. The heatmap
 425 in Figure 4 shows low correlation between coun-
 426 selors and clients, confirming rating disparities
 427 based on self-reports. Counselors tend to rate sig-
 428 nificantly higher than their clients in 21% of coun-
 429 seling sessions (details in Appendix C), indicating
 430 an overly positive bias in self-assessment (Walfish
 431 et al., 2012; Lambert, 2013b). On the client side,
 432 13 out of 81 clients consistently assign the highest
 433 ratings in over half of their counseling sessions, po-
 434 tentially due to social desirability and dissonance-
 435 reducing responses, making their self-reports unre-
 436 liable (Shick Tryon et al., 2007; Tryon et al., 2008).
 437 Therefore, an impartial and professional third-party
 438 assessment becomes essential for unbiased and ac-
 439 curate evaluation (Tichenor and Hill, 1989; Gold-
 440 berg et al., 2020).

441 In contract, consensus among third-party ob-
 442 servers surpasses that between observers and either
 443 counselors or clients. This emphasizes the robust-
 444 ness and feasibility of evaluating the therapeutic
 445 relationship when multiple observers contribute as-
 446 sements from a third-party perspective.

447 6.2 RQ2: Can LLMs Assess the Reliably of 448 Working Alliance ?

449 **Model Self-Consistency.** The reliability of a
 450 model as an annotator depends on its consistency
 451 in multiple independent evaluations of the same
 452 samples. We evaluated ChatGPT and GPT-4 by
 453 assessing their consistency in evidence extraction
 454 using fine-grained guidelines. The summarized re-
 455 sults are shown in Table 3, and detailed results can

456 be found in Table 8 in the Appendix.

457 Our analysis shows that both ChatGPT and GPT-
 458 4 demonstrate a moderate level of self-consistency,
 459 indicating their reliability in evaluations. Notably,
 460 GPT-4 consistently outperforms ChatGPT in all
 461 settings, leading to a significantly higher over-
 462 all agreement with detailed guidelines and CoT
 463 (0.7205 compared to 0.5209). However, it’s im-
 464 portant to note that while general guidelines have
 465 higher correlation scores across different runs, they
 466 have lower correlation scores compared to human
 467 annotators in our later experiments.

468 These findings highlight GPT-4’s superior abil-
 469 ity to comprehend conversation content and consis-
 470 tently follow guidelines, making it a more reliable
 471 evaluation tool.

472 **Alignment with Human Evaluations.** The mod-
 473 els’ capability on the evaluation task is defined as
 474 the extent to which its assessments align with those
 475 of human experts. we calculate the Pearson corre-
 476 lation coefficients (Lee Rodgers and Nicewander,
 477 1988) between human and model ratings on the
 478 fine-grained rating scale ranging from 1 to 5 points.
 479 The results in Table 4 highlight GPT-4’s consis-
 480 tently superior performance compared to ChatGPT
 481 across all settings. The results also show that when
 482 given detailed guidelines and CoT, the LLM eval-
 483 uations has a correlation of 0.5 with human eval-
 484 uation, even higher than the correlation between
 485 counselors and clients.

486 6.3 RQ3: How to Improve LLMs’ Evaluation 487 Capabilities?

488 To investigate the impact of prompt tuning on
 489 the consistency of LLMs with human evaluations,
 490 we focus on two factors: the level of detail in
 491 guidelines and the setting of the Chain-of-Thought.
 492 Given the superior evaluation capabilities of GPT-4
 493 compared to ChatGPT (as shown in section 6.2),
 494 our main objective is to enhance GPT-4’s perfor-
 495 mance. The alignment between GPT-4 and human
 496 evaluations across different experimental settings
 497 is summarized in Table 4.

498 **Guidelines.** We find that GPT-4 maintains a mod-
 499 erate level of self-consistency across all guideline
 500 types, ensuring the validity of its annotated results.
 501 We further analyze the influence of guidelines on
 502 the alignment between GPT-4 and human evalua-
 503 tions.

504 As shown in Table 4, the results consistently
 505 demonstrate that increasing the level of detail in

	ChatGPT		GPT4		
	Detailed Guideline + CoT	No Guideline	General Guideline	Detailed Guideline	Detailed Guideline + CoT
<i>Overall Agreement</i>	0.5209	0.6687	0.7482	0.6854	0.7205

Table 3: The overall agreement of three independent runs of ChatGPT and GPT-4 in evaluating all question and dimension across different experimental settings.

Models		Goal	Approach	Affective Bond	Overall
ChatGPT	Detailed Guidelines + CoT	0.2004	0.3612	0.4122	0.3246
	No Guidelines	0.3591	0.4288	0.3693	0.3857
GPT-4	General Guidelines	0.3320	0.4517	0.3961	0.3933
	Detailed Guidelines	0.4979	0.5481	0.4416	0.4959
	Detailed Guidelines + CoT	0.4938	0.5448	0.4667	0.5018

Table 4: The overall Pearson correlation results of ChatGPT and GPT-4 with human evaluation on the working alliance dimensions across different experimental settings.

guidelines improves the alignment. This improvement is particularly significant when transitioning from general guidelines to more detailed ones, resulting in a notable 26.09% increase in correlation. Detailed guidelines are particularly effective in enhancing GPT-4’s performance on challenging questions. For instance, in the case of discerning whether counselors and clients like each other (Q9), GPT-4 performs poorly without guidelines or with general guidelines. However, when detailed guidelines are provided, there is a remarkable 76% increase in correlation (Detailed results can be found in Table 9 in the Appendix).

These findings highlight the potential to improve the alignment of LLM evaluations with human assessments by refining the guidelines. Ensuring high self-agreement in LLMs is a crucial prerequisite for them to be qualified evaluators.

Chain-of-Thought Prompting. Table 4 demonstrates that integrating CoT improves the alignment of GPT-4 evaluations with human assessments, particularly in forming emotional connections. CoT significantly enhances GPT-4’s performance on the challenging question Q9, resulting in a remarkable 32.05% increase in the Pearson correlation with human evaluations. Thus, facilitating evidence extraction and explanation generation prior to scoring proves to be an effective strategy for enhancing GPT-4’s comprehension of dialogue content and improving assessment accuracy.

6.4 RQ4: Can GPT-4’s Explanations Help Human Annotators?

To explore how GPT-4 assessments can support human annotation, we conducted a qualitative investigation. We assessed how evidence generated

Annotator	Weak Annotation	Modified annotation	Modified ratio
A	44	36	81.81%
B	40	40	100%
C	31	27	87.10%

Table 5: The amount of modified annotations of each annotator after reading the reasons generated by GPT-4.

through Chain-of-Thought prompting could help less experienced annotators refine their evaluations. We identified annotators who disagreed with two others as "weak annotators" and tasked them with re-evaluating samples where GPT-4 consistently aligned with over half of the annotators. To avoid bias, we didn’t disclose GPT-4’s scores to weak annotators. In Table 5, all three annotators revised their scores with a correction rate exceeding 80%, resulting in improved human agreement from 0.6888 to 0.7087 (detailed improvements in Table 6 in Appendix). This suggests that GPT-4’s ability to process extensive textual information may help humans capture crucial evidence that could be overlooked.

7 Conclusion

We developed detailed guidelines, dataset, and LLM-based approaches for evaluating therapeutic working alliance between counselors and clients in text-based counseling from the perspective of observers. Our demonstration suggests that the integration of detailed guidelines and chain-of-thought prompting empower LLMs to assess the working alliance effectively with underlying rationales. Moreover, the identified evidence proves helpful in improving the mutual understanding of working alliance for human annotators.

8 Limitations

As this is the first LLM-based approach to automatically evaluate counseling quality, there is huge room for future improvement. First, our experiments focus exclusively on assessing the performance of two preeminent general-purpose LLMs: GPT-4 and ChatGPT. Because of the high costs, unavailability of open APIs, and inherent limitations in the capabilities of LLMs, we limit our testing scope. As our approach can be generalized to any other language models, we envision future endeavors expanding into broader LLMs, particularly those fine-tuned for psychology applications. Second, despite our efforts to enhance guidelines through consensual qualitative research, improving the alignment between LLMs and human ratings, challenges persist in evaluating certain questions. We will continue to explore how to systematically design guidelines to target improvements in the ability of specific LLMs to effectively assess the working alliance of counseling .

9 Ethics Statement

Data Privacy. This study is granted ethics approval from the Institutional Ethics Committee. All the counselors and clients provided their consent to participate and received reasonable fee for participation. All the participants were notified that the conversations collected on the platform would be utilized for scientific research purposes and potentially shared with third parties for the same purpose. Participants were also informed that they could discontinue counseling and withdraw from the research at any time. The detailed consent form for clients and user services agreement are presented in Appendix D.

Throughout the annotation process, we devoted meticulous attention to manually de-identifying and anonymizing the data, ensuring the utmost protection of the privacy of both clients and counselors. Additionally, our guidelines developers and annotators, prior to accessing the conversation data, formally committed to data confidentiality agreements and adhered to ethical guidelines, underscoring our commitment to upholding the highest standards of privacy and ethical conduct.

LLM Evaluation. To avoid potential privacy concerns during LLMs evaluations, we utilize LLMs through the official API and provide them with the anonymized counseling data. Moreover,

considering that the LLMs used in our study do not achieve perfect alignment with human assessment, and may exhibit inherent biases stemming from their extensive training data, this work only aims to offer an alternative option to human evaluation rather than seeking to replace human judgements at this time. We hope the positive results in this work can provide NLP and psychology researchers with an alternative approach for applying LLMs in the automatic evaluation of counseling conversations, fostering further discussions on this topic. This can facilitate the future research of AI psychology and sociology.

References

- David C Atkins, Mark Steyvers, Zac E Imel, and Padhraic Smyth. 2014. Scaling up the evaluation of psychotherapy: evaluating motivational interviewing fidelity via statistical text classification. *Implementation Science*, 9(1):1–11.
- Edward S Bordin. 1979. The generalizability of the psychoanalytic concept of the working alliance. *Psychotherapy: Theory, research & practice*, 16(3):252.
- Edward S Bordin. 1983. A working alliance based model of supervision. *The counseling psychologist*, 11(1):35–42.
- Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort, Jaques Grobler, Robert Layton, Jake VanderPlas, Arnaud Joly, Brian Holt, and Gaël Varoquaux. 2013. API design for machine learning software: experiences from the scikit-learn project. In *ECML PKDD Workshop: Languages for Data Mining and Machine Learning*, pages 108–122.
- Jie Cao, Michael Tanana, Zac Imel, Eric Poitras, David Atkins, and Vivek Srikumar. 2019. [Observing dialogue in therapy: Categorizing and forecasting behavioral codes](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5599–5611, Florence, Italy. Association for Computational Linguistics.
- MJ Constantino, LG Castonguay, and AJ Schut. 2002. The working alliance: A flagship for the “scientist-practitioner” model in psychotherapy. *Counseling based on process research: Applying what we know*, pages 81–131.
- Andrew Darchuk, Victor Wang, David Weibel, Jennifer Fende, Timothy Anderson, and Adam Horvath. 2000. Department of psychology ohio university december 11, 2000.
- Gunther Eysenbach, John Powell, Marina Englesakis, Carlos Rizo, and Anita Stern. 2004. Health related

669	virtual communities and electronic support groups: systematic review of the effects of online peer to peer interactions. <i>Bmj</i> , 328(7449):1166.	Michael J Lambert. 2013a. <i>Bergin and Garfield's handbook of psychotherapy and behavior change</i> . John Wiley & Sons.	724 725 726
672	Fredrik Falkenström, Fredrik Granström, and Rolf Holmqvist. 2014. Working alliance predicts psychotherapy outcome even while controlling for prior symptom improvement. <i>Psychotherapy Research</i> , 24(2):146–159.	Michael J Lambert. 2013b. Outcome in psychotherapy: the past and important advances.	727 728
673		Joseph Lee Rodgers and W Alan Nicewander. 1988. Thirteen ways to look at the correlation coefficient. <i>The American Statistician</i> , 42(1):59–66.	729 730 731
674		Anqi Li, Jingsong Ma, Lizhi Ma, Pengfei Fang, Hongliang He, and Zhenzhong Lan. 2022. Towards automated real-time evaluation in text-based counseling. <i>arXiv preprint arXiv:2203.03442</i> .	732 733 734 735
675		Anqi Li, Lizhi Ma, Yaling Mei, Hongliang He, Shuai Zhang, Huachuan Qiu, and Zhenzhong Lan. 2023. Understanding client reactions in online mental health counseling . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10358–10376, Toronto, Canada. Association for Computational Linguistics.	736 737 738 739 740 741 742 743
676		Baihan Lin, Guillermo Cecchi, and Djallel Bouneffouf. 2022. Working alliance transformer for psychotherapy dialogue classification. <i>arXiv preprint arXiv:2210.15603</i> .	744 745 746 747
677	Lisa R Fenton, John J Cecero, Charla Nich, Tami L Frankforter, and Kathleen M Carroll. 2001. Perspective is everything: The predictive validity of six working alliance instruments. <i>The Journal of psychotherapy practice and research</i> , 10(4):262.	Baihan Lin, Guillermo Cecchi, and Djallel Bouneffouf. 2023. Deep annotation of therapeutic working alliance in psychotherapy. In <i>International Workshop on Health Intelligence</i> , pages 193–207. Springer.	748 749 750 751
678		Yen-Ting Lin and Yun-Nung Chen. 2023. Llm-eval: Unified multi-dimensional automatic evaluation for open-domain conversations with large language models .	752 753 754 755
679		Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. Towards emotional support dialog systems . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3469–3483, Online. Association for Computational Linguistics.	756 757 758 759 760 761 762 763 764
680		Victor R Martinez, Nikolaos Flemotomos, Victor Ardulov, Krishna Somandepalli, Simon B Goldberg, Zac E Imel, David C Atkins, and Shrikanth Narayanan. 2019. Identifying therapist and client personae for therapeutic alliance estimation. In <i>Inter-speech</i> , volume 2019, page 1901. NIH Public Access.	765 766 767 768 769 770
681		Shikib Mehri and Maxine Eskenazi. 2020. Unsupervised evaluation of interactive dialog with DialoGPT . In <i>Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 225–235, 1st virtual meeting. Association for Computational Linguistics.	771 772 773 774 775 776
682	James Gibson, Dogan Can, Bo Xiao, Zac E Imel, David C Atkins, Panayiotis Georgiou, and Shrikanth Narayanan. 2016. A deep learning approach to modeling empathy in addiction counseling. <i>Commitment</i> , 111:21.	John C Norcross. 2010. <i>The therapeutic relationship. The heart and soul of change: Delivering what works in therapy</i> , pages 113–141.	777 778 779
683			
684			
685			
686			
687	Simon B Goldberg, Nikolaos Flemotomos, Victor R Martinez, Michael J Tanana, Patty B Kuo, Brian T Pace, Jennifer L Villatte, Panayiotis G Georgiou, Jake Van Epps, Zac E Imel, et al. 2020. Machine learning and natural language processing in psychotherapy research: Alliance as example use case. <i>Journal of counseling psychology</i> , 67(4):438.		
688			
689			
690			
691			
692			
693			
694	Peiyuan Gong and Jiaxin Mao. 2023. Coascore: Chain-of-aspects prompting for nlg evaluation. <i>arXiv preprint arXiv:2312.10355</i> .		
695			
696			
697	Robert L. Hatcher and J. Arthur Gillaspay. 2006. Development and validation of a revised short version of the working alliance inventory . <i>Psychotherapy Research</i> , 16(1):12–25.		
698			
699			
700			
701	Emily A Holmes, Ata Ghaderi, Catherine J Harmer, Paul G Ramchandani, Pim Cuijpers, Anthony P Morrison, Jonathan P Roiser, Claudi LH Bockting, Rory C O'Connor, Roz Shafran, et al. 2018. The lancet psychiatry commission on psychological treatments research in tomorrow's science. <i>The Lancet Psychiatry</i> , 5(3):237–286.		
702			
703			
704			
705			
706			
707			
708	Adam O Horvath and Leslie S Greenberg. 1989. Development and validation of the working alliance inventory. <i>Journal of counseling psychology</i> , 36(2):223.		
709			
710			
711	Adam O Horvath and Leslie S Greenberg. 1994. <i>The working alliance: Theory, research, and practice</i> , volume 173. John Wiley & Sons.		
712			
713			
714	Adam O Horvath and Ronald W Marx. 1990. The development and decay of the working alliance during time-limited counselling. <i>Canadian Journal of Counselling and Psychotherapy</i> , 24(4).		
715			
716			
717			
718	Tae Kyun Kim. 2015. T test as a parametric statistic. <i>Korean journal of anesthesiology</i> , 68(6):540–546.		
719			
720	Terry K Koo and Mae Y Li. 2016. A guideline of selecting and reporting intraclass correlation coefficients for reliability research. <i>Journal of chiropractic medicine</i> , 15(2):155–163.		
721			
722			
723			

780	OpenAI. 2023a. [link] .		
781	OpenAI. 2023b. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .		
782			
783	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. 2022. Training language models to follow instructions with human feedback . <i>ArXiv</i> , abs/2203.02155.		
784			
785			
786			
787			
788			
789			
790			
791			
792	Huachuan Qiu, Hongliang He, Shuai Zhang, Anqi Li, and Zhenzhong Lan. 2023. Smile: Single-turn to multi-turn inclusive language expansion via chatgpt for mental health support .		
793			
794			
795			
796	Eugénia Ribeiro, António P Ribeiro, Miguel M Gonçalves, Adam O Horvath, and William B Stiles. 2013. How collaboration in therapy becomes therapeutic: The therapeutic collaboration coding system. <i>Psychology and Psychotherapy: Theory, Research and Practice</i> , 86(3):294–314.		
797			
798			
799			
800			
801			
802	Nathália Soares Ribeiro, Fernando Antonio Basile Colugnati, Nikolaos Kazantzis, and Laisa Marcocorela Andreoli Sartes. 2021. Observing the working alliance in videoconferencing psychotherapy for alcohol addiction: Reliability and validity of the working alliance inventory short revised observer. <i>Frontiers in Psychology</i> , 12:647814.		
803			
804			
805			
806			
807			
808			
809	Faraj A Santirso, Manuel Martín-Fernández, Marisol Lila, Enrique Gracia, and Elena Terreros. 2018. Validation of the working alliance inventory–observer short version with male intimate partner violence offenders. <i>International journal of clinical and health psychology</i> , 18(2):152–161.		
810			
811			
812			
813			
814			
815	Georgiana Shick Tryon, Sasha Collins Blackwell, and Elizabeth Felleman Hammel. 2007. A meta-analytic examination of client–therapist perspectives of the working alliance. <i>Psychotherapy Research</i> , 17(6):629–642.		
816			
817			
818			
819			
820	Patrick E Shrouf and Joseph L Fleiss. 1979. Intraclass correlations: uses in assessing rater reliability. <i>Psychological bulletin</i> , 86(2):420.		
821			
822			
823	Zachary Steel, Claire Marnane, Changiz Iranpour, Tien Chey, John W Jackson, Vikram Patel, and Derrick Silove. 2014. The global prevalence of common mental disorders: a systematic review and meta-analysis 1980–2013. <i>International journal of epidemiology</i> , 43(2):476–493.		
824			
825			
826			
827			
828			
829	Michael Tanana, Kevin Hallgren, Zac Imel, David Atkins, Padhraic Smyth, and Vivek Srikumar. 2015. Recursive neural networks for coding therapist and patient behavior in motivational interviewing . In <i>Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality</i> , pages 71–79, Denver, Colorado. Association for Computational Linguistics.		
830			
831			
832			
833			
834			
835			
836			
	Victoria Tichenor and Clara E Hill. 1989. A comparison of six measures of working alliance. <i>Psychotherapy: Theory, Research, Practice, Training</i> , 26(2):195.	837	838
			839
	Georgiana Shick Tryon, Sasha Collins Blackwell, and Elizabeth Felleman Hammel. 2008. The magnitude of client and therapist working alliance ratings. <i>Psychotherapy: Theory, Research, Practice, Training</i> , 45(4):546.	840	841
			842
			843
			844
	Raphael Vallat. 2018. Pingouin: statistics in python . <i>Journal of Open Source Software</i> , 3(31):1026.	845	846
	Steven Walfish, Brian McAlister, Paul O’Donnell, and Michael J Lambert. 2012. An investigation of self-assessment bias in mental health providers. <i>Psychological reports</i> , 110(2):639–644.	847	848
			849
			850
	Xiangdong Wang, Xilin Wang, and Hong Ma. 1999. Manual for the mental health rating scale. <i>Chinese Mental Health Journal</i> , 13(1):31–35.	851	852
			853
	Yaqing Wang, Jiepu Jiang, Mingyang Zhang, Cheng Li, Yi Liang, Qiaozhu Mei, and Michael Bendersky. 2023. Automated evaluation of personalized text generation using large language models. <i>arXiv preprint arXiv:2310.11593</i> .	854	855
			856
			857
			858
	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022a. Emergent abilities of large language models. <i>arXiv preprint arXiv:2206.07682</i> .	859	860
			861
			862
			863
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.	864	865
			866
			867
			868
	Zixiu Wu, Rim Helaoui, Diego Reforgiato Recupero, and Daniele Riboni. 2023. Towards effective automatic evaluation of generated reflections for motivational interviewing. In <i>Companion Publication of the 25th International Conference on Multimodal Interaction</i> , pages 368–373.	869	870
			871
			872
			873
			874
	Jing Xu, Arthur Szlam, and Jason Weston. 2022. Beyond goldfish memory: Long-term open-domain conversation . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5180–5197, Dublin, Ireland. Association for Computational Linguistics.	875	876
			877
			878
			879
			880
	A Guidelines		881
	A.1 Guidelines Development Process		882
	To facilitate the understanding of questions and the differentiation of scores by observers, we have four developers (One is a postdoctoral fellow in psychology and a State-Certificated Class 3 Psycho-counselor with 4 years of experience; another is a master student in psychology; and the remaining two are a doctoral student and a postdoctoral	883	884
			885
			886
			887
			888
			889

fellow in the intersection of computer science and psychology.) to carefully design specific guidelines for each score associated with each question. Following Darchuk et al. (2000)'s work, we employ the amount of evidence present in counseling conversations as anchor labels for scores, using the middle point (i.e., 3) as the start point representing "no evidence". The higher score denotes more positive evidence, and vice versa. As a result, each question is scored from 1 to 5.

Expanding on the general guidelines, we further design specific descriptions for each score of every question. Here, we introduce the detailed descriptions by taking the question "There is agreement about the usefulness of the current activity in therapy (i.e., the client is seeing new ways to look at his/her problem)" as an example. Firstly, we anchor the extreme scores of the scale with bipolar adjective relevant to this question, resulting in "open claim of useless" at a rating of 1 and "overt statements of usefulness" at a rating of 5. Secondly, we outline counselors and clients' behavioral indicators at each score level, along with the corresponding extent and frequencies. For the exemplar question, the descriptions are formulated based on clients' frequency (always or sometimes) and attitude (actively or passively) towards participating in tasks proposed by counselors.

The resulting guidelines establish conceptual boundaries among questions within the same dimension and provide clear distinctions among the points on the scale, allowing raters to discern subtle changes in the working alliance with greater reliability.

A.2 Detailed Guidelines

— Goal —

Q1: There are doubts or a lack of understanding about what participants are trying to accomplish in therapy.

1 = The counselor or the client explicitly mentions the counseling goals and works around the established objectives, such as understanding information related to the goals and methods to achieve them. The relevance of the dialogue to the goals is evident for both the counselor and the client. They may discuss the goals to acknowledge or comment on the usefulness of the therapeutic process.

2 = The counselor and the client do not explicitly mention the goals but are working towards a common objective. The counselor addresses the client's

concerns immediately and adjusts the therapeutic process to meet the client's needs. The client is satisfied with the progress made.

3 = There is no evidence to suggest that the counselor and the client have established consistent counseling goals, or there is an equal level of confusion and understanding regarding the goals.

4 = There is disagreement between the counselor and the client regarding counseling goals. While there may be some communication between both parties, the counselor's specific tasks or interventions may be questioned or resisted by the client. The counseling may need to be paused multiple times to adjust the goals. The client may express overall dissatisfaction with the counseling. At this stage, the counselor may take on an "expert" role, sometimes overlooking the client's opinions or therapeutic ideas, and instances where the counselor guides but the client is not engaged may occur. The client may become less emotionally invested.

5 = The counselor and the client have clearly identified different goals, and there are disagreements in the order of issues and solutions in therapy. This inconsistency may lead the client to express strong dissatisfaction with the overall counseling process and goals, possibly mentioning the reasons for participating in therapy. This could further trigger a negative reaction from the counselor. At this stage, it seems challenging for both parties to find common ground, making the therapeutic process difficult.

Q2: The client and therapist are working on mutually agreed upon goals.

1 = The shift of topics often occurs abruptly, usually without mutual agreement from both parties. This frequent topic shift may result from one party interrupting or disregarding the other's statements. At this stage, significant conflicts exist between the counselor and the client regarding the appropriateness, definition, and boundaries of the goals, leading to confusion in the rhythm and content of the conversation.

2 = Topics may shift before resolution or conclusion, but the transition typically moves from one relevant topic to another related or less related one. This shift can be initiated by either the counselor or the client. At this stage, both parties may express dissatisfaction with the frequent shift of topics or the overall pace of therapy, but friction is relatively minor and has not escalated into apparent conflict.

3 = There may be some ambiguity or uncertainty between the counselor and the client regarding ses-

992	sion goals. The current stage of communication	1044
993	lacks clear evidence that both parties have reached	1045
994	a common understanding or collaboration, but there	1046
995	is also no explicit conflict or disagreement. Further	1047
996	communication and discussion may be necessary	1048
997	to clarify expectations and goals to ensure the ef-	1049
998	fectiveness of therapy.	
999	4 = The counselor and the client have made some	1050
1000	progress through discussing relevant topics, but	1051
1001	there may still be a small amount of disagreement	1052
1002	or areas that need further exploration. At this stage,	1053
1003	although both parties generally agree on the current	1054
1004	direction and topics of therapy, more communica-	1055
1005	tion and consensus may be needed to ensure the	1056
1006	achievement of goals.	1057
1007	5 = The counselor and the client have achieved	1058
1008	complete agreement on goals through in-depth, tar-	1059
1009	geted discussions, and have had highly productive	1060
1010	discussions on multiple related topics. At this stage,	1061
1011	both parties almost always reach consensus on the	1062
1012	current topic identified by the client as a goal and	1063
1013	then smoothly transition to another relevant topic.	1064
1014	The overall session and communication are very	1065
1015	smooth and efficient.	1066
1016	Q3: The client and therapist have different ideas	1067
1017	about what the client's real problems are.	1068
1018	1 = The counselor and the client have a very clear	1069
1019	and consistent understanding of the client's issues	1070
1020	and goals. At this stage, there is a strong consen-	1071
1021	sus on problem resolution, with both parties often	1072
1022	identifying the same issues and considering therapy	1073
1023	sessions highly effective. This indicates that they	1074
1024	have formed a close collaborative relationship in	1075
1025	the session.	1076
1026	2 = The counselor and the client have a certain	1077
1027	level of consensus on the client's issues and goals.	1078
1028	While not fully synchronized like the first category,	1079
1029	both parties are making efforts to understand each	1080
1030	other and demonstrate open and cooperative atti-	1081
1031	tudes in discussions. This indicates that they are	1082
1032	working towards establishing a common therapeu-	1083
1033	tic direction and goals.	1084
1034	3 = In the communication between the counselor	1085
1035	and the client regarding the client's issues, there	1086
1036	is no clear evidence of agreement or disagreement.	1087
1037	In the current interaction, there may be neither a	1088
1038	clear consensus nor explicit conflict in opinions and	1089
1039	feelings on both sides. Further communication and	1090
1040	discussion may be needed to clarify the positions	1091
1041	and expectations of both parties.	1092
1042	4 = There is some disagreement between the	1093
1043	counselor and the client regarding the client's is-	1094
	ues. This disagreement may manifest as contro-	1095
	versy in response to certain topics or differences	
	in the relevance of counseling goals. At this stage,	
	although there may be occasional confrontations in	
	the interaction between the two, it has not escalated	
	to strong opposition or sustained conflict.	
	5 = There is evident conflict and disagreement	
	between the counselor and the client in defining	
	and addressing the client's issues. The client may	
	strongly oppose the counselor's viewpoints, and	
	the counselor may shift topics, frequently inter-	
	rupt, and express disagreement with the client's	
	perspectives. At this stage, there may be clear con-	
	frontations in the interaction between both parties,	
	leading to a compromised effectiveness of the ses-	
	sion.	
	Q4: The client and therapist have established a	
	good understanding of the changes that would	
	be good for the client.	
	1 = There are clear misunderstandings and dis-	
	agreements between the counselor and the client in	
	the process of change. The client may express con-	
	cerns or doubts about the direction of their change,	
	the expected outcomes of the change, or the meth-	
	ods of change suggested by the counselor. At this	
	stage, more communication and guidance may be	
	needed to build trust and understanding.	
	2 = The client may have doubts or uncertainties	
	in the process of change. Although they may be	
	taking some actions or practices, it is not clear	
	how to achieve the expected change or the actual	
	effectiveness of these practices. The counselor and	
	the client need to further explore and clarify the	
	path and expectations of change.	
	3 = The counselor and the client have a neutral	
	attitude towards the process and goals of change	
	in the conversation. Both parties may not have ex-	
	PLICITLY expressed their understanding or misunder-	
	standing of the change. Expectations and methods	
	of change are neither emphasized nor overlooked	
	in the discussion, resulting in an overall lack of	
	clear consensus or disagreement on the goals and	
	process of counseling.	
	4 = Both the counselor and the client in the con-	
	versation are aware of changes that would benefit	
	the client. This understanding may be reflected in	
	the client's compromise on counseling goals, ex-	
	pressions, or discussions about the client's current	
	situation and future expectations. Both parties are	
	working to clarify the path and direction of change.	
	5 = In the counseling process, there is strong	
	consistency and clarity between the counselor and	

1096 the client regarding the client's goals and how to
1097 achieve them. They not only discuss these goals
1098 frequently and explicitly during the session but also
1099 summarize and confirm the progress and outcomes
1100 achieved at the end. The interaction and discussion
1101 at this stage align completely with the therapeutic
1102 plan.

1103 — Approach —

1104 **Q5: There is agreement about the steps taken to**
1105 **help improve the client's situation.**

1106 1 = The client directly expresses that the tasks
1107 and goals are inappropriate and generally disagrees
1108 with homework or tasks during the session. There
1109 is a disagreement between the client and the coun-
1110 selor regarding the approach to be taken. The client
1111 refuses to engage in tasks.

1112 2 = The client hesitates to explore and does not
1113 follow the counselor's guidance in the change pro-
1114 cess. The client withdraws from the counselor,
1115 seeming to just "go through the motions," not en-
1116 gaging or focusing on the counselor or tasks. Even
1117 after some clarification by the counselor, the client
1118 still seems uncertain about the relevance of the
1119 tasks to their goals. The client appears conflicted
1120 or indifferent towards tasks in therapy and passively
1121 resists them (e.g., limited participation).

1122 3 = There is no clear consensus or disagreement
1123 between the counselor and the client regarding ther-
1124 apy tasks. Both may have vague views on the sig-
1125 nificance and purpose of tasks, resulting in a neutral
1126 attitude towards participation and involvement in
1127 tasks during the session.

1128 4 = The client shows a clear interest and involve-
1129 ment in therapy tasks. Whether occasional clar-
1130 ification is needed or not, the client participates
1131 and follows the exploration process. There is an
1132 unspoken understanding behind the tasks, leading
1133 the client to gradually acknowledge and engage in
1134 the tasks.

1135 5 = The counselor and client strongly agree on
1136 different goals, and there is a clear disagreement
1137 on the order and solutions to issues in therapy. This
1138 inconsistency may lead the client to express strong
1139 dissatisfaction with the overall therapy process and
1140 goals, possibly mentioning the reasons for attend-
1141 ing therapy, which may further trigger a negative
1142 reaction from the counselor. At this stage, finding
1143 common ground seems challenging, making the
1144 therapy process difficult.

1145 **Q6: There is agreement about the usefulness of**
1146 **the current activity in therapy (i.e., the client is**

seeing new ways to look at his/her problem).

1147 1 = The client repeatedly argues against tasks.
1148 The client refuses to participate, claiming that it
1149 is pointless for their goals. Tension exists in the
1150 relationship between the counselor and the client,
1151 and issues are not explored. 1152

1153 2 = The client does not actively engage in the
1154 session tasks, although he/she may not openly ques-
1155 tion the usefulness of the tasks. The client fails to
1156 openly discuss the issues. The client may hesi-
1157 tate to participate in tasks but eventually engages
1158 in them. The counselor accurately conveys the
1159 reasons behind the tasks, enabling the client to un-
1160 derstand the relevance of the tasks to their current
1161 concerns. 1162

1163 3 = There is no clear evidence in the communi-
1164 cation between the counselor and the client about
1165 whether they have reached an agreement or dis-
1166 agreement on the client's issues. In the current
1167 interaction, there is neither a clear consensus nor
1168 an explicit conflict in opinions and feelings. Fur-
1169 ther communication and discussion may be needed
1170 to clarify their positions and expectations. 1171

1172 4 = The client actively participates in and is
1173 committed to therapy tasks, showing no skepticism
1174 about their effectiveness. Regardless of occasional
1175 resistance, the client engages and follows the ex-
1176 ploration process. Both parties share a common
1177 understanding of the tasks' principles, allowing
1178 the client to gradually accept and participate in the
1179 tasks. 1180

1181 5 = In the counseling process, the counselor and
1182 the client have a strong and clear agreement on
1183 the client's goals and how to achieve them. They
1184 not only frequently and explicitly discuss these
1185 goals during the session but also summarize and
1186 confirm the progress and achievements at the end.
1187 The interaction and discussion at this stage align
1188 completely with the therapeutic plan. 1189

1190 **Q7: There is agreement on what is important**
1191 **for the client to work on.**

1192 1 = There is a clear disagreement and opposi-
1193 tion between the counselor and the client regarding
1194 the current focus. This difference may manifest
1195 as the counselor not allowing the client to shift to
1196 different topics or the client showing strong oppo-
1197 sition during the therapy process. Their views on
1198 the direction and outcomes of therapy are entirely
1199 different. 1200

1201 2 = The counselor and the client have some dis-
1202 agreement about the content and direction of ther-
1203 apy, differing in the themes and time allocation to
1204 1205

1199	focus on during therapy.		1251
1200	3 = There are no clear signs of agreement or	pects of therapy tasks, although this agreement may	1252
1201	disagreement in the interaction between the coun-	not always be explicitly expressed. His/her level of	1253
1202	selor and the client regarding the themes or issues	involvement in the therapy process falls between	1254
1203	of therapy. Although they may engage in some	simple compliance and actively providing sugges-	1255
1204	exploration and communication, it is challenging	tions. The client shows a certain level of agreement	1256
1205	to determine whether they share views on therapy	with the collaboration with the counselor, possibly	1257
1206	themes or issues. Their reactions seem neither par-	being more actively involved in certain aspects of	1258
1207	ticularly synchronized nor explicitly conflicting.	therapy.	1259
1208	4 = The client and the counselor respond to each	5 = The client is satisfied and excited about	1260
1209	other's focus and needs to some extent. They ex-	the counselor's methods and approach to problem-	1261
1210	plore and accept each other's views and intentions	solving. His/her performance in therapy is highly	1262
1211	to some degree. Although there may be some dif-	positive, possibly suggesting suggestions to fur-	1263
1212	ferences, they both strive to seek a common under-	ther advance therapy tasks. Overall, the client is	1264
1213	standing and progress the therapy process.	content with therapy work, and their interaction	1265
1214	5 = The counselor and the client are highly ac-	demonstrates a high level of cooperation and enthu-	1266
1215	tively engaged in the therapy process, thoroughly	siasm.	
1216	exploring each other's issues and responding ex-		
1217	PLICITLY and continuously to each other's views and	— Affective Bond —	1267
1218	intentions. They approach therapy themes and is-	Q9: There is a mutual liking between the client	1268
1219	issues with an open mindset, working together, re-	and therapist.	1269
1220	fecting flexibility, and demonstrating a cooperative	1 = There is evident animosity, hostility, or in-	1270
1221	spirit.	difference between the counselor and the client.	1271
1222	Q8: The client believes that the way they are	This may manifest in arguments, derogatory com-	1272
1223	working with his/her problem is correct.	ments, or open hostility. The counselor fails to	1273
1224	1 = The client holds evident doubts and aver-	demonstrate concern for the client and may either	1274
1225	sions towards the counseling process, frequently	forget important details of their life or completely	1275
1226	engaging in arguments with the counselor. Progress	disregard the client.	1276
1227	between the counselor and the client is very lim-	2 = Although there is no direct hostility between	1277
1228	ited, and the time spent arguing may exceed the	both parties, there is noticeable tension and dis-	1278
1229	time dedicated to therapy. This inconsistency and	tance in the relationship. The counselor appears	1279
1230	questioning impact the overall therapy process.	indifferent or mechanical in response to the client,	1280
1231	2 = The counselor and the client sometimes have	lacking enthusiasm. While there may not be ex-	1281
1232	conflicting opinions, but they seem to cooperate	PLICIT negative language, there is a lack of positive	1282
1233	in certain parts of the therapy process. The client	feedback and reinforcement in their interactions.	1283
1234	expresses doubts about the therapy process or oc-	3 = There are no clear signs of warmth or cold-	1284
1235	casionally expresses concerns about certain tech-	ness in the relationship between the counselor and	1285
1236	niques, finding other things to do during most of	the client. Communication lacks strong emotional	1286
1237	the counseling time.	feedback, and both parties seem to maintain a neu-	1287
1238	3 = The client maintains a neutral stance toward	tral stance. Despite engaging in communication,	1288
1239	the therapy process and methods. He/she neither	there is no clear expression or implication of liking	1289
1240	explicitly expresses satisfaction nor dissatisfaction	or disliking each other. The relationship appears	1290
1241	with therapy, nor does he/she clearly indicate agree-	balanced without significant signs of warmth or	1291
1242	ment or disagreement with the therapeutic methods.	indifference.	1292
1243	During the therapy process, the client may com-	4 = In the majority of the sessions, the coun-	1293
1244	ply at certain moments and show reservations at	selor and the client have positive interactions. The	1294
1245	other times, without providing a clear evaluation	counselor shows enthusiasm and care for the client,	1295
1246	of the therapy's effectiveness. This neutral attitude	frequently communicating with empathy and en-	1296
1247	may stem from the client's ongoing assessment of	couragement, exploring and understanding impor-	1297
1248	therapy effectiveness or uncertainty about how to	tant details of the client's life.	1298
1249	evaluate therapy progress.	5 = Throughout the therapy process, the coun-	1299
1250	4 = The client partially agrees with certain as-	selor and the client consistently demonstrate a deep	1300
		care for each other and provide positive feedback.	1301

The counselor not only encourages and reinforces the client's healthy behaviors but also deeply understands and cares about various aspects of the client's life, including their interests and hobbies. This profound care may lead to the client explicitly expressing gratitude and trust in the counselor. The client may also show appreciation for the counselor's care.

Q10: The client feels confident in the therapist's ability to help the client.

1 = The client expresses minimal or no hope for the therapy outcomes. The client significantly questions the therapist's capabilities and may directly challenge the therapist's qualifications or understanding of the client's experiences. The client resists the therapist's suggestions, attempts at assistance, or expresses discouragement and pessimism.

2 = The client harbors doubts about the therapist, the therapy process, or the anticipated outcomes. The client may question whether the therapist truly understands their issues or doubt the interventions/homework provided during the problem-solving stages. These doubts do not come with strong opposition or hostility but noticeably impact the progress of the therapy process.

3 = The client holds a neutral stance regarding the therapist's capabilities. Throughout the therapy process, there is no clear evidence suggesting that the client has high confidence in the therapist, nor is there evidence indicating skepticism about the therapist's abilities. The client's responses and comments neither explicitly appreciate nor question the therapist's skills and capabilities.

4 = The client expresses a certain level of confidence in the therapist's abilities. This confidence may be reflected in the client's in-depth discussions on therapy topics, positive responses to the therapist's guidance, or an optimistic attitude towards resolving current counseling issues. Additionally, the client has substantial trust in the therapist's competency, possibly expressing appreciation for the effectiveness of the therapy or the therapist's abilities.

5 = The client consistently agrees with the therapist's reflections and interventions/guidance, expressing high satisfaction and appreciation for certain aspects of the therapy process or the therapist themselves. There may be multiple discussions during the therapy process highlighting the strengths of the therapy and/or the therapist.

Q11: The client feels that the therapist appreciates him/her.

1 = The client feels that the therapist is indifferent, inattentive, and unconcerned about his/her issues. This is expressed through explicit accusations, disdain, or other negative reactions, indicating a sense of being disregarded or misunderstood by the therapist.

2 = The client harbors some doubts about whether the therapist genuinely cares. These doubts might be indirectly expressed, such as subtle mentions or manifestations of emotions like withdrawal, displeasure, or frustration.

3 = Throughout the therapy process, there is no clear evidence of strong positive or negative reactions from the client regarding the therapist's care and support. The client neither explicitly appreciates nor expresses dissatisfaction or disregard for the therapist's sensitivity and empathetic abilities. The emotional tone of the relationship is neutral, with no apparent strong connection or distance.

4 = The therapist demonstrates a level of acceptance, warmth, and empathy towards the client, and the client perceives and responds to this caring attitude. During the therapy process, the client acknowledges to some extent the therapist's warmth and understanding.

5 = The client strongly senses the therapist's care and support, expressing gratitude for the relationship. They may praise the therapist's sensitivity and empathetic abilities, feeling comfortable and at ease for most of the therapy process.

Q12: There is mutual trust between the client and therapist.

1 = The client has significant mistrust towards the therapist, demonstrated by avoiding discussions on critical issues or directly expressing distrust. This mistrust hinders open communication, and the therapist may also show concerns and discomfort about the therapeutic process.

2 = There is a moderate level of mistrust between both parties, though not as intense as in the first category. The client may hesitate to share private content, and the therapist may feel a sense of uncertainty or slight discomfort regarding the therapeutic situation.

3 = There are no clear signs of trust between the therapist and client, but there are also no apparent behaviors indicating mistrust. There is a balance between trust and mistrust in their interactions, with no explicit demonstration of reliance on each other, nor clear signs of doubt or guardedness.

4 = The client is willing to disclose some personal concerns, and the therapist accepts the

	Before	After
Q1	0.6785	0.6835
Q2	0.8297	0.8341
Q3	0.7337	0.7381
Q4	0.7906	0.8061
Goal	<i>0.7581</i>	0.7655
Q5	0.6034	0.6034
Q6	0.6645	0.6750
Q7	0.6055	0.6398
Q8	0.7612	0.7612
Approach	<i>0.6587</i>	0.6699
Q9	0.6455	0.6906
Q10	0.7124	0.7396
Q11	0.617	0.6357
Q12	0.6241	0.6970
Affective Bond	<i>0.6498</i>	0.6907
Overall	<i>0.6888</i>	0.7087

Table 6: Human agreement on evaluating the working alliance across all dimensions and questions before and after the refinement of weak annotators.

client’s surface statements. The therapist does not overturn or interrupt the client’s thoughts and maintains focus.

5 = The trust between both parties is deep enough that the client not only willingly shares deeper layers of privacy and issues but also accepts and responds to the therapist’s feedback and suggestions. This level of trust enhances the overall smoothness and efficiency of the therapeutic process.

B Human Evaluation

B.1 Human Agreement

Table 6 shows human agreement in evaluating working alliance across all dimensions and questions during the initial annotation phase and after refinement based on evidence generated by GPT-4. Given that we plan to generalize our reliability results to any annotators with similar characteristics as the selected raters in this work, focus on the absolute agreement instead of consistency between annotators, and use the mean value of three annotators as an assessment basis, we adopt the ICC(2, k) form with two-way random effects, absolute agreement, and multiple raters. We use Pingouin package (Vallat, 2018) to calculate the ICC metric.

Besides, Table 5 presents the proportion of revisions made by each annotator during the process of refining labels with evidence from GPT-4.

Dimension	Avg. Score	Question	Avg. Score
Goal	3.57(0.56)	Q1	3.56(0.63)
		Q2	3.69(0.60)
		Q3	3.56(0.67)
		Q4	3.47(0.64)
Approach	3.52(0.56)	Q5	3.46(0.61)
		Q6	3.32(0.64)
		Q7	3.75(0.63)
		Q8	3.57(0.55)
Affective Bond	3.60(0.48)	Q9	3.67(0.55)
		Q10	3.37(0.63)
		Q11	3.39(0.42)
		Q12	3.97(0.52)

Table 7: The average scores annotated on each question and dimension, with standard deviations presented in parentheses. The highest average score in each column is shown in bold.

B.2 Data Characteristics

Based on the annotated data, we analyze the score distribution. Table 7 presents the average scores per dimension and questions along with their standard deviations in parentheses.

C Experimental Results and Analysis

C.1 Biases of Counselors’ Retrospective Self-reports

We employ a paired t-test (Kim, 2015) implemented in scikit-learn package (Buitinck et al., 2013) to examine the relationship between the working alliance scores rated by counselors and clients within each counselor-client pair. We observe that counselors tend to assign significantly higher scores than their clients in nearly 21% of counseling sessions. In these instances, counselors’ average ratings stand at 4.38 ± 0.57 , contrasting sharply with clients’ average rating of only 2.84 ± 0.60 . This disparities demonstrate significant distinctions in how counselors perceive the overall relationship with their clients compared to the assessments provided by the clients themselves.

C.2 Model Self-Agreement

As the final annotation is determined by the average of the model’s three independent annotations, we adopt the intraclass correlation coefficient with the 2-way mixed-effects model, absolute agreement definition, and the mean of k measurements type as the measure of the model’s self-reliability (Koo and Li, 2016; Shrout and Fleiss, 1979). Table 8 presents models’ intra-rater agreement on evaluating all the questions.

	ChatGPT		GPT4		
	Detailed Guideline + CoT	No Guideline	General Guideline	Detailed Guideline	Detailed Guideline + CoT
Q1	0.2921	0.5359	0.4136	0.7210	0.7111
Q2	0.5972	0.4327	0.6193	0.6884	0.6978
Q3	0.0195	0.5935	0.5174	0.5368	0.6432
Q4	0.6179	0.7516	0.8716	0.8500	0.8086
<i>Goal</i>	<i>0.3817</i>	<i>0.5784</i>	<i>0.6055</i>	<i>0.6991</i>	<i>0.7152</i>
Q5	0.7828	0.8674	0.8806	0.7648	0.7424
Q6	0.4448	0.6768	0.8188	0.7137	0.6580
Q7	0.5755	0.4903	0.7278	0.4286	0.7158
Q8	0.6710	0.8279	0.8218	0.8115	0.7885
<i>Approach</i>	<i>0.6185</i>	<i>0.7156</i>	<i>0.8123</i>	<i>0.6797</i>	<i>0.7262</i>
Q9	0.7148	0.8439	0.9232	0.5222	0.5449
Q10	0.6225	0.6476	0.7942	0.7920	0.7786
Q11	0.4708	0.6716	0.8913	0.7175	0.8117
Q12	0.4418	0.6849	0.6992	0.6781	0.7456
<i>Affective Bond</i>	<i>0.5625</i>	<i>0.7120</i>	<i>0.8270</i>	<i>0.6775</i>	<i>0.7202</i>
<i>Total</i>	<i>0.5209</i>	<i>0.6687</i>	<i>0.7482</i>	<i>0.6854</i>	<i>0.7205</i>

Table 8: The intrarater reliability of models in evaluating each question and dimension across different experimental settings.

C.3 Alignment with Human Evaluations

The alignment between LLMs and human evaluations are presented in Table 9.

D The Consent Form and User Services Agreement

Below are the English translation of consent forms and user services agreement used in the current work, the original documents are in Mandarin Chinese. Every client gave their consent to attend the online text-based psycho-counseling on our counseling platform and agreed to data usage for the current work.

D.1 Consent Form

Dear clients,

Thank you for your trust. Before we formally begin the counselings, there are some relevant matters that need to be communicated to you, so that the consultation can proceed smoothly and effectively. This agreement is the basic framework to ensure the normal conduct of the psychological consultation process. Please read it carefully and tick the box at the bottom to indicate your agreement. If you have any questions, please raise them with your counselor after the counselings.

1. Duration and Frequency of Consultation: Psychological consultations require regular sessions, each typically lasting 50 minutes. The frequency and total duration of the consultations will be jointly determined by you and your counselor based on the nature of your psychological distress and personal needs.

2. Confidentiality and Exceptions to Confidentiality: In general, your counselor will keep the information you provide confidential, including case records, test materials, letters, recordings, videos, and other materials, all of which are considered professional information and are stored under strict confidentiality to prevent public disclosure in any public setting. However, there are exceptions to confidentiality in the following cases, and relevant individuals and institutions will be notified:

1) Violation of relevant laws (e.g., if you pose a danger to others; suspicion of child or elder abuse or abuse of someone dependent on you for care, etc.)

2) If your situation endangers your own safety (e.g., suicide, self-harm, mental illness, severe depression, etc.), we will notify your relatives or guardians when necessary and consult your opinion to ensure your safety.

3) Counselors need to receive supervision during their work. Counselors will discuss parts of the consultation content and visitor information in personal supervision and case discussions. Privacy information unrelated to the consultation, such as personal names and regions, will be anonymized; supervisors and case discussion members are also bound by the aforementioned confidentiality rules. If there is a need to publicly release or publish consultation details, the visitor's written consent must be obtained first.

3. Adjusting Consultation Times: If you wish to adjust your consultation time, please do so at least 24 hours in advance on the platform. Adjustments cannot be made if the time limit is exceeded.

	ChatGPT	GPT4			
	Detailed Guideline + CoT	No Guideline	General Guideline	Detailed Guideline	Detailed Guideline + CoT
Q1	0.1139	0.2406*	0.3012**	0.5379***	0.4292***
Q2	0.2877*	0.3423**	0.3698***	0.4712***	0.5379***
Q3	0.2430*	0.3869***	0.2920**	0.4907***	0.4510***
Q4	0.1570	0.4667***	0.3651***	0.4919***	0.5569***
<i>Goal</i>	<i>0.2004</i>	<i>0.3591</i>	<i>0.3320</i>	<i>0.4979</i>	<i>0.4938</i>
Q5	0.4222***	0.5710***	0.6423***	0.5618***	0.6025***
Q6	0.3599**	0.5237***	0.6190***	0.5065***	0.5371***
Q7	0.3392**	0.3764***	0.2921**	0.5341***	0.4924***
Q8	0.3233**	0.2439*	0.2532*	0.5898***	0.5472
<i>Approach</i>	<i>0.3612</i>	<i>0.4288</i>	<i>0.4517</i>	<i>0.5481</i>	<i>0.5448</i>
Q9	0.3752***	0.0106	0.1325	0.2337*	0.3086**
Q10	0.4273***	0.5164***	0.6339***	0.5114***	0.4520***
Q11	0.4570***	0.4994***	0.3874***	0.6113***	0.6103***
Q12	0.3892***	0.4506***	0.4305***	0.4101***	0.4960***
<i>Affective Bond</i>	<i>0.4122</i>	<i>0.3693</i>	<i>0.3961</i>	<i>0.4416</i>	<i>0.4667</i>
<i>Total</i>	<i>0.3246</i>	<i>0.3857</i>	<i>0.3933</i>	<i>0.4959</i>	<i>0.5018</i>

Table 9: Pearson correlation between human and model annotations on each dimension and question. Statistic significance levels for individual question correlations are denoted by *** $p < 0.001$, ** $p < 0.01$, and * $p < 0.05$. The overall and dimension-specific correlations are calculated as the averages of the correlations on corresponding questions.

4. Handling of Lateness: You may enter the counseling from the start of the scheduled appointment until it ends, but the end time of the consultation will not be extended due to your lateness. If you are late and do not log in to start the consultation by the service end time, the consultation will be considered expired, and the consultation fee will not be refunded.

5. Responsibilities of the Clients: During the consultation process, visitors need to:

1) Attend and participate in the consultation sessions;

2) Express and share their thoughts and feelings as much as possible during the consultation;

3) Seriously reflect on their own expressions, the counselor's responses, and the interaction process between the two.

6. Responsibilities of the Counselor: Counselors need to:

1) Arrange a suitable consultation schedule for both parties;

2) Strive to guide visitors towards an understanding of themselves and their current situation, and help them better deal with the various difficulties and life events they are facing;

3) Regularly participate in professional learning and case discussions to ensure their competence in counseling work with visitors;

4) Be aware of their limitations as a counselor and discuss ending the consultation or referrals with the visitor in a timely manner if the consultation is ineffective or unsuccessful.

7. Duration and Frequency of Consultation:

1) Psychological consultations are regular sessions, typically 50 minutes each, once a week. Changes to the interval and frequency will be determined based on the nature of your psychological issues and personal needs.

2) Consultation sessions will start and end on time. Flexibility in timing will not exceed 5 minutes.

8. Emergency Consultation: In urgent situations, you may make a temporary appointment or call the local crisis intervention hotline.

9. Crisis Intervention Measures: In the event that you are experiencing severe psychological stress with thoughts of suicide and impulses, it is necessary to discuss potential risks and coping strategies with a counselor. This includes how to access local support resources and techniques for self-regulation. Due to the limitations of remote counseling, counselors may be unable to work with visitors at high risk of suicide. In cases of intense suicidal urges or self-destructive behavior, counselors are obligated to discuss referral to appropriate assistance agencies. (National 24-Hour Suicide Intervention Hotline: 4001619995)

10. Physical symptoms and psychological symptoms often interact, and if necessary, we may discuss the need for consultation and treatment in medical institutions during counseling. Additionally, medication can be beneficial at the appropriate time in alleviating both physical and mental issues. Throughout the treatment process, based on your

1595 specific situation, the counselor may recommend
1596 relevant laboratory and instrumental examinations,
1597 providing detailed explanations as needed.

1598 11. Psychological counseling and therapy are
1599 complex processes that may require coordination,
1600 continuous goal adjustment, or referrals and other
1601 interventions during the course.

1602 12. Voluntary Withdrawal: You have the right
1603 to terminate your counseling at any time, but it is
1604 recommended to discuss and carefully conclude
1605 with your counselor before formal withdrawal.

1606 13. If there are other research and teaching mat-
1607 ters that require your participation, your counselor
1608 will inform you and negotiate with you to sign an
1609 additional written agreement.

1610 14. During the period of the consultation work,
1611 if there is a need to adjust or modify the agreement,
1612 both parties can propose it during the consultation.
1613 After thorough discussion and agreement, corre-
1614 sponding changes will be made.

1615 **Remote/Online Counseling Additional Matters:**

1616 When conducting online counseling, identity ver-
1617 ification is required. For this purpose, you need to
1618 provide some materials (such as personal infor-
1619 mation, current situation, etc.) to complete this
1620 process.

1621 For situations not suitable for online counsel-
1622 ing, such as suicidal or homicidal thoughts, life-
1623 threatening circumstances, a history of suicidal,
1624 abusive, or violent tendencies, hallucinations, and
1625 substance or alcohol abuse, it is recommended to
1626 consider face-to-face counseling or alternative in-
1627 tervention methods.

1628 Considering the potential impact on the counsel-
1629 ing relationship, please refrain from recording au-
1630 dio or video during the counseling process. If there
1631 is a genuine need for such recordings, it should
1632 be discussed thoroughly and agreed upon by both
1633 parties.

1634 The smooth conduct of online counseling de-
1635 pends on stable network conditions, communica-
1636 tion devices, and a disturbance-free room. Please
1637 ensure that you are adequately prepared before
1638 starting online counseling. Additionally, be psy-
1639 chologically prepared for unforeseen events such
1640 as network interruptions during online counseling.

1641 [] I fully understand and agree to the above
1642 terms.

1643 **D.2 Informed Consent Form in the User** 1644 **Services Agreement**

1645 VI. Informed Consent

1646 6.1 To protect your rights, please read and agree
1647 before activating the dialogue service of this appli-
1648 cation: Users agree to accept the online text coun-
1649 seling or venting services (hereinafter referred to as
1650 the service) provided by this application based on
1651 my confusions. Users understand that the current
1652 service provided by this application is AI-assisted
1653 psychological counseling/venting, with real human
1654 counselors also providing services. Users need to
1655 understand that the online text venting/counseling
1656 service is an internet-based form of instant psy-
1657 chological confusion resolution and psychological
1658 knowledge popularization service. This service is
1659 provided in Chinese. Users need to understand
1660 that the service content includes support and help
1661 for psychological confusions (including, but not
1662 limited to: emotional issues, relationship issues,
1663 family relations, interpersonal relationships, per-
1664 sonal growth, career development, etc.). Although
1665 it is difficult to guarantee a complete improvement
1666 in psychological conditions and resolution of con-
1667 fusions, we serve you with the attitude of "some-
1668 times curing, often helping, always comforting".
1669 Users need to understand that during the service
1670 process: conversations will involve the user's psy-
1671 chological/psychological health and emotional state
1672 among other related information. Users have the
1673 right to privacy in the venting/counseling service,
1674 and the personal information disclosed by users
1675 will, in principle, be kept strictly confidential. At
1676 the same time, the user's right to privacy is pro-
1677 tected and restricted by national laws in terms of
1678 content and scope. Users need to understand, based
1679 on national laws, there are exceptions to the princi-
1680 ple of confidentiality, including but not limited to
1681 the following situations:

1682 1) When the service seeker or others are prepar-
1683 ing or in the process of engaging in actions that
1684 endanger the safety of themselves or others' person
1685 or property;

1686 2) When the service seeker may endanger others
1687 (such as in cases of contagious diseases);

1688 3) When the information disclosed by the ser-
1689 vice seeker involves a minor being or about to be
1690 sexually abused;

1691 4) When the service seeker or others are prepar-
1692 ing or in the process of engaging in actions that
1693 endanger national security or public safety;

1694 5) In cases where data is anonymized for discus-
1695 sions, consultations, or when receiving supervision
1696 and training among consulting members;

1697 6) In cases where data is anonymized for scien-

1698	tific research.	
1699	7) When disclosure is required by law.	1750
1700	6.2 Users must agree that for the aforementioned	1751
1701	non-confidential situations, for the fundamental	1752
1702	reason of protecting the rights of the user or re-	1753
1703	lated individuals, we may disclose information to	1754
1704	the minimal extent necessary and only within the	1755
1705	necessary scope of personnel. Furthermore, users	1756
1706	must understand that since the counseling service	1757
1707	is conducted over the internet, although we strive	1758
1708	to protect users' privacy to the greatest extent, it is	1759
1709	difficult to avoid the possibility of personal infor-	1760
1710	mation being leaked due to internet security vulner-	1761
1711	abilities, technical failures, or unauthorized access.	1762
1712	Users must understand that under the following	1763
1713	conditions, we are unable to provide effective vent-	1764
1714	ing/counseling services, and it is necessary to seek	1765
1715	professional offline treatment or counseling ser-	1766
1716	vices:	1767
1717	1. Having thoughts or plans of suicide;	1768
1718	2. Having thoughts or plans of harming oneself	1769
1719	or others;	1770
1720	3. Having any psychiatric disorder diagnosed by	1771
1721	a hospital;	1772
1722	4. Meeting the diagnostic criteria for any psychi-	1773
1723	atric disorder.	1774
1724	Users need to understand that if the physiologi-	1775
1725	cal, psychological, mental state, and behavior plans	1776
1726	described or reflected in their information meet any	1777
1727	of the above criteria, we cannot continue to pro-	1778
1728	vide services to them, and may suggest seeking	1779
1729	professional offline treatment or counseling ser-	1780
1730	vices. Users must understand that this application	1781
1731	provides support and help for psychological con-	1782
1732	fusions (including but not limited to: emotional	1783
1733	issues, relationship issues, family relations, inter-	1784
1734	personal relationships, personal growth, career de-	1785
1735	velopment, etc.), but there still exist some services	1786
1736	that are difficult to provide:	1787
1737	1) Crisis intervention for suicide or other harmful	
1738	behaviors;	
1739	2) Diagnosis and treatment of psychiatric disor-	
1740	ders;	
1741	3) Specific advice on the use of psychiatric med-	
1742	ications;	
1743	4) Dealing with severe psychological trauma;	
1744	5) Providing specific resources or information	
1745	for careers, academics, etc.;	
1746	6) Providing views on social phenomena and	
1747	interpretations of policies;	
1748	7) Interpretation of dreams (e.g., explaining the	
1749	meaning of dreams, why certain people or things	
	appear in dreams, etc.).	
	8) To answer psychological confusions not re-	
	lated to myself (for example, those of my friends,	
	family, online friends, etc.).	
	Users need to understand that when the de-	
	scribed situation exceeds our service scope (which	
	does not include the aforementioned 8 types), we	
	cannot meet their needs. Users need to understand	
	the potential benefits and risks of internet-based	
	text venting/counseling services. The benefits in-	
	clude, but are not limited to, being able to access	
	services more conveniently without the need to	
	travel to a designated location. And, although the	
	risks are small, users still understand that there may	
	be potential risks. These risks include, but are not	
	limited to: due to possibly insufficient information	
	provided by the user, the services received may not	
	fully resolve the user's confusions or improve the	
	user's psychological state; due to possible techni-	
	cal failures or other unforeseen reasons, the user	
	may not receive timely analysis and advice for their	
	psychological confusions. Users must agree that	
	when the application provides services, it follows	
	the laws and regulations of mainland China, not	
	the laws and regulations of the user's location. The	
	above informed consent remains effective during	
	the user's single or multiple uses of the service.	
	6.3. I agree to convert the collected psychologi-	
	cal counseling dialogue text data into digital and	
	graphical forms for use in non-profit academic co-	
	operation, academic conferences, journal publica-	
	tions, and other academic activities by certified	
	third-party academic institutions (*1).	
	(*1) Certified third-party academic institutions	
	refer to universities and research institutes officially	
	recognized by the state, and researchers working	
	within them have undergone formal academic train-	
	ing.	