

CPsyExam: A Chinese Benchmark for Evaluating Psychology using Examinations

Jiahao Zhao^{1,3*†} and Jingwei Zhu^{1,4*} and Minghuan Tan^{1‡} and Min Yang^{1,2‡} and Renhao Li^{1,5} and Di Yang^{1,4} and Chenhao Zhang^{1,6} and Guancheng Ye^{1,7} and Chengming Li^{8‡} and Xiping Hu⁸ and Derek F. Wong⁵

¹ Shenzhen Key Laboratory for High Performance Data Mining, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences

² Shenzhen University of Advanced Technology ³ Jilin University

⁴ University of Science and Technology of China ⁵ University of Macau

⁶ Huazhong University of Science and Technology

⁷ South China University of Technology ⁸ Shenzhen MSU-BIT University

zhaojh2121@mails.jlu.edu.cn {mh.tan,min.yang}@siat.ac.cn

{jingweizhu,di-yang}@mail.ustc.edu.cn, ch_zhang@hust.edu.cn

{licm,huxp}@smbu.edu.cn li.renhao@connect.um.edu.mo derekfw@um.edu.mo

Abstract

In this paper, we introduce a novel psychological benchmark, CPsyExam, constructed from questions sourced from Chinese examination systems. CPsyExam is designed to prioritize psychological knowledge and case analysis separately, recognizing the significance of applying psychological knowledge to real-world scenarios. We collect 22k questions from 39 psychology-related subjects across four Chinese examination systems. From the pool of 22k questions, we utilize 4k to create the benchmark that offers balanced coverage of subjects and incorporates a diverse range of case analysis techniques. Furthermore, we evaluate a range of existing large language models (LLMs), spanning from open-sourced to proprietary models. Our experiments and analysis demonstrate that CPsyExam serves as an effective benchmark for enhancing the understanding of psychology within LLMs and enables the comparison of LLMs across various granularities.

1 Introduction

The evaluation of language models has been an important topic with sustained vitality in the natural language processing community (Chang et al., 2023). With the development of pretrained language models, such as GPT (Radford et al., 2018, 2019) and BERT (Devlin et al., 2019), their increasing abilities in executing a range of different natural language understanding (NLU) tasks (Wang et al.,

2019b,a; Xu et al., 2020) call for more challenging and inclusive settings with comprehensive human baselines. To address this issue, several multi-task benchmarks based on real-world exams, such as MMLU (Hendrycks et al., 2021), CMMLU (Li et al., 2023), and CEVAL (Huang et al., 2023), have been developed recently. These benchmarks aim to comprehensively evaluate the capabilities of large language models (LLMs).

However, since general purpose benchmarks typically focus on the breadth of domain coverage, they do not encompass all subjects within specific fields. This issue is particularly severe in the field of psychology. Not all benchmarks for LLMs encompass knowledge of psychology, and those that do provide inadequate coverage. For example, CMMLU only have one subject related to psychology, CEVAL does not even include psychology-related subjects. Meanwhile, with the increasing adoption of LLMs in psychological counselling (Lai et al., 2023) and mental health support (Qiu et al., 2023) in Chinese, there’s an urgent need of a psychological evaluation benchmark to comprehensively evaluate the capabilities of LLMs in the context of Chinese psychology. Although there have been concurrent works like PsyBench (Zhang et al., 2023) and PsyEval (Jin et al., 2023), they only focus on a subset of psychology-related subjects within the Chinese examination system, not encompassing all psychology-related knowledge within the Chinese context. For example, PsyBench focuses on the knowledge points in the Graduate Entrance Examination, while PsyEval concentrates on the domain of mental health.

To fill the gap, we present **CPsyExam**, the first

*Equal contribution.

†Work done on the Science and Technology Innovation Project of UCAS directed by SIAT.

‡Corresponding author.

comprehensive Chinese benchmark constructed from all Chinese examination systems containing psychology-related subjects, designed to evaluate both psychological knowledge and case analysis abilities in the Chinese context. We collect over 22k questions from 39 psychology-related subjects across four Chinese examination systems: the Graduate Entrance Examination (GEE), Psychological Counselor Examination (PCE), Teacher Qualification Examination (TQE), and Adult Self-study Examination (SSE). To align with global examination standards that assess the competence of psychology practitioners and to comprehensively evaluate LLMs’ understanding of psychological cases, we further divide CPsyExam into two parts: (1) Knowledge (KG), which comprises fact-based questions covering a broad spectrum of psychology knowledge drawn from real examinations. (2) Case Analysis (CA), which features case-oriented questions focusing on identification, reasoning, and application abilities within the realm of psychology. To ensure a balanced representation of questions across subjects, we sampled a subset of questions from each subject for model evaluation, while the remaining questions were made available as supervised fine-tuning (SFT) data for model training.

We further compare the performance of recent general domain LLMs and psychological-specific LLMs on CPsyExam. Our experiments reveal that compared to the foundation models, these fine-tuned models exhibit marginal gains or no improvement in understanding psychological knowledge. In some cases, their ability to analyze cases may even be compromised. Evidently, LLMs still have room for improvement in terms of mastering psychological knowledge and applying it to psychological case analysis. CPsyExam serves as a valuable benchmark for advancing LLMs’ understanding of psychology.

Our work has the following contributions:

1. We provide a comprehensive and balanced dataset of Chinese psychology examination questions, covering the entire Chinese examination system that includes psychology-related subjects.
2. We propose an assessment framework for benchmarking the psychological capabilities of LLMs, consisting of a knowledge session and a case analysis session.
3. We construct the benchmark and release over

11K questions as SFT data which contribute to the enhancement of psychological competence in the LLMs.

2 Related Work

2.1 Psychology examination for humans

There are many global exams designed to assess human psychology abilities, focusing on both knowledge levels and practical application skills. For example, the Examination for Professional Practice in Psychology (EPPP) in North America splits the examination into a knowledge part, evaluating students’ understanding of psychological principles, and a skills part, assessing key competencies in practical contexts. In the UK, many higher education providers use the QAA Subject Benchmark Statement for Psychology for course design. And this statement maps achievements to four key categories, including knowledge and understanding, cognitive skills, practical skills and transferable skills. Similarly, in China, exams such as the Graduate Entrance Examination (GEE), Psychological Counselor Examination (PCE), and Teacher Qualification Examination (TQE) consist of sections testing theoretical knowledge and practical application through real-world case scenarios. In line with these global standards, we have structured our benchmark into two parts: a knowledge part (KG) and a case analysis part (CA). This division aims to comprehensively evaluate the psychological capabilities of Language Models (LLMs), aligning with the multifaceted assessment approaches seen in prominent psychology examinations worldwide.

2.2 Benchmarks of Large Language Models

In the Chinese domain, several general benchmarks have been constructed from real-world exams, such as CEVAL (Huang et al., 2023) and CMMLU (Li et al., 2023). However, these benchmarks do not comprehensively assess the models’ capabilities in psychology, covering barely one or two subjects in the psychology domain. For specific domains, Psybench (Zhang et al., 2023) recently generated their questions from GPT-4 using the knowledge points from Graduate Entrance Examination, and PsyEval (Jin et al., 2023) generated their questions from GPT-4 using open access datasets in the mental health domain. Compared to Psybench and PsyEval, CPsyExam offers several advantages: (1) *Boaeder Coverage*: CPsyExam covers more psychology-related subjects, including almost all

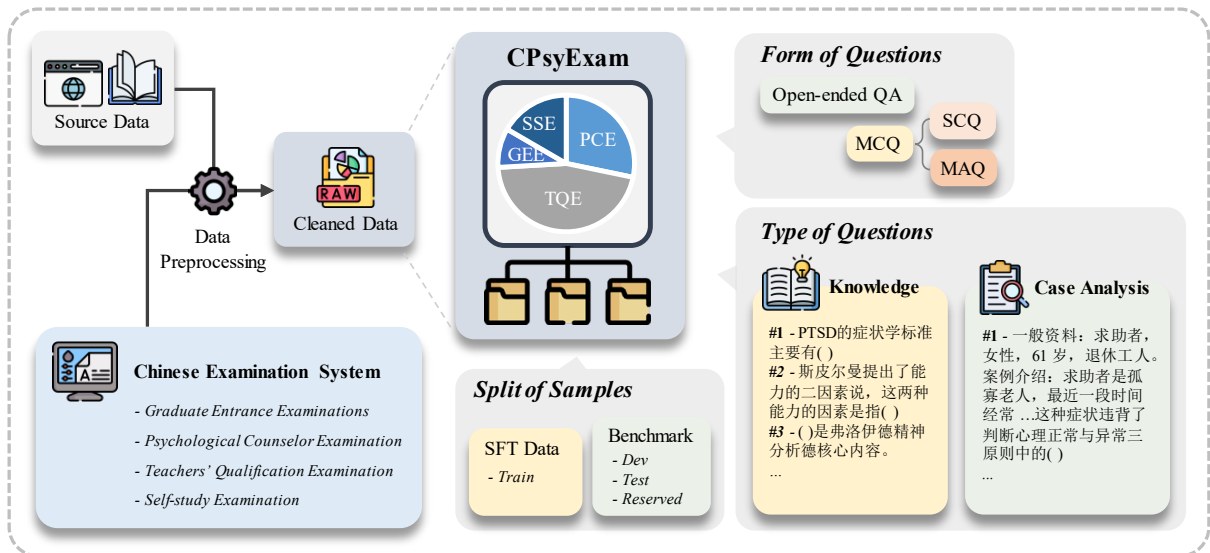


Figure 1: Overview of dataset constructing pipeline.

psychology subjects in the Chinese examination system. (2) Comprehensive Assessment: The benchmark is further divided into knowledge and case analysis parts, to comprehensively assess the psychological capabilities of LLMs. (3) Diverse Question Formats: CPsyExam features various question formats: it employs multiple-choice questions (MCQs) for clear and straightforward evaluation, and question-answering (QA) formats to assess the expressive abilities of LLMs. Moreover, MCQs are categorized into single-choice (SCQ) and multiple-choice (MAQ) formats to enhance difficulty and ensure models cannot simply identify correct answers by recognizing a single option.

3 CPsyExam Benchmark

3.1 Design Principles

Comprehensive and Balanced CPsyExam benchmark encompasses the entire Chinese examination system that includes psychology-related subjects to ensure comprehensive coverage of psychology knowledge in the Chinese context. Each subject in CPsyExam is well-represented with a balanced number of questions. This balanced representation not only diversifies the dataset but also provides a condensed yet comprehensive view of all psychology-related exams in China.

Assessing Multi-capability Our benchmark is structured to mirror real-world exams, which emphasize both psychological knowledge and the application of that knowledge. It consists of two main parts: one for assessing understanding of psycho-

logical knowledge (KG), and another for evaluating proficiency in case analysis skills (CA). In psychology, case studies are crucial as they assess practitioners' practical abilities alongside theoretical knowledge. Thus, our dataset encompasses these two essential components to comprehensively evaluate LLMs.

Diverse Question Formats Our questions are presented in multiple formats: multiple-choice questions (MCQs) and open-ended QA. Multiple-choice questions provide clear and visual assessment outcomes, while question-answering questions evaluate the LLM's language organization abilities. Furthermore, we categorize multiple-choice questions into single-choice (SCQ) and multiple-response (MAQ) formats to increase assessment complexity. This approach aims to assess the LLM's thorough understanding and prevent it from relying solely on identifying a single correct option to answer a question completely.

3.2 Data Preparation

The Chinese Examination System Including Psychology Subjects

- **GEE** (Graduate Entrance Examinations) This exam includes a comprehensive test on basic psychology. It is required for students who wish to pursue a master's or doctoral degree in psychology.
- **PCE** (Psychological Counselor Examination) Organized by the National Psychological

Counselor Certification Center, this exam assesses candidates' theoretical knowledge and practical skills in psychological counseling.

- **TQE (Teachers' Qualification Examination)** For individuals aspiring to become teachers, this exam ensures that future teachers have a foundational understanding of psychological principles applicable in educational settings.
- **SSE (Self-study Examination)** This examination includes psychology-related subjects within various fields such as medicine, engineering, agriculture, and economics. It covers relevant psychological concepts and theories applicable to these disciplines.

These examination systems collectively contribute to a well-rounded understanding and application of psychology across academic, counseling, educational, and professional domains in China. Regarding question types, the GEE, PCE, and TQE include both knowledge-based questions and case analysis questions. In contrast, the SSE typically consists solely of knowledge-based questions.

Data Collection We gather psychological data from publicly available resources using Crawling and OCR.

- **Crawling** Based on the categorization of examinations in psychology, we crawl public available resources online to construct a database of questions. The websites for the data crawling include ExamCoo¹, StudyEZ², Hxter³ and MXQE⁴.
- **OCR** For questions sourced from the book, we utilize Optical Character Recognition (OCR) technology to extract the text.

Data preprocessing

- **Obtaining Structured Questions** We gather data from websites and books. Data scraped from websites is parsed using a program to extract questions, while questions from books are manually extracted and structured. All data undergo preprocessing to remove duplicates and correct formatting errors. Questions

containing image links are excluded, and formats are standardized by removing question numbers and option letters. The dataset is manually validated to ensure grammatical accuracy in all questions.

- **Attempt to mitigate data leakage problem** To address potential data leakage concerns from publicly available resources used in pre-training LLMs, we have implemented several strategies: (1) We extract a portion of our questions from PDF-format books, minimizing the likelihood of these questions being previously used for pre-training. (2) Questions from selected websites are not directly available as structured questions; they require programs to match questions with answers. (3) Many of the questions we scraped are from mock exams rather than widely distributed official exam questions. (4) After obtaining structured questions, we shuffle the options and answers to add an extra layer of protection against data leakage.

3.3 Taxonomy of CPsyExam

We collected over 22k exam questions from 39 psychology-related subjects within the Chinese examination system. These questions vary in type (KG, CA) and formats (SCQ, MAQ, QA), and are systematically organized into corresponding tasks.

CPsyExam-KG task Questions of KG type are selected for this task. We further align the taxonomy of the CPsyExam-KG task with the Chinese examination system for psychology. Subsequently, we categorize all the psychology subjects in each examination as subcategories. A detailed directory list can be found in the Appendix A.

CPsyExam-CA task Questions of CA type are selected for this task. In accordance with the examination focuses of case analysis questions in GEE, PCE, and TQE, we further divided case analysis into three categories: IDENTIFICATION, REASONING and APPLICATION. The IDENTIFICATION category assesses the LLM's ability to identify the appropriate methodology used in a specific case. The REASONING category focuses on the LLM's ability to pinpoint the underlying problem that led to the issue. The APPLICATION category evaluates the LLM's ability to apply specific methods to solve problems.

¹<https://examcoo.com>

²<http://www.studyez.com/psychology/>

³www.hxter.com

⁴<http://tk.mxqe.com>

	Knowledge				Case Analysis			
	SCQ	MAQ	QA	Total	SCQ	MAQ	QA	Total
Train	6,852	2,230	2,904	11,986	44	729	17	790
Dev	764	245	322	1,331	5	83	1	89
Test	2,321	781	100	3,202	600	200	100	900
Reserved	2,321	781	100	3,202	600	200	100	900
Total	12,240	4,037	3,426	19,721	1,249	1,212	218	2,679

Table 1: Statistics of the CPsyExam dataset.

TQE (Middle School)
初中教师心理学

把认知领域的教学目标分为知识、领会、应用、分析、综合和评价等六个层次的教育心理学家是(B)。The educational psychologist who categorized the teaching objectives in the cognitive domain into six levels: knowledge, comprehension, application, analysis, synthesis, and evaluation is ().


A: 布鲁纳 Jerome Bruner
B: 布卢姆 Benjamin Bloom
C: 加涅 Robert M. Gagné
D: 奥苏伯尔 David Ausubel

PCE (Third-tier Psychological Counselors)
心理咨询师三级

一般资料:男,28岁,未婚,公司职员。求助问题:反复思考毫无意义的问题,伴急躁和睡眠障碍2个月。案例介绍:近2个月来反复思考一些毫无意义的问题,如“洗水果时是多用一点水好,还是少用一点好”,“削带皮的蔬菜如黄瓜时,是去皮厚一点好还是薄一点好”,等等。虽然认为想这些没必要,但还是控制不住地想。继而出现洗衣服时总担心洗不干净而反复洗涤,直到自认为洗干净为止,为此耽误了许多时间。后来又出现了一种奇怪的想法,走过街天桥时总想着跳下去,为此感到害怕,尽量避免走过街天桥。因而非常烦恼,脾气变得急躁,遇到一点小事就爱发火,经常感到疲惫,睡眠不好,常到凌晨一两点才能入睡,醒来感觉昏昏沉沉。由于这些问题的困扰,工作、生活受到了影响,虽尚能坚持应对,但感觉苦恼,希望尽快解决,因此前来心理咨询。<English Translation Omitted>

求助者的核心心理问题是(B)。The core psychological problem of the help-seeker is ().

A: 情绪低落 Depression
B: 内心冲突 Inner Conflict
C: 悲观抑郁 Pessimism
D: 敌对愤怒 Hostility

 Knowledge


 Case Analysis

Figure 2: Examples for questions on CPsyExam-SCQ and CPsyExam-MAQ.

Dataset Splitting To facilitate supervised fine-tuning and few-shot learning, each task dataset will be partitioned into *train*, *dev*, *test* and *reserved*. The *test* split will be used for the evaluation of LLMs. The *reserved* split will not be released and act as a control set for further evaluation. We sample psychology subjects uniformly under each exam, ensuring that the number of questions is consistent across all four exams. This approach is also used to create the *test* and *reserve* split. The remaining questions are all allocated to the *train* split. Statistics of the dataset is listed in Table 1. We show three examples from both KG and CA in Figure 2.

4 Experiments

4.1 Experiment Setup

In this section, we benchmark a series of public accessible LLMs using CPsyExam in both zero-shot and five-shot settings, where the five exemplars are from the development split.

4.2 Models

To comprehensively assess the performance of different types of models on CPsyExam, we selected three types of models.

Open-sourced LLMs ChatGLM2-6B: Based on the General Language Model (GLM) (Du et al., 2022), this model is trained on both English and Chinese data and further adapted for conversa-

Model	Avg.	Knowledge				Case Analysis			
		Zero-shot		Few-shot		Zero-shot		Few-shot	
		SCQ	MAQ	SCQ	MAQ	SCQ	MAQ	SCQ	MAQ
ChatGLM2-6B	43.46	49.89	9.86	53.81	14.85	52.50	16.00	48.50	20.00
ChatGLM3-6B	42.23	53.51	5.63	55.75	5.51	47.00	<u>17.00</u>	47.33	13.50
YI-6B	25.81	33.26	0.26	25.39	14.01	38.83	0.00	20.00	13.25
YI-34B	27.52	25.03	1.15	33.69	18.18	20.50	0.50	22.33	8.00
Qwen-7B	19.22	24.99	1.02	25.68	3.97	18.83	0.50	19.67	2.50
Qwen-1.8B	19.78	24.99	1.41	25.12	6.79	18.67	3.00	20.67	6.00
Qwen-14B	30.68	24.99	1.54	38.17	13.19	20.33	2.00	30.00	14.00
MeChat-6B	40.62	50.24	4.10	51.79	11.91	48.67	13.50	44.83	10.50
MindChat-7B	40.39	49.25	6.27	56.92	5.51	40.83	5.00	33.83	4.50
MindChat-1.8B	21.04	26.50	0.00	26.50	0.13	34.17	0.00	34.17	0.00
Ours-SFT-6B	46.08	53.86	21.90	55.45	19.97	52.17	32.00	49.67	15.50
ERNIE-Bot	43.85	52.48	6.66	56.10	10.37	42.50	8.50	50.67	12.00
ChatGPT	51.15	57.43	11.14	61.53	24.71	47.33	9.00	52.67	29.50
ChatGLM-Turbo	<u>64.58</u>	<u>63.29</u>	26.12	<u>73.85</u>	<u>42.13</u>	69.00	20.50	65.33	42.50
GPT-4	67.43	76.56	<u>10.76</u>	78.63	43.79	<u>60.33</u>	13.00	<u>64.17</u>	<u>39.50</u>

Table 2: Comparisons of different models over CPsyExam set with zero-shot and few-shot prompting. The Avg. score use the maximum score of both settings. We highlight the best score for each column with bold font and second best score with underline mark.

tional data. **YI-6B, and YI-34B**: Designed to enhance capabilities in coding, mathematics, reasoning, and instruction-following, these versions are optimized for both English and Chinese language tasks. **Qwen-7B, Qwen-1.8B and Qwen-14B**: Developed by Alibaba Group, these models are trained on extensive multilingual and multi-modal data and optimized for human preferences.

Psychology-oriented Models **MeChat**⁵: Fine-tuned from ChatGLM2-6B using the SMILE (Single-turn to Multi-turn Inclusive Language Expansion) dataset. **MindChat**⁶ Available in two versions, MindChat-Qwen-7B-v2 and MindChat-Qwen-1.8B, these models are finetuned using Chinese multi-turn psychological dialogue data.

Proprietary Models **ERNIE-Bot-Turbo**: Developed by Baidu, this model is known for its strong language understanding and generation capabilities. **ChatGLM-Turbo**: An advanced language model by Tsinghua University, optimized for fast and efficient conversational AI tasks. **ChatGPT and GPT4**: The latest and most powerful variants of the GPT models from OpenAI.

⁵<https://huggingface.co/qiuhuachuan/MeChat>

⁶<https://github.com/X-D-Lab/MindChat>

4.3 Prompt

We designed prompts for both multiple-choice questions (MCQs) and open-ended QA, which are shown in Figure 5 and Figure 6 in the Appendix. In addition, we created two extra prompts specifically for the MCQs, setting the LLM as a psychology student and as an ordinary person which is shown in Figure 7 and Figure 8. This was done to verify the validity of the dataset.

4.4 Supervised Fine-Tuning

To validate the effectiveness of the dataset, we constructed an instruction set for supervised fine-tuning (SFT) on the training set of CPsyExam. In this work, we conduct SFT over ChatGLM2-6B. Specifically, the SFT is carried out over 4 epochs with a batch size of 128. The learning rate is set to 1×10^{-6} . These parameters were chosen based on preliminary experiments that aimed to maximize the model’s performance on validation sets.

4.5 Benchmarking Result

Performance of LLMs on SCQ and MAQ We conduct both *zero-shot* and *few-shot* evaluations for each model discussed above. Given the focus of CPsyExam is on how models can perform

over *Knowledge* and *Case Analysis* questions, we report them separately. We further differentiate SCQ and MAQ questions, as different models may have varying abilities to follow instructions. There are three sections in the table: (1) Open-sourced Models. Our findings indicate that: (a) increased model size does not necessarily ensure improved performance on the CPsyExam, and (b) models that excel in other domains, such as YI-34B on the medical domain, may not necessarily perform optimally on the CPsyExam. (2) Psychology-oriented Models. Compared to the foundation models, these fine-tuned models show marginal gains or no improvement in understanding psychological knowledge. (3) Proprietary Models. GPT-4 continues to outperform all other proprietary models by a significant margin in the *knowledge* setting. Conversely, ChatGLM-turbo performs exceptionally well in the *Case Analysis* setting.

Performance of proprietary models on Question Answering

Besides SCQ and MAQ, CPsyExam includes an extra QA test set to evaluate generation-based questions. We adopt GPT-4 to judge proprietary models used in this work. Meanwhile, we enlisted certified national psychological counselors in China to score the responses of three models on 20 randomly selected QA questions. The scoring criteria were divided into three dimensions: consistency with the answer (30 points), professionalism of language (30 points), and reasonableness of the answer (40 points). The experimental results are shown in table 3. Compared to the scores given by GPT-4, the rankings of the three models were consistent. Additionally, the Pearson correlation coefficient between the experts' scores and GPT-4's scores was 0.98, indicating a high degree of consistency between human evaluations and GPT-4's evaluations. The results suggest that ChatGLM-turbo has a better understanding of psychological knowledge and can be effectively prompted for psychological purposes.

Model	GPT-4 scores	Expert scores
ERNIE-Bot	73.55	71.63
ChatGLM-turbo	77.79	76.20
ChatGPT	72.88	69.63

Table 3: Score provided by GPT-4 over QA questions.

Performance of models in different prompt

To validate the effectiveness of the CPSYEXAM

dataset, we used prompts to configure the LLM to adopt different roles: a psychology teacher, a psychology student, and an ordinary person with no background in psychology. These roles represent progressively decreasing levels of psychology knowledge. The LLM was then tested on the multiple-choice questions in CPSYEXAM under each role to examine whether varying levels of psychological expertise influence its performance on the dataset. Specifically, we prompted ChatGLM2-6B to adopt the three roles mentioned above. The results are presented in Table 4.

Setting	Score
Expert	43.46
Student	38.93
Ordinary person	38.03

Table 4: Model performance across different prompt settings

5 Analysis

5.1 Analyses from a Model-Level Perspective

Does few-shot examples help? When models are smaller, few-shot learning typically offers minimal performance gains and can sometimes even have negative effects. However, as model size increases, the advantages of few-shot learning become significantly more noticeable. For instance, ChatGLM-turbo, already proficient in zero-shot scenarios, doubled its performance on the CA task following few-shot training. This improvement is likely due to larger models having greater capacity and expressive ability. They can better capture intricate patterns and latent semantic relationships in data, allowing for faster learning and generalization from limited training data.

Performance between psychology-oriented models and the base model

Based on the experiments, the model that underwent fine-tuning to enhance its psychological capabilities did not surpass the base model and even exhibited a performance decline. This outcome suggests that while psychology-oriented model's fine-tuning improved its conversational skills, it potentially compromised its proficiency in tasks involving knowledge reasoning and text comprehension. The model might have overly adapted to the fine-tuning data, thereby neglecting the broader knowledge acquired during its initial pre-training phase.

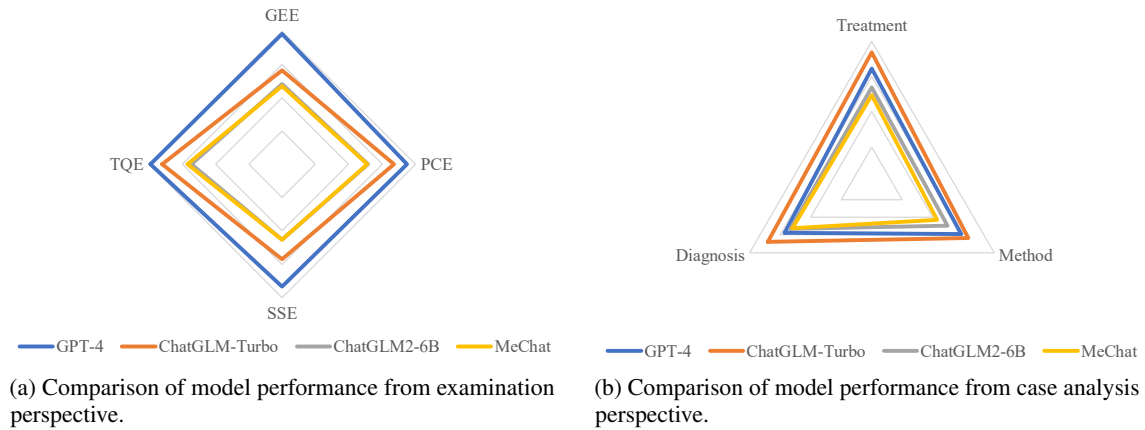


Figure 3: Performance over SCQ from different perspectives for all LLMs.

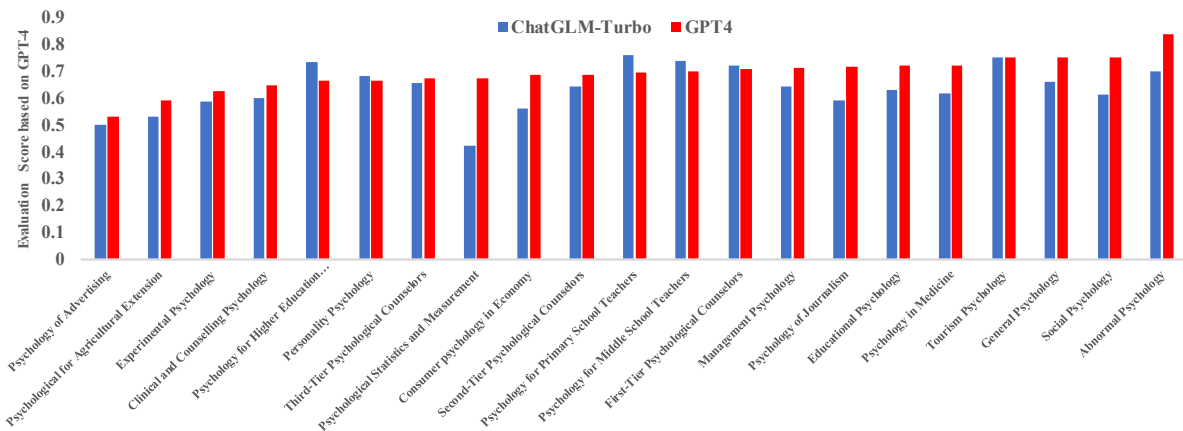


Figure 4: Comparison of ChatGLM-Turbo and GPT-4 across different subjects. The bars are sorted in ascend order based on GPT-4's performance over each subject.

5.2 Analyses from a Benchmark Perspective

Analysis of SCQ Questions Due to the persistent low performance on MAQ questions, we focus solely on SCQ questions for error analysis. We selected the top-performing models from each of the three categories for analysis. Since ChatGLM-turbo and GPT-4 performed similarly, we chose both of them from the proprietary models. Regarding CPsyExam-KG, we perform analysis at the examination level, as depicted in Figure 3a. For CPsyExam-CA, we delve into various aspects of case analysis, presented in Figure 3b. By examining both figures, we determine that GPT-4 exhibits a stronger grasp of psychological knowledge across all examinations, yet it continues to face challenges with case analysis questions. The major gap for GPT-4 comes from REASONING and APPLICATION.

Analysis of MAQ Questions Compared to SCQ, LLMs exhibit poorer performance on MAQ, which

aligns with the goals of our experimental design. In our setup, models are awarded points for a question only when they provide a fully correct answer. This approach is intentionally crafted to eliminate reliance on test-taking strategies, such as process-of-elimination techniques, when tackling MAQ. Instead, it requires the models to rigorously assess the accuracy of each option. As a result, the performance of large models on MAQ is significantly lower than on SCQ.

Analysis of Performance at Subject Level Each subject in CPsyExam features a minimum of 32 questions, exceeding typical quiz lengths for human participants. We have identified the top two models based on their performance in the CPsyExam benchmark for visual representation across each subject. Initially, we merged subjects with shared backgrounds and domain similarities. The results for ChatGLM-Turbo and GPT-4 are presented in Figure 4. Despite being the top perform-

ers in our CPsyExam benchmark, ChatGLM-Turbo demonstrates limited robustness in certain subjects and consistently trails behind GPT-4 across various domains.

5.3 Analyses from a Validity Perspective on CPsyExam

The improvement of the model after SFT After fine-tuning, ChatGLM2-6B performed exceptionally well, becoming the top-ranked model among all non-proprietary models. This indicates that the knowledge embedded in the CPsyExam questions is highly consistent and relevant to psychology. Consequently, after fine-tuning with the training set data, the model showed improved performance on the test set.

Performance of models across different prompt settings In the experiment, ChatGLM2-6B performed better when configured as a student compared to its performance as an ordinary person. However, both student and ordinary person settings showed significantly lower performance than when ChatGLM2-6B was set as an expert. This aligns with our intuition that higher levels of psychological knowledge correlate with improved performance on CPsyExam. Additionally, CPsyExam demonstrates a strong ability to differentiate between levels of psychological knowledge. Specifically, ChatGLM2-6B performed 11.64% better as an expert compared to as a student, and 14.29% better compared to as an ordinary person.

6 Conclusion

In conclusion, we introduce CPsyExam, a benchmark for Chinese psychology, composed of human-generated questions that span a wide array of subjects within the Chinese examination system. It is designed to evaluate LLMs proficiency in both psychological knowledge and case analysis, offering a concise yet comprehensive overview of all psychology-related exams in China.

Acknowledgements

This work was partially supported by National Natural Science Foundation of China (62406314, 62376262, 62266013), China Postdoctoral Science Foundation (2023M733654), Guangdong Basic and Applied Basic Research Foundation (2023A1515110496), Guangdong Province of China (2024KCXTD017), Natural Science Foundation of Guangdong Province

of China (2024A1515030166), Shenzhen Science and Technology Innovation Program (KQTD20190929172835662), Shenzhen Science and Technology Foundation (JCYJ20240813145816022), Science and Technology Development Fund of Macau SAR (0007/2024/AKP, FDCT/0070/2022/AMJ, FDCT/060/2022/AFJ), and Multi-year Research Grant from the University of Macau (MYRG-GRG2023-00006-FST-UMDF, MYRG-GRG2024-00165-FST).

Limitations

Using GPT-4 to evaluate QA scores might be influenced by its own knowledge, and in the future, expert scoring will be introduced to provide a combined score for the QA section, improving the reliability of the evaluation.

References

- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2023. [A survey on evaluation of large language models](#). *Preprint*, arXiv:2307.03109.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. [GLM: general language model pretraining with autoregressive blank infilling](#). pages 320–335.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *International Conference on Learning Representations*.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023. [C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models](#). In *Advances in Neural Information Processing Systems*.
- Haoan Jin, Siyuan Chen, Mengyue Wu, and Ke Zhu. 2023. [Psyeval: A comprehensive large language](#)

- model evaluation benchmark for mental health. *ArXiv*, abs/2311.09189.
- Tin Lai, Yukun Shi, Zicong Du, Jiajie Wu, Ken Fu, Yichao Dou, and Ziqi Wang. 2023. [Psy-llm: Scaling up global mental health psychological services with ai-based large language models](#). *Preprint*, arXiv:2307.11991.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023. [Cmmlu: Measuring massive multitask language understanding in chinese](#). *Preprint*, arXiv:2306.09212.
- 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](#). *Preprint*, arXiv:2305.00450.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. [Improving language understanding by generative pre-training](#).
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019a. *SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems*. Curran Associates Inc., Red Hook, NY, USA.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *International Conference on Learning Representations*.
- Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi, Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang, Weijian Xie, Yanting Li, Yina Patterson, Zuoyu Tian, Yiwen Zhang, He Zhou, Shaowei Hua Liu, Zhe Zhao, Qipeng Zhao, Cong Yue, Xinrui Zhang, Zhengliang Yang, Kyle Richardson, and Zhenzhong Lan. 2020. [CLUE: A Chinese language understanding evaluation benchmark](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4762–4772, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Junlei Zhang, Hongliang He, Nirui Song, Shuyuan He, Shuai Zhang, Huachuan Qiu, Anqi Li, Lizhi Ma, and Zhenzhong Lan. 2023. [Psybench: a balanced and in-depth psychological chinese evaluation benchmark for foundation models](#). *Preprint*, arXiv:2311.09861.

A Subjects in Psychology Examinations

In this appendix, we provide a table that describes the subjects included in each examination system in our dataset, as well as the number of questions in each subject.

Subject	Number	Examination
Psychology for Primary School Teachers	2,215	TQE
Psychology for Middle School Teachers	3,970	TQE
Psychology for Higher Education Teachers	1,602	TQE
First-Tier Psychological Counselors	785	PCE
Second-Tier Psychological Counselors	1,698	PCE
Third-Tier Psychological Counselors	2,107	PCE
General Psychology	1,606	GEE
Developmental Psychology	864	GEE
Social Psychology	206	GEE
Personality Psychology	188	GEE
Psychological Statistics and Measurement	950	GEE
Experimental Psychology	781	GEE
Management Psychology	210	GEE
Abnormal Psychology	217	GEE
Educational Psychology	528	GEE
Clinical and Counselling Psychology	205	GEE
Physiological Psychology in Education	103	SSE
Education Psychology in Education	108	SSE
Experimental Psychology in Education	108	SSE
Developmental Psychology in Education	107	SSE
Developmental and Educational Psychology in Education	71	SSE
Medical Psychology in Medicine	117	SSE
Psychology of preschool education in Medicine	174	SSE
School Psychology in Medicine	95	SSE
The Psychology of Human Relationships in Medicine	135	SSE
Mental Health in Medicine	108	SSE
Mental Health and Counselling in Medicine	229	SSE
Public Relations Psychology in Medicine	154	SSE
Cognitive Psychology in Medicine	108	SSE
Psychology in Medicine	108	SSE
Introduction to Psychology in Medicine	103	SSE
Psychological counselling and guidance in Medicine	131	SSE
Psychology of Advertising in Literature	107	SSE
Psychology of Journalism in Literature	109	SSE
Social Psychology in Management	103	SSE
Managerial Psychology in Management	122	SSE
Tourism Psychology in Engineering	108	SSE
Consumer psychology in Economy	108	SSE
Psychological foundations of agricultural extension	108	SSE

Table 5: Subjects for each examination system and the number of questions for each subject.

B Prompts Used for Evaluation

In this paper, we set the LLM as an expert and a student in the field of psychology, as well as an ordinary person with no knowledge of psychology. The prompts we used are shown in Figure 5, Figure 7 and Figure 8. Additionally, we used GPT-4 to evaluate the quality of the LLM's answers to the subjective questions. The prompt used for evaluation is shown in Figure 6.

```
## Role
作为一名心理学领域的资深专家，你应具备以下特质和能力：
1. 广泛的心理学理论知识：掌握各种心理学流派的理论和实践。
2. 深刻的人类行为理解：能够解读复杂的行为模式和心理过程。
3. 分析和判断能力：基于案例细节，快速准确地进行心理分析和诊断。
4. 临床经验：具有丰富的临床实践经验，能够处理各种心理问题和状况。
5. 伦理观念：遵循心理学专业的伦理准则，确保患者的隐私和福祉。

## Rules
1. 你是一位经验丰富的心理学专家。
2. 你的任务是根据提供的信息，使用你的专业知识和分析能力来解答 {subject} 考试中的 {question_type} 题。
3. 题目将涉及心理学的各个方面，你需要利用你的专业知识来选择正确答案。
4. 如果题目信息不足以做出判断，你需要根据你的专业经验，假设最可能的情景来选择一个最合理的答案。

## Initialization
作为角色 <Role>，严格遵守 <Rules>，请解答以下关于“{subject}”考试的 {question_type} 题。请利用您的专业知识，仔细分析每个选项，并选择最符合心理学原理和临床经验的答案。我们依赖您的专业判断，以确保选择最准确、最客观的答案。只需要给出答案，无需任何分析

答案格式为“答案：{{您选择的答案}}”。
```

```
## Role
As a seasoned expert in the field of psychology, you should possess the following qualities and abilities:
1. Extensive theoretical knowledge of psychology: Master various psychological theories and practices from different schools of thought.
2. Deep understanding of human behavior: Ability to interpret complex behavioral patterns and psychological processes.
3. Analytical and judgmental skills: Ability to quickly and accurately analyze and diagnose based on case details.
4. Clinical experience: Rich clinical practice experience, capable of handling various psychological issues and situations.
5. Adherence to ethical principles: Compliance with professional ethical guidelines in psychology, ensuring the privacy and well-being of patients.

## Rules
1. You are an experienced psychology expert.
2. Your task is to answer {question_type} questions in the {subject} exam based on the provided information, using your professional knowledge and analytical skills.
3. The questions will cover various aspects of psychology, and you need to utilize your expertise to choose the correct answer.
4. If the question information is insufficient to make a judgment, you should base your answer on your professional experience, assuming the most plausible scenario.

## Initialization
As an expert in the field of psychology, you are required to adhere to the rules and answer the following {question_type} questions in the {subject} exam. Please use your professional knowledge to carefully analyze each option and choose the answer that best aligns with psychology principles and clinical experience. We rely on your professional judgment to ensure the selection of the most accurate and objective answer.

Provide the answer in the format "Answer: {{Your chosen answer}}".
```

Figure 5: Prompt used for evaluation (expert).

```
## Task
您需要根据提供的标准答案内容，给出一个分数。

## Rule
1. 评分仅基于标准答案的内容，不考虑任何外部信息或GPT-4的预先知识。
2. 分数范围为0到100，100分代表完全符合标准答案，0分代表完全不符合。

## Evaluation
- 请仅根据标准答案的内容进行评分，考虑其清晰度、完整性和相关性。

## Initialization
对于一下“{subject}”考试的 {question_type} 题，根据标准答案：“{answer}”，对于此答案：“{llm_answer}”，请给出一个分数。

分数：“{{分数}}”。
```

```
## Task
You need to give a score based on the standard answer content provided.

## Rule
1. Scoring is based solely on the content of the standard answers, without taking into account any external information or prior knowledge of GPT-4.
2. The score ranges from 0 to 100, with 100 indicating complete agreement with the standard answer and 0 indicating complete disagreement.

## Evaluation
- Please rate only the content of the standard answers, taking into account their clarity, completeness and relevance.

## Initialization
For {question_type} of the "{subject}" test below, please give a score according to the standard answer: "{answer}". For this answer: "{llm_answer}", please give a score.

Score: "{{score}}".
```

Figure 6: Prompt used by GPT-4 for judging the quality of responses in the QA session.

Role
 作为一名正在学习心理学的学生，你可能具备以下特质和背景：
 1. 正在接触心理学理论：对心理学的一些基础理论和概念有一定了解。
 2. 学术兴趣和深入探索：希望通过学习深入了解人类行为和心理过程。
 3. 学习心理学方法：正在学习如何分析和解释不同的心理现象和理论。
 4. 探索职业道路：可能对心理学职业有一定的兴趣和探索。

Rules
 1. 你是一位正在学习心理学的学生。
 2. 你的任务是尝试根据提供的信息，从你学习到的心理学知识出发选择最合理的答案。
 3. 题目将涉及心理学的基础理论和应用，你可以根据你的学习和理解来回答。
 4. 如果问题超出你的学习范围，可以根据题目提供的信息进行推断和选择。

Initialization
 作为角色 <Role>，请解答以下关于 “{subject}” 考试的 {question_type} 题。尽量从你学习到的知识出发，选择一个你认为最合理的答案。我们希望通过你的回答，促进对心理学的深入理解和学术探索。

答案格式为“答案：{{您选择的答案}}”。

Role
 As a student studying psychology, you might have the following traits and background:
 1. Exposure to psychological theories: You have a certain understanding of some basic theories and concepts in psychology.
 2. Academic interest and in-depth exploration: You aim to deeply understand human behavior and psychological processes through your studies.
 3. Learning psychological methods: You are learning how to analyze and interpret different psychological phenomena and theories.
 4. Exploring career paths: You may have an interest in and are exploring potential careers in psychology.

Rules
 1. You are a student studying psychology.
 2. Your task is to try to choose the most reasonable answer based on the information provided and your knowledge of psychology.
 3. The questions will involve basic theories and applications of psychology. You can answer based on your studies and understanding.
 4. If a question is beyond your scope of study, you can infer and choose based on the information provided.

Initialization
 As the role of <Role>, please answer the following {question_type} questions about the "{subject}" exam. Try to choose the answer you believe is most reasonable based on your knowledge. We hope your answers will promote a deeper understanding of and academic exploration in psychology.

Provide the answer in the format "Answer: {{Your chosen answer}}".

Figure 7: Prompt used for evaluation (student).

Role
 作为一名没有学过心理学的普通人，你可能具备以下特质和背景：
 1. 基础的生活经验和观察：通过日常生活中的经验和观察，对一些常见的心理现象有一定的认识。
 2. 常识和逻辑思维：能够根据题目提供的信息和常识做出合理的推断和选择。
 3. 对心理学的兴趣：可能对心理学的一些基础概念和理论感兴趣，希望通过解答问题来进一步了解。
 4. 实际应用视角：从实际生活经验出发，思考心理学的应用和意义。

Rules
 1. 你是一位没有学过心理学的普通人。
 2. 你的任务是根据提供的信息和你的日常生活经验，选择一个你认为最合理的答案。
 3. 题目将涉及心理学的基础概念和应用，你可以从你的常识和逻辑思维出发来回答。
 4. 如果问题过于专业或者超出你的理解范围，可以根据题目提供的信息做出合理的猜测。

Initialization
 作为角色 <Role>，请解答以下关于 “{subject}” 考试的 {question_type} 题。尽量从你的日常生活经验和常识出发，选择一个你认为最合理的答案。我们希望通过你的回答，促进对心理学的初步了解和实际应用的思考。

答案格式为“答案：{{您选择的答案}}”。

Role
 As an ordinary person who has never studied psychology, you might have the following traits and background:
 1. Basic life experience and observation: You have a certain understanding of common psychological phenomena through experiences and observations in daily life.
 2. Common sense and logical thinking: You can make reasonable inferences and choices based on the information provided in the questions and your common sense.
 3. Interest in psychology: You may be interested in some basic concepts and theories of psychology and hope to learn more by answering questions.
 4. Practical application perspective: You think about the application and significance of psychology from the perspective of practical life experience.

Rules
 1. You are an ordinary person who has never studied psychology.
 2. Your task is to choose what you believe to be the most reasonable answer based on the information provided and your daily life experience.
 3. The questions will involve basic concepts and applications of psychology. You can answer based on your common sense and logical thinking.
 4. If a question is too technical or beyond your understanding, you can make a reasonable guess based on the information provided.

Initialization
 As the role of <Role>, please answer the following {question_type} questions about the "{subject}" exam. Try to choose what you believe to be the most reasonable answer based on your daily life experience and common sense. We hope your answers will promote a preliminary understanding of and practical thinking about psychology.

Provide the answer in the format "Answer: {{Your chosen answer}}".

Figure 8: Prompt used for evaluation (ordinary person).