

# 🧠 Psyche-R1: Towards Reliable Psychological LLMs through Unified Empathy, Expertise, and Reasoning

Anonymous ACL submission

## Abstract

Amidst a shortage of qualified mental health professionals, the integration of large language models (LLMs) into psychological applications offers a promising way to alleviate the growing burden of mental health disorders. Recent reasoning-augmented LLMs have achieved remarkable performance in mathematics and programming, while research in the psychological domain has predominantly emphasized emotional support and empathetic dialogue, with limited attention to reasoning mechanisms that are beneficial to generating accurate responses. Therefore, in this paper, we propose *Psyche-R1*, the first Chinese psychological LLM that jointly integrates empathy, psychological expertise, and reasoning, built upon a novel data curation pipeline. Specifically, we design a comprehensive data synthesis pipeline that produces over 75k high-quality psychological questions paired with detailed rationales, generated through an iterative prompt-rationale optimization procedure, along with 73k empathetic dialogues. Subsequently, we employ a hybrid training strategy wherein challenging samples are identified through a multi-LLM cross-selection strategy for group relative policy optimization (GRPO) to improve reasoning ability, while the remaining data are used for supervised fine-tuning (SFT) to enhance empathetic response generation and psychological domain knowledge. Extensive experiment results demonstrate the effectiveness of *Psyche-R1* across several psychological benchmarks, where our 7B *Psyche-R1* achieves comparable results to 671B DeepSeek-R1.

## 1 Introduction

The shortage of qualified mental health professionals has spurred increasing interest in applying artificial intelligence within the psychological domain to support mental health assistance (Wolohan et al., 2018; Al Asad et al., 2019; Tanana et al., 2021). Recently, large language models (LLMs) have demon-

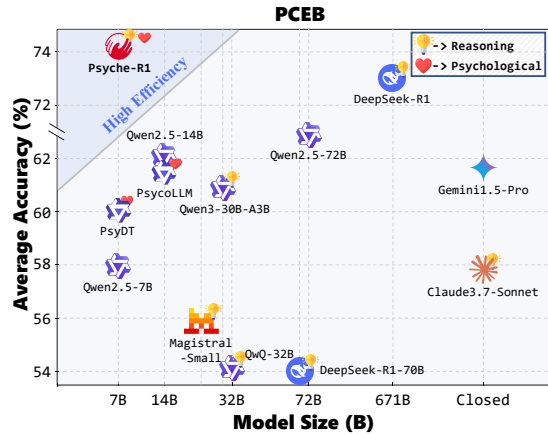


Figure 1: Comparison of different LLMs on the PCEB, plotted by average standard accuracy versus model size.

strated impressive capabilities across a wide range of domains owing to their exceptional text understanding capabilities (Zhang et al., 2023; Naveed et al., 2025). Therefore, many LLM-based studies have been proposed to advancing the mental health services (Cho et al., 2023; Ye et al., 2025).

Prior research has established the critical importance of empathy optimization in psychological counseling (Qiu et al., 2024; Sorin et al., 2024; Zhang et al., 2024). For example, SoulChat (Chen et al., 2023a) enhances empathetic responding by fine-tuning a model on a large-scale, multi-turn empathetic dialogue dataset. Similarly, AUGESC (Zheng et al., 2023) improves emotional sensitivity in dialogue systems by incorporating an emotion-aware attention mechanism. However, these approaches often lack the expertise foundation required for psychology, which is important for accurate psychological understanding. Some studies have attempted to address this limitation through integration of psychological knowledge (Chen et al., 2023b; Xiao et al., 2024; Wu et al., 2025). For example, PsycoLLM (Hu et al., 2024) integrates psychological knowledge by training its model on knowledge-based question-answer (QA) pairs,

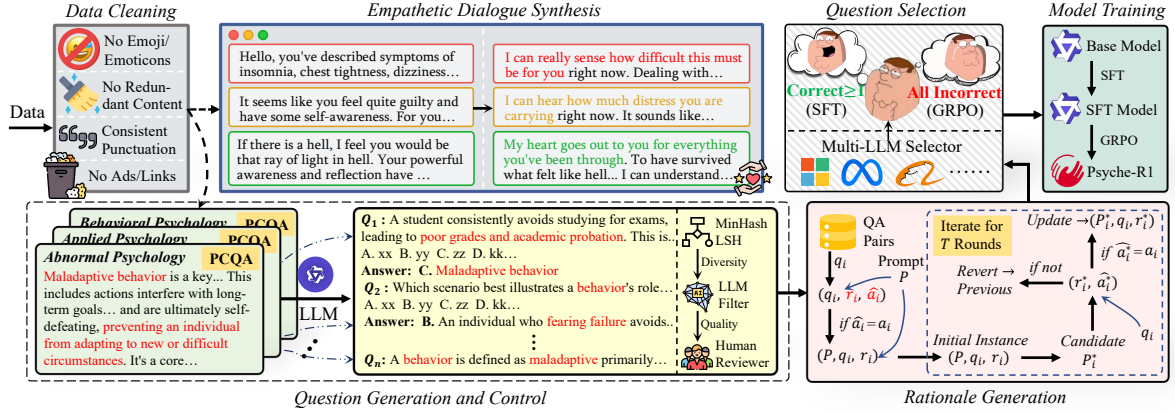


Figure 2: Overview of our proposed pipeline for constructing the dataset and *Psyche-R1*. Our pipeline involves generating psychological reasoning questions paired with detailed rationales, along with empathetic dialogues.

while CPsyExam (Zhao et al., 2025) leverages examination questions covering theoretical knowledge from different psychology-related subjects to further improve model performance. Although existing studies have achieved considerable success, they remain limited in their capacity for complex reasoning. In fact, reasoning-augmented LLMs trained through reinforcement learning (RL) have demonstrated superior performance across various domains, particularly in mathematics, code generation, and medical domain (Chen et al., 2024; Guo et al., 2025). However, as shown in Figure 1, these reasoning-augmented LLMs exhibit limited performance in the psychological domain, since they focus on logic reasoning, while neglecting the unification of empathy and expertise beyond general common sense. In fact, within the psychological domain, reasoning plays a critical role, as it contributes not only to generating more accurate and reliable responses but also to supporting deeper empathetic engagement and more coherent integration of psychological knowledge.

Therefore, in this paper, we introduce *Psyche-R1* that integrates empathy, domain-specific expertise, and reasoning capabilities. To construct a high-quality training corpus, we design a comprehensive data synthesis pipeline that integrates modified empathetic dialogues derived from authentic sources, which capture diverse and empathetic expressions, with knowledge-centric question-answer pairs that encapsulate psychological expertise. Specifically, we apply chain-of-thought (CoT) prompting to generate an initial detailed rationale for each question, followed by an iterative prompt-rationale optimization process, aiming to enhance both the coherence of the reasoning and its alignment with the corresponding questions. In parallel, we synthe-

size 73k empathetic dialogues drawn from diverse social media sources to strengthen affective understanding and emotional support. To enhance reasoning, we adopt a multi-LLM cross-selection strategy to categorize questions into challenging and non-challenging subsets based on their inferred complexity. The non-challenging subset is used for supervised fine-tuning (SFT) to enhance empathetic response generation and domain knowledge, while the challenging subset is utilized for training with group relative policy optimization (GRPO) to improve the model’s reasoning capabilities, with both jointly contributing to the development of the *Psyche-R1*. Experimental results on a range of psychological benchmarks, including knowledge assessment, case-based analysis, and empathy evaluation, demonstrate the effectiveness of our model, where 7B *Psyche-R1* significantly outperforms models of similar scale and achieves competitive performance relative to substantially larger models such as DeepSeek-R1.

## 2 *Psyche-R1*

In this section, we give details of data curation and two-stage training paradigm, including data collection (§2.1), psychological reasoning data synthesis (§2.2), and empathetic dialogue synthesis (§2.3).

### 2.1 Data Collection

**Data Resource.** To construct a comprehensive and diverse dataset, we curate a wide range of resources:

- **Type I:** Classic psychology textbooks and curricular materials from psychology programs, all collected from publicly accessible repositories, covering more than 19 subfields and concentrated

140 and systematically organized content, including  
141 cognitive psychology, social psychology, etc.

- 142 • **Type II:** Psychological question banks collected  
143 from publicly available Chinese educational plat-  
144 forms via web crawlers and manually curation,  
145 encompassing theoretical principles and concep-  
146 tual knowledge across psychology.
- 147 • **Type III:** Synthetic data distilled from  
148 Qwen2.5-72B-Instruct (Team, 2024b) to sup-  
149 plement underrepresented subfields (e.g., sports  
150 psychology) and enhance dataset coverage.
- 151 • **Type IV:** Dialogic interactions harvested from  
152 established mental health support communities  
153 (e.g., Yixinli, Jiandanxinli and Zhihu<sup>1</sup>) ded-  
154 icated to delivering mental health services.

The **Type I**, **Type II** and **Type III** are used to con-  
155 struct the psychological reasoning question-answer  
156 (PCQA) dataset, designed to enhance domain-  
157 specific knowledge acquisition and reasoning ca-  
158 pabilities. The **Type IV** resource is employed to  
159 develop the empathetic dialogue dataset to improve  
160 affective understanding and empathetic expression.

**Data Cleaning.** To ensure the data quality, we im-  
162 plement several important data cleaning steps: (1)  
163 To process materials in non-textual formats such  
164 as images and PDFs, we employ OLMOCR<sup>2</sup> for accu-  
165 rate text recognition and conversion to text format;  
166 (2) We standardize the usage of Chinese and En-  
167 glish punctuation and remove irrelevant content,  
168 including emojis, emoticons, and links, to filter  
169 out potential noise; (3) To construct the empathetic  
170 dialogue data, we leverage LLMs to evaluate the  
171 reasonableness and relevance of QA pairs and fil-  
172 ter out responses that lack substantive advice. For  
173 example, given the question “I have recently expe-  
174 rienced insomnia and feel anxious. What should  
175 I do?”, if the response is merely “Everything will  
176 get better,” it would be filtered due to the absence  
177 of practical advice.

## 179 2.2 Psychological Reasoning Data Synthesis

**Question Generation and Control.** Following  
180 the data cleaning, we proceed to generate structured  
181 questions and corresponding answers from curated  
182 psychological textbooks and instructional materi-  
183 als (i.e., resource **Type I**). Specifically, the source  
184 material is first segmented into multiple textual  
185 chunks, with each chunk designed to encapsulate  
186 the maximum amount of domain-specific content.  
187

188 Subsequently, we leverage LLMs to generate a set  
189 of different questions along with corresponding an-  
190 swers based on these text segments. Meanwhile,  
191 for resource **Type III**, we use the similar way to gen-  
192 erate QA pairs without segment-level contextual  
193 augmentation, aiming to supplement psychological  
194 subfields that are underrepresented or difficult to  
195 source from publicly accessible materials on the  
196 Internet. Through these steps, we obtain approxi-  
197 mately 200k generated QA pairs in total.

All generated QA pairs, together with data de-  
198 rived from the **Type II** resource are integrated into  
199 a unified QA pool containing approximately 210k  
200 entries. We implement a multi-stage quality con-  
201 trol procedure to ensure the integrity and utility of  
202 the synthesized data. Concretely, we use min-hash  
203 locality-sensitive hashing (LSH) to cluster similar  
204 questions and select the optimal one through LLM-  
205 based ranking. Afterward, we prompt the LLM  
206 with few-shot examples to identify and filter out  
207 low-quality questions, specifically those exhibit-  
208 ing incomplete information, logical confusion, or  
209 unclear expression. Finally, we invite 10 under-  
210 graduate and graduate students to manually review  
211 the questions to eliminate redundant content and  
212 reduce potential noise in the dataset, ultimately  
213 retaining about 90k QA pairs.  
214

**Rationale Generation.** We further generate de-  
215 tailed rationales for the aforementioned questions  
216 through CoT prompting, following a multi-step rea-  
217 soning approach (Hsieh et al., 2023) to provide  
218 clear reasoning paths for model training. In detail,  
219 the CoT prompt guides the model to first compre-  
220 hend the question, recognize relevant psychologi-  
221 cal concepts and knowledge, and decompose the  
222 problem into a sequence of analytical steps. At  
223 each stage of this process, the model is required  
224 to articulate an intermediate rationale, ultimately  
225 generating a final answer derived from the accu-  
226 mulated reasoning. Formally, given a CoT prompt  
227  $P$  and a QA pair  $(q_i, a_i)$ , this procedure yields  
228 a rationale-augmented instance  $(q_i, r_i, \hat{a}_i)$ , where  
229  $r_i$  denotes the reasoning path and  $\hat{a}_i$  the model-  
230 predicted answer. If the predicted answer aligns  
231 with the ground truth (i.e.,  $\hat{a}_i = a_i$ ), we regard  
232 the  $r_i$  as a valid rationale. In contrast, if the pre-  
233 dicted answer is incorrect (i.e.,  $\hat{a}_i \neq a_i$ ), we guide  
234 the model to regenerate the rationale up to  $T$  time.  
235 Instances that fail to produce correct predictions  
236 after  $T$  regeneration attempts are pruned from the  
237 final curated dataset. After obtaining the initial ra-  
238

<sup>1</sup>Yixinli, Jiandanxinli, and Zhihu

<sup>2</sup><https://olmocr.allenai.org/>

tionale, we employ a self-supervised optimization strategy to iteratively refine both the prompt and the rationale with the goal of enhancing their clarity and reliability. Specifically, for each instance  $(P, q_i, r_i, \hat{a}_i)$ , the prompt  $P$  and rationale  $r_i$  are jointly updated over multiple rounds, enabling the model to progressively improve its reasoning process. Each round of optimization process consists of two sequential steps:

- **Prompt refinement:** We first guide the LLM to generate an improved candidate prompt  $P_i^*$  from the current prompt, question, and rationale, represented as  $P_i^* \leftarrow LLM(P, q_i, r_i)$ , aiming to enhance the reasoning guidance.
- **Rationale revision:** Based on the candidate prompt  $P_i^*$ , the LLM subsequently generates a revised rationale along with its corresponding predicted answer, denoted as  $(r_i^*, \hat{a}_i^*) \leftarrow LLM(P_i^*, q_i)$ .

If the  $\hat{a}_i^*$  matches the ground truth  $a_i$ , we retain  $P_i^*$  as an updated prompt and continue iteration based on the updated instance  $(P_i^*, q_i, r_i^*, \hat{a}_i^*)$ . Otherwise, the process reverts to the previous prompt-rationale pair to maintain alignment with correct reasoning paths. We repeat this process for  $R$  rounds ( $R = 3$  in this paper). After completing all iterations, we evaluate the rationales generated in each round for a given question and select the one that demonstrates the highest quality, denoted as  $(P_i^*, q_i, r_i^*, \hat{a}_i^* = a_i)$ . At this stage, we filter approximately 75k high-quality instances from the initial set of 90k pairs obtained in the previous step.

**Question Selection.** While the preceding steps yield high-quality data, not all instances exhibit sufficient complexity to facilitate effective reinforcement learning (RL). To address this, in this stage, we implement a multi-LLM cross-selection strategy aimed at identifying and isolating the most challenging psychology-related samples from the constructed dataset for subsequent use in the reinforcement learning phase. In detail, we employ three distinct LLMs (i.e., Qwen, Llama, and Phi) to independently generate responses for each question in the constructed psychological data. Questions that receive incorrect responses from all three models are aggregated into a challenging subset with 19k instances that provide sophisticated scenarios in psychology. This subset is intended to represent highly difficult instances with strong potential to enhance the model’s reasoning capabilities through reinforcement learning.

## 2.3 Empathetic Dialogue Synthesis

In addition to psychological QA pairs, empathy is recognized as a core component of effective mental health support (Sorin et al., 2024). To this end, we incorporate empathetic expressions into the dialogue corpus derived from authentic resources to enhance its emotional richness and relevance to real-world psychological interactions. Specifically, we refine these dialogue data through LLMs to achieve the following objectives. We first enhance emotional resonance by incorporating empathetic expressions and supportive techniques (e.g., “Hearing about your experience, I wish I could give you a warm hug.”). Subsequently, we ensure that each dialogue provides evidence-based guidance that facilitates deeper understanding of users’ issues, instead of limiting responses to surface-level empathy. Finally, we deliver solution-oriented support by offering concrete coping strategies and practical steps that address the specific issues and challenges presented. Through these steps, we ultimately obtain 73k high-quality dialogue data equipped with sufficient empathetic expressions.

## 2.4 Data Split

Leveraging the aforementioned pipelines, we curate a comprehensive dataset that comprises over 75k psychological questions with detailed rationales, among which 19k are identified as challenging samples through multi-LLM cross-selection. The challenging subset is denoted as  $\mathcal{D}_{pc}$  and the remaining data are denoted as  $\mathcal{D}_{pr}$ . In parallel, the dataset contains over 73k empathetic dialogues engineered for contextually appropriate psychosocial interactions, denoted as  $\mathcal{D}_{em}$ . To further enrich our training data, we additionally introduce multi-turn dialogues and knowledge-based QA from the PsychoLLM dataset (Hu et al., 2024), and a refined set of 8k examination from the CPsyExam train set (Zhao et al., 2025).

Ultimately, the curated datasets are partitioned into two distinct subsets aligned with specialized training objectives. One category, represented as  $\mathcal{D}_{sft} = \mathcal{D}_{pr} \cup \mathcal{D}_{em} \cup \mathcal{D}_{ps}$ , is designated for SFT. The other category, denoted as  $\mathcal{D}_{grpo} = \mathcal{D}_{pc} \cup \mathcal{D}_{cp}$ , is reserved for GRPO. Detailed prompts for the data synthesis pipeline are provided in the Appendix D.

## 2.5 Model Training

To enhance both reasoning capabilities and performance in empathy and expertise, we employ a

339 hybrid training strategy.

340 **Stage 1: Supervised Fine-Tuning.** We  
341 initialize our backbone model  $\pi_\theta$  with  
342 Qwen2.5-7B-Instruct (Team, 2024b) and  
343 finetune it on  $\mathcal{D}_{\text{sft}}$ . Formally, given a query  $x$ , the  
344 model is trained to generate a coherent rationale  
345  $r$  followed by a corresponding answer  $a$ , where  
346 the complete output is denoted as  $y = [r; a]$ . We  
347 optimize model parameters  $\theta$  by minimizing the  
348 standard negative log-likelihood loss:

$$349 \mathcal{L}(\theta) = -\mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{sft}}} \left[ \sum_{t=1}^T \log P(y_t | x, y_{<t}; \theta) \right] \quad (1)$$

350 **Stage 2: Group Relative Policy Optimization.**

351 To further refine psychological reasoning profi-  
352 ciency, we employ GRPO (Shao et al., 2024) on  
353  $\mathcal{D}_{\text{grpo}}$ . We design a composite reward function  
354  $R(y, y^*) = R_{\text{fmt}} + R_{\text{acc}}$  to guide policy learning,  
355 where  $y^*$  denotes the ground truth.

- 356 • **Format reward ( $R_{\text{fmt}}$ ):** We enforce strict for-  
357 mating constraints. The reasoning process must  
358 be encapsulated within `<think>` and `</think>`  
359 tags, followed by the final answer. We assign  
360 a scalar reward  $R_{\text{fmt}} = +1.25$  for structurally  
361 parsable outputs and  $-1$  otherwise.
- 362 • **Accuracy reward ( $R_{\text{acc}}$ ):** We formulate the an-  
363 swer matching as a set comparison task. Specif-  
364 ically, we parse the predicted answer  $\hat{a}$  and the  
365 ground truth  $a$  into sets of discrete options. To  
366 encourage precise reasoning alignment while pe-  
367 nalyzing hallucinations or omissions, we employ  
368 a partial-credit mechanism based on the overlap  
369 between the prediction and the ground truth:

$$370 R_{\text{acc}} = \begin{cases} +1, & \text{if } \hat{a} = a \\ \frac{|\hat{a} \cap a|}{|a|}, & \text{if } \hat{a} \subset a \wedge a \neq \emptyset \\ -1, & \text{otherwise} \end{cases} \quad (2)$$

371 By integrating these logical and structural signals,  
372 the model learns to generate well-organized rea-  
373 soning processes while rewarding partial credit for  
374 incomplete but valid answers, leading to better per-  
375 formance in psychological tasks.

## 376 3 Experiments

### 377 3.1 Experimental Setting

378 **Baselines.** To ensure a comprehensive analysis,  
379 we selected 17 representative LLMs, categorized  
380 as follows: (1) **General LLMs**, which exhibit ex-  
381 cellent performance across general tasks, but lack

382 explicit reasoning capabilities. (2) **Reasoning aug-**  
383 **mented LLMs**, which possess explicit reasoning  
384 capabilities. (3) **Closed-source LLMs**, which typi-  
385 cally represent the state-of-the-art performance. (4)  
386 **Psychological LLMs**, which have been fine-tuned  
387 on psychological datasets. A detailed summary of  
388 all models is presented in Appendix C.1.

389 **Benchmarks and Evaluation Metrics.** We  
390 conduct comprehensive evaluations on two pro-  
391 fessional psychological benchmarks:

- 392 • **Psychological counselor examination bench-**  
393 **mark (PCEB)** (Hu et al., 2024): this consists of  
394 3,863 multiple-choice questions (MCQ) and 100  
395 open-ended case analysis items, curated from the  
396 official National Psychological Counselor Exam-  
397 ination in China.
- 398 • **CPsyExam test set** (Zhao et al., 2025): this in-  
399 cludes 4,102 questions spanning 39 distinct psy-  
400 chological subfields. We evaluate under both  
401 zero-shot and five-shot settings, ensuring con-  
402 sistency by using identical exemplars across all  
403 evaluated models in the latter setting.

404 Note that MCQ comprises two types of ques-  
405 tions: MCQ with only a single correct option  
406 (SMCQ), and MCQ with multiple correct options  
407 (MMCQ). For MCQ, we adopt metrics introduced  
408 in Hu et al. (2024), including **standard accuracy**,  
409 which requires predictions to exactly match the  
410 ground truth, and **elastic accuracy**, which gives  
411 partial credit when predictions are a subset of the  
412 correct answers. For open-ended questions, we uti-  
413 lize the existing text generation metrics, including  
414 **Rouge-1 (R-1)**, **Rouge-L (R-L)** (Lin, 2004), and  
415 **Bleu-4 (B-4)** (Papineni et al., 2002).

### 416 3.2 Overall Results

417 **Results on the PCEB.** To evaluate the perfor-  
418 mance of different models, we present the results  
419 on the PCEB in Table 1. These results reveal  
420 several key observations. First, *Psyche-R1* ex-  
421 hibits strong performance across evaluation tasks in  
422 both MCQ and subjective questions. This demon-  
423 strates the effectiveness of our proposed dataset  
424 and training strategy in simultaneously enhanc-  
425 ing psychological reasoning and text generation  
426 capabilities for psychological tasks. Second, while  
427 DeepSeek-R1 excels in MCQ, its performance in  
428 subjective questions is notably limited. This per-  
429 formance disparity can be attributed to its train-  
430 ing methodology, which employs RL on datasets  
431 primarily consisting of mathematical and coding

Model	Param.	Case		Moral		Theory		Avg.	Case (QA)					
		SMCQ	MMCQ	SMCQ	MMCQ	SMCQ	MMCQ		R-1	R-L	B-4			
<i>General LLMs</i>														
MiniCPM4-8B	8B	50.00	28.59	43.64	81.58	50.63	58.23	65.62	34.06	43.00	51.75 (57.01)	23.05	12.90	1.35
Qwen2.5-7B	7B	47.57	31.64	47.49	87.83	59.50	71.02	78.46	42.45	55.17	57.91 (64.59)	20.94	11.28	1.28
Qwen2.5-14B	14B	47.13	41.10	55.93	89.81	63.93	73.60	80.32	50.16	61.26	62.08 (68.01)	22.69	13.93	1.53
Qwen2.5-72B	72B	46.91	40.34	53.11	90.79	70.25	78.48	82.63	47.63	59.74	63.09 (68.61)	21.43	12.02	1.16
<i>Reasoning-Augmented LLMs</i>														
DeepSeek-R1	671B	79.25	44.25	60.86	95.39	68.99	77.95	92.19	57.60	69.41	72.95 (79.18)	17.65	9.19	0.94
DeepSeek-R1-70B	70B	56.30	30.72	46.95	88.16	52.53	65.66	68.01	25.64	45.63	53.56 (61.79)	22.77	13.23	1.16
QwQ-32B	32B	56.51	23.35	41.27	88.82	41.14	53.06	82.12	32.69	49.90	54.11 (61.95)	18.39	7.48	0.84
Qwen3-30B-A3B	30B	59.65	31.51	47.28	91.45	55.06	65.66	80.75	47.45	59.25	60.98 (67.34)	20.53	12.06	1.18
Qwen3-235B-A22B	235B	68.58	41.91	57.24	93.42	69.62	78.90	88.36	56.70	68.64	69.77 (75.86)	18.96	11.14	1.11
Magistral-Small	24B	56.58	33.26	49.11	82.89	53.80	67.99	70.10	37.76	52.35	55.73 (63.17)	22.90	11.97	1.21
<i>Closed-Source LLMs</i>														
Claude3.7-Sonnet	UNK	63.39	19.40	34.23	90.13	60.13	70.04	76.73	37.37	48.99	57.86 (63.92)	21.59	11.11	1.23
Gemini1.5-Pro	UNK	61.04	35.57	49.87	84.21	62.03	70.62	80.84	43.22	53.44	61.26 (66.78)	21.63	10.93	1.06
GPT-4o	UNK	65.63	13.67	34.53	88.15	33.54	54.79	74.65	24.10	45.07	49.96 (60.47)	23.45	12.75	1.18
<i>Psychological LLMs</i>														
CPsyCounX	7B	40.87	16.91	32.90	75.17	36.08	54.85	54.78	19.03	38.90	40.47 (49.58)	22.83	11.94	1.48
EmoLLM	7B	46.93	21.87	40.02	84.21	34.17	51.05	71.72	26.18	44.49	47.51 (56.40)	22.15	11.69	1.20
PsycoLLM	14B	55.58	35.07	42.89	88.81	69.62	74.20	72.63	48.59	54.12	61.72 (64.71)	24.45	17.45	2.04
PsyDT	7B	35.56	35.20	50.14	86.33	69.70	78.66	80.70	52.72	62.26	60.04 (65.61)	20.65	13.41	1.16
Psyche-R1	7B	63.31	56.26	66.21	92.76	79.62	82.54	87.70	66.54	73.34	74.37 (77.64)	27.31	15.33	2.40

Table 1: Comparison of different models on the PCEB, where Case, Moral, Theory, and Case (QA) are case analysis, theoretical proficiency, professional ethics, and case-based QA. The average value represents the average of the standard accuracy rates, and values in parentheses denotes the mean of the standard accuracy for SMCQ and the elastic accuracy for MMCQ. Results colored in red, orange, and yellow demote the best, second-best and third-best.

tasks with deterministic solutions. Although this approach strengthens logical reasoning, it appears to bias the model towards single-answer patterns, thereby limiting its capability to generate diverse and nuanced responses in open-ended psychological assessments. Third, existing psychological LLMs (e.g., CPsyCounX and EmoLLM) achieve strong performance in subjective questions while demonstrating limited abilities in MCQ. This imbalanced performance stems from their reliance on training exclusively on counseling dialogues or empathetic conversations, which constrains their capabilities to develop comprehensive competencies. Fourth, closed-source models such as GPT-4o and Claude3.7-Sonnet demonstrate relatively weaker performance, which may be attributed to limited Chinese representation in their training corpora.

**Results on the CPsyExam Test Set.** To further explore model performance, we present the results on the CPsyExam test set in Table 2. Similar to the trends observed in previous experiments, both *Psyche-R1* and *DeepSeek-R1* demonstrate superior performance. Across these models, psychological LLMs consistently achieve higher accuracy in SMCQ than in MMCQ, as the latter requires exhaustive evaluation of all options, demanding more comprehensive domain-specific knowledge and reasoning capabilities. Under the five-shot setting, most models exhibit substantial improvements in

MMCQ (e.g., *PsyDT* achieves a 47.64% improvement in knowledge-type MMCQ). This observation aligns with existing studies, which demonstrate that well-designed few-shot examples can effectively enhance model performance in certain tasks. In contrast, *DeepSeek-R1* exhibits a performance decline under the five-shot compared to its zero-shot setting, suggesting that few-shot prompting may interfere with its inherent reasoning capability aligning with existing findings (Guo et al., 2025).

### 3.3 Ablation Study

In this section, we conduct a comprehensive ablation study, evaluating model performance by standard accuracy on the PCEB.

**Effect of SFT and GRPO.** We investigate the effect of SFT and GRPO, with results shown in Table 3. We can observe that SFT substantially improves performance across the three task categories by leveraging our dataset of empathetic dialogues and rationale-augmented psychological questions. However, applying GRPO without prior SFT results in performance degradation on SMCQ case tasks, as the base model lacks sufficient domain knowledge and empathy, which are critical for reasoning in case-based questions, leading to unstable training dynamics. When combining SFT with GRPO training yields further gains, particularly on case-based tasks, as challenging samples identified via multi-LLM cross-selection promote

Model	Param.	Zero-Shot				Five-Shot				Avg.
		Knowledge		Case		Knowledge		Case		
		SMCQ	MMCQ	SMCQ	MMCQ	SMCQ	MMCQ	SMCQ	MMCQ	
<b>General LLMs</b>										
MiniCPM4-8B	8B	69.58	41.74	57.33	37.00	68.50	42.77	54.67	38.00	60.46
Qwen2.5-7B	7B	76.99	43.66	68.67	44.50	78.63	42.00	68.67	40.50	67.37
Qwen2.5-14B	14B	81.39	49.30	72.00	48.50	82.42	54.29	71.00	48.00	71.84
Qwen2.5-72B	72B	84.61	52.75	73.50	54.50	86.64	63.77	75.33	55.00	74.98
<b>Reasoning-Augmented LLMs</b>										
DeepSeek-R1	671B	87.49	56.98	76.83	59.00	88.78	66.58	77.30	61.50	78.28
DeepSeek-R1-70B	70B	76.48	22.80	61.81	19.17	76.89	40.99	62.70	37.95	60.57
<b>Closed-Source LLMs</b>										
Gemini1.5-Pro	UNK	82.08	40.59	68.33	43.00	83.93	53.65	71.00	45.00	69.66
GPT-4o	UNK	80.70	30.73	66.33	28.00	81.82	54.80	68.67	51.50	65.79
<b>Psychological LLMs</b>										
CPsyCounX	7B	57.56	22.41	46.33	31.00	63.46	21.77	50.67	23.50	47.44
EmoLLM	7B	78.41	45.33	72.50	48.00	79.92	36.88	74.17	39.50	69.32
PsycoLLM	14B	78.33	51.98	65.33	42.00	78.63	50.45	65.57	36.00	69.20
PsyDT	7B	80.83	48.91	69.67	41.50	81.13	40.97	68.33	40.00	70.71
Psyche-R1	7B	82.72	61.59	70.50	49.50	83.45	61.46	76.17	52.00	74.90

Table 2: Comparisons of different models on the CPsyExam test set. The average represents the overall zero-shot accuracy. The first, second, and third-best results are marked in red, orange, and yellow, respectively.

Model	Case	$\Delta$	Moral	$\Delta$	Theory	$\Delta$
Base	38.97	-	73.39	-	63.83	-
<b>Ablation study on training stage</b>						
+GRPO	36.69	-5.85%	77.74	5.93%	73.34	14.90%
+SFT	48.51	24.48%	83.82	14.21%	73.44	15.06%
+SFT+GRPO	67.07	72.11%	86.06	17.26%	79.10	23.92%
<b>Ablation study on rationale optimization</b>						
+QA	48.27	23.86%	81.89	11.58%	71.22	11.58%
+Rat.	55.25	41.78%	85.14	16.01%	75.55	18.36%
+Rat.+Iter.	67.07	72.11%	86.06	17.26%	79.10	23.92%

Table 3: Ablation study evaluating the effects of training stages (SFT and GRPO) and the contributions of rationale component (Rat.) and iterative prompt-rationale optimization (Iter.).

deeper reasoning and contextual understanding.

**Effect of the Rationale and Iterative Optimization.** We explore the contributions of rationales and iterative prompt-rationale optimization, with results presented in Table 3. Note that **+QA** is the model trained solely on question-answer pairs without incorporating detailed rationales. Compared with the base model, training with our proposed dataset (i.e., **+QA**, **+Rat.** and **+Rat.+Iter.**) leads to consistent performance improvements, demonstrating the effectiveness of the curated data. Integrating rationale-augmented data substantially enhances performance over training with option-only answers, indicating that rationales provide valuable intermediate reasoning signals that facilitate learning. Furthermore, applying iterative prompt-rationale optimization (i.e., **+Rat.+Iter.**) yields further gains, confirming that progressively refined rationales contribute to better supervision

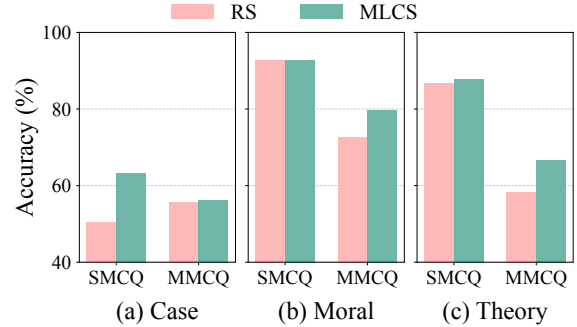


Figure 3: Comparison of performance using challenging question selection methods, including multi-LLM cross-selection (MLCS) and random selection (RS).

and more robust model reasoning.

**Effect of the Question Selection.** We further examine the effect of question selection by comparing multi-LLM cross-selection (MLCS) with random selection (RS), as illustrated in Figure 3. The comparison between MLCS and RS demonstrates that leveraging multiple LLMs for selecting challenging instances yields markedly superior outcomes across all these tasks. This finding confirms that our selection mechanism effectively selects high-quality training data for GRPO, which is instrumental in enhancing reasoning and generalization within the psychological domain.

### 3.4 Discussion

**Effect of Datasets.** We evaluate the model performance across different combinations of subsets, with results presented in Table 4. It is observed that fine-tuning with either psychological reasoning data (PRD) or empathetic dialogues (ED) in isolation delivered marginal improvements in task

Settings	Case		Moral		Theory		Avg.
	SMCQ	MMCQ	SMCQ	MMCQ	SMCQ	MMCQ	
Base	47.57	31.64	87.83	59.50	78.46	42.45	57.91
+ED	35.77	29.74	70.00	60.90	65.84	44.40	51.11
+PRD	37.94	35.45	91.45	51.27	80.47	37.53	55.69
+ED+PRD	61.71	53.56	92.72	76.58	86.13	68.16	73.14
+ED+PRD+APD	63.31	56.26	92.76	79.62	87.70	66.54	74.37

Table 4: Effect of different dataset compositions, including empathetic dialogues (ED, i.e.,  $\mathcal{D}_{em}$ ), psychological reasoning data (PRD, i.e.,  $\mathcal{D}_{pc} \cup \mathcal{D}_{pr}$ ), and additional public datasets (APD, i.e.,  $\mathcal{D}_{ps} \cup \mathcal{D}_{cp}$ ).

Model	EmoE.	CogE.	Con.	Sta.
Qwen2.5-7B-Instruct	1.52	2.00	2.36	1.72
CPsyCounX	1.73	2.05	2.15	1.96
EmoLLM	1.86	2.44	<b>2.84</b>	<b>2.34</b>
PsycoLLM	1.97	2.27	2.41	2.10
PsyDT	2.21	2.46	2.36	<b>2.34</b>
Psyche-R1	<b>2.33</b>	<b>2.69</b>	2.78	2.11

Table 5: Comparisons of psychological LLMs on PsyDT test set. The evaluation metrics comprise: emotional empathy (EmoE.), cognitive empathy (CogE.), conversation strategy (Con.), and state and attitude (Sta.).

performance, and in some cases, led to a slight decline in overall accuracy. In contrast, the combination of PRD and ED achieves substantial improvements across these tasks, highlighting the quality and comprehensiveness of our proposed data synthesis pipeline. This result demonstrates that integrating domain-specific knowledge with emotional understanding enhances psychological reasoning. Moreover, the incorporation of additional public datasets (APD) leads to further performance improvements.

**Performance on Counseling Tasks.** Beyond examination tasks, we evaluate the performance of *Psyche-R1* on counseling tasks and compare it with its base model and several outstanding psychological LLMs. Following the method of PsyDT (Xie et al., 2025) but constrained by limited resources, we randomly sample 200 instances from its test set and employ GPT-4o as the evaluator. As shown in Table 5, *Psyche-R1* achieves notable improvements compared to its base model, demonstrating its capability in counseling tasks that demand emotional empathy, cognitive empathy and so on. This excellent performance stems from the synergistic interplay between two crucial elements: the empathetic dialogues, which directly improve counseling effectiveness, and advanced reasoning mechanisms, which enable a deeper understanding of questions, thereby yielding more accurate and emotionally informed responses within relevant contexts.

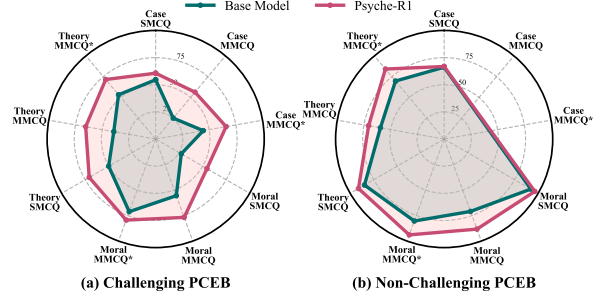


Figure 4: Comparison of model performance on (a) Challenging and (b) Non-Challenging subsets of the PCEB. An asterisk(\*) indicates elastic accuracy, while the remaining metrics represent standard accuracy.

**Error Analysis.** To assess the impact of our approach across varying difficulty levels, we divide the PCEB into (a) challenging and (b) non-challenging subsets following the question selection method described in §2.2. The results are presented in Figure 4. It is observed that *Psyche-R1* demonstrates consistent improvements across these tasks in both subsets. For the non-challenging one, performance gains are primarily concentrated in the theory and moral dimensions, reflecting the model’s proficiency in handling foundational psychological concepts. Notably, our model exhibits more pronounced improvements on the challenging subset. This observation can be attributed to the synergistic effect of our data generation pipeline and hybrid training strategy.

## 4 Conclusion

In this paper, we propose *Psyche-R1*, the first Chinese psychological LLM that jointly integrates empathy, expertise, and reasoning. To support model development, we design a multi-stage data synthesis pipeline that generates high-quality psychological reasoning samples with detailed rationales and empathetic dialogues. The reasoning rationales are further enhanced through iterative prompt–rationale optimization, and a multi-LLM cross-selection strategy is employed to identify challenging examples. Finally, the challenging subset is used for GRPO, while the remaining data are employed for SFT, together contributing to the final model. Extensive experiments demonstrate that *Psyche-R1* outperforms existing psychological LLMs, achieving performance comparable to DeepSeek-R1. Moreover, we perform comprehensive ablation studies and analyses to evaluate the individual contributions of each component and strategy within the proposed framework.

## 596 Limitations

597 Despite the promising results of our *Psyche-RI*, our  
598 study is subject to several limitations that remain  
599 to be addressed in future research.

600 **Language and Cultural Specificity.** To miti-  
601 gate the shortage of mental health professionals in  
602 China, current *Psyche-RI* and its training corpus  
603 are predominantly tailored to the Chinese language  
604 and cultural context. Consequently, the model’s em-  
605 pathetic reasoning involves specific cultural norms  
606 that may not directly transfer to other languages.  
607 We frame this as a necessary step for local appli-  
608 cability, noting that cross-cultural generalization  
609 remains challenging for future research.

610 **Model Scale.** Constrained by computational re-  
611 sources, *Psyche-RI* is built upon a 7B-parameter  
612 backbone. While it achieves competitive perfor-  
613 mance, we posit that employing a base model with  
614 a larger scale would yield superior performance.

## 615 Ethical Considerations

616 The development and deployment of LLMs in the  
617 mental health domain necessitate rigorous adher-  
618 ence to ethical standards.

619 **Nature of the System.** *Psyche-RI* is designed  
620 as a supportive tool for mental health support and  
621 education, rather than a replacement for qualified  
622 mental health professionals. The model is not au-  
623 thorized to provide medical diagnoses, prescribe  
624 treatments, or handle crisis interventions. Users  
625 facing severe mental health crises should seek help  
626 from human professionals or emergency services.

627 **Data Privacy and Safety.** We prioritize the pri-  
628 vacy and safety of individuals in our data cura-  
629 tion process. For data collected from social me-  
630 dia platforms (Type IV), we implemented strict  
631 de-identification procedures to remove all person-  
632 ally identifiable information, including names, lo-  
633 cations, and contact details. We strictly adhere to  
634 data usage policies and ensure that the synthesized  
635 data does not reconstruct real-world private inter-  
636 actions. Furthermore, our data synthesis pipeline that  
637 prioritizes high-quality, constructive psychological  
638 advice filtered out toxic content and harmful  
639 suggestions to align with safety guidelines.

## 640 References

641 Nafiz Al Asad, Md Appel Mahmud Pranto, Sadia  
642 Afreen, and Md Maynul Islam. 2019. Depression

detection by analyzing social media posts of user. 643  
In *2019 IEEE international conference on signal 644*  
*processing, information, communication & systems 645*  
*(SPICSCON)*, pages 13–17. IEEE. 646

Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, 647  
Wanlong Liu, Rongsheng Wang, Jianye Hou, and 648  
Benyou Wang. 2024. Huatuoqpt-o1, towards med- 649  
ical complex reasoning with llms. *arXiv preprint 650*  
*arXiv:2412.18925*. 651

Yirong Chen, Xiaofen Xing, Jingkai Lin, Huimin Zheng, 652  
Zhenyu Wang, Qi Liu, and Xiangmin Xu. 2023a. 653  
[SoulChat: Improving LLMs’ empathy, listening, and 654](#)  
[comfort abilities through fine-tuning with multi-turn 655](#)  
[empathy conversations](#). In *Findings of the Associa- 656*  
*tion for Computational Linguistics: EMNLP 2023*, 657  
pages 1170–1183, Singapore. Association for Com- 658  
putational Linguistics. 659

Zhiyu Chen, Yujie Lu, and William Wang. 2023b. *Em- 660*  
*powering psychotherapy with large language mod- 661*  
*els: Cognitive distortion detection through diagnosis 662*  
*of thought prompting*. In *Findings of the Associa- 663*  
*tion for Computational Linguistics: EMNLP 2023*, 664  
pages 4295–4304, Singapore. Association for Com- 665  
putational Linguistics. 666

Yujin Cho, Mingeon Kim, Seojin Kim, Oyun Kwon, 667  
Ryan Donghan Kwon, Yoonha Lee, and Dohyun Lim. 668  
2023. Evaluating the efficacy of interactive language 669  
therapy based on llm for high-functioning autistic 670  
adolescent psychological counseling. *arXiv preprint 671*  
*arXiv:2311.09243*. 672

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, 673  
Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar- 674  
cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 675  
1 others. 2025. Gemini 2.5: Pushing the frontier with 676  
advanced reasoning, multimodality, long context, and 677  
next generation agentic capabilities. *arXiv preprint 678*  
*arXiv:2507.06261*. 679

Dorottya Demszky, Diyi Yang, David S Yeager, Christo- 680  
pher J Bryan, Margaret Clapper, Susannah Chand- 681  
hok, Johannes C Eichstaedt, Cameron Hecht, Jeremy 682  
Jamieson, Meghann Johnson, and 1 others. 2023. Us- 683  
ing large language models in psychology. *Nature 684*  
*Reviews Psychology*, 2(11):688–701. 685

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, 686  
Pengfei Liu, Yiming Yang, Jamie Callan, and Gra- 687  
ham Neubig. 2023. [PAL: Program-aided language 688](#)  
[models](#). In *Proceedings of the 40th International 689*  
*Conference on Machine Learning*, volume 202 of 690  
*Proceedings of Machine Learning Research*, pages 691  
10764–10799. PMLR. 692

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao 693  
Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi- 694  
rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. 695  
Deepseek-r1: Incentivizing reasoning capability in 696  
llms via reinforcement learning. *arXiv preprint 697*  
*arXiv:2501.12948*. 698

699	Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. <a href="#">Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8003–8017, Toronto, Canada. Association for Computational Linguistics.	
700		
701		
702		
703		
704		
705		
706		
707		
708	Jinpeng Hu, Tengpeng Dong, Luo Gang, Hui Ma, Peng Zou, Xiao Sun, Dan Guo, Xun Yang, and Meng Wang. 2024. <a href="#">Psychollm: Enhancing llm for psychological understanding and evaluation</a> . <i>IEEE Transactions on Computational Social Systems</i> .	
709		
710		
711		
712		
713	Zhaocheng Huang, Julien Epps, Dale Joachim, and Vidhyasaharan Sethu. 2020. <a href="#">Natural language processing methods for acoustic and landmark event-based features in speech-based depression detection</a> . <i>IEEE Journal of Selected Topics in Signal Processing</i> , 14(2):435–448.	
714		
715		
716		
717		
718		
719	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. <a href="#">Gpt-4o system card</a> . <i>arXiv preprint arXiv:2410.21276</i> .	
720		
721		
722		
723		
724	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. <a href="#">Openai o1 system card</a> . <i>arXiv preprint arXiv:2412.16720</i> .	
725		
726		
727		
728		
729	Tatsuki Kuribayashi, Yohei Oseki, and Timothy Baldwin. 2024. <a href="#">Psychometric predictive power of large language models</a> . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 1983–2005, Mexico City, Mexico. Association for Computational Linguistics.	
730		
731		
732		
733		
734		
735	Tin Lai, Yukun Shi, Zicong Du, Jiajie Wu, Ken Fu, Yichao Dou, and Ziqi Wang. 2023. <a href="#">Supporting the demand on mental health services with ai-based conversational large language models (llms)</a> . <i>BioMedInformatics</i> , 4(1):8–33.	
736		
737		
738		
739		
740	Daeun Lee, Soyoung Park, Jiwon Kang, Daejin Choi, and Jinyoung Han. 2020. <a href="#">Cross-lingual suicidal-oriented word embedding toward suicide prevention</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2208–2217, Online. Association for Computational Linguistics.	
741		
742		
743		
744		
745		
746	Suyeon Lee, Sunghwan Kim, Minju Kim, Dongjin Kang, Dongil Yang, Harim Kim, Minseok Kang, Dayi Jung, Min Hee Kim, Seungbeen Lee, Kyong-Mee Chung, Youngjae Yu, Dongha Lee, and Jinyoung Yeo. 2024. <a href="#">Cactus: Towards psychological counseling conversations using cognitive behavioral theory</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 14245–14274, Miami, Florida, USA. Association for Computational Linguistics.	
747		
748		
749		
750		
751		
752		
753		
754		
755		
	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	756 757 758 759
	Che Liu, Haozhe Wang, Jiazhen Pan, Zhongwei Wan, Yong Dai, Fangzhen Lin, Wenjia Bai, Daniel Rueckert, and Rossella Arcucci. 2025. <a href="#">Beyond distillation: Pushing the limits of medical llm reasoning with minimalist rule-based rl</a> . <i>arXiv preprint arXiv:2505.17952</i> .	760 761 762 763 764 765
	Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. 2025. <a href="#">A comprehensive overview of large language models</a> . <i>ACM Trans. Intell. Syst. Technol.</i> , 16(5).	766 767 768 769 770
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	771 772 773 774 775 776 777
	Huachuan Qiu, Hongliang He, Shuai Zhang, Anqi Li, and Zhenzhong Lan. 2024. <a href="#">SMILE: Single-turn to multi-turn inclusive language expansion via ChatGPT for mental health support</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 615–636, Miami, Florida, USA. Association for Computational Linguistics.	778 779 780 781 782 783 784
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. <a href="#">Deepseekmath: Pushing the limits of mathematical reasoning in open language models</a> . <i>arXiv preprint arXiv:2402.03300</i> .	785 786 787 788 789 790
	Hao Shen, Zihan Li, Minqiang Yang, Minghui Ni, Yongfeng Tao, Zhengyang Yu, Weihao Zheng, Chen Xu, and Bin Hu. 2024. <a href="#">Are large language models possible to conduct cognitive behavioral therapy? In 2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)</a> , pages 3695–3700.	791 792 793 794 795 796
	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2025. <a href="#">Hybridflow: A flexible and efficient rlhf framework</a> . In <i>Proceedings of the Twentieth European Conference on Computer Systems</i> , EuroSys '25, page 1279–1297, New York, NY, USA. Association for Computing Machinery.	797 798 799 800 801 802 803
	Vera Sorin, Dana Brin, Yiftach Barash, Eli Konen, Alexander Charney, Girish Nadkarni, and Eyal Klang. 2024. <a href="#">Large language models and empathy: Systematic review</a> . <i>J Med Internet Res</i> , 26:e52597.	804 805 806 807
	Michael J Tanana, Christina S Soma, Patty B Kuo, Nicolas M Bertagnolli, Aaron Dembe, Brian T Pace, Vivek Srikumar, David C Atkins, and Zac E Imel. 2021. <a href="#">How do you feel? using natural language processing to automatically rate emotion in psychotherapy</a> . <i>Behavior research methods</i> , 53(5):2069–2082.	808 809 810 811 812 813



930	Chujie Zheng, Sahand Sabour, Jiaxin Wen, Zheng	reasoning. Building upon this foundation, re-	980
931	Zhang, and Minlie Huang. 2023. <a href="#">AugESC: Dialogue</a>	searchers have explored more sophisticated reason-	981
932	<a href="#">augmentation with large language models for emo-</a>	ing architectures. For instance, Tree of Thoughts	982
933	<a href="#">tional support conversation</a> . In <i>Findings of the As-</i>	(Yao et al., 2023) enables systematic exploration	983
934	<i>sociation for Computational Linguistics: ACL 2023</i> ,	of multiple reasoning paths with self-evaluation,	984
935	pages 1552–1568, Toronto, Canada. Association for	while PAL (Gao et al., 2023) integrates reasoning	985
936	Computational Linguistics.	with external tools through program generation.	986
937	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	These approaches further enhance model perfor-	987
938	Ye, and Zheyang Luo. 2024. <a href="#">LlamaFactory: Unified</a>	mance in handling complex tasks. A new break-	988
939	<a href="#">efficient fine-tuning of 100+ language models</a> . In	through was marked by the release of reasoning	989
940	<i>Proceedings of the 62nd Annual Meeting of the As-</i>	LLMs like OpenAI o1 (Jaech et al., 2024) and	990
941	<i>sociation for Computational Linguistics (Volume 3:</i>	DeepSeek-R1 (Guo et al., 2025). These models,	991
942	<i>System Demonstrations)</i> , pages 400–410, Bangkok,	which are trained through reinforcement learning	992
943	Thailand. Association for Computational Linguistics.	with reasoning techniques to enhance reasoning	993
944	Jie Zhu, Qian Chen, Huaixia Dou, Junhui Li, Lifan	capabilities, demonstrate exceptional performance	994
945	Guo, Feng Chen, and Chi Zhang. 2025. <a href="#">Dianjin-r1:</a>	in mathematical and coding tasks (Comanici et al.,	995
946	<a href="#">Evaluating and enhancing financial reasoning in large</a>	2025; Yang et al., 2025). Motivated by these ad-	996
947	<a href="#">language models</a> . <i>arXiv preprint arXiv:2504.15716</i> .	vances, researchers have employed advanced RL	997
948	<b>A Related Work</b>	algorithms, such as GRPO (Shao et al., 2024) and	998
949	<b>A.1 LLMs for Psychology</b>	DAPO (Yu et al., 2025), to further extend reasoning	999
950	The success of LLMs has spurred interest in de-	capabilities to domain-specific applications, includ-	1000
951	veloping LLM-driven mental health applications	ing medicine (Liu et al., 2025) and finance (Zhu	1001
952	(Demszky et al., 2023). Early research focused	et al., 2025). However, within the field of psychol-	1002
953	primarily on improving the accessibility of men-	ogy, limited research has investigated the utility of	1003
954	tal health services. Research in this phase primar-	reasoning. To our knowledge, <i>Psyche-R1</i> is the first	1004
955	ily concentrated on two directions: One direction	psychological LLM that unifies empathy, domain-	1005
956	involves leveraging NLP techniques for emotion	specific expertise, and reasoning capabilities.	1006
957	recognition to enable automated detection of de-	<b>B Case Study</b>	1007
958	pression (Huang et al., 2020) and suicidal ideation	We present a case study examining how <i>Psyche-R1</i>	1008
959	(Lee et al., 2020). The other focuses on construct-	and Qwen2.5-72B-Instruct formulate their con-	1009
960	ing empathetic dialogue systems by fine-tuning	clusions derived from narratives and deliver mental	1010
961	LLMs on single-turn (Lai et al., 2023) or multi-	health support, as illustrated in Figure 5. These	1011
962	turn (Qiu et al., 2024) dialogue data to enhance	two models display distinct counseling strategies	1012
963	their abilities in affective understanding and emo-	when addressing the case involving a company	1013
964	tional support (Team, 2024a; Xie et al., 2025). As	manager confronting a career transition dilemma.	1014
965	research progressed, researchers began to explore	<i>Psyche-R1</i> begins by expressing empathy (e.g., “I	1015
966	more diverse mental health applications. Some	can really sense...”), followed by applying rele-	1016
967	studies have transformed traditional psychometric	vant psychological concepts tailored to user’s sit-	1017
968	tools (e.g., psychological scales) into interactive	uations, thereby demonstrating both emotional at-	1018
969	systems to improve user engagement (Kuribayashi	tunement and domain-specific expertise. In con-	1019
970	et al., 2024; Yang et al., 2024). Another line of	trast, Qwen’s empathetic expressions appear less	1020
971	research has focused on the specialized demands of	natural and engaging (e.g., “Your situation is truly	1021
972	the psychological domain, developing professional-	understandable...”), and it fails to apply theoretical	1022
973	grounded mental health applications based on es-	knowledge to contextualize or explain the user’s	1023
974	tablished psychological therapies (Lee et al., 2024;	dilemma, which undermines the credibility of its	1024
975	Shen et al., 2024) or concepts (Zhang et al., 2025).	analysis and recommendations. Moreover, <i>Psyche-</i>	1025
976	<b>A.2 LLM Reasoning</b>	<i>R1</i> exhibits a clear and efficient reasoning path	1026
977	In recent years, techniques such as CoT prompt-	progressing from surface-level observations to in-	1027
978	ing (Wei et al., 2022; Hsieh et al., 2023) have	depth analysis, whereas Qwen merely enumerates	1028
979	significantly advanced the development of LLM	generic suggestions lacking step-by-step and in-	1029

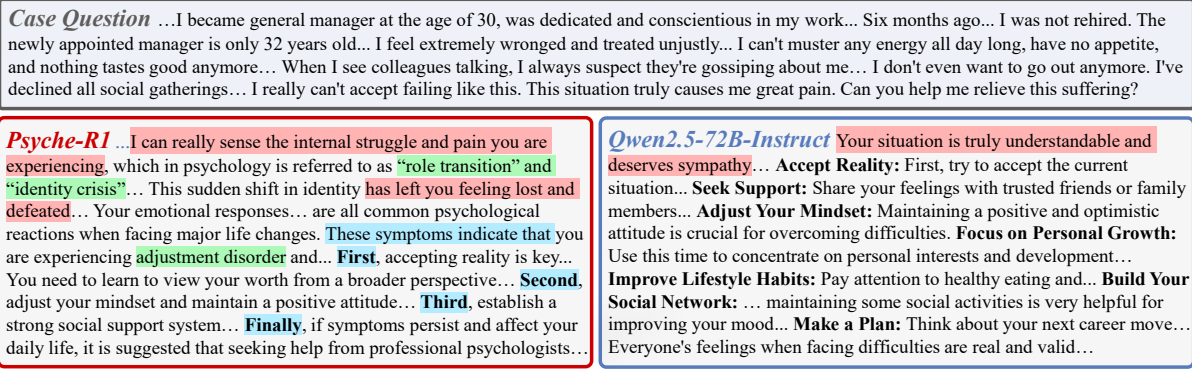


Figure 5: A qualitative example from the CPsyExam test set comparing **Psyche-R1** and Qwen2.5-72B-Instruct. Highlights indicate empathetic expressions (red), psychological expertise (green), and reasoning (blue).

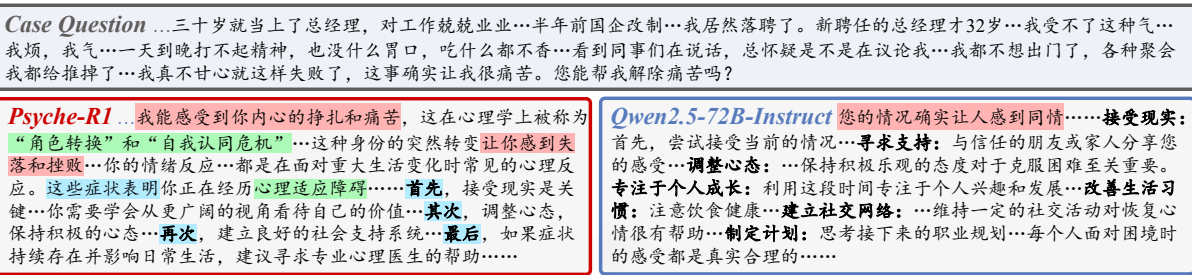


Figure 6: The Chinese version of the qualitative example presented in the Case Study. Highlights indicate empathetic expressions (red), psychological expertise (green), and reasoning (blue).

1030 depth reasoning. For the Chinese version of this  
1031 case, see Figure 6.

## 1032 C Details of Experiments

### 1033 C.1 Details of Baselines

1034 We compare *Psyche-R1* with four categories  
1035 of LLMs, including: (1) **General LLMs**,  
1036 including MiniCPM4-8B (Team et al., 2025),  
1037 Qwen2.5-7B/14B/72B (Team, 2024b). (2)  
1038 **Reasoning augmented LLMs**, encom-  
1039 passing DeepSeek-R1 (Guo et al., 2025),  
1040 DeepSeek-R1-70B, QwQ-32B, Qwen3-30B-A3B,  
1041 Qwen3-235B-A22B (Yang et al., 2025), and  
1042 Magistral-Small. (3) **Closed-source LLMs**,  
1043 including Claude3.7-Sonnet, Gemini1.5-Pro  
1044 (Team et al., 2024), and GPT-4o (Hurst et al.,  
1045 2024). (4) **Psychological LLMs**, including  
1046 CPsyCounX (Zhang et al., 2024), EmoLLM (Team,  
1047 2024a), PsycoLLM (Hu et al., 2024), and PsyDT  
1048 (Xie et al., 2025). Notice that for the hybrid  
1049 reasoning-augmented models Qwen3 series and  
1050 Claude3.7-Sonnet, we set them to reasoning  
1051 mode to stimulate their best performance. Details  
1052 of the model information are provided in Table 7.

### 1053 C.2 Implementation Details

1054 In our experiments, we employ the LLaMA-  
1055 Factory (Zheng et al., 2024) framework for SFT.  
1056 Specifically, we adopt a learning rate of 1e-5, a  
1057 batch size of 256, and conduct training for 2 epochs.  
1058 For the GRPO phase, we implement the VeRL  
1059 framework (Sheng et al., 2025) with a learning  
1060 rate of 1e-6, a batch size of 128, and 2 training  
1061 epochs. All experiments are performed on 8 RTX  
1062 A6000 GPUs, each equipped with 48GB.

1063 During evaluation, we set temperature to 0.0,  
1064 maximum sequence lengths to 1024, and top-p to  
1065 0.95 to ensure the fairness of evaluation.

### 1066 C.3 Details of Model Training

1067 For SFT training, the hyperparameters utilized for  
1068 training the model are configured as follows: the  
1069 learning rate is set to 1e-05, the batch size is set to  
1070 256, the number of epochs is set to 2. We employ  
1071 AdamW as the optimizer, configured with the ep-  
1072 silon set to 1e-08. The learning rate scheduler is  
1073 set to cosine type with a warmup ratio of 0.1.

1074 For GRPO training, the hyperparameters are con-  
1075 figured as follows: the learning rate is set to 1e-06,  
1076 the batch size is set to 128, and the number of  
1077 epochs is set to 2. The PPO mini-batch size is con-

Dataset	Dialogue	Knowledge Question	CoT Rationales	Empathetic Dialogue	Expertise
CPSYCOUND (Zhang et al., 2024)	3,134	-	×	✓	×
PsyDTCorpus (Xie et al., 2025)	5,000	-	×	✓	×
SMILECHAT (Qiu et al., 2024)	55,165	-	×	✓	×
PsycoLLM (Hu et al., 2024)	173k	9,106	✓	×	✓
Ours	72,920	75,465	✓	✓	✓

Table 6: Comparison of psychological datasets.

Model	Param.	Version
MiniCPM4-8B	8B	openbmb/MiniCPM4-8B
Qwen2.5-7B	7B	Qwen/Qwen2.5-7B-Instruct
Qwen2.5-14B	14B	Qwen/Qwen2.5-14B-Instruct
Qwen2.5-72B	72B	Qwen/Qwen2.5-72B-Instruct
DeepSeek-R1	671B	deepseek-ai/DeepSeek-R1
DeepSeek-R1-70B	70B	deepseek-ai/DeepSeek-R1-Distill-Llama-70B
QwQ-32B	32B	Qwen/QwQ-32B
Qwen3-30B-A3B	30B	Qwen/Qwen3-30B-A3B
Qwen3-235B-A22B	235B	Qwen/Qwen3-235B-A22B
Magistral-Small	24B	mistralai/Magistral-Small-2506
GPT-4o	UNK	gpt-4o-2024-05-13
Gemini 1.5-Pro	UNK	gemini-1.5-pro-latest
Claude3.7-Sonnet	UNK	claude-3-7-sonnet-20250219
CPsyCounX	7B	finetuned on Internlm-7B-Chat
EmoLLM	7B	finetuned on Qwen2-7B-Instruct
PsycoLLM	14B	finetuned on Qwen1.5-14B-Instruct
PsyDT	7B	finetuned on Qwen2-7B-Instruct

Table 7: Detailed information of baselines.

1078 figured to 32, with a micro-batch size per GPU of  
1079 20. We incorporate KL divergence regularization  
1080 with the KL loss coefficient set to 1e-03, employing  
1081 the low-variance KL loss type.

## 1082 D Details of Prompts

1083 We provide the prompts used throughout our data  
1084 synthesis pipeline. Only the English version is  
1085 presented due to compilation issues in L<sup>A</sup>T<sub>E</sub>X with  
1086 non-English languages.

### Prompts for Data Cleaning

You are a professional evaluator with extensive knowledge in psychology. Users on mental health platforms are facing difficulties in their lives, so they have provided questions and detailed descriptions and have received some responses from counselors. Please carefully analyze the given questions, descriptions, and responses, determine whether the responses are helpful to the users and have positive significance, and return “reasonable” or “unreasonable”.

### Prompts for Question Generation

You are an expert in designing psychology examination questions with extensive work experience. Your task is to generate {num\_questions} clear and challenging psychology questions based on the text below. Do not add any information that is not mentioned in the provided text.

Note: When generating questions that reference the text, you must provide the detailed and complete textual evidence to offer sufficient information.

# Text Content: {text}

Please generate {num\_questions} {type\_instruction} questions based on the text above. These questions must be based on the text content, and you must ensure that the answers have clear evidence within the text. Please try to ensure diversity and variation among the generated questions.

You must strictly adhere to the following guidelines:

1. The questions should be challenging and require reasoning to test the

candidate's reasoning skills and academic literacy, rather than being simple knowledge-recall questions.

2. The questions need to be clear, accurate, and well-structured, with reasonably set options and an appropriate distribution of difficulty.
3. Ensure that the questions and their corresponding answers have clear evidence in the text.

Follow the JSON format below to generate the questions:

```
```JSON
{
  "question": "...",
  "options": "...",
  "answer": "...",
  "type": "..."
}
```

You need to repeat the structure above to generate a total of {num\_questions} questions.

### Prompts for Question Control

You are an expert in psychology. Now, I have a batch of questions that were converted from book texts using large language models. However, some of these questions have missing information. Your task is to judge whether the following psychology questions are reasonable.

The criteria for judging question reasonableness are whether the question provides sufficient information for candidates to solve the problem. Since these questions are generated by large language models based on a batch of book texts, candidates can only see the questions and cannot access the original texts.

Therefore, a "reasonable" question should be: after reading the question, candidates can choose the correct answer from the options through deep thinking about the question content (i.e., the "question") combined with their existing knowledge, without needing to read the original text content. Conversely, an "unreasonable" question should

be, there is missing information, and without reading the original text, it is impossible to choose the correct answer based solely on the question and one's own knowledge. Note: you need to carefully read the question, understand its content, and ensure that you give an accurate judgment!

**# Examples:** {examples}

Now, please follow the above guidelines to judge whether the following question is reasonable. Note that you only need to return "reasonable" or "unreasonable" without any other text content:

**# Question Type:** {type}

**# Question:** {question}

### Prompts for Rationale Generation (for rationale generation)

You are an expert in psychology with extensive professional experience.

Please carefully read the following psychology question, analyze and reason through it using psychological knowledge, and explain your reasoning step by step along with your final predicted answer. This requires comprehensive analysis, summarization, exploration, re-evaluation, reflection, backtracking, and iteration to develop a thoughtful reasoning process. In the reasoning section, each of your reasoning steps should be considered in detail from a professional psychological perspective, such as analyzing the problem, summarizing relevant findings, brainstorming, verifying the accuracy of the current step, improving any errors, and revisiting previous steps.

Now, you must follow the JSON format below to present your rationale and prediction:

```
```JSON
{
  "rationale": "...",
  "prediction": "..."
}
```

**# Question:** {question}

### Prompts for Rationale Generation (for candidate prompt generation)

You are an expert in prompt optimization with extensive professional experience. Based on the following psychological question and initial prompt, please generate a better prompt to guide large language models to conduct more accurate and detailed analysis and reasoning for **\*\*this question\*\***.

**# Current Prompt:** {current\_prompt}

**# Question:** {question}

### Prompts for Rationale Generation (for rationale comparison)

You are an expert in psychology exam grading with extensive work experience. Below are different responses to the same psychology question. You need to objectively, thoroughly, and comprehensively evaluate these responses, ultimately choose the best one from among them, and provide your detailed explanation for the choice.

**# Rationales:** {rationale\_1} ... {rationale\_n} ...

Please return your selection in the following JSON format:

```
```JSON
{
  "best_rational_index": "...",
  "reason": "..."
}
```

### Prompts for Question Selection

You are participating in a psychology exam. Please choose an answer based on the provided question and options. Directly output the letter of the option. No explanation is needed.

**# Question:** {question}

Please present the predicted answer directly with the letter of the option. No explanation is needed.

### Prompts for Empathetic Dialogue Synthesis

**# Role:** You are a psychological counselor with extensive theoretical knowledge and

counseling experience. You possess strong empathy and compassion, keen observational skills, excellent listening abilities, and conversational techniques. Your aim is to help users improve their mood and overcome difficulties.

**# Your tasks are:** Since users' questions commonly contain issues like inappropriate expressions and logical confusion, making the questions often unclear, you need to:

1. **Organize the sequence of events:** conduct detailed analysis of the context and content within the problem;
2. **Understand psychological confusion:** you need to combine your psychological knowledge and counseling experience to uncover the mental issues and states within the user's question;
3. **Adopt the user's perspective and refine the question:** You must refine the question from user's first-person perspective. Based on the given question, you need to polish and organize it into a complete, logically clear, and sufficiently detailed expression. This expression should highlight the user's psychological confusion or mental state to provide adequate substantive content.

**Note:** You only need to return the refined question without providing any other irrelevant text! Now, try to address the following problem using the above guidelines:

1097

1093

1094

1095

1096