# PsychBench: A comprehensive and professional benchmark for evaluating the performance of LLM-assisted psychiatric clinical practice

Ruoxi Wang[1†], Shuyu Liu[1†], Ling Zhang[2,3†], Xuequan Zhu[2,3†], Rui Yang[2,3],

Xinzhu Zhou[2], Fei Wu[1], Zhi Yang[2,3], Cheng Jin[1,2,4*], Gang Wang[2,3*]

[1*]School of Biomedical Engineering, Shanghai Jiao Tong University, Shanghai, 200240, China.

[2*]Beijing Key Laboratory of Mental Disorders, National Clinical Research Center for Mental Disorders National Center for Mental Disorders, Beijing Anding Hospital, Capital Medical University, Beijing, 100088, China.

[3*]Advanced Innovation Center for Human Brain Protection, Capital Medical University, Beijing, 100088, China.

[4*]Shanghai Artificial Intelligence Laboratory, Shanghai, 200232, China

## Abstract

The advent of Large Language Models (LLMs) offers potential solutions to address problems such as shortage of medical resources and low diagnostic consistency in psychiatric clinical practice. Despite this potential, a robust and comprehensive benchmarking framework to assess the efficacy of LLMs in authentic psychiatric clinical environments is absent. This has impeded the advancement of specialized LLMs tailored to psychiatric applications. In response to this gap, by incorporating clinical demands in psychiatry and clinical data, we proposed a benchmarking system, PsychBench, to evaluate the practical performance of LLMs in psychiatric clinical settings. The PsychBench is composed of a comprehensive dataset and an evaluation framework. The dataset includes 300 real-world patient cases sourced from three geographically diverse medical centers across northern, central, and southern China, ensuring broad regional and cultural representation. The evaluation framework encompasses five key clinical tasks—clinical text understanding and generation, principal diagnosis, differential analysis, medication recommendation, and long-term course management—each supported by psychiatry-specific quantitative evaluation metrics to ensure rigorous performance assessment. Together, the dataset and framework provide a robust system for evaluating the application of LLMs in psychiatric clinical tasks. We conducted a comprehensive quantitative evaluation of 16 LLMs using PsychBench, and investigated the impact of prompt design, chain-of-thought reasoning, input text length, and domain-specific knowledge fine-tuning on model

performance. Through detailed error analysis, we identified strengths and potential limitations of the existing models and suggested directions for improvement. Subsequently, a clinical reader study involving 60 psychiatrists of varying seniority was conducted to further explore the practical benefits of existing LLMs as supportive tools for psychiatrists of varying seniority. Through the quantitative and reader evaluation, we show that while existing models demonstrate significant potential, they are not yet adequate as decision-making tools in psychiatric clinical practice. The reader study further indicates that, as an auxiliary tool, LLM could provide particularly notable support for junior psychiatrists, effectively enhancing their work efficiency and overall clinical quality. To promote research in this area, we will make the dataset and evaluation framework publicly available, with the hope of advancing the application of LLMs in psychiatric clinical settings.

# Introduction

In recent years, the prevalence of mental disorders has been steadily increasing, becoming a major global public health challenge[1,2]. However, this rising number of patients contrasts sharply with the relative scarcity of mental health resources, particularly in terms of the availability of psychiatrists and access to specialized care[3,4]. This imbalance has driven the exploration of new technologies in psychiatric practice. Against this backdrop, the emergence of LLMs presents a new potential solution to this issue. Given the heavy reliance on verbal communication and text analysis in psychiatric care, LLMs demonstrate a greater application advantage in supporting diagnosis, treatment, and patient management in psychiatry compared to other fields focused on organic diseases. By rapidly analyzing and interpreting patients' emotional expressions, thought patterns, and linguistic features, LLMs can offer real-time, intelligent decision support for psychiatrists[5-12]. However, to implement LLMs effectively in psychiatric practice, it is essential to ensure their comprehensive and reliable performance, which necessitates systematic and scientific evaluation. Currently, research on evaluating the performance of LLMs in psychiatric applications is still in its early stages, lacking sufficient empirical evidence and evaluation frameworks. This underscores the urgent need of an evaluation framework to explore and validate the feasibility and effectiveness of LLMs in psychiatric clinical practice.

At present, the evaluation of LLMs mainly revolves around standardized exams and simulated clinical data, where models are presented with straightforward information and multiple-choice options, requiring little in-depth analysis to reach an answer. Multiple studies have shown that LLMs perform exceptionally well on these tests, achieving results comparable to, or even surpassing, those of human doctors in medical knowledge and diagnostic reasoning, suggesting strong capabilities in processing medical information[6,7,13-21]. However, these evaluation methods mainly emphasize static and general knowledge assessment, fails to fully capture the model's response when faced with complex patient scenarios in real-world medical settings[22,23], especially in psychiatry, which requires long-term, multi-turn interactions and the ability to handle diverse patient backgrounds. Some studies have assessed the performance of LLMs in making clinical decisions for bipolar disorder and offering treatment recommendations for mild depression using hypothetical clinical vignettes[8,10]. There is still a need for comprehensive evaluations and standardized guidelines to assess the effectiveness and safety of using LLMs in real-world psychiatric practice.

To develop a comprehensive, professional, and reliable evaluation system for LLMs in the clinical field of psychiatry, it is essential to base the framework on real clinical data, strictly adhere to clinical guidelines, and fully account for the diversity and complexity of actual clinical needs[24-26]. Based on discussions within our clinical committee, the key requirements for LLMs to assist psychiatric practice can be summarized into five aspects. Firstly, LLMs should reduce the burden of clinical documentation, as psychiatrists spend considerable time drafting medical records[12]. By automating document generation, LLMs can allow psychiatrists to save time in writing medical records and thus spend more time with patients. Secondly, LLMs must provide robust diagnostic support, as diagnosing psychiatric disorders often involves interpreting complex symptoms and subjective descriptions. Accurate diagnostic capabilities are crucial for helping psychiatrists analyze patient symptoms effectively. Additionally, LLMs should offer comprehensive differential analysis, assisting psychiatrists in ruling out misdiagnosis in the context of overlapping psychiatric symptoms. Moreover, LLMs need to deliver scientific medication recommendations. Given the complexities of psychiatric pharmacotherapy, LLMs should provide personalized medication plans tailored to individual patient needs. Finally, LLMs should support long-term course management by rapid and reliable information retrievaling[27]. By quickly analyzing historical patient data, LLMs can aid psychiatrists in developing effective long-term treatment plans and providing real-time support during consultations.

In this study, a comprehensive evaluation system, PsychBench, was proposed for assessing the performance of LLMs in the clinical psychiatric field. The PsychBench system includes a dataset, and an evaluation framework built around this dataset. The dataset comprises 300 real cases of psychiatric disorders from three specialized psychiatric medical centers. The evaluation framework addresses five key clinical tasks: clinical text generation, primary diagnosis, differential analysis, medication recommendations, and long-term course management. For each clinical task, specific quantitative evaluation metrics have been developed to ensure scientific rigor and accuracy. During the evaluation, relevant patient information, such as history of present illness and psychiatric examination results, is provided to the LLMs. Detailed instructions, derived from clinical guidelines, are provided to the model, requiring adherence to these prompts when executing the designated clinical tasks. Fig. 1 presents the primary pipeline of our study. A comprehensive quantitative evaluation of 16 LLMs was conducted using PsychBench. Further assessment examined how factors such as prompt design, chain-of-thought reasoning, input text length, and domain-specific fine-tuning influenced the models' overall performance. Additionally, to identify potential limitations and areas for improvement, a detailed error analysis was performed for each clinical task.

To further explore the practical benefits of existing LLMs as supportive tools for psychiatrists of varying levels of experience, a clinical reader study was conducted. The study included 60 psychiatrists, divided equally into three groups: 20 junior, 20 intermediate, and 20 senior psychiatrists. The experimental design featured two scenarios: one without LLM assistance and one with LLM support. In both scenarios, participants were tasked with completing a set of specific clinical tasks. The study primarily measured the time taken by the clinicians to complete these tasks under each condition, as well as the task performance of each group. Two specialist psychiatrists from a review committee evaluated the participants' task performance based on a predefined scoring criterion.

The establishment of PsychBench offers a scientific foundation for evaluating the practical application of LLMs in psychiatric clinical work. And we make the dataset and evaluation framework publicly available to provide insights for future research and technological advancements.

# Results

## Creating the PsychBench dataset and evaluation framework

The construction of the PsychBench dataset and evaluation framework was meticulously designed to ensure scientific rigor and comprehensiveness. For the dataset, a power analysis was conducted to determine the necessary sample size, ensuring statistical significance and reliability of the results. Based on the power analysis and insights from relevant research[12,28], 300 de-identified clinical cases were collected from three geographically diverse and representative medical centers across northern, central, and southern China, ensuring broad regional and cultural representation. These cases incorporated comprehensive information including patients' history of present illness, past treatments, family history, physical and mental status examinations, and ancillary test results, etc. This data provided a realistic and detailed clinical context for LLM assessment. An independent expert committee was also established to audit and validate the dataset, ensuring data accuracy and consistency. For detailed procedures on power analysis and dataset construction, refer to the "Dataset" section in Methods.

In building the evaluation framework, five independent clinical tasks were designed based on the dataset: clinical text generation, primary diagnosis, differential analysis, medication recommendation, and long-term course management. Each task was paired with standard answer given by expert committee, and specific quantitative evaluation metrics to precisely measure LLM performance across different clinical scenarios. Specifically, patient information, along with task instructions, were input into the LLMs, which were then tasked with completing specific objectives and generating outputs. The models' outputs were analyzed using the general and psychiatry-specific metrics to compare their performance across tasks. For detailed design on each clinical task and associated prompts and evaluation metrics, see the "Evaluation Framework" and "Quantitative metrics" section in Methods.

## Quantitative evaluation of LLMs using PsychBench

This section presents a detailed report on the quantitative evaluation of 16 mainstream LLMs using PsychBench, focusing on their performance across five psychiatric clinical tasks: clinical text understanding and generation, principal diagnosis, differential analysis, medication recommendation, and long-term course course management. Integrated performance of each model across the five tasks is illustrated in Fig. 2. To gain a deeper understanding of the models' real-world performance, comprehensive error analysis was conducted for each task to identify and explain potential issues in their outputs. Further, the evaluation examined how factors such as prompt design, chain-of-thought reasoning, input text length, and domain-specific fine-tuning influenced the models' outputs. By adjusting these variables, we could identify effective strategies for enhancing LLM performance in psychiatric tasks and pinpoint elements that might contribute to errors. These detailed analyses offer valuable insights for future optimization and application of LLMs in psychiatric clinical settings.

**Clinical Text Understanding and Generation.** This task assessed the capabilities of LLMs in comprehending and generating clinical text. The experimental results revealed that GPT-3.5-Turbo achieved a *BLEU* score of 19.52±14.76, a *ROUGE-L* score of 42.31±12.18, and a *BERTScore* of 77.34±4.86; GPT-3.5-Turbo showcased a *BLEU* score of 20.72±11.60, a *ROUGE-L* score of 39.20±11.65, and a *BERTScore* of 76.09±4.46; Doubao-pro-32k obtained a *BLEU* score of 19.90±10.56, a *ROUGE-L* score of 44.08±10.73, and a *BERTScore* of 78.31±4.62, all indicating a significant advantage in understanding and summarizing clinical text. Additionally, most models achieved 100% of *Diagnostic Criteria Completeness Index (DCCI)*, suggesting a robust ability to understand and follow instructions and being capable of outputting the required medical record content modules in full based on prompts. However, the Spark4-Ultra model scored only 95.51±11.28%, which highlights its deficiencies in instruction following and the completeness of generated content.

In terms of the generating accurate and standardized chief complaints and diagnostic criteria, we defined the indicators *MNER-F1* and *MNER-BERTScore* for quantitative evaluation, which measures the correctness of the generated key information such as the year of the course of diseases, the form of onset, the description of symptom etc. The Doubao-pro-32k model achieved the highest *MNER-F1*, at 28.36±13.46%, with GPT-3.5-Turbo closely following with an *MNER-F1* of 27.24±15.74%. This indicates that Doubao-pro-32k excels in summarizing and articulating key information such as the course of illness and severity, aligning clinical manifestations with psychiatric terminology more precisely and professionally.

In the clinical text understanding and generation task, we systematically analyzed model errors to identify common error types and their underlying causes. We classified the errors into four primary categories: course summary errors, onset pattern summary errors, symptom summary errors, and clinical standardization errors. Our analysis revealed that onset pattern summary errors accounted for

40% of all errors, often manifesting as incorrect judgments of disease onset patterns, such as misclassifying an intermittent course as a continuous one. Course summary errors made up 25%, frequently occurring in cases with subtle symptoms, where the model overlooked the patient's long-term mild symptoms, leading to inaccurate course summarization. Clinical standardization errors comprised 27%, mainly reflected in the extracted chief complaints, which either lacked explicit course and symptom details or exceeded the 20-word limit. Errors in symptom summary were relatively infrequent, accounting for 6% of total errors, indicating that the model can summarize patient symptoms. However, for cases with more complex medical histories or fluctuating conditions, there remains a risk of error. The statistics of the results and examples of each error type were presented in Extended Data Fig. 4. As shown in Extended Data Fig. 9-A and 9-B, the MNER-F1 Score and MNER-BERTScore for cases correctly answered by the LLMs are significantly higher than those in the Course and Symptom Summary Errors group (independent t-test p-value <0.05). This demonstrates that the two metrics based on psychiatric named entity recognition designed in PsychBench effectively measure the professionalism and accuracy of the vocabulary used to describe key terms such as patient symptoms and disease course in the outputs of LLMs.

**Principal Diagnosis.** In this task, we evaluated the LLMs' ability to process complex patient information to provide primary and comorbid diagnoses. The precision of the diagnoses given by the models was evaluated at a fine-grained level using *ICD-10 guided Primary Diagnosis Accuracy (ICD10-PDA)*. As shown in Supplementary Table S2, there were significant differences in the diagnostic accuracy among the models. GPT-4 achieved the highest *ICD10-PDA* at 56.06±40.41%, followed by GLM-4 and Gemini-1.5-pro with an *ICD10-PDA* of 53.14±41.41% and 53.14±41.41%, which was a small gap behind GPT-4. Hunyuan-lite had the lowest *ICD10-PDA*, at only 31.50±38.07%.

We found that most models achieved an *ICD10-PDA* rate of around 50%, with an average *ICD10-PDA* of 48.54%. While this result may seem relatively low, it is influenced by several key factors. Firstly, the evaluation criteria were designed to assess diagnostic accuracy at the level of psychiatric subtype, which significantly increased the difficulty of the task. In clinical practice, different subtypes of psychiatric disorders often exhibit complex and subtle differences in symptoms, posing a challenge to the models' ability to distinguish and categorize them accurately. Additionally, in the task's prompt design, the models were presented with a list of 77 ICD-10 diagnostic categories and required to provide a precise diagnosis accurate to the third decimal place of the ICD-10 code. This design, while adding to the difficulty, closely mirrors real-world clinical scenarios where psychiatrists must make accurate judgments from a range of similar diagnostic options. Therefore, although the current models have not yet achieved an ideal level of accuracy under these conditions, this evaluation method more effectively tests their applicability and reliability in practical clinical settings.

In this task, diagnostic errors were categorized by specific causes into course assessment errors, symptom assessment errors, severity assessment errors, and unclear diagnosis. Symptom assessment errors accounted for 64% of all errors, primarily manifesting as misdiagnosis or missed diagnosis of patients' psychotic symptoms, as well as errors in identifying depressive and manic symptoms in bipolar disorder. The symptom assessment in this task differs from the symptom summary in Task 1. In Task 1, the focus is primarily on summarizing symptoms presented in the patient's medical history. In contrast, symptom evaluation in this task emphasizes identifying specific symptoms from the patient's current clinical presentation and examination results that support the diagnosis. Evidently, symptom analysis in this task is more challenging. Errors due to course assessment made up 18%, mainly involving cases of recurrent depressive disorder and first-episode depressive disorder. The remaining errors were due to severity assessment, and unclear diagnosis, accounting for 11%, and 6%, respectively. Further analysis of error sources reveals several key contributors: (1) Complexity of input information: When patients have a complex medical history or variable disease progression, the model often struggles to identify and extract essential information such as symptoms, disease course, and severity, leading to the omission of critical symptoms or historical details. This finding aligns with following experimental results on the impact of input length. (2) Insufficient clinical knowledge: When diagnoses require adherence to clinical guidelines, such as the ICD-10 diagnostic standards, the model's limited understanding of distinctions among subtypes can lead to an increased rate of misdiagnosis. More detailed results were shown in Extended Data Fig. 5.

**Differential Analysis.** In Task 3, we evaluated the LLMs' ability to perform a thorough differential analysis to assist psychiatrists in ruling out misdiagnosis and minimizing diagnostic errors. The results reveal significant differences in performance among the LLMs on Task 3.

Doubao-pro-32k achieved the highest *Accuracy_main*, at 53.33±48.36%. Meanwhile, it excelled on the *BLEU* and *BERTScore*, securing the top position, and secured the second-highest scores on *ROUGE-L*. Qwen-max, on the other hand, performed well on generating differential diagnosis, with an *Accuracy_diff* of 34.56±40.92% , but underperformed in text generation metrics (*BLEU*, *ROUGE-L*, and *BERTScore*), suggesting that the model may be more adept at making precise differential diagnosis decisions but struggles with generating coherent and semantically rich text that aligns with the writing habits of differential diagnostic analysis. In the differential diagnosis subtask, the top-performing models show relatively small differences in their metrics. The Qwen-max model stood out with the highest *Accuracy_diff* rate at 34.56±40.92%. Doubao-pro-32k, which had the highest accuracy rate for principal diagnosis, secured the fourth position in differential diagnosis accuracy with a rate of 32.58±41.77%, slightly trailing behind Qwen-max, Hunyuan-pro and Yi-large.

Most models achieved an $Acc_{main}$ rate of over 48% and an $Acc_{diff}$ rate of over 25%, with an average $Acc_{main}$ rate of 48.42% and an average $Acc_{diff}$ rate of 29.35%. In this task, models were primarily

focused on differentiating among 26 broader disease categories (roughly corresponding to the part of the ICD-10 code before the decimal point), rather than precisely differentiating between disease subtypes. Although the number of diagnostic options is reduced in Task3, the primary diagnostic accuracy achieved by the model is not improved compared with Task2.

In the process of conducting differential diagnostic analysis, there was a significant variance in the performance of the models evaluated with respect to the professionalism and accuracy in describing critical information such as key symptoms and test results. Among them, ERNIE-4-8k emerged as the top performer, achieving an *MNER-F1* score of 28.89±12.93%. In contrast, GPT-4 fared the worst with an *MNER-F1* score of 11.73±9.92%. Deepseek, on the other hand, secured the highest *MNER-BERTScore* and the third-highest *MNER-F1* score, with respective values of 89.83±2.41% and 28.26±12.58%.

As shown in Extended Data Fig. 6, the model's errors in this task can be categorized into several types: symptom judgment errors, disease course judgment errors, misunderstanding of diagnostic criteria, lack of specificity in differential analysis, and omission of medical history information. Similar to Tasks 1 and 2, the primary errors are concentrated in the misjudgment of disease course and symptoms, accounting for 33% and 26% of the errors, respectively. Additionally, 21% of the errors stem from the model's misunderstanding of diagnostic criteria. For example, despite clearly identifying recurrent depressive episodes in the patient's history, the model still fails to accurately diagnose recurrent major depressive disorder. Notably, while 14% of errors are attributed to a lack of specificity in differential diagnosis, this also indicates that the model has already developed a certain level of ability to provide targeted differential analysis based on the patient's specific clinical situation. We specified in the prompt that the model should perform differential analysis between the primary diagnostic result and the two most commonly confusable diseases. Therefore, when the model is tasked with differentiating between a broader range of diseases, the LLM prediction may hit more differential diagnosis in the reference answer and the error rate is expected to decrease. As depicted in Extended Data Fig. 9-C and 9-E, the BERTScore and MNER-BERTScore for cases correctly answered by the model are significantly higher than those in the Misunderstanding of Diagnostic Criteria group, affirming that the two metrics based on text similarity and psychiatric named entity recognition-related BERTScore, as designed within PsychBench, effectively gauge the semantic similarity between the outputs of LLMs and the reference differential diagnostic analyses. Additionally, as illustrated in Extended Data Fig. 9-D, the MNER-F1 Score for cases correctly answered by the model is significantly higher than the scores for all other error type groups (independent t-test p-value <0.05), with a particularly pronounced difference observed in comparison to the Symptom, Disease Course, and Medical History Judgment Errors groups (independent t-test p-value <0.01). This evidence validates that the MNER-F1 metric, based on psychiatric named entity recognition as designed in PsychBench, effectively measures whether LLMs have correctly attended to key information in cases such as symptoms, disease course, and medication history, thereby guiding medication decision-making.

**Medication Recommendation.** This task evaluates the capability of LLMs in recommending medications within the context of psychiatric clinical practice. Specifically, the task requires the models to provide medication recommendations from the candidate drugs, ranked by recommendation priority from highest to lowest, based on the patient's medical records and various test results. Additionally, the models must articulate the reasons for their recommendations, as well as the co-medication situations if it applies. It is worth mentioning that the reference answer for each case is the medication that has proved effective for the patient in actual clinical treatment, but it is not a unique answer to this task. Therefore, in addition to the quantitative evaluation, we will ask specialist psychiatrists to conduct in-depth analysis of the medication recommendations given by the model in the reader study.

Metrics *Top Choice Alignment Score (TCAS)*, *Medication Match Score (MMS)*, and *Recommendation Coverage Rate (RCR)* calculate the degree of match between the answer given by LLMs and the reference answer from the hit rate perspective, and the specific definitions and calculation formulas are detailed in the "Quantitative metrics" section of Methods. The 16 LLMs involved in the leaderboard achieved an average *TCAS* of 10.64%, among which Hunyuan-lite, Moonshot-v1-32k, and GPT-4 all have *TCASs* above 13%, reaching 15.67±36.35%, 13.15±33.79%, and 13.00±33.63% respectively. In terms of *MMS*, the 16 models involved in the evaluation achieved an average of 33.34%, with Hunyuan-pro and GPT-4 achieving the top 2 *MMS* performance, reaching 37.71±27.22% and 37.11±32.92% respectively. In terms of *RCR*, the 16 models involved in the evaluation achieved an average of 36.32%, with Moonshot-v1-32k and Hunyuan-pro achieving the top 2 *RCR*, reaching 43.72±30.99% and 41.63±30.85% respectively. The experimental outcomes indicate that Moonshot-v1-32k achieved the highest average min-max normalized scores across *TCAS*, *MMS*, and *RCR* metrics. Notably, it attained the top score in *RCR* with 43.72 $\pm$ 30.99%. It also ranked second in *TCAS* at 13.15±33.79% with a disadvantage of 2.52% and third in *MMS* at 35.92±25.48% with a disadvantage of 1.79%. This suggests that Moonshot-v1-32k's recommendations are not only aligned with the reference answers at a high rate but also cover a broad range of potential medications. Hunyuan-lite has obvious advantages in *TCAS* indicators, indicating its recommended primary medications align more closely with clinical practitioners' choices among all the evaluated models. From a clinical perspective, this task presents significant challenges. Most models participating in the test achieved an *TCAS* rate of around 10% for the first-choice medication recommendation, as well as an average *MMS* and *RCR* of around 35%.

In our error analysis of medication recommendations, we assessed the clinical feasibility of the model's suggestions and identified six major categories of errors or inappropriate recommendations: basic medication usage errors, inadequate consideration of adverse drug reactions, errors in combined medication use, overtreatment, lack of reference to the patient's treatment history, and treatment

plans conflicting with the current condition. Among these, treatment conflicts with the current condition were the most common, comprising 33% of all errors. This type of error often involved recommending antidepressants to patients with pronounced manic or psychotic symptoms, potentially exacerbating mania or inducing psychotic symptoms. Overtreatment and failure to consider the patient's previous treatment history made up 24% and 14% of errors, respectively, indicating that the model can still improve in accurately incorporating patient history and assessing symptoms. Errors related to basic medication usage, insufficient caution regarding adverse reactions, and combined medication errors accounted for 4%, 18%, and 6% of total errors, respectively. These errors suggest that while the model generally demonstrates strong knowledge in psychiatric pharmacology, there are areas where further refinement could enhance the clinical safety and efficacy of its recommendations. The examples of each error type were illustrated in Extended Data Fig. 7. The evaluation of the medication recommendation task underscored practical challenges in psychiatric pharmacotherapy. The selection of psychiatric medications often requires iterative adjustments to identify the optimal treatment for a patient, resulting in a slower, costlier treatment process. If the model can reliably assist in this area, it has the potential to expedite the medication selection process, reduce clinical costs, and improve the overall treatment experience for patients.

**Long-term Disease Course Management.** The objective of this task is to evaluate the LLMs' abilities in extracting and comprehending information within psychiatric clinical settings, particularly in terms of precision when dealing with lengthy patient medical records spanning multiple phases. In the analysis of our experimental results, we observed that models such as Doubao-pro-32k and GPT-3.5-Turbo demonstrated high capabilities in handling long texts and multi-phase medical records. These models effectively captured a wide range of contextual information, including medical history, symptoms, and treatment processes, and achieved commendable performance on the evaluation metrics, obtaining average scores of 20.88 for *BLEU*, 49.91 for *ROUGE-L*, 75.92 for *BERTScore*, 53.66% for *MNER-F1*, 83.61% for *MNER-BERTScore* and 81.45% for *Accuracy*, respectively. Doubao-pro-32k achieved the first result with obvious advantages in both question-answering(QA) and multiple choice(MC) forms of long-term medical record reading and analysis tasks. It is worth mentioning that the model participating in the evaluation achieved an average accuracy rate of 82.11% on the MC task, and more than half of the LLMs had an accuracy rate exceeding 85%, which demonstrates the model's ability to extract and analyze key information from long-term medical records. This performance indicates a high level of proficiency in understanding and processing complex medical data, which is crucial for applications in healthcare where accurate diagnosis and treatment rely heavily on the precise extraction and interpretation of detailed patient histories.

Given the nature of psychiatric disorders, the medical records of patients with these conditions tend to be lengthy and complex, involving changes in symptoms and treatment strategies over time. The clinical language in these records is intricate and often features ambiguous terminology, high levels of

ellipses, and non-standard grammatical structures, which pose challenges for accurate interpretation. In clinical practice, errors in information retrieval, such as false recall and misinterpretations of medication dosages, can have severe consequences for patient treatment. Moreover, enhancing diagnostic efficiency and reducing misdiagnosis in clinical practice directly rely on the strong information extraction and comprehension abilities. By rapidly integrating and accurately comprehending medical record information, LLMs can assist clinicians in working more efficiently and reduce the risk of misdiagnosis due to omitted or inaccurate information. Finally, a profound understanding of medical records also supports clinical research and data mining, providing a foundation for large models to conduct disease trend analysis and offer reliable treatment recommendations.

Based on question type, we categorized model errors into variation information judgment errors, summary information judgement errors, comparative information judgement errors, locational information judgement errors, and knowledge deficiency caused errors. The results were illustrated in Extended Data Fig. 8. Variation information judgment errors and summary information judgement errors accounted for 42% and 39% of total errors, respectively. Answering these types of questions requires the model to accurately identify and extract all key information across the input and perform logical analysis. A complex patient history further increases the difficulty in accurately addressing these tasks. In contrast, comparative information judgement errors and locational information judgement errors were less frequent, making up 3% and 5% of total errors. Given the simpler logic and positioning demands of these tasks, the model was more likely to produce correct answers. Finally, errors due to knowledge deficiency accounted for 11% of all errors, often occurring when the model failed to recognize if a biochemical test value exceeded reference ranges due to limited domain knowledge.

## The influence of different prompt strategies

In our study, we delved into the effects of few-shot and Chain of Thought (CoT) prompting techniques on the performance of LLMs. Specifically, for Task 1 and Task 3, we investigated the impact of including or excluding examples in the prompts on model performance. For Task 2 and Task 4, we examined the influence of employing versus not employing the CoT prompting technique on model performance. The rationale behind this experimental setup is that, in clinical scenarios involving Task 1 (distillation of chief complaints and structured summary) and Task 3 (differential diagnostic analysis), clinicians are required to write and organize medical records in a standardized format. In these instances, the large model must not only ensure the accuracy of concepts and semantics but also pay attention to the standardization of the format. We conducted a few-shot experiments for these two tasks to explore whether the large model can learn the formatting standards for writing mental health clinical records based on a limited number of examples.

As shown in Fig. 3-A to Fig. 3-D, our results indicate that providing a single exemplar in the prompts can significantly improve the overall performance of the large model on task1, particularly those related to formatting standards and language expression habits such as *BLEU*, *ROUGE-L*, and *BERTScore*. This reflects that there is indeed a unique set of formatting standards for psychiatric clinical record-keeping, and the LLMs can effectively master these standards with a small number of examples, thereby aligning its output format with that presented in the examples. Additionally, we observed a marked improvement in the *MNER-F1* score and *MNER-BERTScore* of extracting chief complaints and structured summary in the 1-shot setting for task 1, indicating that the LLMs can standardize output format, and learn professional terms and expressive habits with minimal example prompting. However, for task3, the use of 1-shot prompts did not yield a substantial overall performance enhancement across the evaluated metrics. While there were noticeable improvements in *BLEU*, *ROUGE-L*, and *BERTScore*, which are indicative of the linguistic quality and semantic similarity between the model's output and the reference answers, the *MNER-F1* and *MNER-BERTScore* metrics, which are critical for assessing the accuracy of key information extraction and expression, experienced a slight decline. In the 0-shot and 1-shot settings, the average $Acc_{main}$ of task3 obtained by the evaluation model is basically unchanged. It is noteworthy that the accuracy rate of differential diagnosis ($Acc_{diff}$) slightly decreased in the 1-shot scenario, which may be due to the exemplar inducing a bias in the model's differential diagnostic choices.

For Task 2 (primary diagnosis) and Task 4 (medication recommendations), in real clinical practice, clinicians engage in multi-step complex and implicit logical reasoning when making diagnoses and developing treatment plans. These reasoning processes are not detailed in a fixed format within the text of medical records. Consequently, we designed CoT comparative experiments for these two tasks to investigate whether the LLM can enhance the quality of diagnostic and treatment suggestions by simulating the thought processes of clinicians.

The results of the CoT prompting comparative experiments are depicted in Fig. 3-E to Fig. 3-G. For task 2 (primary diagnosis), the use of CoT-style prompts led to a decrease in the *ICD10-PDA* of primary diagnosis. For Task 4, similar to the situation in Task 2, the average performance of LLMs on the metrics of *TCAS*, *RCR*, and *MMS* decreased after using the CoT-form prompt. The *TCAS* dropped from 10.65% to 9.15%, *RCR* from 38.11% to 33.62%, and *MMS* from 35.27% to 28.53%. Upon detailed analysis of the responses provided by the model under the CoT-form prompt, we found that LLMs were capable of comprehensively and accurately summarizing patient information and historical medication usage and efficacy from medical records. However, when suggesting medication, LLMs tended to recommend drugs that had appeared in the medical records, despite their previous suboptimal treatment outcomes, showcase in Supplementary Table S8. In contrast, the recommended medications by psychiatric clinicians in the standard answers showed a lower overlap with previously used drugs, with a preference for adjusting medications to achieve better therapeutic effects.

Additionally, we observed that even when LLMs identified a drug's poor efficacy in the analysis phase, they failed to make correct and reasonable adjustments in the final medication recommendation, which highlights the current LLMs' insufficient reasoning ability in transitioning from past medication and efficacy analysis to drug adjustment plans. The phenomenon that the CoT-based prompting strategy did not yield the expected performance improvement in Task 2 and Task 4 also reflects the complexity and uniqueness of clinical diagnostic thinking in psychiatry and psychology.

## The influence of input length

The distribution of input lengths for the 5 tasks in PsychBench is shown in Fig. 4-A. Compared to common nature language processing (NLP) tasks, psychiatric clinical diagnosis and treatment tasks involve input texts that are lengthy and complex, necessitating LLMs to possess fine-grained extraction and analysis capabilities. Fig. 4-B illustrates the mean values of various evaluation metrics for LLM groups with different context window lengths across 5 tasks. Based on the model's context length, the LLMs tested were categorized into four groups: $8k$, $32k$, $128k$, and $> 128k$. Since GPT-3.5-turbo has a context length of $16k$, it was not included in this analysis. It is evident that on Task 1, as input length increases, the performance of models with various context lengths shows a trend of first slightly decreasing, then rising and then decreasing. This suggests that the relationship between input length and performance may not be linear, and certain models may perform better with specific input lengths. This could be attributed to the fact that longer inputs might introduce additional complexity or noise, which affects the model's ability to generate concise and accurate outputs. On the other hand, shorter inputs may not provide enough context for the model to generate a comprehensive and accurate chief complaints or diagnosis criteria. The task demands both extraction of relevant details and the ability to generate a coherent response within the constraints of clinical standards, making it essential for the model to balance brevity with completeness.

On Task 2, as the input length increases, LLMs with a context length of $8k$ exhibit a consistent decline in performance, whereas models with a context length of $32k$ or more show a trend of decreasing performance followed by an increase. This phenomenon reflects the trade-off between the difficulty of information extraction and analysis in long texts and the richer information provided by more detailed patient information and medical records for diagnosis. For LLMs with shorter context lengths, the increase in input length results in a richer set of diagnostic information, but the model lacks the capacity to extract it effectively; for models with longer context lengths, when the input exceeds 6000 words, the positive impact of the additional information outweighs the negative effects of analyzing longer texts, leading to an upturn in overall performance. It can be observed from the polyline graph of Task2 that LLMs with longer context show a more significant performance improvement after the inflection point.

For Task 3, as input length increases, the performance of all models showed an upward trend regardless of context length. For LLMs with a context length of $32k$ or more, the performance bound is more significant after the point of 4500-5000 input length, with the 128k LLM group with the longest context window achieving the best average performance on this task for input lengths >4500.

For task 4, contrary to task 3, the performance of all four groups of context length models shows a decreasing trend as the input length increases. This downward trend suggests that LLMs struggle to fully analyze and understand more detailed and complex historical medication and disease progression records. Despite longer inputs providing richer information, the participating LLMs are unable to effectively utilize this information to assist in reasoning and make better medication recommendations. At the same time, in clinical practice, developing the next step in the treatment plan for patients with chronic, recurrent conditions and extensive historical medication records is indeed a more challenging task. Additionally, it should be noted that the input length distribution of Task3 and Task4 is different, the input length of Task4 is mostly distributed between 2500 and 3500, while the input length of Task3 is mostly distributed between 3500 and 4500, which may also explain the opposite relationship between the performance of the model and the input length on these two tasks.

In Task 5, for the question answering (QA) task, as input length increases, the performance of LLMs with context length of 32k or more trends upwards, with models having a context length greater than $128k$ showing a more rapid performance improvement before the input length reaches 4000 words. In contrast, the group of models with an 8k context window exhibit a performance trend that first decreases and then increases, with the turning point also occurring in the 3000-4000 words range. This could be due to the model's difficulty in managing excessively large inputs, leading to the degradation of performance as the input length surpasses a certain threshold. As the context window expands beyond 4000 words, the model may struggle to maintain focus on the most relevant information, causing a decline in its ability to generate accurate and meaningful responses. Conversely, models with smaller context windows (such as 8k) exhibit a performance trend characterized by an initial decline, likely due to the limited context they can process, followed by an improvement when the input length approaches a manageable size. This pattern suggests that an optimal context window exists, where the model can balance information retrieval and reasoning without being overwhelmed by excessive details.

For the multiple-choice (MC) task, the performance of all groups of models with different context lengths generally shows a downward trend as the input length increases. This suggests that the ability of the LLMs to handle complex, lengthy input may be limited in certain contexts, leading to a diminished ability to extract precise information required for answering multiple-choice questions. Notably, when comparing models with different context lengths within task5, the longer-context models do not demonstrate a clear advantage and even, in some cases, perform worse than those with

shorter context windows. As shown in Extended Data Fig. 10, for models with a context window greater than or equal to 32k, the accuracy on the multiple-choice questions of Task 5 fluctuates with the position of the correct answer in the medical records, showing a trend of first decreasing and then increasing. Notably, when the answer is located at a relative position of 0.2-0.4 in the medical records, the accuracy decreases most significantly. This phenomenon is closely related to the "lost in the middle"[29] effect, suggesting that models with longer context windows tend to lose focus on key information in the middle part of the text when processing long documents, leading to a noticeable decline in accuracy in the middle section. LLMs with context windows of >128k and 32k exhibit a more pronounced "lost in the middle" effect when the answer is located at the 0.2-0.4 relative position in the medical records. The model with a context length of 128k shows a "lost in the middle" effect when the answer is in the 0.4-0.6 range, but the accuracy fluctuation is less severe in comparison. On the other hand, the model with an 8k context window shows a different trend: as the answer's position in the long medical records shifts later, its accuracy increases. This may indicate that shorter context windows help the model focus more on the key information in the medical record, preventing the occurrence of the "lost in the middle" phenomenon. Therefore, the model can maintain relatively high accuracy over long records, particularly when the answer appears in the later positions.

This reflects that the current strategies for extending LLM context length may impair their analytical and reasoning abilities, as the 5 tasks designed by PsychBench require not only the extraction of key information but also a certain level of understanding and analysis of the input content combined with psychiatric expertise. Other studies have also found that after extending the context window, LLMs do not necessarily "understand" the content better[30], and model performance is influenced by the position of the answer within the input[31]. These results alert us to reconsider the current strategies for extending LLM context length and the methods of evaluation.

## The comparation between general-purpose LLMs and LLMs finetuned on medical domain

To enhance the capabilities and adaptability of LLMs in the medical field, numerous efforts have been made to fine-tune general-purpose LLMs using medical literature, medical encyclopedias, or consultation records from internet hospitals, thereby constructing medical-specific LLMs. For instance, HuaTuoGPT2 is a medical large model fine-tuned based on the general-purpose model Baichuan2. In this evaluation, we conducted a performance comparison between HuaTuoGPT2 and Baichuan2 under both 0-shot and 1-shot scenarios. The relative capabilities of the two models across five tasks are illustrated in Fig. 4-D. In the heatmap presented, colors are determined based on the comparative ratio of HuaTuoGPT2 to Baichuan2 on specific performance metrics, with red hues indicate that HuaTuoGPT2 outperforms Baichuan2 in terms of the specific metric, while blue hues suggest that Baihcuan2 has the advantage. The depth of coloration corresponds to the magnitude of the performance differential. It is evident that for the five tasks designed for PsychBench, the fine-tuned

HuaTuoGPT2 in the medical domain demonstrates a nuanced superior or comparable performance on most metrics compared to the general-purpose model Baichuan2. This advantage is more pronounced in terms of the $Acc_{main}$ in differential diagnosis tasks and the *MMS* of medication recommendations in supportive treatment decision-making tasks. These results indicate that fine-tuning in the medical domain can bring about a subtle improvement in the overall performance of LLMs in psychiatric clinical diagnosis and treatment tasks. Moreover, the experimental results also reveal that in both 0-shot and 1-shot scenarios, the fine-tuned HuaTuoGPT2 in the medical domain exhibits slightly inferior performance than the general model Baichuan2 or shows no advantage over it in terms of the *ICD10-PDA* of principal diagnosis and the *TCAS* and *RCR* rate of medication recommendations, compared to the general-purpose model Baichuan2. Fig. 4-E presents a comparison of the performance evaluated by *BLEU*, *ROUGE-L*, and *BERTScore* of the two models on tasks 1, 3, and 5, which involve the composition of summaries and analytical texts. The figure is structured such that the vertical axis denotes the scores achieved by Baichuan2 for the respective metrics, while the horizontal axis represents the corresponding scores for HuaTuoGPT2. Each data point within the plot corresponds to an individual test case from the benchmark. The distribution of points within the upper half of the quadrant would signify that Baichuan2 attains superior scores to HuaTuoGPT2 across a greater number of test cases, and conversely, a concentration in the lower half would imply a superior showing by HuaTuoGPT2. The results reveal that the data points are almost evenly distributed on both sides of the dash line, which represents equivalent performance between the two models. This observation is further supported by the heatmap, where the colors corresponding to these metrics are relatively light, trending towards white, suggesting a lack of clear distinction in performance between HuaTuoGPT2 and the general-purpose model Baichuan2. In other words, despite the targeted fine-tuning of HuaTuoGPT2 in the medical domain, it does not demonstrate a clear advantage over the general-purpose large model in the psychiatric clinical diagnostic tasks designed by PsychBench.

## Reader study

Psychiatric clinical work heavily relies on clinical experience, leading to differences in performance among psychiatrists with varying levels of experience when completing clinical tasks. Therefore, to more thoroughly examine the effectiveness of LLMs in assisting psychiatrists at different experience levels, and to further analyze the potential strengths and limitations of LLMs to provide directions for future research, we designed and conducted a clinical reader study. We recruited 60 psychiatrists with varying levels of experience: 20 junior, 20 intermediate, and 20 senior psychiatrists. Extended Data Fig. 1 illustrates the detailed design of the reader study. Participants were asked to complete a series of clinical tasks (including diagnosis, differential analysis, and medication recommendations) under two conditions: with and without LLM assistance. Subsequently, specialist psychiatrists evaluated their responses to compare the performance between the two scenarios, as well as across different experience levels. The scoring criteria specifically for the reader study were developed based on ICD-

10 guidelines, as shown in Extended Data Tab. 5. The reader study user interface is presented in Extended Data Fig. 2.

As depicted in Fig. 5, the assistance of existing LLM had varying effects on psychiatrists with different levels of experience. A substantial improvement was observed in the overall performance of junior psychiatrists, with average overall scores increasing from 22.85 to 26.25 (p-value = 0.013). Psychiatrists with intermediate and higher levels of seniority demonstrated slight performance enhancement, with average overall scores rising from 26.35 to 27.9 (p-value = 0.276) and from 29.0 to 30.2 (p-value = 0.242), respectively.

In the diagnostic task, the results of the reader study indicated that physician groups with different levels of experience performed well in completing the task, consistently providing correct diagnoses (scoring 5 points), as shown in Fig. 5-B. However, analysis of the violin plot shapes revealed that the lower half of the LLM-assisted group was narrower compared to the group without LLM assistance. This change suggests that the assistance of LLMs has, to some extent, reduced the likelihood of incorrect diagnoses in the diagnostic task. Notably, the effect of LLM assistance was more pronounced in the lower- and mid-experience groups. This finding indicates that for less experienced psychiatrists, LLMs can provide valuable support and reference, helping to reduce errors and biases during the diagnostic process.

The goal of differential diagnosis is to analyze and differentiate potential similar diseases based on the patient's clinical condition. Differential accuracy mainly measures the hit rate of identifying important potential diseases after considering the specific clinical context. As shown in Fig. 5-B, the accuracy of differential diagnosis was relatively lower in the low-experience group. This is primarily because less experienced psychiatrists often lack specificity when performing differential analysis. For example, in the case of depressed patients with delusions or anxiety, psychiatrists should consider differential diagnoses such as delusional disorder and generalized anxiety disorder. Psychiatrists in the low-experience group sometimes overlooked these possible similar conditions. With the assistance of LLM, the lower bound of the differential accuracy in the junior group improved, although it did not reach a statistically significant difference (P-value = 0.16). The differential completeness mainly measures the ability to conduct a thorough analysis of potential diseases. The results in Fig. 5-B indicate that LLM assistance significantly improved the comprehensiveness of differential diagnosis in junior and intermediate groups. The effect of LLM assistance is primarily reflected in its ability to provide a detailed analysis based on the patient's clinical condition. Psychiatrists can quickly reference this content to capture the patient's condition more effectively and develop a clear differential thought process, thereby delivering accurate and comprehensive differential analyses and reducing the risk of missing potential diseases.

In psychiatric practice, there is often no single correct treatment plan; multiple feasible pharmacological options may exist for the same patient. Therefore, we evaluated medication recommendations from physicians with varying levels of experience from multiple dimensions. First, in terms of medication accuracy, physicians in the intermediate and senior groups performed significantly better than those in the junior group. The deficiencies in the junior group primarily stemmed from insufficient analysis of the patient's symptoms progression and treatment history. For example, as presented in Supplementary Table S5, a junior psychiatrist failed to recognize that a patient with treatment-resistant depression and anxiety had been on an adequate dose of venlafaxine for six months with poor efficacy, and did not offer any medication options to address the patient's anxiety symptoms. While benefiting from LLM's detailed analysis of the patient's condition and treatment history, along with its provided medication suggestions, the average score of medication accuracy of junior group increased by 14%. However, in terms of medication adherence to clinical guideline, the effect of LLM assistance was negligible across all experience groups. This result reflects the limitations of current LLMs in adhering to clinical medication guidelines. The primary issue is the model's lack of reliable, real-world clinical guideline knowledge and practical experience, which impedes its performance in strictly following specific treatment protocols. Thus, while LLMs can provide valuable medication recommendations, there remains room for improvement in ensuring that these recommendations fully comply with clinical standards and treatment guidelines. In the evaluation of medication contraindication accuracy and comprehensiveness, junior group showed improvement with the assistance of LLM, with accuracy and comprehensiveness increasing by 14% (p-value = 0.18) and 19% (p-value = 0.08), respectively. The primary contribution of LLMs was providing detailed interpretations of the patient's condition and relevant test results. Furthermore, LLMs assisted physicians by offering knowledge about drug interactions and contraindications, helping to reduce the risk of prescribing medications that are contraindicated. This is especially valuable for junior psychiatrists, as they may lack sufficient experience when managing complex cases. However, for more experienced physicians, the effect of LLM assistance was negligible. Nevertheless, LLM still contributed by offering a rapid analysis of the patient's condition, which can enhance efficiency in clinical decision-making. Therefore, we performed a statistical analysis of the efficiency across the different groups.

In terms of productivity, as presented in Fig. 5-C, the LLM has been shown to markedly reduce the time psychiatrists require to formulate primary diagnoses, conduct differential diagnoses, and devise medication regimens. For the junior group, the average time to process a case was 535.7 seconds, which was significantly reduced to 292.6 seconds with LLM assistance, indicating the most pronounced efficiency gains. The intermediate group demonstrated the highest efficiency, with an average case processing time of 337.0 seconds, further reduced to 217.4 seconds with the aid of LLM. Conversely, the senior group exhibited a less pronounced reduction in average case processing time,

decreasing from 524.2 seconds to 399.6 seconds with LLM assistance. Notably, the efficiency gains attributed to LLM for the senior group were not statistically significant, with a p-value of 0.705, suggesting that the impact of current LLMs on the workflow of senior psychiatrists may be limited.

## Discussion

This study conducts an in-depth analysis of the application of LLMs in the field of psychiatric clinical practice. As the prevalence of mental disorders rises, the traditional psychiatric clinical practice faces increasingly evident challenges. The emergence of LLMs presents new possibilities for addressing these critical issues. However, the actual effectiveness of LLMs in psychiatric clinical practice has yet to be thoroughly validated, which limits their practical application and hinders further research on LLMs tailored for psychiatric applications. To address this gap, we have developed a benchmarking framework—PsychBench—grounded in real clinical data, standardized clinical guidelines, and the actual demands of clinical practice. The PsychBench is designed to comprehensively evaluate the performance of LLMs in psychiatric clinical settings, providing robust evidence for the reliable assessment of their efficacy in real-world applications, and guide future research in this area.

The PsychBench framework stands out for several reasons. Firstly, its design acknowledges the distinctiveness of mental health by decoupling and defining clear, clinically significant sub-tasks with customized evaluation indicators. This approach ensures that the evaluation of LLMs is aligned with the practical demands of psychiatric clinical practice, something that general medical benchmarks have failed to achieve. Secondly, the framework is grounded in high-quality, annotated data from real-world clinical scenarios. This ensures that the evaluation indicators can effectively and objectively measure LLM performance, capturing the subtleties of mental health assessments that are often missed by more generalized benchmarks. Thirdly, PsychBench offers an easy way to comprehensively evaluate the capacity of LLMs in psychiatric clinical practice. By defining task-specific prompts and quantitative evaluation metrics for each clinical task, PsychBench enables multidimensional assessments that are both efficient and thorough. This structured approach facilitates a nuanced understanding of model performance across various aspects of psychiatric care. Utilizing PsychBench, we evaluated the psychiatric clinical performance of 16 LLMs varying with respect to open-source properties, manufacturers, number of parameters, and specific domains, obtaining an advanced and holistic view of LLM's strengths and challenges in the field of psychiatric clinical practice.

LLMs have demonstrated advantages in clinical text comprehension and generation tasks, particularly in structured summarization. This Transformer architecture allows LLMs to effectively capture long-range dependencies within text, leading to outstanding performance in both generation and comprehension tasks, which is very helpful for analyzing complex patient histories in psychiatric clinical practice. Furthermore, leveraging few-shot learning, LLMs can quickly adapt to new task

requirements with minimal example support. This learning paradigm reduces training costs while enhancing task execution flexibility. For instance, by providing one single example, most LLMs can rapidly learn to generate structured summaries adhering to four specific dimensions: symptom criteria, course criteria, severity criteria, and exclusion criteria. This capability of LLMs provides a potential approach to assist psychiatrists in saving time on tedious documentation while ensuring that the generated text meets clinical needs. However, we observed that LLMs may occasionally exhibit errors in summarizing clinical information, particularly when dealing with complex cases, such as inaccuracies in assessing disease progression or misjudgment of symptoms. Therefore, further optimization of LLMs is necessary to ensure they meet the higher accuracy standards required for clinical applications.

LLMs exhibit inadequate performance in diagnostic tasks and currently do not meet the clinical demands for accurate diagnoses. The underlying reason is that most models evaluated are general-purpose, lacking in-depth training specific to psychiatric clinical expertise. This results in a deficiency in the models' mastery of specialized knowledge, adversely affecting diagnostic accuracy. Furthermore, the requirement for the models to provide definitive diagnoses from among 77 possible subtypes of mental disorders undoubtedly complicates the diagnostic process. Although this requirement poses challenges for the models, it more accurately reflects the complexities and diversity inherent in psychiatric diagnosis in real clinical settings. The error analysis revealed that most diagnostic errors made by LLMs are concentrated on the misinterpretation of patient symptoms. This limitation stems from the complexity and variability of symptoms in psychiatric patients, where accurate diagnosis requires not only precise identification of the patient's current symptoms, but also thorough analysis and comprehensive consideration of the patient's past symptom and diagnosis changes. This highlights the importance of fine-tuning LLMs with clinical data, as current models primarily possess guideline-based knowledge but lack extensive clinical experience. The authenticity and complexity of clinical data can bridge this gap, enhancing the models' performance in tackling real-world clinical tasks.

LLMs can provide psychiatrists with helpful differential analyses assistance. Although existing models still exhibit limitations in their decision-making abilities for accurate diagnoses, their robust text comprehension and analytical capabilities enable them to generate detailed differential analyses based on patient information and specific instructions. The reader study indicated that LLMs offer particularly significant support to junior psychiatrists, helping them access more thorough differential analysis references and thereby improving diagnostic accuracy. This auxiliary function not only enhances their work efficiency but also bolsters their capacity to handle complex cases. Error analysis showed that the primary issues still stem from inaccuracies in assessing patients' medical history and symptomatology.

Existing LLMs have acquired extensive knowledge about medications and related diseases, allowing them to provide valuable medication recommendations for psychiatrists. The LLM can first conduct a thorough analysis of the patient's condition and relevant auxiliary test results, and then, by integrating knowledge of relevant medications and diseases, generate an appropriate medication recommendation. For instance, LLMs can suggest appropriate combinations of antidepressants and antipsychotics to a patient with slight psychotic symptoms, while some junior psychiatrists tend to focus solely on the patient's depressive symptoms, neglecting potential psychotic symptoms. In the reader study, the LLM demonstrated strong supportive effects, providing substantial references for doctors in formulating medication plans. The errors in medication recommendations made by the LLM mainly stem from unfamiliarity with clinical medication guidelines and, as mentioned earlier, incorrect judgments regarding patient symptoms. This highlights the need for fine-tuning the model with real clinical data to enhance its adherence to clinical guidelines and experiential knowledge. By incorporating further specialized knowledge training and employing techniques such as retrieval-augmented generation (RAG)[32,33], the practicality and safety of LLMs in individualized medication recommendations can be enhanced.

LLMs possess the capability to rapidly retrieve target information from lengthy texts, a feature that holds significant value in the long-term management of patients with mental disorders. This ability enables models to swiftly integrate patients' historical medical records, treatment histories, and symptom changes, thus aiding psychiatrists in extracting key information from complex datasets to optimize clinical decision-making processes. The LLM demonstrated a high accuracy rate in this task. However, this capability has certain limitations. First, it depends on the model's context window size; if a patient's medical history exceeds the model's processing capacity, critical information may be overlooked or inadequately utilized. Given the "lost in the middle" phenomenon, simply increasing the context window length does not effectively address this challenge. Additionally, the model's understanding of specialized terminology directly impacts the accuracy of its information retrieval. If the model lacks sufficient comprehension of specific medical terms or clinical jargon, it may lead to incomplete or inaccurate information retrieval, ultimately affecting the effectiveness of long-term patient management.

We explored multiple factors that may influence LLM performance to provide guidance for subsequent work. First, in-context-learning (ICL) can significantly improve model performance, even if only one example is used for model adaptation. However, the impact of more examples on performance was not investigated due to the long context of the clinical tasks involved in this study exceeds the limits of some models' context windows. Moreover, we investigated the effect of input context length on model performance. A key observation is that the relationship between input length and model performance is not linear, and the optimal model context window length appears to vary by task and input length. The results highlight the challenges LLMs face when dealing with long and complex clinical notes,

requiring fine-grained extraction and inference capabilities that are often affected by the context window of the model. This complex interplay between input length and LLM performance across psychiatric clinical tasks reveals that longer context window length or inputs do not necessarily lead to improved outcomes. While extended context windows allow models to process more comprehensive information, they also introduce challenges related to information extraction and coherence, particularly in tasks that require fine-grained analysis and reasoning.

The CoT-based prompting strategy did not necessarily improve model performance on specific tasks. Our experiments suggest that the use of CoT-style prompts even led to a decrease in the performance of primary diagnosis and treatment planning. A similar situation was found in a study by Yang et al. on ChatGPT's ability to perform mental health analysis and emotional reasoning tasks[34]. The reason is that in the clinical domain of psychiatry, the physician's thought-decision process, which is constructed over a long period of extensive clinical practice, is multifaceted, nonlinear, and somewhat personalized. It is infeasible to construct chains of thought to boost model's performance by simply adjusting input prompt to improve model performance for the specialized domains and complex tasks. Conversely, its low-quality or erroneous analyses may lead to greater biases. A more feasible way to inject clinical decision-making thought into a model is to fine-tune the model using real clinical data. However, it is important to note that fine-tuning needs to be targeted. As shown in Fig. 4-D and Fig. 4-E, HuatuoGPT2, which has been fine-tuned with medical data, does not perform significantly better than the generic model Baichuan2 on PsychBench. This discrepancy may arise due to the unique terminology norms and diagnostic decision-making processes in psychiatric clinical practice. These knowledge and logic cannot be acquired through low-quality, broad medical internet corpora. Instead, it necessitates the collection and organization of high-quality clinical corpora, including real-world clinical case records, authoritative guidelines, and cutting-edge academic papers, etc.

In order to further assess the effectiveness of existing LLMs as auxiliary tools for psychiatrists of different experience levels, we conducted a clinical reader study. The results revealed that, in terms of work quality, LLMs did not significantly improve the performance of senior and intermediate psychiatrists. However, they notably enhanced the clinical performance of junior psychiatrists, particularly in the comprehensiveness of differential diagnosis analysis and medication recommendations. Regarding work efficiency, LLMs demonstrated significant improvements for both junior and mid-level psychiatrists, boosting their ability to complete clinical tasks more efficiently. These findings highlight the potential value of LLMs in supporting psychiatric clinical work and underscore the differing needs of psychiatrists at various stages of their careers. This differences suggest that future LLM development should consider tailoring LLM assistance based on the doctor's level of expertise and area of specialization, to more comprehensively support the development of psychiatric practice. For example, for junior psychiatrists, the LLMs should focus on supporting

foundational knowledge and disease management, while for more experienced psychiatrists, the assistance can be more centered on the latest research findings and updates to clinical guidelines.

We also identified the following limitations of our study. Firstly, we only investigated the performance of the original model and the model adapted using ICL, not the performance of the model after adaptation using fine-tuning methods such as LoRA. This is because fine-tuning models specializing in psychiatric clinical tasks requires a large amount of high-quality clinical data, which is a very complex task, and there is no existing related work. This study demonstrates that the existing generic models already can assist psychiatrists with clinical tasks. In the future, we will build on this study to further explore the development of models specialized in psychiatric clinical tasks that serve psychiatrists. Besides, we constructed the evaluation dataset using clinical patient data from multiple authoritative mental health centers across China, enhancing the study's practicality and relevance to some extent. However, the challenges associated with data acquisition prevented us from incorporating clinical data from other countries and regions. As a result, the generalizability of our findings in non-Chinese contexts remains to be further explored. Future work should consider expanding the data sources to validate the model's applicability across diverse cultural and healthcare settings, thereby enhancing the broader impact of the research.

In conclusion, this study proposes a benchmark for evaluating the performance of LLMs in assisting psychiatric clinical practice, known as PsychBench, which includes a dataset and an evaluation framework. Through quantitative assessments of existing LLMs and a clinical reader study, we identify the potential of these models to assist psychiatric clinician. However, despite demonstrating certain potential advantages, LLMs exhibit significant shortcomings and do not fully meet clinical application needs. These deficiencies primarily manifest in diagnostic accuracy, application of specialized knowledge, and handling of complex cases. Overall, this research provides a systematic evaluation framework and reference for the future development of LLMs in the psychiatric field, underscoring the importance of further optimizing models to achieve greater clinical adaptability and effectiveness.

# Methods

## Dataset

To ensure the validity and reliability of our research findings, we first determined the required sample size through power analysis. In this study, we established the following parameters: Effect Size: Based on literature and previous studies[35,36], we assumed an effect size of 0.5, defined using Cohen's d, which is considered a medium effect. Alpha Level ($\alpha$): We selected 0.05 as the significance level to control the risk of Type I errors. Statistical Power ($1 - \beta$): We set this to 0.90, indicating our aim to have a 90% chance of detecting a true effect, thereby reducing the risk of Type II errors. We utilized the 'statsmodels' library in Python to perform the power analysis calculations, which yielded a required

sample size of 85.03. This means we need to collect data from at least 86 patients for the assessment in this study. A power analysis curve is presented in Extended Data Fig. 3. This systematic approach to sample size calculation ensures that our research possesses adequate statistical power, thereby enhancing the credibility and generalizability of the results.

During the data collection phase, we collaborated with three prestigious psychiatric medical centers to ensure diversity and representativeness in our sample, including Beijing Anding Hospital affiliated with Capital Medical University, Fourth People's Hospital of Wuhu, and Second People's Hospital of Dali. Based on prior sample size calculations, we determined that at least 86 patient data points were necessary. Ultimately, we collected clinical patient data from each center, totaling 300 cases, with 100 cases from each facility. This approach not only exceeds the requirements established by our power analysis but also provides a broader context and richer data characteristics for the study. The Beijing Anding Hospital affiliated with Capital Medical University, is a leading psychiatric institution in China. Located in the northern region of the country, it boasts a strong clinical foundation and abundant research resources, representing the highest level of psychiatric diagnosis and treatment in China. Fourth People's Hospital of Wuhu, located in central China, serves a diverse patient population in the central region, reflecting the actual medical resources of the area. Second People's Hospital of Dali, an important medical institution in southwest China, represents the mental health service situation and patient characteristics of more remote areas. The selection of these three medical centers covers different geographical regions and medical backgrounds, providing a solid foundation for the representativeness of the data. In addition, patients of multiple ethnic minorities were collected, whose living habits, cultural backgrounds, and health beliefs vary significantly. These differences may not only have a potential impact on the pathogenesis and clinical presentation of psychiatric disorders but also influence communication methods and medication adherence during the treatment process. By including data from ethnic minority patients, this study offers a more comprehensive representation of psychiatric patient characteristics across different cultural backgrounds. This diversity of culture and habits adds richer dimensions to the construction of evaluation benchmarks, ensuring the applicability of the research findings across various ethnic and cultural contexts. The study was reviewed and approved by the ethics committees of all participating institutions according to the data verification and annotation guidelines (Extended Data Fig. 12), ensuring compliance with local cultural norms and ethical standards.

We implemented several key screening standards to ensure data quality and the validity of the study. First, we focused on inpatients admitted after 2022. This criterion was designed to exclude clinical data biases introduced by the COVID-19. Second, to reflect the actual clinical scenarios of different medical centers, we selected case data by referencing the statistical distribution of patient diagnoses at each center, focusing primarily on schizophrenia and mood disorder spectrum conditions. This selection process ensured not only the authenticity of the dataset but also the diversity of disease types,

providing a robust foundation for evaluating the performance of LLMs in psychiatric clinical tasks. The patient data we collected encompasses essential information, including demographic details (such as age, gender, and occupation), history of present illness, past medical history, personal history, family history, and treatment history. Additionally, complete clinical records from the hospitalization process were included, which contain multiple physical examinations, psychiatric assessments, results of auxiliary examinations, medical orders, and physician ward round notes. To protect patient privacy, all personal information underwent rigorous de-identification to ensure that no identifiable data is disclosed during the analysis. An independent expert committee in psychiatry reviewed and validated the collected data to ensure its accuracy and reliability. Detailed statistics of the dataset are presented in Extended Data Tab. 1.

## Evaluation framework

The construction of the evaluation framework comprises three components: the design of evaluation tasks, the formulation of quantitative metrics for each task, and the design of prompts for each task. We have designed five clinical tasks: clinical text understanding and generation, principal diagnosis, differential analysis, medication recommendations, and long-term disease course management. The following subsections will provide detailed descriptions of each task.

**Clinical text understanding and generation.** This task requires the model to extract and generate chief complaints and structured summaries that adhere to clinical standards from detailed patient information. Specifically, the input patient information includes age, gender, present history of illness, past medical history, personal history, family history, treatment history, and results from the first physical examination, mental examination, and auxiliary tests conducted after hospitalization. The model's output for the chief complaint should summarize the patient's symptoms and clinical course in no more than 20 words, while the structured summary must cover four dimensions: 'symptoms', 'clinical course', 'severity', and 'exclusion'. Unlike conventional text summarization, this task necessitates that the LLM not only extracts relevant patient information from lengthy input but also complies with clinical norms and conveys the information in an exceedingly concise and clear manner. This task allows for the assessment of the model's understanding of clinical information, adherence to instructions, and adaptability to specific standards.

For Task 1, the prompt delineated the standards for writing clinical chief complaints, adhering to clinical medical record conventions. We required the LLM to limit the chief complaint to no more than 20 words. For the structured summary, the prompt states detailed explanation of the four aspects: symptoms, course of disease, severity, and exclusion. The prompts of all of the five tasks are illustrated in Supplementary Table S7.

**Principal diagnosis.** This task focuses on fine-grained psychiatric diagnosis, requiring the model to provide precise diagnoses based on detailed patient information. The input patient information is identical to that used in Task 1. The model must adhere to the ICD-10 diagnostic criteria and provide diagnostic results specified to the fourth character of the ICD code, such as F31.4 (bipolar disorder, currently in a major depressive episode without psychotic features). This requirement necessitates that the model not only classify the primary psychiatric disorders but also conduct more nuanced subtype diagnoses based on the patient's specific condition. This diagnostic process is more complex than standard diagnostic tasks and closely reflects actual clinical scenarios, which is significant for the subsequent development of treatment plans. This task can assess the model's analytical ability regarding detailed patient conditions and its understanding and application of clinical diagnostic standards, ensuring that the generated diagnostic results possess practical clinical applicability and provide robust support for subsequent treatment. Through this design, we can comprehensively evaluate the model's capability in psychiatric diagnosis.

For Task 2, the prompt specified that the primary diagnosis should be based on the ICD-10 diagnostic criteria and required the model to refine the diagnosis to the disease subtype. The ICD-10 diagnostic codebook will be embedded in the prompt, thereby guiding and constraining the outputs of LLMs to align with standardized medical nomenclature. Specifically, we listed 77 common psychiatric disorders along with their corresponding ICD-10 codes in the prompt, thereby standardizing and limiting the model's diagnostic outputs.

**Differential Analysis.** The focus of this task is the differential analysis of potential psychiatric disorders. Given the complexity of psychiatric disorders, patients may exhibit similar symptoms yet suffer from different diseases, necessitating distinct treatment approaches. In such scenarios, the accuracy of differential diagnosis plays a pivotal role. For instance, the differentiation between bipolar affective disorder, characterized by depressive phases, and major depressive disorder, which is typified solely by depressive episodes, is imperative for planning the therapeutic trajectory. LLM is prompted to provide a primary diagnosis and two differential diagnoses for each patient case: the primary diagnosis is intended to accurately reflect the patient's actual condition, while the differential diagnoses involve distinguishing and ruling out several similar diseases through comparative analysis of the patient's symptoms, signs, medical history, and examination results to ascertain the most likely diagnosis. Through this task, we can evaluate the model's ability to meticulously identify and distinguish symptoms in clinical reasoning and diagnostic processes, and verify the model's capability to integrate patient information, generate coherent differential diagnosis suggestions, and provide clinical decision support for psychiatrists.

For Task 3, the prompt directed the model to analyze the provided patient information and offer primary diagnoses and differential diagnoses. We outlined the main content that should be included

in the differential diagnosis analysis according to clinical guidance and specified the number of primary and differential diagnoses to be given. Like in Task 2, the prompt also listed 26 common psychiatric disorders with their ICD-10 codes, with a wider range from F00 to F98, to restrict the model's selection of primary and differential diagnoses within the given range. The difference in this task is that the ICD-10 codes and names are broader, only precise to the decimal point before the ICD-10 code, and in some cases even provided as a range (e.g., F70-F79 Intellectual disabilities). The inclusion of a broader range of ICD-10 codes in this task, compared to Task 2, allows the model to accommodate the varying levels of specificity that can arise in clinical practice, where an initial assessment may not immediately point to a single, clear diagnosis. Instead, clinicians often begin with a set of potential differential diagnoses that need further investigation or refinement based on additional patient information, clinical tests, or observations. By prompting the model to generate a range of differential diagnoses, we aim to simulate this process, ensuring that the model can explore multiple possibilities and assess how different disorders might explain a patient's presentation.

**Medication recommendation.** This task necessitates LLMs to prescribe the optimal psychiatric therapeutic medication based on the medical context and disease progression of the patients. Specifically, the medications prescribed by the physician at the time of the patient's discharge will be considered as the appropriate treatment for that patient. These medications are treated as the reference answers to this task. LLMs are tasked with providing optimal therapeutic medications and the reasons, informed by analyzing patient's present history of illness, past medical history, personal history, family history, treatment history, and results from the first physical examination, mental examination, and auxiliary tests conducted after hospitalization. The specific medications should be listed according to their recommendation priority, and the reason for choosing the medication should be illustrated. For example, for a patient with a history of arrhythmias and major depressive disorder, LLMs should identify the contraindication for Amitriptyline due to their potential to exacerbate cardiac conduction abnormalities. Through this task, the model's ability to understand disease characteristics and treatment needs, as well as its depth of pharmacological knowledge, can be evaluated. The model must also accurately grasp the patient's specific condition to ensure that the medication recommendations are personalized and targeted, thereby enhancing their practical value in clinical decision-making.

For Task 4, the prompt required the model to analyze the provided patient information and offer medication recommendations in order of preference. The prompt, based on clinical practice guidelines, outlined the factors and strategies that the model should consider during medication recommendation, with a particular emphasis on the careful evaluation of drug interactions and adverse reactions. To reduce task complexity and standardize the output of medication names, we provided a list of 34 commonly used medications for psychiatric disorders in the prompt, limiting the range of

recommended drugs. The list of alternative disease names and drug names provided in the prompts of task2-4 contains all the diseases and drugs covered in the standard answer.

**Long-term course management.** This task involves retrieving information from long-term disease courses to complete question answering (QA) and multiple-choice (MC) tasks. Accurate knowledge of the patient's longitudinal medical trajectory and history is instrumental in enhancing the efficiency of therapeutic interventions. We collected the complete course records of each patient in our dataset. These records are structured in terms of dates of hospitalization, documenting the patients' clinical presentations, physical examination findings, psychiatric assessments, and outcomes of ancillary tests conducted on each respective day. For each patient profile, we extracted three tailored questions from their extensive longitudinal medical records that pertain to their condition, focusing on details such as changes in the patient's condition, adjustments to treatment plans, and examination outcomes. For instance, "How did the patient's psychiatric state manifest following the initial MECT session on the tenth day of their hospital stay?" or "In the recent auxiliary examination, which indicator was higher than the reference value but was not mentioned in the previous examination? A. Glutamate aminotransferase B. Aspartate aminotransferase C. Triglycerides D. High-density lipoprotein". All responses to the queries were sourced from authentic textual data within the medical records. LLMs are challenged to answer these questions or make accurate choices utilizing the precise text extracted from the disease courses records. Through this task, the model's ability to retrieve and analyze long-term hospitalization information in real time can be evaluated, which enables clinicians to quickly identify information like key medication responses and changes in the patient's condition during hospitalization, allowing for more effective adjustments to treatment plans.

For Task 5, the prompt instructed the model to thoroughly examine the patient's multi-phase course records, inspection results, and other information. We then require the model to provide answers to these questions in one session.

## Quantitative metrics

We have designed detailed evaluation criteria for the five tasks within the PsychBench. The guiding principle for the design of these metrics is to enable quantitative assessment of the performance of specific tasks based on their characteristics, and these metrics can be calculated automatically.

*BLEU*, *ROUGE-L*, and *BERTScore* are traditional metrics more commonly used in machine translation and summarization tasks. However, they are limited in their capacity for the medical domain. These scores reflect the degree of structural and lexical similarity between the generated text and the provided reference, but they do not specifically assess critical information such as symptom descriptions, medication usage, disease names, anatomical and physiological terms, and laboratory tests in diagnostic and treatment contexts. The challenge is particularly significant in the psychiatric

medical domain, where generating diagnosis and treatment plans often involves navigating abstract concepts, precisely grasping and defining symptoms, paying particular attention to past medication and complications, and dealing with omissions and hallucinations (fabrication, falsification, and plagiarism). In response to this issue, PsychBench further proposes evaluation metrics based on medical named entity recognition: the *Medical NER F1 Score (MNER-F1)* and the *Medical NER BERTScore (MNER-BERTScore)*, which assess the quality of key information in LLM outputs from the perspectives of strict keyword matching and keyword semantic similarity, respectively. Specifically, in line with the approach outlined by Bureaux Tao et al. (https://github.com/Bureaux-Tao/ccksyidu4k-ner), we developed a specialized medical named entity recognition (NER) model, termed $M_{NER}$, optimized for the analysis of medical electronic health record notes. This model was trained on the CHIP2020 dataset, which encompasses 2.2 million characters, 47,194 sentences, and 938 documents, with an average document length of 2,355 characters. The dataset includes a diverse range of medical entities across nine major categories, such as 504 common diseases, 7,085 anatomical names, 12,907 clinical manifestations, and 4,354 medical procedures. To enhance the model's applicability to the psychiatric and psychological domain, we assembled a separate collection of psychiatric clinical electronic medical records, distinct from the 300 cases in PsychBench. We meticulously annotated a training set of psychiatric and psychological electronic medical record notes, with a focus on domain-specific entities and expressions. This annotation process involved the identification and labeling of entities such as mental disorders, psychiatric treatments, and psychological terminology. Subsequent to this annotation, we fine-tuned the NER model on this specialized dataset. The resulting $M_{NER}$ model, based on this dataset, demonstrates an F1-score of 0.66 for the identification of medical entity keywords, indicating a robust performance in recognizing relevant entities in the psychiatric and psychological domain. If we denote the reference answer for the $i$th case in the benchmark as $r_i$ and the LLM's output as $o_i$, we have:

$$M_{NER}(r_i) = \{en_{r_i1}, en_{r_i2}, \ldots, en_{r_im}\}$$

$$M_{NER}(o_i) = \{en_{o_i1}, en_{o_i2}, \ldots, en_{o_in}\}$$

$$Medical\ NER\ Precision = \frac{|M_{NER}(r_i) \cap M_{NER}(o_i)|}{|M_{NER}(o_i)|}$$

$$Medical\ NER\ Recall = \frac{|M_{NER}(r_i) \cap M_{NER}(o_i)|}{|M_{NER}(r_i)|}$$

$$Medical\ NER\ F1 = 2 \times \frac{Medical\ NER\ Precision\ \times\ Medical\ NER\ Recall}{Medical\ NER\ Precision\ +\ Medical\ NER\ Recall}$$

Where $\{en_{r_i1}, en_{r_i2}, \ldots, en_{r_im}\}$ represents the m medical entities predicted after inputting $r_i$ into $M_{NER}$. || denotes the number of elements in the set. Considering that the reference answers and LLM outputs

may have slightly different descriptions for the same symptoms, we use MedicalBERT to further compare the semantic similarity between the two sets of identified named entities at the semantic level. Specifically, we have:

$$MedicalBERT\left(en_{r_ij}\right) = v_{r_ij} \in R^{1 \times d}$$

$$MedicalBERT\left(M_{NER}(r_i)\right) = V_{r_i} = \left[v_{r_i1}, v_{r_i2}, \dots, v_{r_im}\right] \in R^{m \times d}$$

$$MedicalBERT\left(M_{NER}(o_i)\right) = V_{o_i} = \left[v_{o_i1}, v_{o2}, \dots, v_{o_im}\right] \in R^{n \times d}$$

$$Medical\ NER\ BERTScore = \frac{1}{m} \Sigma_{k=1}^{m} \max_{j} \left(cosine\ similarity\left(V_{r_i}, V_{o_i}\right)\right)$$

In task1 (Clinical text Understanding and Generation), task3 (Differential Analysis), and task5 (Management of Long-Term Disease Progression), which involve summary, information extraction and analysis, we use these two metrics to provide a more comprehensive quantitative evaluation of model performance.

Below, we will introduce the evaluation metrics for each of the five tasks in detail:

**Clinical text Understanding and Generation.** This task encompasses two key components: the abstraction of the patient's chief complaint and the synthesis of structured summary through comprehensive analysis of the patient's information and medical history. The chief complaint serves as a concise abstraction of the patient's narrative, while deriving the structured summary necessitates a thorough analysis and synthesis of the patient's medical data. Therefore, in this benchmark, we use the commonly employed *BLUE*[37], *ROUGE-L*[38], and *BERTScore* metrics for the distillation of the chief complaint and the generation of the structured summary. Additionally, the description of the disease courses in the chief complaint and the analysis and summary of the structured summary involves the calculation of time and the mapping of symptoms to psychiatric professional descriptions. Hence, in addition to evaluating the summarization ability through *BLUE*, *ROUGE-L*, and *BERTScore*, this benchmark also includes *Diagnostic Criteria Completeness Index (DCCI)*, *MNER-F1* and *MNER-BERTScore* to assess the integrity and accuracy of the generated structured summary, respectively. Specifically, the percentage of answers generated by the LLM that cover all four diagnostic criteria—'symptom criteria', 'disease course criteria', 'severity criteria', and 'exclusion criteria'—is used as the *Diagnostic Criteria Completeness Index (DCCI)* metric. In addition to the aforementioned metrics, this task also employs the *MNER-F1* and the *MNER-BERTScore* as evaluation indicators. These metrics are specifically designed to assess the quality of medical named entity recognition in the outputs of LLMs.

**Principal Diagnosis.** This task aims to evaluate the LLM's ability to process complex patient information and provide a primary diagnosis. The model's output includes the International

Classification of Diseases, 10th Edition (ICD-10)[39] codes and their corresponding disease names. To comprehensively measure the diagnostic accuracy of the model, we use *ICD-10 guided Primary Diagnosis Accuracy (ICD10-PDA)* as the evaluation metric, calculated based on the overlap between the model's predictions and the reference ICD-10 codes. The method is as follows: if the LLM's predicted ICD-10 code exactly matches the reference answer, it indicates that the model successfully predicted the disease category and sub-type, and the case's *Accuracy* is scored as 1. If the first three digits of the LLM's predicted ICD-10 code match the reference answer but differ from the fourth digit onwards, it shows that the model correctly predicted the disease category but failed to precisely identify the sub-type, scoring 0.5 for that case. If the first three digits of the LLM's predicted ICD-10 code do not match the reference answer, the model is considered to have failed in diagnosing the disease category, and the *Accuracy* is scored as 0. This fine-grained accuracy calculation method allows us to evaluate the model's performance more comprehensively in disease diagnosis tasks, reflecting its strengths and weaknesses in recognizing different disease categories and sub-types.

**Differential Analysis.** The objective of this task is to evaluate the LLM's capability in distinguishing between potential psychiatric diagnoses. The model is tasked with accurately identifying the primary diagnosis and suggesting two most probable differential diagnoses, supported by a comprehensive analysis and rationale derived from the patient's information. Additionally, the model should highlight the key differential aspects of the primary and differential diagnoses under consideration. To quantify the LLM's efficacy in pinpointing the primary diagnosis from the choices presented in the prompt, we employ the $Acc_{main}$ metric. If the LLM's identified primary diagnosis aligns with the reference diagnosis, the case is assigned an $Acc_{main}$ score of 1; otherwise, it receives a score of 0. For evaluating the LLM's performance in identifying the differential diagnoses, we utilize the $Acc_{diff}$ metric. The LLM earns a 0.5 increment on the $Acc_{diff}$ score for each of the two differential diagnoses it proposes that matches the reference list, up to a maximum of 1. To gauge the depth and quality of the LLM's differential diagnosis analysis, we compute the *BLEU*, *ROUGE-L*, and *BERTScore* metrics; to specifically measure the LLM's grasp of key information and the accuracy of its analysis of key symptoms, we calculated the *MNER-F1* and *MNER-BERTScore*. These metrics are used to compare the LLM's analytical output with the reference analysis provided by clinical experts in the fields of psychiatry and psychology. This comparative analysis ensures a comprehensive assessment of the LLM's diagnostic reasoning capabilities.

**Medication recommendation.** This task requires the LLM to provide medication recommendations based on the patient's medical history and various test results, in order of recommended priority, along with an explanation of the reasons. To rigorously assess the concordance between the LLM's medication suggestions and the reference answer, we have developed a set of evaluation metrics grounded in the hit rate concept. Let $D_{LLM} = \{d_{l1}, d_{l2}, \ldots, d_{ln}\}$ be the set of medications recommended

by the LLM and $D_{ref} = \{d_{r1}, d_{r2}, \ldots, d_{rn}\}$ be the set of medications recommended by the reference answer. The medications recommended in both sets are sorted in order of recommendation priority from high to low. The metrics to evaluate LLMs performance on this task are defined as follows:

*Recommendation Coverage Rate (RCR)*: This metric is calculated as the ratio of the number of medications recommended by the LLM that are also present in the benchmark response's recommended medications list.

$$Recommendation\ Coverage\ Rate = \frac{|D_{LLM} \cap D_{ref}|}{|D_{ref}|}$$

Recommendation Coverage Rate quantifies the exhaustiveness of the LLM's recommendations.

*Medication Match Score (MMS)*: Medication Match Score is determined by dividing the number of medications from the benchmark response that are correctly identified by the LLM by the total number of medications suggested by the LLM. The formula for Medication Match Score is:

$$Medication\ Match\ Score = \frac{|D_{LLM} \cap D_{ref}|}{|D_{LLM}|}$$

Medication Match Score gauges the exactness or appropriateness of the LLM's medication suggestions.

*Top Choice Alignment Score (TCAS)*: Accuracy is a measure of whether the LLM's highest-ranked medication corresponds with the top medication in the benchmark response. The formula for Top Choice Alignment Score is:

$$Top\ Choice\ Alignment\ Score = \begin{cases} 1\ if\ d_{l1} = d_{r1} \\ 0\ if\ d_{l1} \neq d_{r1} \end{cases}$$

This metric evaluates whether the LLM's top-ranked medication aligns with the top choice in the reference answer, showcasing the reliability of the LLM in critical decision-making.

**Long-Term course management.** This task necessitates that the LLM accurately extract and comprehend information from patients' longitudinal medical record texts to complete custom-designed reading comprehension and multiple-choice questions.

In the reading comprehension subtask, three sets of reading comprehension questions and answers were generated for each medical record, with a focus on details such as changes in the patient's condition, adjustments to treatment plans, and examination results. This task leans more towards comprehending summaries, where the model needs to read and analyze the long texts of medical records produced during a patient's prolonged hospitalization. The goal is to accurately capture the specific information queried in the questions, interpret the information within the medical records using its own psychological and psychiatric domain knowledge, consider clinical examination

indicators, and analyze the medication and its effects. The evaluation metrics for this task include *BLEU*, *ROUGE-L*, and *BERTScore*, to measure the accuracy and fluency of the LLM's responses against the reference answers, as well as *MNER-F1* and *MNER-BERTScore* to measure the precision in identifying key entities in the responses. These two sets of evaluation metrics together provide standardized and precise quantitative assessment.

In the multiple-choice subtask, 4 multiple-choice questions were generated for each medical record, with each question having only one correct option among the four to five options provided. These questions target specific information such as medication dosages on particular days within the long-term case records and specific numerical values of certain indicators from an examination, thereby examining the LLM's ability to accurately extract specific information from lengthy texts. The evaluation metric for this task is the average *Accuracy* rate of the LLM's responses to the multiple-choice questions.

## LLMs

We investigated a variety of LLMs through PsychBench, with a particular focus on those utilizing autoregressive architectures, known for their ability to produce adaptive and contextually aware outputs. These models, which rely solely on transformer decoders, sequentially generate tokens in a manner that each new token is contingent upon its predecessors, thereby capturing contextual and long-range dependencies. Typically trained on unannotated data, they have proven adept at various NLP applications, such as text generation, question-answering, and dialogue management, which align with our task requirements. The chosen LLMs varied with respect to open-source properties, manufacturers, number of parameters, and specific domains for a comprehensive comparison. Considering the extensive textual inputs involved, we selected models with a minimum context length of 4k tokens or more, ensuring they can process lengthy clinical texts effectively. We organized the LLMs into two categories:

**General-Purpose Models.** In this benchmark, we have included the most advanced multilingual models and Chinese general-purpose large models. The multilingual models comprise the GPT series (GPT-3.5-turbo[40], GPT-4o-mini[41], GPT-4[42]) and Gemini-1.5-pro[43]. For Chinese general-purpose large models, we have incorporated the latest products from various manufacturers to broadly test their capabilities in the field of psychiatric practice. The Chinese general-purpose large models included in this evaluation are GLM4[44], Qwen2.5 (Qwen-max)[45], Doubao-pro-32k[46], Moonshot-v1-32k[47], Spark4-Ultra[48], ERNIE4-8k[49], Baichuan4[50], Yi-large[51], Hunyuan-pro[52], Huanyuan-lite[52], Deepseek-chat-v2[53], and MiniMax[54].

**Medical Domain Fine-tuned Large Models.** To compare the capabilities of large models fine-tuned in the medical domain with those of general-purpose large models in assisting psychiatric diagnosis

and treatment, this evaluation also includes state-of-the-art medical domain fine-tuned large models. To conduct a fairer comparison, we selected Baichuan2-7B-chat[55], which was fine-tuned on general data, and HuaTuoGPT2[56], which was fine-tuned on medical data, both based on the Baichuan2-7B-base model.

The names, parameter sizes, and context lengths of the large models involved in the evaluation are shown in Extended Data Tab. 2. We conduct experiments of Baichuan2-7b-base and HuaTuoGPT2 on a single Nvidia Tesla A100 GPU with 80GB of memory. The results of all the rest of the LLM experiments are obtained by calling the corresponding API. Each prompt is fed independently to avoid the effects of dialogue history.

## Prompt strategy

Based on the prompts specifically designed for each clinical task, we employed three proven prompt strategies to guide the model in completing psychiatric clinical tasks and compared their final performance across different tasks.

**Zero-shot learning** does not rely on any examples but instead directly depends on the task description and contextual information for reasoning[57]. In this study, we used the answers generated by the model using the zero-shot learning strategy as a baseline to assess its ability to handle tasks without any prior examples. Additionally, zero-shot learning was also used as a reference standard for evaluating the effectiveness of other prompt strategies.

**Few-shot learning** helps the model better understand the context and task requirements by providing a small number of examples, without adjusting the model's weights, thereby improving its performance on specific tasks[57]. In this study, we applied few-shot learning to Task 1-4 (Clinical Text Understanding and Generation, Principal Diagnosis, Differential Diagnosis and Medication Recommendation), and detailed analyzed the performance of the LLMs on task1 and task3 using 0-shot and 1-shot prompting strategies. These tasks have high requirements for the content and format of the model's output. By providing a small number of examples, the model could more precisely understand the specific requirements of the task and effectively capture the relationship between input and output, thereby enhancing task completion and accuracy. It is important to note that due to the length of the patient information, including multiple examples in the prompt could exceed the context window limit of some models. Therefore, in this study, we only used one example for few-shot learning to ensure the model's context window limit was not exceeded. Regarding the selection of examples, research indicates that choosing relevant examples can effectively enhance model performance[58]. However, to ensure fairness in evaluation, we selected random samples as examples for testing.

**Chain of Thought** (CoT) strategy is an approach designed to guide the model through step-by-step reasoning, helping it draw more logical conclusions when facing complex tasks[57]. In this study, we applied the CoT strategy to Task 2 (Mental Illness Diagnosis) and Task 4 (Medication Recommendations). These tasks require the model to perform diagnostic reasoning and medication suggestions based on the clinical information provided. This not only demands strong information processing capabilities but also requires detailed and rigorous reasoning. By explicitly guiding the model through a step-by-step reasoning process, the CoT strategy theoretically enhances the model's accuracy and rationality when handling complex information, enabling it to make better diagnostic and medication recommendations. However, although the CoT strategy can provide some level of guidance in the form of instructions within the prompt, the actual effectiveness still largely depends on the model's inherent reasoning ability. Therefore, in some complex clinical tasks, the application of the CoT strategy may not significantly improve the model's performance, particularly when handling intricate clinical decision-making scenarios. In Task 2, we constructed the reasoning chain based on the ICD-10 diagnostic criteria, guiding the model to first analyze the patient according to the ICD-10 standards and then provide a diagnosis. In Task 4, we developed a reasoning chain based on several factors, including the patient's condition, symptoms, diagnostic and examination results, adverse drug reactions, drug interactions, and adherence to treatment protocols. The model was required to perform a thorough analysis of these aspects before providing a final medication recommendation.

## Reader study

After quantitative evaluating the LLMs performance in completing psychiatric clinical tasks using automated metrics, we designed and conducted a reader study to thoroughly assess the application of LLMs as assistive tools for doctors with varying levels of experience, thereby providing more insight for further development of related research. The specific design of the reader study is illustrated in Extended Data Fig. 1. We primarily analyzed the study results from two perspectives: work quality and efficiency. Unlike purely quantitative evaluations, the reader study offers a more realistic reflection of the practical utility and potential of LLMs in clinical practice.

During the preparation phase of the study, we began by recruiting participants. A total of 60 psychiatric psychiatrists were recruited for the study, including 20 psychiatrists each from three experience levels: junior (less than 5 years of experience), intermediate (5-10 years of experience), and senior (more than 10 years of experience). Additionally, we invited two specialist psychiatrists from an independent expert committee to serve as evaluators. One expert was responsible for scoring the participants' responses, while the other conducted a review of the scores to minimize potential biases and maintain reliability of assessment. This experimental design not only allows us to analyze the auxiliary effects of LLMs across psychiatrists with varying levels of experience but also helps us understand the potential development directions for LLMs in real-world clinical applications.

The development of scoring criteria was a critical component of the reader study. To ensure scientific rigor, objectivity, and reproducibility, we designed detailed scoring standards based on the ICD-10 diagnostic guidelines, with reference to frameworks SaferDx[24]. These standards cover multiple dimensions, including (1) Diagnostic Accuracy: The correctness of the diagnosis provided; (2) Differential Accuracy: The precision in differentiating between similar conditions; (3) Differential Completeness: The thoroughness of the differential diagnosis process; (4) Medication Accuracy: The correctness of prescribed medications; (5) Medication Adherence to Guidelines: Adherence to standard protocols in medication prescription; (6) Contraindication Accuracy: The avoidance of contraindicated medications in the prescription; (7) Contraindication Completeness: The comprehensiveness of identifying and avoiding contraindications. Each dimension was accompanied by clear scoring guidelines and corresponding point definitions, as illustrated in Extended Data Tab. 5. The scoring standards underwent review and revision by an independent expert committee to ensure comprehensiveness and consistency. This thorough validation process guarantees that the criteria can be used reliably in future studies, facilitating comparison and reproducibility across different research efforts. Additionally, we will create a public leaderboard to showcase the performance of different models to encourage further research and advancements in this area.

The specific execution process of the reader study is detailed in Extended Data Fig. 1. A subset of 100 patient cases were randomly selected from the entire dataset for reader study. Psychiatrists were divided into three groups based on their experience levels (junior, intermediate, and senior) and completed tasks under two conditions (without LLM assistance and with LLM assistance), resulting in a total of six groups, each comprising 20 psychiatrists. Each psychiatrist was required to complete 10 cases out of the 100, ensuring that each case was repeated twice in each group and appeared a total of 12 times across all groups. For each case, psychiatrists needed to accomplish three primary tasks: primary diagnosis, differential analysis, and medication recommendation. We excluded tasks related to text summarization and case management, as these primarily involved text processing and did not effectively assess psychiatric expertise. In the no-LLM assistance condition, psychiatrists provided answers independently; in the LLM assistance condition, they modified their answers by referencing the responses generated by the LLM before finalizing their responses. After collecting responses from the six groups, the expert psychiatrists would score the responses, and the scoring results would be aggregated by group. Additionally, psychiatrists were required to record the time spent on each task to analyze the impact of LLM assistance on work efficiency. This design ensures the rigor and the reproducibility of the study.

## Statistics analysis

In determining the sample size, we referred to similar studies and incorporated statistical power analysis. We collected 100 clinical cases from each of the three medical centers, with a uniform

distribution of samples from each center, totaling 300 cases for the evaluation dataset. During the power analysis, we followed relevant guidelines and literature recommendations to set the effect size and statistical power parameters, avoiding reliance on determining the effect size through pilot studies, as this approach is considered unreliable and may waste data[59]. We believe the sample size is sufficient as it enabled reproducible and highly credible results when conducting the same experiment with a different set of samples.

For all the selected LLMs, we adopt their default hyper-parameters to maintain consistency with standard operational settings. To ensure the generation of deterministic responses, the temperature parameter was configured to 0.1. Additionally, to prevent premature termination of responses, the maximum token limit for new generations, denoted as $max\_new\_tokens$, was set to 4096, thereby ensuring the integrity of the generated text.

We applied min-max normalization to revalue each evaluation metric for every task and then calculated the mean of all metrics for each task as the overall performance indicator. Subsequently, we computed the mean of the overall performance indicators across the five tasks to serve as the comprehensive evaluation metric for the large models in the field of psychiatric care.

In the quantitative evaluation, we found that GPT-4 exhibited the highest diagnostic accuracy in the diagnostic tasks. This highlights GPT-4's ability to accurately analyze patient conditions and its deep understanding and application of psychiatric clinical knowledge. Accurate diagnosis is the cornerstone of psychiatric clinical practice and serves as the foundation for developing subsequent treatment plans. Considering its overall performance across all tasks, we selected GPT-4 as the LLM-assisted tool for the reader study to comprehensively evaluate its effectiveness in supporting psychiatrists with varying levels of experience in real-world clinical tasks.

In the reader study, we randomly selected 100 cases for comparative analysis. The study was conducted by 60 psychiatrists with varying levels of experience, who were evenly distributed into three groups based on their experience levels: junior, intermediate, and senior psychiatrists, with 20 psychiatrists in each group. To ensure the fairness and scientific integrity of the experiment and adhere to the principle of repetition, each of the 100 cases was completed and evaluated by 6 psychiatrists (2 psychiatrists from each of the 3 experience-based groups). Specifically, each psychiatrist was assigned 10 cases, following a structured allocation: cases 1–10 were assigned to the first psychiatrist in each group, cases 11–20 to the second psychiatrist, and so on. This ensured that each group had two psychiatrists reviewing the same set of 10 cases, and across the three groups, a total of six psychiatrists evaluated each set of 10 cases. This design minimized the potential influence of individual differences on the results. During the expert evaluation phase, the six responses from each group were randomized and presented to experts in a blinded manner to ensure the objectivity and reliability of the assessment process.

## Ethics approval

This study adhered to the principles outlined in the Declaration of Helsinki. Informed consent was obtained from each psychiatrist before their participation. This study used only retrospective, de-identified data that fell outside the scope of institutional review board oversight.

## Data and code availability

All the data and code used in this study are accessible at https://github.com/wangrx33/PsychBench.

## Supplementary information

Supplementary Table. 1 | The example outputs of evaluated models on a random case in terms of 5 tasks.

Supplementary Table. 2 | The example answers given by each group in reader study.

Supplementary Table. 3 | The example of junior group misclassifying depressive episodes as recurrent depressive disorder.

Supplementary Table. 4 | Example answer given by LLM for differential diagnosis.

Supplementary Table. 5 | The example of LLM assisting medication recommendation.

Supplementary Table. 6 | Reference answers of reader study.

Supplementary Table. 7 | The prompts and sample input/output for tasks 1-5.

Supplementary Table. 8 | The example of LLMs tended to recommend drugs that had appeared in the medical records.

Supplementary Table. 9 | The detailed quantitative results of evaluated models on PsychBench in terms of metrics (1-shot)

Supplementary Table. 10 | The detailed quantitative results of evaluated models on PsychBench in terms of metrics (0-shot)

## Acknowledgements

## Author contributions

R.W., S.L., L.Z., and X.Z. contributed equally to this study. R.W. and S.L. designed the pipeline of the study, preprocessed data, ran experiments, designed reader study, analyzed results, created figures and wrote the manuscript. All authors reviewed the manuscript and provided meaningful feedbacks. X.Z. collected clinical and analyzed the results. L.Z. and R.Y. organized the reader study and participated in the reader study as specialist psychiatrist. L.Z. provided the clinical guidance of the whole study and participated in the reader study as specialist psychiatrist. X.Z. and F.W. processed data and ran experiments. Z.Y. provided technical advice and assisted in quantitative evaluation. G.W. and C.J. guided the whole study.

## Competing interests

The authors declare no competing interests.

# References

1       Vigo, D., Thornicroft, G. & Atun, R. Estimating the true global burden of mental illness. *Lancet Psychiatry* **3**, 171-178 (2016). https://doi.org:10.1016/S2215-0366(15)00505-2

2       Patel, V. *et al.* The Lancet Commission on global mental health and sustainable development. *Lancet* **392**, 1553-1598 (2018). https://doi.org:10.1016/S0140-6736(18)31612-X

3       Evans-Lacko, S. *et al.* Socio-economic variations in the mental health treatment gap for people with anxiety, mood, and substance use disorders: results from the WHO World Mental Health (WMH) surveys. *Psychol Med* **48**, 1560-1571 (2018). https://doi.org:10.1017/S0033291717003336

4       van Ginneken, N. *et al.* Primary-level worker interventions for the care of people living with mental disorders and distress in low- and middle-income countries. *Cochrane Database Syst Rev* **8**, CD009149 (2021). https://doi.org:10.1002/14651858.CD009149.pub3

5       Shatte, A. B. R., Hutchinson, D. M. & Teague, S. J. Machine learning in mental health: a scoping review of methods and applications. *Psychol Med* **49**, 1426-1448 (2019). https://doi.org:10.1017/S0033291719000151

6       Kanjee, Z., Crowe, B. & Rodman, A. Accuracy of a Generative Artificial Intelligence Model in a Complex Diagnostic Challenge. *JAMA* **330**, 78-80 (2023). https://doi.org:10.1001/jama.2023.8288

7       Singhal, K. *et al.* Large language models encode clinical knowledge. *Nature* **620**, 172-180 (2023). https://doi.org:10.1038/s41586-023-06291-2

8       Inbar, L. & Zohar, E. Identifying depression and its determinants upon initiating treatment: ChatGPT versus primary care physicians. *Family Medicine and Community Health* **11**, e002391 (2023). https://doi.org:10.1136/fmch-2023-002391

9       Obradovich, N. *et al.* Opportunities and risks of large language models in psychiatry. *NPP—Digital Psychiatry and Neuroscience* **2**, 8 (2024). https://doi.org:10.1038/s44277-024-00010-z

10      Perlis, R. H., Goldberg, J. F., Ostacher, M. J. & Schneck, C. D. Clinical decision support for bipolar depression using large language models. *Neuropsychopharmacology* **49**, 1412-1416 (2024). https://doi.org:10.1038/s41386-024-01841-2

11      Jiang, L. Y. *et al.* Health system-scale language models are all-purpose prediction engines. *Nature* **619**, 357-362 (2023). https://doi.org:10.1038/s41586-023-06160-y

12      Van Veen, D. *et al.* Adapted large language models can outperform medical experts in clinical text summarization. *Nature Medicine* **30**, 1134-1142 (2024). https://doi.org:10.1038/s41591-024-02855-5
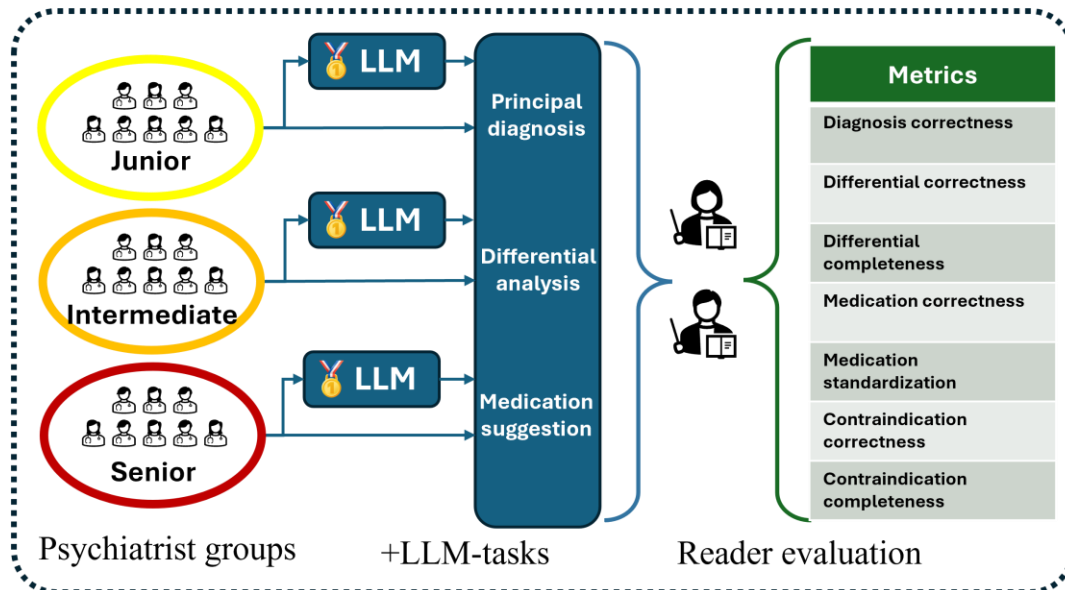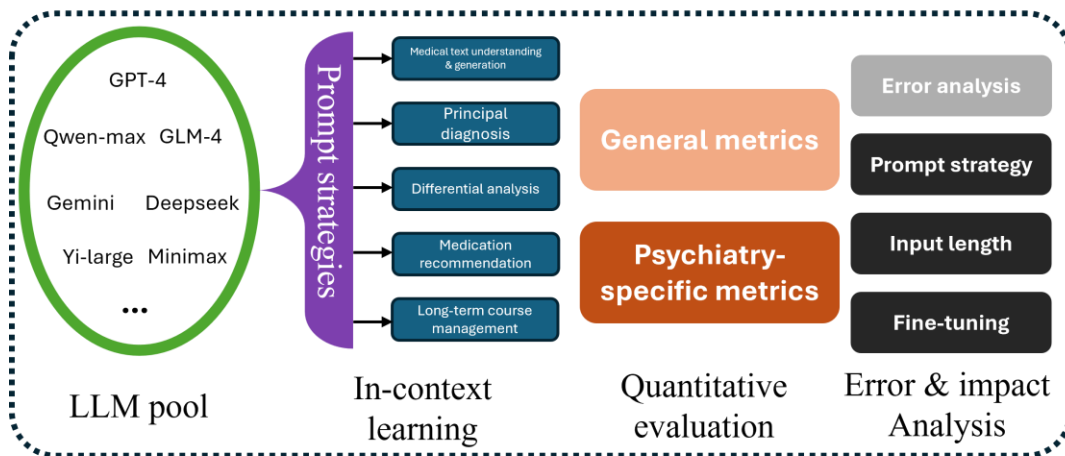
13      Nori, H., King, N., McKinney, S. M., Carignan, D. & Horvitz, E. Capabilities of GPT-4 on Medical Challenge Problems. *ArXiv* **abs/2303.13375** (2023).

14      Singhal, K. *et al.* Towards Expert-Level Medical Question Answering with Large Language Models. *ArXiv* **abs/2305.09617** (2023).

15      Thirunavukarasu, A. J. *et al.* Trialling a Large Language Model (ChatGPT) in General Practice With the Applied Knowledge Test: Observational Study Demonstrating Opportunities and Limitations in Primary Care. *JMIR Med Educ* **9**, e46599 (2023). https://doi.org:10.2196/46599

16      Kung, T. H. *et al.* Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models. *PLOS Digit Health* **2**, e0000198 (2023). https://doi.org:10.1371/journal.pdig.0000198

17      Toma, A. *et al.* Clinical Camel: An Open-Source Expert-Level Medical Language Model with Dialogue-Based Knowledge Encoding. *ArXiv* **abs/2305.12031** (2023).

18      Nori, H. *et al.* Can Generalist Foundation Models Outcompete Special-Purpose Tuning? Case Study in Medicine. *ArXiv* **abs/2311.16452** (2023).

19      McDuff, D. *et al.* Towards Accurate Differential Diagnosis with Large Language Models. *ArXiv* **abs/2312.00164** (2023).

20      Eriksen, A. V., Möller, S. & Ryg, J. Use of GPT-4 to Diagnose Complex Clinical Cases. *NEJM AI* (2023).

21      Pal, A., Umapathi, L. K. & Sankarasubbu, M. MedMCQA : A Large-scale Multi-Subject Multi-Choice Dataset for Medical domain Question Answering. in *ACM Conference on Health, Inference, and Learning*.

22      Tiffen, J., Corbridge, S. J. & Slimmer, L. Enhancing clinical decision making: development of a contiguous definition and conceptual framework. *J Prof Nurs* **30**, 399-405 (2014). https://doi.org:10.1016/j.profnurs.2014.01.006

23      Popatia, R. Berman's Pediatric Decision Making. *JAMA* **307**, 617-618 (2012). https://doi.org:10.1001/jama.2012.104

24      Singh, H., Khanna, A., Spitzmueller, C. & Meyer, A. N. D. Recommendations for using the Revised Safer Dx Instrument to help measure and improve diagnostic safety. *Diagnosis (Berl)* **6**, 315-323 (2019). https://doi.org:10.1515/dx-2019-0012

25      in *Depression in adults: treatment and management  National Institute for Health and Care Excellence: Guidelines*   (2022).

26      Feng, Y. *et al.* Guidelines for the diagnosis and treatment of depressive disorders in China: The second edition. *J Affect Disord* **253**, 352-356 (2019). https://doi.org:10.1016/j.jad.2019.04.104

27      Kessler, R. C. *et al.* Lifetime prevalence and age-of-onset distributions of DSM-IV disorders in the National Comorbidity Survey Replication. *Arch Gen Psychiatry* **62**, 593-602 (2005). https://doi.org:10.1001/archpsyc.62.6.593

28      Hager, P. *et al.* Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nature Medicine* **30**, 2613 - 2622 (2024).

29     Liu, N. F. *et al.* Lost in the Middle: How Language Models Use Long Contexts. *Transactions of the Association for Computational Linguistics* **12**, 157-173 (2024). https://doi.org:10.1162/tacl_a_00638

30     Iyyer, M. K. a. K. T. a. K. L. a. T. G. a. M. One Thousand and One Pairs: A "novel" challenge for long-context language models. *arXiv* (2024).

31     Liu, N. F. a. *et al.* Lost in the Middle: How Language Models Use Long Contexts. *Transactions of the Association for Computational Linguistics* **12**, 157--173 (2024). https://doi.org:10.1162/tacl_a_00638

32     Gao, Y. *et al.* Retrieval-Augmented Generation for Large Language Models: A Survey. *ArXiv* **abs/2312.10997** (2023).

33     Lewis, P. *et al.* Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *ArXiv* **abs/2005.11401** (2020).

34     Ananiadou, K. Y. a. S. J. a. T. Z. a. Q. X. a. Z. K. a. S. Towards Interpretable Mental Health Analysis with Large Language Models. *arXiv* (2023).

35     Fritz, C. O. D., Morris, P. E. & Richler, J. J. Effect size estimates: current use, calculations, and interpretation. *Journal of experimental psychology. General* **141 1**, 2-18 (2012).

36     Funder, D. C. & Ozer, D. J. Evaluating Effect Size in Psychological Research: Sense and Nonsense. *Advances in Methods and Practices in Psychological Science* **2**, 156 - 168 (2019).

37     Papineni, K., Roukos, S., Ward, T. & Zhu, W.-J. in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*   311–318 (Association for Computational Linguistics, Philadelphia, Pennsylvania, 2002).

38     Lin, C.-Y. in *Annual Meeting of the Association for Computational Linguistics*.

39     Bramer, G. R. International statistical classification of diseases and related health problems. Tenth revision. *World Health Stat Q* **41**, 32-36 (1988).

40     Brown, T. B. *et al.* in *Proceedings of the 34th International Conference on Neural Information Processing Systems*   Article 159 (Curran Associates Inc., Vancouver, BC, Canada, 2020).

41     OpenAI. *GPT-4o-mini*, <https://platform.openai.com/docs/models/gpt-4o-mini> (2024).

42     Anadkat and Red Avila and Igor Babuschkin and Suchir Balaji and Valerie Balcom and Paul Baltescu and Haiming Bao and Mohammad Bavarian and Jeff Belgum and Irwan Bello and Jake Berdine and Gabriel Bernadett-Shapiro and Christopher Berner and Lenny Bogdonoff and Oleg Boiko and Madelaine Boyd and Anna-Luisa Brakman and Greg Brockman and Tim Brooks and Miles Brundage and Kevin Button and Trevor Cai and Rosie Campbell and Andrew Cann and Brittany Carey and Chelsea Carlson and Rory Carmichael and Brooke Chan and Che Chang and Fotis Chantzis and Derek Chen and Sully Chen and Ruby Chen and Jason Chen and Mark Chen and Ben Chess and Chester Cho and Casey Chu and Hyung Won Chung and Dave Cummings and Jeremiah Currier and Yunxing Dai and Cory Decareaux and Thomas Degry and Noah Deutsch and Damien Deville and Arka Dhar and
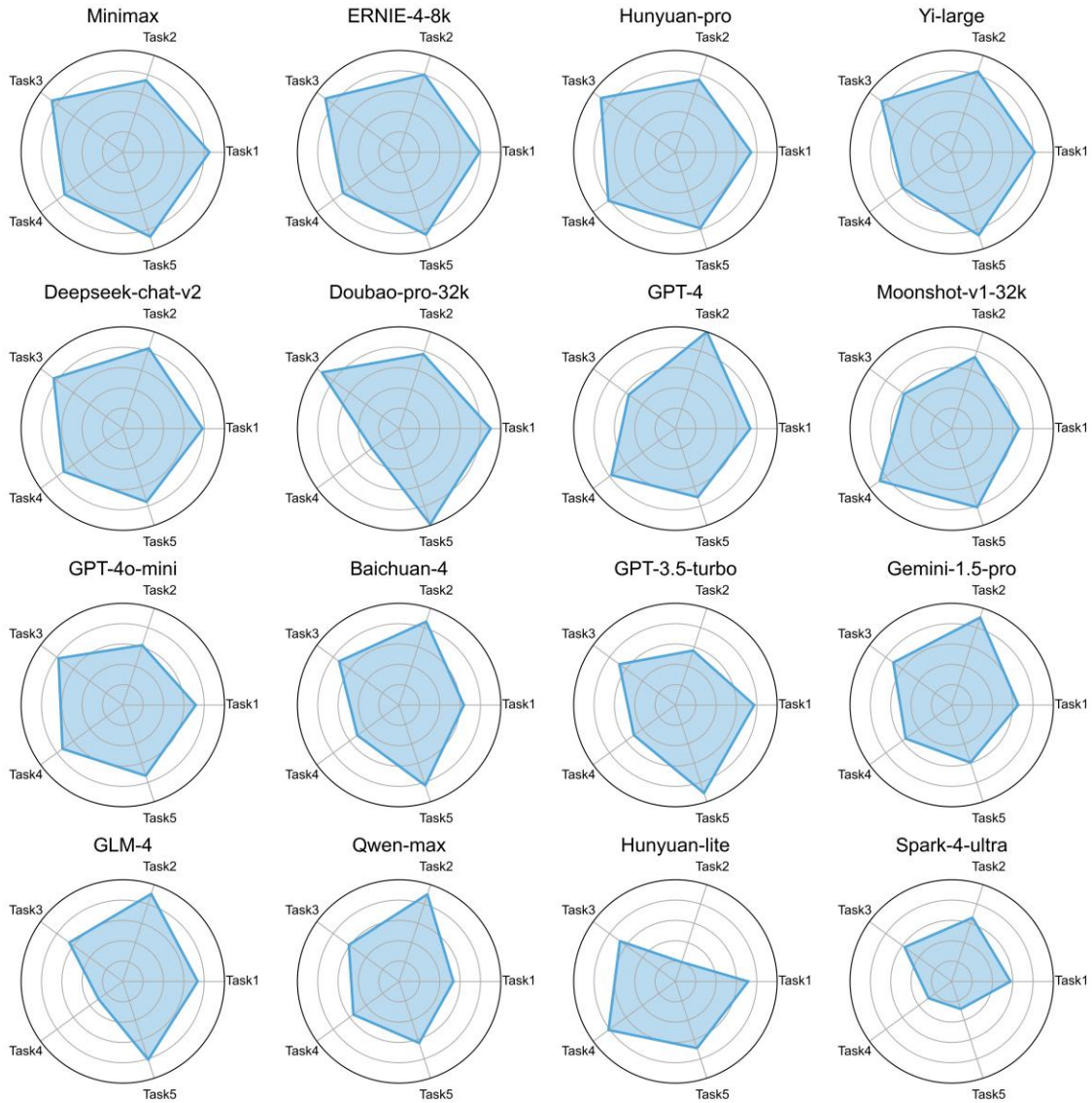
David Dohan and Steve Dowling and Sheila Dunning and Adrien Ecoffet and Atty Eleti and Tyna Eloundou and David Farhi and Liam Fedus and Niko Felix and Simón Posada Fishman and Juston Forte and Isabella Fulford and Leo Gao and Elie Georges and Christian Gibson and Vik Goel and Tarun Gogineni and Gabriel Goh and Rapha Gontijo-Lopes and Jonathan Gordon and Morgan Grafstein and Scott Gray and Ryan Greene and Joshua Gross and Shixiang Shane Gu and Yufei Guo and Chris Hallacy and Jesse Han and Jeff Harris and Yuchen He and Mike Heaton and Johannes Heidecke and Chris Hesse and Alan Hickey and Wade Hickey and Peter Hoeschele and Brandon Houghton and Kenny Hsu and Shengli Hu and Xin Hu and Joost Huizinga and Shantanu Jain and Shawn Jain and Joanne Jang and Angela Jiang and Roger Jiang and Haozhun Jin and Denny Jin and Shino Jomoto and Billie Jonn and Heewoo Jun and Tomer Kaftan and Łukasz Kaiser and Ali Kamali and Ingmar Kanitscheider and Nitish Shirish Keskar and Tabarak Khan and Logan Kilpatrick and Jong Wook Kim and Christina Kim and Yongjik Kim and Jan Hendrik Kirchner and Jamie Kiros and Matt Knight and Daniel Kokotajlo and Łukasz Kondraciuk and Andrew Kondrich and Aris Konstantinidis and Kyle Kosic and Gretchen Krueger and Vishal Kuo and Michael Lampe and Ikai Lan and Teddy Lee and Jan Leike and Jade Leung and Daniel Levy and Chak Ming Li and Rachel Lim and Molly Lin and Stephanie Lin and Mateusz Litwin and Theresa Lopez and Ryan Lowe and Patricia Lue and Anna Makanju and Kim Malfacini and Sam Manning and Todor Markov and Yaniv Markovski and Bianca Martin and Katie Mayer and Andrew Mayne and Bob McGrew and Scott Mayer McKinney and Christine McLeavey and Paul McMillan and Jake McNeil and David Medina and Aalok Mehta and Jacob Menick and Luke Metz and Andrey Mishchenko and Pamela Mishkin and Vinnie Monaco and Evan Morikawa and Daniel Mossing and Tong Mu and Mira Murati and Oleg Murk and David Mély and Ashvin Nair and Reiichiro Nakano and Rajeev Nayak and Arvind Neelakantan and Richard Ngo and Hyeonwoo Noh and Long Ouyang and Cullen O'Keefe and Jakub Pachocki and Alex Paino and Joe Palermo and Ashley Pantuliano and Giambattista Parascandolo and Joel Parish and Emy Parparita and Alex Passos and Mikhail Pavlov and Andrew Peng and Adam Perelman and Filipe de Avila Belbute Peres and Michael Petrov and Henrique Ponde de Oliveira Pinto and Michael and Pokorny and Michelle Pokrass and Vitchyr H. Pong and Tolly Powell and Alethea Power and Boris Power and Elizabeth Proehl and Raul Puri and Alec Radford and Jack Rae and Aditya Ramesh and Cameron Raymond and Francis Real and Kendra Rimbach and Carl Ross and Bob Rotsted and Henri Roussez and Nick Ryder and Mario Saltarelli and Ted Sanders and Shibani Santurkar and Girish Sastry and Heather Schmidt and David Schnurr and John Schulman and Daniel Selsam and Kyla Sheppard and Toki Sherbakov and Jessica Shieh and Sarah Shoker and Pranav Shyam and Szymon Sidor and Eric Sigler and Maddie Simens and Jordan Sitkin and Katarina Slama and Ian Sohl and Benjamin Sokolowsky and Yang Song and Natalie Staudacher and Felipe Petroski Such and Natalie

Summers and Ilya Sutskever and Jie Tang and Nikolas Tezak and Madeleine B. Thompson and Phil Tillet and Amin Tootoonchian and Elizabeth Tseng and Preston Tuggle and Nick Turley and Jerry Tworek and Juan Felipe Cerón Uribe and Andrea Vallone and Arun Vijayvergiya and Chelsea Voss and Carroll Wainwright and Justin Jay Wang and Alvin Wang and Ben Wang and Jonathan Ward and Jason Wei and CJ Weinmann and Akila Welihinda and Peter Welinder and Jiayi Weng and Lilian Weng and Matt Wiethoff and Dave Willner and Clemens Winter and Samuel Wolrich and Hannah Wong and Lauren Workman and Sherwin Wu and Jeff Wu and Michael Wu and Kai Xiao and Tao Xu and Sarah Yoo and Kevin Yu and Qiming Yuan and Wojciech Zaremba and Rowan Zellers and Chong Zhang and Marvin Zhang and Shengjia Zhao and Tianhao Zheng and Juntang Zhuang and William Zhuk and Barret Zoph, O. a. J. A. a. S. A. a. S. A. a. L. A. a. I. A. a. F. L. A. a. D. A. a. J. A. a. S. A. a. S. GPT-4 Technical Report. (arXiv, 2024).

43      Team, G. *et al*. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv:2403.05530 (2024). <https://ui.adsabs.harvard.edu/abs/2024arXiv240305530G>.

44      and Hanyu Lai and Hao Yu and Hongning Wang and Jiadai Sun and Jiajie Zhang and Jiale Cheng and Jiayi Gui and Jie Tang and Jing Zhang and Jingyu Sun and Juanzi Li and Lei Zhao and Lindong Wu and Lucen Zhong and Mingdao Liu and Minlie Huang and Peng Zhang and Qinkai Zheng and Rui Lu and Shuaiqi Duan and Shudan Zhang and Shulin Cao and Shuxun Yang and Weng Lam Tam and Wenyi Zhao and Xiao Liu and Xiao Xia and Xiaohan Zhang and Xiaotao Gu and Xin Lv and Xinghan Liu and Xinyi Liu and Xinyue Yang and Xixuan Song and Xunkai Zhang and Yifan An and Yifan Xu and Yilin Niu and Yuantao Yang and Yueyan Li and Yushi Bai and Yuxiao Dong and Zehan Qi and Zhaoyu Wang and Zhen Yang and Zhengxiao Du and Zhenyu Hou and Zihan Wang, T. G. a. a. A. Z. a. B. X. a. B. W. a. C. Z. a. D. Y. a. D. Z. a. D. R. a. G. F. a. H. Z. ChatGLM: A Family of Large Language Models from GLM-130B to GLM-4 All Tools. (arXiv, 2024).

45      Team, Q. Qwen2.5: A Party of Foundation Models. (2024).

46      Team, D. *Doubao-pro*, <https://www.volcengine.com/> (2024).

47      Team, M. *Moonshot-v1*, <https://platform.moonshot.cn/> (2024).

48      Team, S. *Spark AI*, <https://xinghuo.xfyun.cn/> (2024).

49      Sun, Y. *et al*. ERNIE 3.0: Large-scale Knowledge Enhanced Pre-training for Language Understanding and Generation. *ArXiv* **abs/2107.02137** (2021).

50      Team, B. *Baichuan4*, <https://www.baichuan-ai.com/> (2024).

51      Chen, A. a. a. A. Y. a. B. C. a. C. L. a. C. H. a. G. Z. a. G. Z. a. H. L. a. J. Z. a. J. Yi: Open Foundation Models by 01.AI. (arxiv, 2024).

52      Team, T. H. *Hunyuan*, <https://cloud.tencent.com/document/product/1729/97731> (2024).

53      Shao, Z., Dai, D., Guo, D., Liu , B. & Wang, Z. DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model. *ArXiv* **abs/2405.04434** (2024).

54 Team, M. *Minimax*, <https://www.minimaxi.com/> (2024).

55 Yang, A. M. *et al*. Baichuan 2: Open Large-scale Language Models. *ArXiv* **abs/2309.10305** (2023).

56 Chen, J. *et al*. HuatuoGPT-II, One-stage Training for Medical Adaption of LLMs. *ArXiv* **abs/2311.09774** (2023).

57 Brown, T. B. *et al*. Language Models are Few-Shot Learners. *ArXiv* **abs/2005.14165** (2020).

58 Nie, F., Chen, M., Zhang, Z. & Cheng, X. Improving Few-Shot Performance of Language Models via Nearest Neighbor Calibration. *ArXiv* **abs/2212.02216** (2022).

59 Albers, C. J. & Lakens, D. When power analyses based on pilot data are biased: Inaccurate effect size estimators and follow-up bias. *Journal of Experimental Social Psychology* **74**, 187-195 (2018).

**Fig. 1 | Overview of the framework in this study.** The proposed PsychBench is composed of a dataset and an evaluation framework. The dataset comprises 300 real patient cases collected from three specialized psychiatric medical centers. The evaluation framework consists of five specifically designed psychiatric clinical tasks and corresponding quantitative metrics tailored for each task. The tasks include clinical text understanding and generation, principal diagnosis, differential analysis, medication recommendation, and long-term course management. In this study, we first quantitatively evaluated 16 existing LLMs using PsychBench. We also performed error analysis and assessed the impact of prompt strategies, input length, and domain-specific fine-tuning on model performance. We then conducted a clinical reader study to further evaluate the effectiveness of LLMs in assisting psychiatrists with different levels of experience. Sixty psychiatrists with varying levels of work experience were recruited to accomplish specific tasks in PsychBench with and without the assistance of LLM respectively. Two specialist psychiatrists then scored the answers given by different groups based on predefined evaluation criteria.
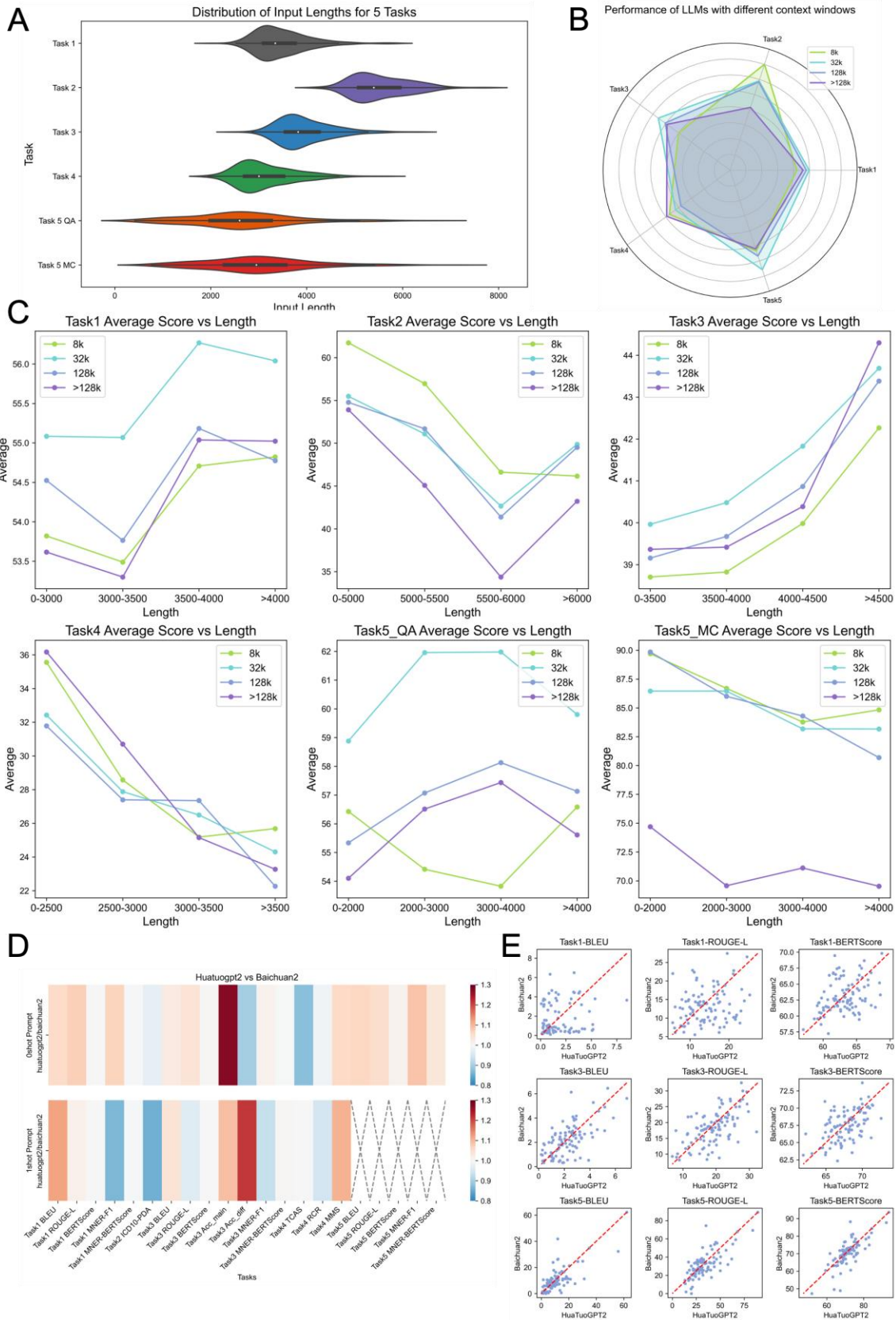
**Fig. 2 | The normalized quantitative results of evaluated LLMs across five tasks.** We applied min-max normalization to revalue each evaluation metric for every task and then calculated the mean of all metrics for each task as the overall performance indicator of the corresponding task. Subsequently, we computed the mean of the overall performance indicators across the five tasks to serve as the comprehensive evaluation metric for the large models in the field of psychiatric care as indicated as the area of each radar map. The radar maps are arranged from left to right and top to bottom in descending order of the comprehensive evaluation metric, reflecting the overall performance of each large model in the psychiatric care domain.
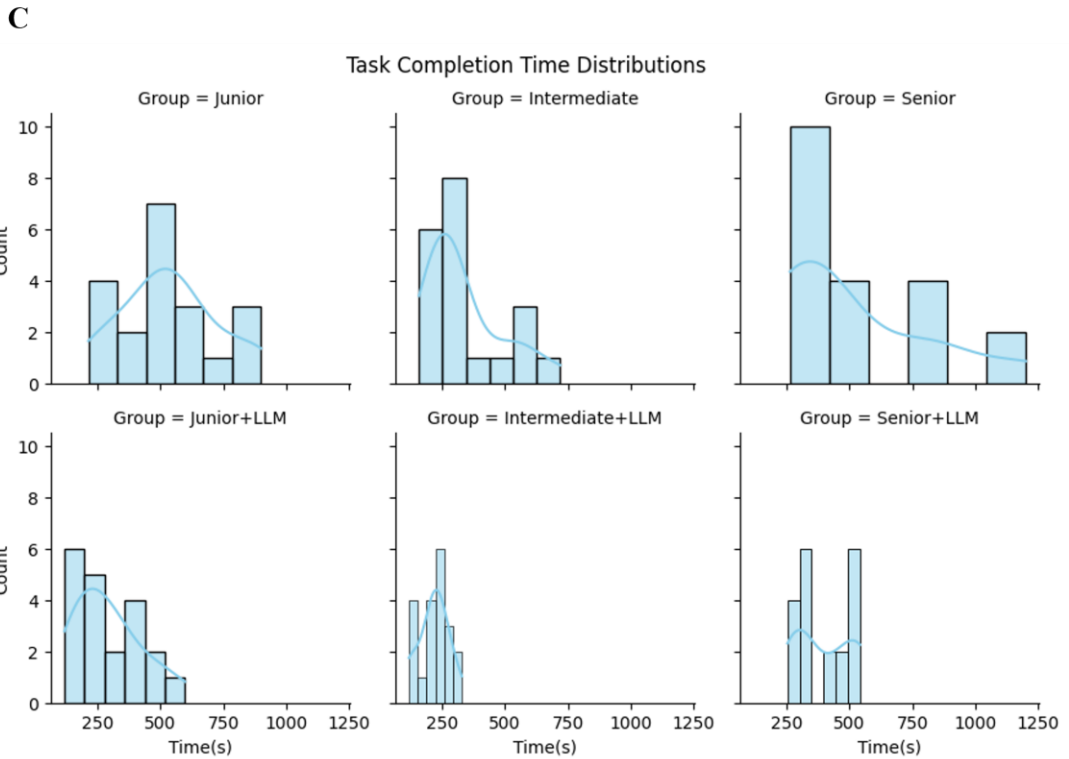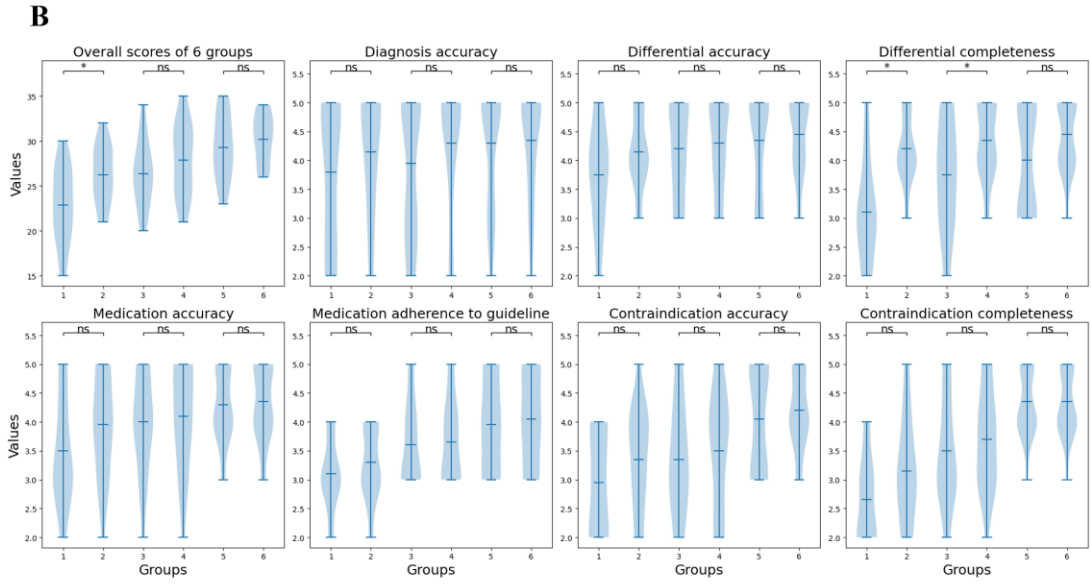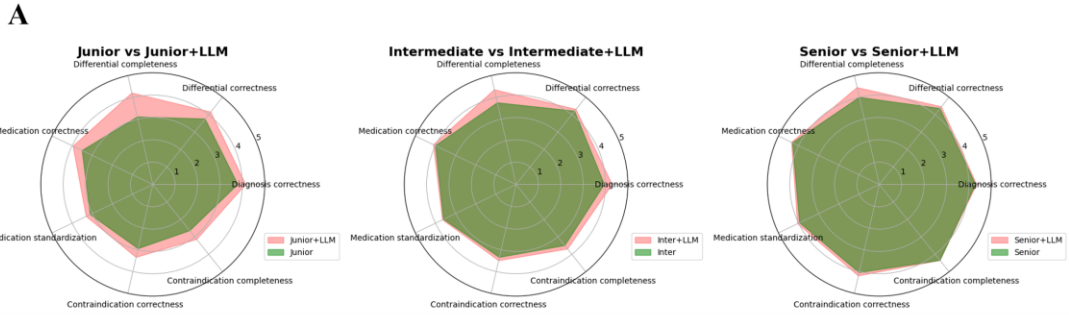
**Fig. 3 | Identifying the influence of prompt strategies: few-shot learning and chain-of-thought (CoT). A**, The performance of models with 0-shot prompt and 1-shot prompt on task1 across *BERTScore*, *ROUGE-L*, *BLEU*, completeness, and accuracy. **B**, The boxplots comparison between models with 0-shot prompt and 1-shot prompt on each metric in task1. **C**, The performance of models with 0-shot prompt and 1-shot prompt on task3 across *BERTScore*, *ROUGE-L*, *BLEU*, completeness, accuracy of principal diagnosis, and accuracy of differential diagnosis. **D**, The boxplots comparison between models with 0-shot prompt and 1-shot prompt on each metric in task3. **E**, The boxplots comparison between models without CoT prompt (0-shot) and with CoT prompt (0-shot) on each metric in task2 and task4. **F**, The heatmap of the predicted ICD codes given by models without CoT prompt against the reference ICD codes in task2. **G**, The heatmap of the predicted ICD codes given by models with CoT prompt against the reference ICD codes in task2.

**A** Distribution of Input Lengths for 5 Tasks

**B** Performance of LLMs with different context windows

**C**
Task1 Average Score vs Length
Task2 Average Score vs Length
Task3 Average Score vs Length
Task4 Average Score vs Length
Task5_QA Average Score vs Length
Task5_MC Average Score vs Length

**D** Huatuogpt2 vs Baichuan2

**E**
Task1-BLEU  Task1-ROUGE-L  Task1-BERTScore
Task3-BLEU  Task3-ROUGE-L  Task3-BERTScore
Task5-BLEU  Task5-ROUGE-L  Task5-BERTScore

51

**Fig. 4 | Identifying the influence of the input context length and the medicine-oriented fine-tuning.** A, the distribution of the input context lengths across 5 tasks in PsychBench. The unit of length is Chinese words. B, The performance of LLMs with different lengths of context windows across five tasks. Since GPT-3.5-turbo has a context length of $16k$, it was not included in this analysis. C, How model performance varies with input context length across five tasks. D-E, Performance comparison between medicine-oriented fine-tuned model HuatuoGPT2 and universal model Baichuan2 across each metric of 5 tasks.

**Fig. 5 | The analysis of clinical reader study.** A, The specialist evaluation of six groups (junior, junior+LLM, intermediate, intermediate+LLM, senior, senior+LLM) across Diagnostic correctness, Differential correctness, Differential completeness, Medication correctness, Medication standardization, Contraindicated correctness, and Contraindicated completeness. The overall scores are indicated as areas of each radar map. B, The comparison of overall scores and scores of each evaluation dimension of six groups. Group 1 to 6 represents group junior, group junior + LLM, group intermediate, group intermediate + LLM, group senior, and group senior + LLM, respectively. '*' indicates a statistically significant difference between the two groups (p-value less than 0.05), while 'ns' indicates no statistically significant difference between the two groups (p-value greater than 0.05). C, The distribution of time taken for each group to complete the three clinical tasks—diagnosis, differential analysis, and medication recommendation.

**Extended Data Tab. 1 | The Statistics of PsychBench dataset.**

| | Center 1 | Center 2 | Center 3 |
|---|---|---|---|
| **Age, M[IQR]** | 35.0 [25.75, 53.0] | 38.5 [23.0, 52.25] | 26.5 [15.0, 49.0] |
| **Gender, N** | | | |
| - Male | 34 | 13 | 40 |
| - Female | 66 | 87 | 60 |
| **Marriage, N** | | | |
| - Married | 51 | 52 | 41 |
| - Single | 31 | 38 | 55 |
| - Divorced or widowed | 18 | 10 | 4 |
| **Career, N** | | | |
| - Student | 12 | - | 44 |
| - Employed | 46 | - | 42 |
| - Unemployed | 33 | - | 10 |
| - Retired | 9 | - | 4 |
| **Ethnic group, N** | | | |
| - Han | 77 | - | 38 |
| - Man | 12 | - | 0 |
| - Hui | 11 | - | 2 |
| - Zang | 0 | - | 3 |
| - Bai | 0 | - | 32 |
| - Yi | 0 | - | 15 |
| - Other | 0 | - | 10 |
| **Family history, N** | | | |
| - **Yes** | 43 | 30 | 13 |
| - **No** | 57 | 70 | 87 |
| **Duration of illness, M[IQR]** | 4.0 [0.42, 10.0] | 4.0 [2.0, 15.0] | 2.0 [1.0, 5.0] |
| **Principal diagnosis (ICD-10)** | | | |
| - F10.x | 0 | 0 | 6 |
| - F20.x | 20 | 45 | 7 |
| - F30.x | 20 | 2 | 0 |
| - F31.x | 20 | 15 | 6 |
| - F32.x | 20 | 23 | 33 |
| - F33.x | 20 | 15 | 11 |
| - F90.x | 0 | 0 | 8 |
| - F98.x | 0 | 0 | 7 |
| - Others | 0 | 0 | 22 |

**Extended Data Tab. 2 | The statistics of evaluated models.**

| Model name | Parameters | Context length |
|---|---|---|
| Baichuan4 | - | 32k |
| Deepseek | 236B | 128k |
| Doubao-pro-32k | - | 32k |
| ERNIE-4-8k | - | 8k |
| Gemini-1.5-pro | 175B (MoE) | 2M |
| GLM-4 | 9B | 128k |
| Hunyuan-lite | MoE | 256k |
| Hunyuan-pro | 1T+ (MoE) | 32k |
| Minimax | 1T+ | 245k |
| Moonshot-v1-32k | 100B+ | 32k |
| Qwen-max | 100B+ | 8k |
| Spark-4ultra | - | 8k |
| Yi-large | 100B+ | 32k |
| GPT-3.5-turbo | 175B | 16k |
| GPT-4o | - | 128k |
| GPT-4 | - | 8k |
| Baichuan2-7b | 7B | 4k |
| HuatuoGPT2-7b | 7B | 4k |

**Extended Data Tab. 3 | The leaderboard of LLMs on PsychBench (1-shot).**

| Rank | Model | task1 | task2 | task3 | task4 | task5 | Overall |
|------|-------|-------|-------|-------|-------|-------|---------|
| 1 | Minimax | 85.18 | 74.30 | 86.22 | 70.75 | 87.43 | 80.78 |
| 2 | ERNIE-4-8k | 79.13 | 80.27 | 89.64 | 68.30 | 85.28 | 80.52 |
| 3 | Hunyuan-pro | 74.53 | 74.84 | 90.56 | 81.43 | 78.71 | 80.01 |
| 4 | Yi-large | 81.59 | 83.54 | 85.34 | 59.74 | 85.91 | 79.22 |
| 5 | Deepseek | 78.50 | 82.99 | 84.00 | 71.88 | 75.79 | 78.63 |
| 6 | Doubao-pro-32k | 90.07 | 77.01 | 93.87 | 33.08 | 98.96 | 78.60 |
| 7 | GPT-4 | 73.66 | 100.00 | 56.71 | 77.76 | 70.96 | 75.82 |
| 8 | Moonshot-v1-32k | 66.12 | 74.01 | 57.99 | 87.48 | 81.38 | 73.39 |
| 9 | GPT-4o-mini | 71.75 | 61.81 | 78.15 | 73.31 | 73.16 | 71.63 |
| 10 | Baichuan4 | 63.69 | 86.24 | 72.97 | 50.73 | 83.11 | 71.35 |
| 11 | GPT-3.5-turbo | 77.45 | 56.30 | 68.23 | 50.46 | 91.06 | 68.70 |
| 12 | Gemini-1.5-pro | 65.39 | 90.49 | 70.99 | 56.02 | 59.50 | 68.48 |
| 13 | GLM4 | 73.61 | 90.49 | 65.07 | 29.53 | 80.87 | 67.92 |
| 14 | Qwen-max | 53.35 | 90.08 | 61.13 | 55.56 | 63.67 | 64.76 |
| 15 | Hunyuan-lite | 71.53 | 20.00 | 67.46 | 81.30 | 68.97 | 61.85 |
| 16 | Spark4-Ultra | 57.86 | 65.88 | 57.31 | 28.10 | 28.49 | 47.53 |

**Extended Data Tab. 4 | The leaderboard of LLMs on PsychBench (0-shot).**

| Rank | Model | task1 | task2 | task3 | task4 | task5 | Overall |
|---|---|---|---|---|---|---|---|
| 1 | ERNIE-4-8k | 84.00 | 87.62 | 89.42 | 72.64 | 85.28 | 83.79 |
| 2 | Hunyuan-pro | 88.59 | 65.52 | 85.34 | 82.57 | 78.71 | 80.15 |
| 3 | Minimax | 86.40 | 64.20 | 87.21 | 71.55 | 87.43 | 79.36 |
| 4 | Yi-large | 70.75 | 82.76 | 77.28 | 76.79 | 85.91 | 78.70 |
| 5 | Doubao-pro-32k | 91.62 | 88.07 | 90.29 | 24.11 | 98.96 | 78.61 |
| 6 | GPT-4o-mini | 76.48 | 73.04 | 78.10 | 80.61 | 73.16 | 76.28 |
| 7 | Deepseek | 77.97 | 80.55 | 82.08 | 50.94 | 75.79 | 73.47 |
| 8 | Moonshot-v1-32k | 60.99 | 63.39 | 61.69 | 93.31 | 81.38 | 72.15 |
| 9 | GPT-3.5-turbo | 76.59 | 71.71 | 79.50 | 33.32 | 91.06 | 70.44 |
| 10 | GPT-4 | 80.99 | 100.00 | 55.57 | 44.22 | 70.96 | 70.35 |
| 11 | Gemini-1.5-pro | 76.95 | 80.63 | 66.32 | 67.36 | 59.50 | 70.15 |
| 12 | GLM4 | 62.27 | 86.44 | 60.85 | 53.66 | 80.87 | 68.82 |
| 13 | Baichuan4 | 51.15 | 77.90 | 64.41 | 57.33 | 83.11 | 66.78 |
| 14 | Qwen-max | 45.56 | 90.28 | 49.33 | 59.23 | 63.67 | 61.61 |
| 15 | Hunyuan-lite | 69.51 | 20.00 | 62.45 | 67.23 | 68.97 | 57.63 |
| 16 | Spark4-Ultra | 60.37 | 69.20 | 53.79 | 39.64 | 28.49 | 50.30 |

**Extended Data Tab. 5 | The evaluation criteria designed for reader study.**

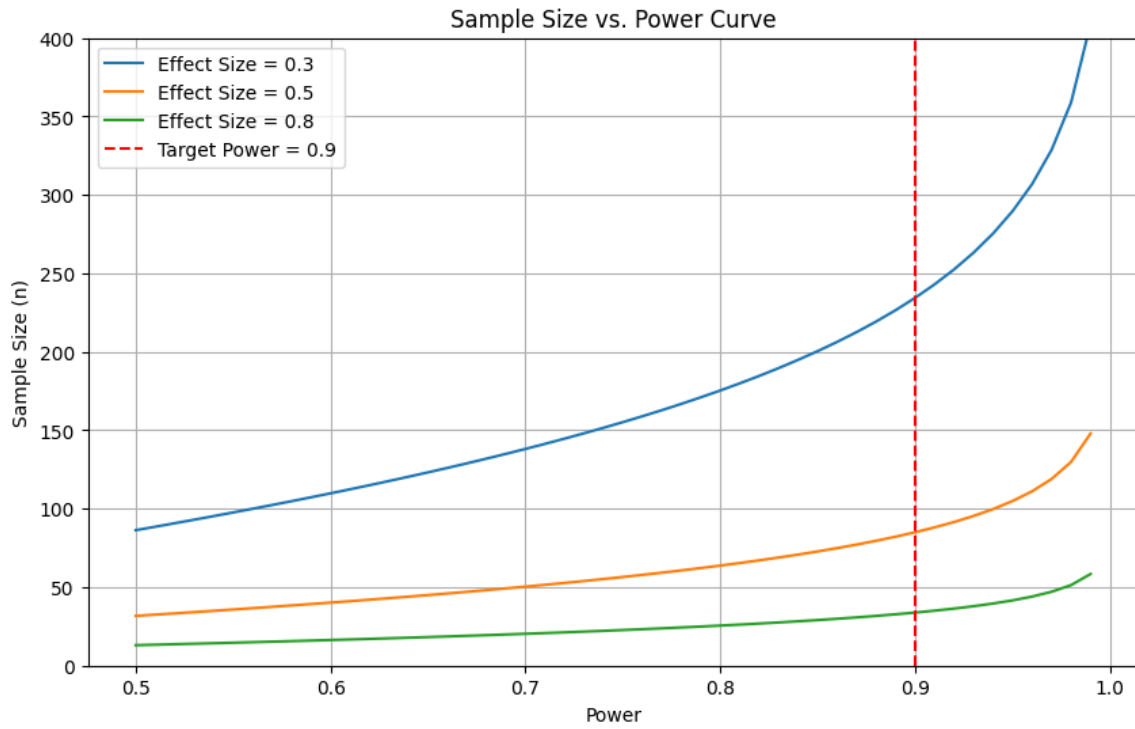| | **The Evaluation Criteria:** 7 dimensions for clinical performance assessment | | |
|---|---|---|---|
| | **Evaluation dimension** | **Example scenarios** | **Rating (1-5)** |
| **1** | **Diagnosis accuracy** | Assess whether the clinician accurately identifies and establishes the primary diagnosis using appropriate tools and methods, avoiding misdiagnosis. | ☐ 5 - Excellent<br>☐ 4 - Good<br>☐ 3 - Fair<br>☐ 2 - Needs Improvement<br>☐ 1 - Poor |
| | | **Scenario 1**: The clinician accurately assesses and confirms the primary psychiatric diagnosis based on symptoms, history, and clinical examination. | |
| | | **Scenario 2**: The clinician reviews prior diagnoses to prevent treatment delays caused by initial diagnostic errors. | |
| **2** | **Differential Accuracy** | Evaluate whether the clinician accurately rules out potential misdiagnoses, ensuring the diagnosis aligns with the patient's clinical presentation. | ☐ 5 - Excellent<br>☐ 4 - Good<br>☐ 3 - Fair<br>☐ 2 - Needs Improvement<br>☐ 1 - Poor |
| | | **Scenario 1**: The clinician confirms the primary diagnosis through further examination, ruling out any misdiagnoses. | |
| | | **Scenario 2**: The clinician uses clinical evidence to eliminate possible misdiagnoses. | |
| **3** | **Differential Completeness** | Assess whether the clinician comprehensively considers alternative diagnoses similar to the primary diagnosis, covering all relevant possibilities. | ☐ 5 - Excellent<br>☐ 4 - Good<br>☐ 3 - Fair<br>☐ 2 - Needs Improvement<br>☐ 1 - Poor |
| | | **Scenario 1**: The clinician lists potential differential diagnoses, covering all major conditions for a thorough assessment. | |
| | | **Scenario 2**: The clinician conducts a comprehensive analysis based on history and symptoms, ruling out all relevant psychiatric conditions. | |
| **4** | **Medication Accuracy** | Evaluate whether the clinician's medication recommendations align with the diagnosis and the patient's specific needs, ensuring the appropriateness of drug selection and dosage. | ☐ 5 - Excellent<br>☐ 4 - Good<br>☐ 3 - Fair<br>☐ 2 - Needs |

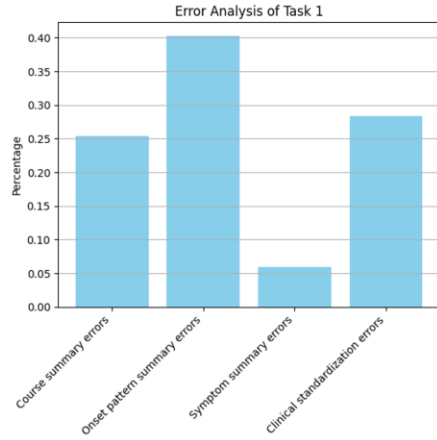| | | | Improvement<br>☐ 1 - Poor |
|---|---|---|---|
| | | **Scenario 1**: The medication was in accordance with the medication specifications, and there were no basic medication errors. | |
| | | **Scenario 2**: The medication was consistent with the patient's symptoms. | |
| 5 | **Medication Adherence to Guidelines** | Assess whether the clinician's medication recommendations follow clinical guidelines and standards, avoiding inappropriate practices. | ☐ 5 - Excellent<br>☐ 4 - Good<br>☐ 3 - Fair<br>☐ 2 - Needs Improvement<br>☐ 1 - Poor |
| | | **Scenario 1**: The clinician prescribes according to the latest clinical guidelines, with no inappropriate medication practices. | |
| | | **Scenario 2**: The clinician consults clinical guidelines before prescribing to ensure an evidence-based decision. | |
| 6 | **Contraindication Accuracy** | Verify whether the clinician accurately identifies and avoids contraindicated medications to ensure safe prescribing. | ☐ 5 - Excellent<br>☐ 4 - Good<br>☐ 3 - Fair<br>☐ 2 - Needs Improvement<br>☐ 1 - Poor |
| | | **Scenario 1**: The clinician identifies contraindications in the patient's profile and selects alternative medication. | |
| | | **Scenario 2**: The clinician thoroughly reviews the patient's history to ensure no contraindicated medications are prescribed. | |
| 7 | **Contraindication Completeness** | Assess whether the clinician thoroughly considers the patient's allergy history, past medical history, and potential drug interactions to avoid contraindications. | ☐ 5 - Excellent<br>☐ 4 - Good<br>☐ 3 - Fair<br>☐ 2 - Needs Improvement<br>☐ 1 - Poor |
| | | **Scenario 1**: The clinician gathers a comprehensive medication and medical history from the patient to avoid drug interaction risks. | |
| | | **Scenario 2**: The clinician assesses contraindications based on past medical and allergy history, avoiding all potential contraindicated medications. | |

**Extended Data Fig. 1 | Study design of the clinical reader study.**

| Input | Group1 | Group2 | Group3 | Group4 | Group5 | Group6 |
|---|---|---|---|---|---|---|
| Patient is a male, 58 years old. Starting in early 2020, without any apparent cause, the patient gradually began to experience disordered thinking, excessive rumination, and increased worry. He developed a low mood, reluctance to engage in activities… | **Task1:** F32.301 | **Task1:** F32.301 | **Task1:** F33.301 | **Task1:** F32.301 | **Task1:** F32.301 | **Task1:** F32.301 |
| | **Task2:** Schizophrenia… | **Task2:** Schizophrenia,GAD.. | **Task2:** BD,Schizophrenia… | **Task2:** BD,GAD,DD… | **Task2:** GAD… | **Task2:** GAD,DD… |
| | **Task3:** Duloxetine… | **Task3:** Mirtazapine … | **Task3:** Escitalopram… | **Task3:** Aripiprazole… | **Task3:** Milnacipran… | **Task3:** Aripiprazole… |
| Diagnosis correctness | 5 | 5 | 2 | 5 | 5 | 5 |
| Differential correctness | 3 | 4 | 2 | 4 | 3 | 4 |
| Differential completeness | 3 | 3 | 2 | 3 | 2 | 3 |
| Medication correctness | 5 | 3 | 2 | 3 | 3 | 3 |
| Medication standardization | 5 | 3 | 2 | 3 | 3 | 3 |
| Contraindication correctness | 1 | 4 | 3 | 4 | 3 | 4 |
| Contraindication completeness | 1 | 3 | 2 | 3 | 2 | 3 |

**Extended Data Fig. 2 | The reader study user interface.** Ellipses indicate that the information is not fully presented.

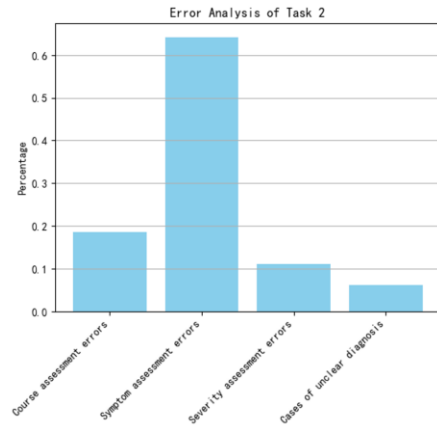**Extended Data Fig. 3 | Power analysis curve.**

**Error Analysis of Task 1**

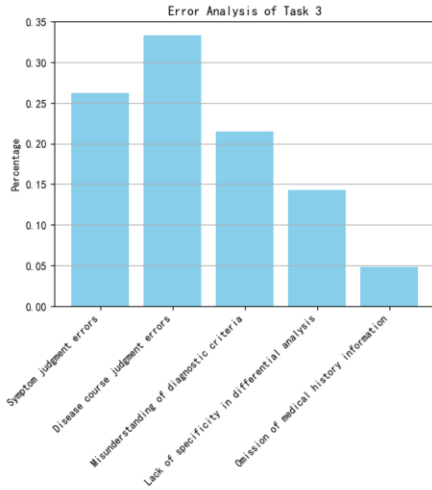| Error type | Example | Green: Correct answer / Red: Wrong answer |
|---|---|---|
| **Course summary errors** | **Chief Complaint:** Low mood accompanied by suicidal thoughts for 3 years (8 years), worsening over the past month.<br>**Case Characteristics:**<br>**Course of Illness:** Persistent depressive symptoms lasting for 3 years (8 years), , with a notable exacerbation over the past month.<br>**Symptomatology:** Primary symptoms include low mood, loss of interest in life, aversion to the idea of a second child, presence of suicidal ideation, poor sleep, and moderate appetite. The patient's mental activity is relatively well-coordinated, and partial insight is observed.<br>**Severity:** The illness has significantly impaired the patient's social functioning, disrupted family relationships, and there is evidence of suicidal behavior.<br>**Exclusion Criteria:** Currently, there is no evidence of a mental disorder induced by psychoactive substances. | |
| **Onset pattern summary errors** | **Chief Complaint:** Insomnia, low mood, and suicidal thoughts lasting for 9 months.<br>**Case Characteristics:**<br>**Course of Illness:** Chronic onset with a persistent course (intermittent course) over the past 9 months.<br>**Symptomatology:** The patient is fully conscious, with a worried expression and passive demeanor, but responds appropriately to questions. Speech volume and amount have increased compared to baseline. The patient spends most of the time lying in bed, exhibiting minimal activity. They report feeling agitated, depressed, concerned about family finances, and worried about both physical and mental health issues. They express a sense of meaninglessness in life, occasionally having suicidal thoughts as a form of escape, but without a clear plan or any suicidal actions. Sleep is shallow at night, and appetite has decreased, citing that food seems overly salty. Partial insight is present.<br>**Severity:** Significant impairment in social functioning, and the patient has ceased working.<br>**Exclusion Criteria:** Thyroid cancer and side effects of medications have been ruled out as causes for the mental health symptoms. | |
| **Symptom summary errors** | **Chief Complaint:** Poor appetite, crying, low mood, and episodic anxiety lasting 7 months following a work mistake.<br>**Case Characteristics:**<br>**Course of Illness:** Chronic, persistent course lasting 7 months.<br>**Symptomatology:** Fully oriented, passive interaction, with low mood and anxiety. Occasional symptoms include excessive excitement, increased activity (there is no these symptoms), feelings of being talked about by others, and concerns about being harmed. Partial insight is present.<br>**Severity:** Significant impairment in social functioning, unable to work.<br>**Exclusion Criteria:** Excludes mental disorders caused by organic brain conditions or psychoactive substance use. | |
| **Clinical standardization errors** | **Chief Complaint:** Low mood, prominent suicidal ideation, and one instance of self-harming behavior. (Poor mood, suicidal ideation for over two years, with exacerbation in the past month.)<br>**Case Characteristics:**<br>**Course of Illness:** Approximately 3 years, with recent exacerbation of symptoms.<br>**Symptomatology:** Depressed mood, loss of interest in life, prominent suicidal ideation, self-harming behavior, poor sleep, and difficulty concentrating.<br>**Severity:** Impaired social functioning, unable to attend school regularly.<br>**Exclusion Criteria:** No symptoms of mania, hypomania, or psychosis; no substance or other medical conditions causing the mental disorder. | |

**Extended Data Fig. 4 | Error analysis of Task 1.**

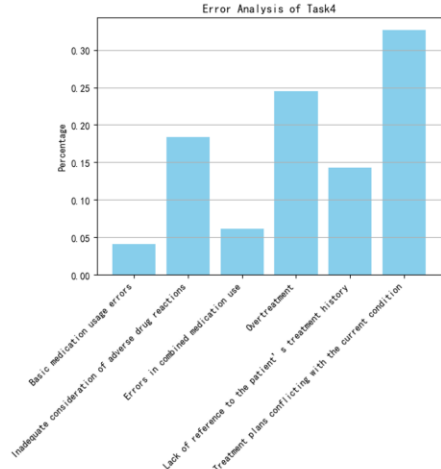| Error type | Example | Green: Correct answer<br>Red: Wrong answer |
|---|---|---|
| Course assessment errors | **Primary Diagnosis:** F32.201 Major Depressive Episode without Psychotic Features (F33.201 Recurrent Depressive Disorder, currently a Major Episode without Psychotic Features)<br>**Comorbid Psychiatric Diagnosis:** F42.- Obsessive-Compulsive Disorder | |
| Symptom assessment errors | **Primary Diagnosis:** F32.201 Major Depressive Episode without Psychotic Features (F32.301 Major Depressive Episode with Psychotic Features)<br>**Comorbid Psychiatric Diagnosis:** None | |
| Severity assessment errors | **Primary Diagnosis:** F32.001 Mild Depressive Episode (F32.201 Major Depressive Episode Without Psychotic Features)<br>**Comorbid Psychiatric Diagnosis:** F90 Attention-Deficit/Hyperactivity Disorder (ADHD) | |
| Case of unclear diagnosis | **Primary Diagnosis:** F30.802 Manic Episode (F31.101 Bipolar Disorder, Currently in a Manic Episode Without Psychotic Features)<br>**Comorbid Psychiatric Diagnosis:** None | |

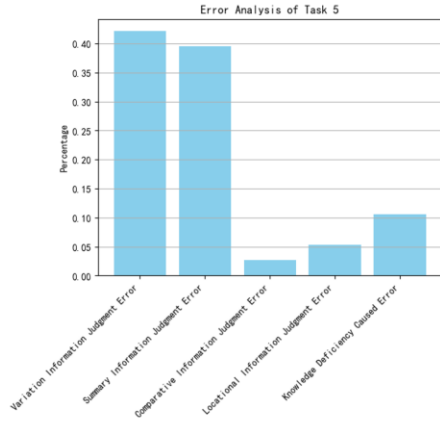**Extended Data Fig. 5 | Error analysis of Task 2.**

**Error Analysis of Task 3**

Percentage (y-axis: 0.00 to 0.35)

Categories (x-axis):
- Symptom judgment errors
- Disease course judgment errors
- Misunderstanding of diagnostic criteria
- Lack of specificity in differential analysis
- Omission of medical history information

| Error type | Example | **Green**: Correct answer  **Red**: Wrong answer |
|---|---|---|
| **Symptom judgement errors** | **Differential Diagnosis Analysis:** The patient is a 34-year-old female, and her pre-morbid personality was not clearly mentioned. The patient developed symptoms of emotional agitation and loss of emotional control following childbirth and under significant work stress. The course of illness has shown marked fluctuations, characterized by intermittent episodes. There is no history of psychiatric disorders in the patient's personal or family history. Symptom analysis reveals that the patient primarily exhibits emotional agitation, uncontrolled emotions, verbal aggression, family arguments, irritability, throwing objects, physical aggression towards family members, talking to herself, disorganized speech, and poor sleep. **Diagnosis: Primary Diagnosis:** [Bipolar Disorder]. The patient's symptoms meet the ICD-10 diagnostic criteria for Bipolar Disorder, with prominent mood swings, alternating between elevated and depressed moods, and a history of fluctuating course. ([Manic Episode without Psychotic Features]. The patient's main clinical symptoms include intermittent mood instability, excitement, disorganized speech, impulsive behavior, irritability, increased activity, and decreased need for sleep, which are indicative of manic symptoms. There are no prominent hallucinations or delusions, nor has there been any previous history of depressive or manic episodes.) **Differential Diagnosis:** <Differential Diagnosis 1> [Acute and Transient Psychotic Disorder], Diagnostic Point: Although the patient's symptoms meet the diagnostic criteria for Acute and Transient Psychotic Disorder, the marked fluctuations in the course of illness and the spontaneous improvement of symptoms without medication more strongly support the diagnosis of Bipolar Disorder. <Differential Diagnosis 2> [Manic Episode], Diagnostic Point: Although the patient's symptoms meet the diagnostic criteria for a Manic Episode, the significant fluctuations in the course of illness and the spontaneous resolution of symptoms without medication are more consistent with the diagnosis of Bipolar Disorder. |
| **Disease course judgement errors** | **Differential Diagnosis Analysis:** The patient is a 21-year-old female, with no clear mention of pre-morbid personality traits. The onset may be related to academic and employment pressures, with a course of about 5 years, presenting as a persistent episodic pattern, with suicidal behaviors. (The medical history is characterized by intermittent episodes lasting more than 2 weeks, primarily presenting with obsessive thoughts, low mood, reduced interest, and sleep disturbances.) Her past medical history includes laser surgery for a facial birthmark and surgery for a breast fibroadenoma. Family history includes hypertension and diabetes in her grandfather. Symptom analysis: The patient primarily exhibits symptoms of low mood, lack of interest in anything, recurrent unwanted thoughts, suicidal ideation and behavior, poor sleep, low self-esteem, pessimism, world-weariness, and a sense of worthlessness. **Diagnosis: Primary Diagnosis:** [Depressive Disorder]. (Recurrent Depressive Disorder, currently a Major Episode without Psychotic Features) The patient's main symptoms include low mood, lack of interest in anything, recurrent unwanted thoughts, suicidal ideation and behavior, poor sleep, low self-esteem, pessimism, world-weariness, and a sense of worthlessness, which meet the ICD-10 diagnostic criteria for depressive disorder. **Differential Diagnosis:** <Differential Diagnosis 1> [Obsessive-Compulsive Disorder], Differential Point: The patient experiences recurrent unwanted thoughts, which is a typical symptom of obsessive-compulsive disorder. However, the patient's primary symptoms are low mood, lack of interest in anything, suicidal ideation and behavior, poor sleep, low self-esteem, pessimism, world-weariness, and a sense of worthlessness, which align more with depressive disorder. <Differential Diagnosis 2> [Bipolar Disorder], Differential Point: The hallmark of bipolar disorder is the alternation of manic and depressive episodes, but the patient's symptoms are primarily depressive without any signs of mania, thus not fitting the criteria for bipolar disorder. |
| **Misunderstanding of diagnostic criteria** | **Differential Diagnosis Analysis:** The patient is a 17-year-old female, with no clear pre-existing personality traits. The onset was triggered by unsatisfactory mock exam results. The course of illness is brief, with a frequency of two episodes, a short duration, and significant fluctuations, presenting as an intermittent episodic pattern. The patient's past medical history includes a head injury at age 4 from falling off a truck, but no abnormalities were found upon examination at the hospital. The family history is negative. Symptom analysis: The patient mainly presents with excessive talking, excitement, increased self-confidence, impulsive verbal aggression, suspicion that family members intend to harm her, and auditory hallucinations (hearing voices). The severity of the symptoms is high, and social functioning is severely impaired. **Diagnosis: Primary Diagnosis:** [Bipolar Disorder]. (Manic Episode without Psychotic Features) The patient displays symptoms of excitement, excessive talking, increased self-confidence, and impulsive verbal aggression, with episodes occurring twice, each lasting a short duration, and significant fluctuations throughout the course. These characteristics align with the diagnostic criteria for bipolar disorder. (The patient exhibits only manic symptoms without any depressive features. This does not meet the diagnostic criteria for Bipolar Affective Disorder, as the condition requires the presence of both manic and depressive episodes. Therefore, the diagnosis should be revised to Manic Episode without Psychotic Features.) **Differential Diagnosis:** <Differential Diagnosis 1> [Schizophrenia], Diagnostic Point: The patient has auditory hallucinations, which are typical of schizophrenia. However, the patient's episodes are brief, with a frequency of only two, and significant fluctuations in the course. There is no obvious thought disorder, which does not meet the diagnostic criteria for schizophrenia. <Differential Diagnosis 2> [Manic Episode], Diagnostic Point: The patient exhibits symptoms of excitement, excessive talking, increased self-confidence, and impulsive verbal aggression, which are consistent with a manic episode. However, the frequency and short duration of the episodes, coupled with significant fluctuations, and the absence of significant thought disorders, do not meet the diagnostic criteria for a manic episode. |
| **Lack of specificity in differential analysis** | **Differential Diagnosis Analysis:** The patient is a 67-year-old female, with no clear description of pre-morbid personality. The onset trigger is unclear, and the course has lasted for about 4 months, presenting with symptoms such as temper outbursts, excessive talking, throwing things, refusal to take medication, constantly sending WeChat messages, disorganized speech, chaotic behavior, and lying on the floor. The past medical history includes chronic conditions such as hypertension, diabetes, hypothyroidism, and sinus bradycardia, with a history of pacemaker implantation. Symptom analysis: The patient mainly presents with symptoms of excitement, excessive talking, rapid speech, pressured speech, grandiosity, emotional instability, irritability, and a high demand for attention. When her demands are unmet, she becomes angry and ignores others. She denies having any illness, requests discharge and has no insight into her condition. **Diagnosis: Primary Diagnosis:** [Manic Episode]. According to the ICD-10 diagnostic criteria, the primary diagnosis is Manic Episode. **Differential Diagnosis:** <Differential Diagnosis 1> Bipolar Disorder, Diagnostic Point: The characteristic of bipolar disorder is the alternation of manic and depressive episodes, but the patient has not presented any obvious depressive symptoms, so this diagnosis is not currently considered. <Differential Diagnosis 2> Schizophrenia, Diagnostic Point: Schizophrenia primarily includes symptoms such as thought disorder, hallucinations, and delusions, but the patient has not exhibited any of these symptoms, so this diagnosis is not currently considered. (<Differential Diagnosis> Organic mental disorder. The patient has comorbidities including hypertension, diabetes, and other physical conditions, all of which are risk factors for cardiovascular and cerebrovascular diseases. Given the patient's emotional instability and excitement, organic causes should be considered. However, the current physical examination does not reveal any significant abnormalities, and the head CT scan shows no obvious infarct lesions. Therefore, this condition does not match the typical characteristics of an organic disorder and is not currently considered.) |
| **Omission of medical history information** | **Differential Diagnosis Analysis:** The patient is a 37-year-old female, with no clear pre-existing personality description. She has marital conflicts, depression, a history of suicidal behavior, a drinking history, and a family history (mother with depressive disorder, sister with bipolar disorder). The onset of symptoms was triggered by marital conflict, with a disease course lasting 4 years. The symptom onset frequency is continuous, with a duration of 4 years, and the course shows fluctuation, with periods of improvement and worsening. Past medical history: Physically healthy. Family history: Mother with depressive disorder, sister with bipolar disorder. Symptom analysis: The patient mainly presents with low mood, insomnia, a sense of meaninglessness, compulsive drinking behavior, and suicidal behavior, with a history of alcohol use. **Diagnosis: Primary diagnosis:** [Depressive Disorder]. The patient has persistent low mood, insomnia, a sense of meaninglessness, suicidal behavior, and a history of alcohol use, which meets the ICD-10 diagnostic criteria for depressive disorder. ([Bipolar Disorder]. The patient's history of depressive episodes was overlooked, so the patient should be diagnosed with Bipolar Disorder.) **Differential Diagnosis:** < Differential Diagnosis 1 > [Bipolar Disorder], Differential point: Bipolar disorder is characterized by alternating periods of manic and depressive episodes. However, the patient primarily presents with persistent low mood and has no significant manic episodes, making the diagnosis of bipolar disorder unlikely. < Differential Diagnosis 2 > [Alcohol-related Mental and Behavioral Disorders], Differential point: Alcohol-related mental and behavioral disorders typically involve cognitive, mood, and behavioral disturbances due to long-term excessive alcohol consumption. Although the patient has a history of alcohol use, her main symptoms are persistent low mood, insomnia, a sense of meaninglessness, and suicidal behavior, which more closely align with depressive disorder. Therefore, the likelihood of an alcohol-related mental and behavioral disorder is low. |

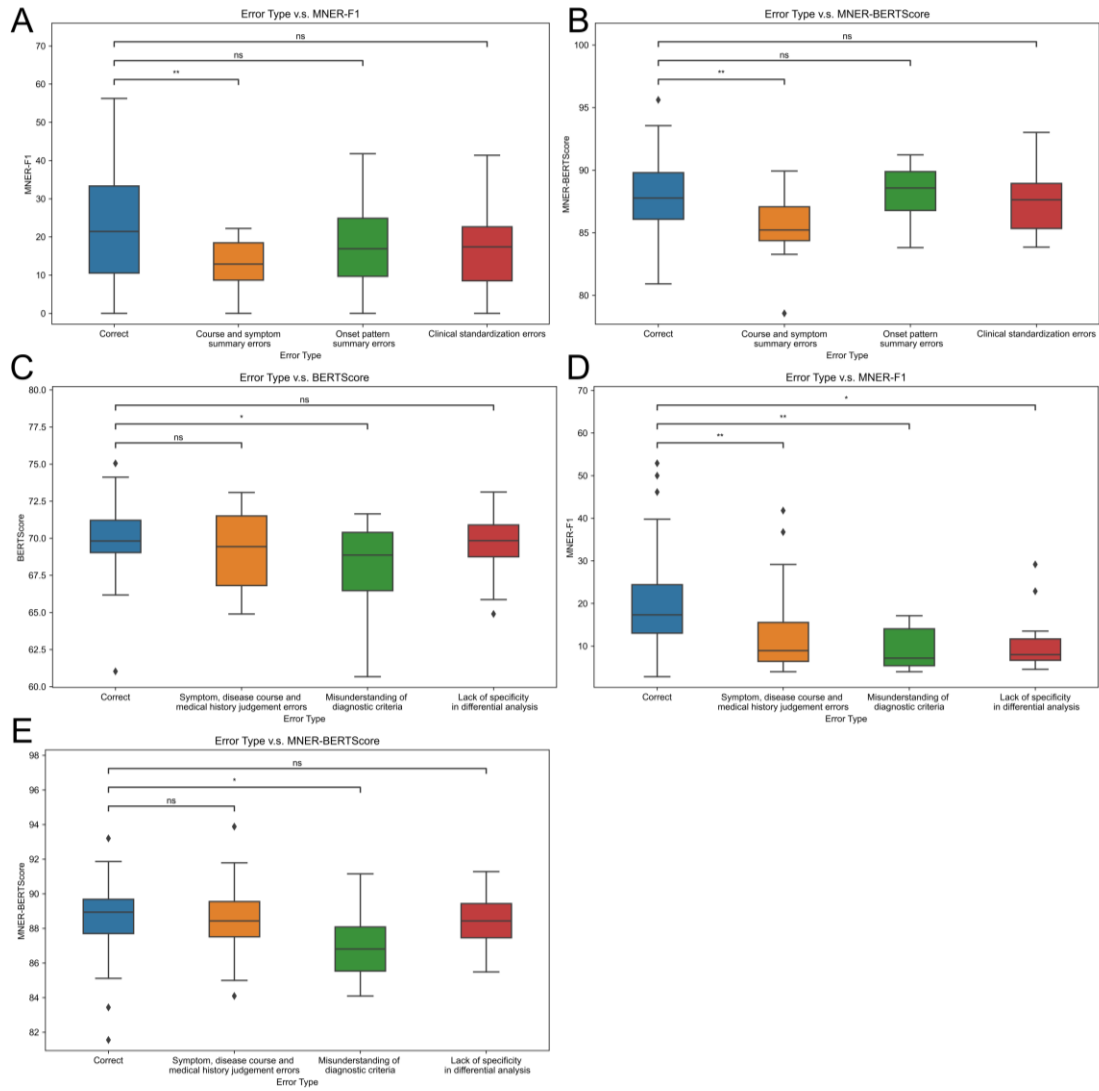**Extended Data Fig. 6 | Error analysis of Task 3.**

Error Analysis of Task4

Percentage (y-axis: 0.00 to 0.30)

Categories:
- Basic medication usage errors
- Inadequate consideration of adverse drug reactions
- Errors in combined medication use
- Overtreatment
- Lack of reference to the patient's treatment history
- Treatment plans conflicting with the current condition

| Error type | Example | Green: Correct answer / Red: Wrong answer |
|---|---|---|
| Basic medication usage errors | **Analysis:** According to the patient's history, they experienced an acute and brief episode of a psychiatric disorder, characterized by speech disorganization, hallucinations, and delusions. Symptoms improved after the use of haloperidol in the emergency department. The patient was subsequently hospitalized at another hospital, where risperidone and lorazepam were administered. The patient reported improvement and was discharged. However, after discharge, the patient discontinued the medication on their own, leading to a recurrence of symptoms, such as depressive mood and chest tightness. Symptoms improved again after resuming medication, but the patient still had ongoing emotional issues. Given the family history of psychiatric disorders and the current symptom profile, it is recommended to continue the use of antipsychotic medications and consider adding an antidepressant. **Recommended medications:** **Risperidone:** The patient responded well to risperidone in the past, so it can be continued. Risperidone is an atypical antipsychotic with good efficacy for both positive symptoms (such as hallucinations and delusions) and negative symptoms (such as emotional flatness and social withdrawal). It is recommended to start at 2 mg (1mg) in the evening, with dose adjustments based on symptoms. **Lorazepam:** The patient responded well to lorazepam in the past, so it can be continued. Lorazepam is a benzodiazepine with anxiolytic and sedative properties that can help alleviate the patient's tension and chest tightness. It is recommended to maintain a dose of 0.5 mg in the evening. **Paroxetine:** Given the patient's current symptoms of low mood and lack of interest, which are suggestive of depression, the addition of paroxetine is recommended. Paroxetine is a selective serotonin reuptake inhibitor (SSRI) that is effective for treating depressive symptoms. It is recommended to start at 10 mg (20mg) daily, with dose adjustments based on symptoms. **Monitoring and follow-up:** It is important to regularly monitor the patient's biochemical indicators, such as liver and kidney function tests, as well as electrocardiogram (ECG) to detect potential side effects of the medications. Additionally, the patient should undergo psychotherapy, such as cognitive behavioral therapy (CBT), to improve emotional and behavioral issues. Family members should also provide support and understanding to help the patient through this difficult period. | |
| Inadequate consideration of adverse drug reactions | **Analysis:** The patient has a long history of depression, anxiety symptoms, and alcohol dependence syndrome. The patient has previously been treated with medications including oxazepam, valproate, fluoxetine, and alprazolam, but adherence to treatment has been poor, with multiple instances of self-discontinuation or dose adjustments. The patient's symptoms worsen after alcohol consumption, including mood swings, reckless behavior, hallucinations, and delusions. These symptoms may be related to alcohol withdrawal or alcohol-induced psychiatric disorders. Additionally, the patient has severe fatty liver disease, so the impact of medications on liver function must be considered. In light of these factors, it is recommended to adjust the medication regimen and enhance psychological support and alcohol cessation guidance. **Recommended Medications:** **Lorazepam:** For short-term relief of anxiety symptoms. Care should be taken with the dosage and duration of use to avoid dependency. **Paroxetine:** For treating depressive and anxiety symptoms. Monitoring for drug response and side effects is essential. **Agomelatine:** An antidepressant with sleep-improving properties, suitable for patients with co-occurring sleep disturbances. (Agomelatine is metabolized by the liver, and the incidence of liver dysfunction is relatively high. Given that the patient has liver-related issues, the use of agomelatine should be avoided.) **Mirtazapine:** An antidepressant that can improve appetite, suitable for patients with appetite loss. **Lithium Carbonate:** Used to stabilize mood and prevent manic episodes. Blood lithium levels must be monitored regularly. **Aripiprazole:** Used to treat alcohol withdrawal symptoms and improve behavior control. Dose adjustments and side effect monitoring are necessary. **Monitoring and Follow-up:** When using the above medications, the patient's specific conditions—such as liver and kidney function and potential drug interactions—must be considered. Dose adjustments and monitoring should be done under medical guidance. Additionally, it is recommended that the patient receive psychotherapy and alcohol cessation support, alongside lifestyle modifications, to improve treatment outcomes and overall quality of life. | |
| Errors in combined medication use | **Analysis:** Based on the patient's history, they have a long-standing history of psychiatric symptoms, including auditory hallucinations, mood swings, irritability, depression, and suicidal ideation. The patient previously used valproate, which had partial efficacy but was discontinued due to side effects. The symptoms are consistent with depression and possible psychotic features, requiring a comprehensive approach that includes both pharmacological and psychological treatments. When selecting medications, factors such as the patient's age, gender, medical history, family history, and potential drug interactions need to be considered. **Recommended Medications:** **Lorazepam:** Used for short-term relief of the patient's anxiety and tension. **Paroxetine:** As an antidepressant, it can help alleviate the patient's depressive symptoms and suicidal ideation. **Quetiapine:** As an atypical antipsychotic, it can be used to treat the patient's auditory hallucinations and mood fluctuations. **Duloxetine:** As an antidepressant, it can be used in combination with paroxetine to enhance the antidepressant effect. (The combination of duloxetine and paroxetine is generally not recommended, as their mechanisms of action overlap, increasing the risk of side effects such as serotonin syndrome.) **Sertraline:** As an antidepressant, it can be used as an alternative or adjunct to paroxetine to improve antidepressant effects. **Aripiprazole:** As an atypical antipsychotic, it can be used to treat the patient's psychotic symptoms. **Lithium carbonate:** As a mood stabilizer, it can be used to prevent mood swings and suicidal behavior. **Monitoring and Follow-up:** During treatment, close monitoring of the patient's symptoms and side effects is required. Adjustments to medication dosage and treatment plans should be made according to the patient's responses. Additionally, the patient should receive psychological therapy, such as cognitive-behavioral therapy, to improve their coping skills and ability to manage emotional fluctuations. Moreover, attention should be given to the patient's family environment and support system to help establish a strong social support network and reduce mental stress. | |
| Overtreatment | **Analysis:** Based on the patient's medical history, she has been diagnosed with "anxiety state, insomnia" and has received treatment with paroxetine, trazodone, and alprazolam. Although her sleep duration has improved after treatment, she still experiences difficulties falling asleep, shallow sleep, low mood, anxiety, and pessimism. Additionally, the patient has developed persecutory delusions and delusions of reference, which may suggest that her condition has progressed to a more severe mental illness, such as depression with psychotic features or schizophrenia. Given the patient's age and comorbidities, including diabetes and hypertension, careful consideration is required when selecting medications to avoid drug interactions and side effects. **Recommended medications:** **Olanzapine:** As an atypical antipsychotic, olanzapine can effectively control the patient's psychotic symptoms, such as delusions and hallucinations, while also providing some mood stabilization. However, it is important to monitor for potential metabolic side effects, such as weight gain, elevated blood glucose, and lipid levels. **Mirtazapine:** As an antidepressant, mirtazapine can improve symptoms of low mood and pessimism. It also has sedative properties, which may help with sleep improvement. However, it may cause side effects such as weight gain and hypotension. **Lorazepam:** As a benzodiazepine, lorazepam can relieve the patient's anxiety symptoms and improve sleep. However, attention should be paid to the risk of dependence and tolerance with prolonged use. **Duloxetine:** As an SNRI (serotonin-norepinephrine reuptake inhibitor), duloxetine can help alleviate depressive symptoms and also provides some anti-anxiety effects. However, potential side effects include nausea, dry mouth, and headaches. **Lithium carbonate:** As a mood stabilizer, lithium carbonate can help stabilize mood fluctuations and reduce episodes of mania and depression. It is important to monitor for side effects such as tremors, nausea, and polyuria, as well as regularly checking lithium blood levels. (Lithium carbonate is a potential medication for bipolar disorder or treatment-resistant depression. However, the likelihood of this patient having bipolar disorder is very low, and she is not considered a treatment-resistant case. Considering the potential side effects of lithium carbonate, it is unnecessary for this patient to use it.) **Monitoring and Follow-up:** During treatment, it is crucial to closely monitor the patient for any changes in symptoms and adverse drug reactions. The dosage and medication regimen should be adjusted according to the patient's response and tolerance. Additionally, psychological therapies, such as cognitive behavioral therapy (CBT), should be recommended to improve the patient's psychological symptoms and coping strategies. It is also essential to address the patient's comorbidities, including diabetes and hypertension, by adjusting medications as necessary to minimize the risk of complications. | |
| Lack of reference to the patient's treatment history | **Analysis:** The patient has experienced symptoms such as low mood, repetitive thoughts, and poor sleep since 2016, and was diagnosed with Obsessive-Compulsive Disorder (OCD) and treated with Olanzapine and Paroxetine, with some improvement. After discontinuing medication in 2018, the patient developed symptoms of low mood and lack of interest again in 2021, likely triggered by stress. The patient was re-treated with Paroxetine and Quetiapine, but no significant improvement was noted. In 2022, the patient self-increased the dosage of Paroxetine, but the condition did not improve, and depressive thoughts and suicidal ideation developed. Currently, the patient is experiencing low mood, suicidal thoughts, and poor sleep. Based on the patient's medical history, it is likely that the patient is suffering from depressive symptoms and a treatment adjustment is needed. **Recommended Medications:** **Paroxetine:** The patient has had some benefit from Paroxetine previously, but the condition has not fully resolved. Consider adjusting the dosage while monitoring changes in the patient's mood. (Although Paroxetine was effective in the past, its full-dose treatment over the past 4 months has shown poor efficacy, and the patient has developed suicidal behavior. Therefore, Paroxetine should not be recommended again.) **Quetiapine:** The patient has previously used Quetiapine, but at a low dose. Consider increasing the dose to help with sleep and mood symptoms. **Mirtazapine:** Mirtazapine has antidepressant effects and also helps with sleep. It may be considered as an adjunctive medication. **Agomelatine:** Agomelatine helps regulate the biological clock and has antidepressant properties. It could be considered as an adjunct medication. **Lithium Carbonate:** The patient has expressed suicidal ideation, and Lithium could be used to stabilize mood and help prevent suicidal behavior. **Monitoring and Follow-up:** In addition to medication adjustments, it is strongly recommended that the patient undergo psychotherapy, such as Cognitive Behavioral Therapy (CBT), to help improve cognitive patterns and emotional regulation. Close monitoring of the patient's condition and timely adjustments to the treatment plan are essential. | |
| Treatment plans conflicting with the current condition | **Analysis:** Based on the patient's history, she has shown symptoms of emotional instability, irritability, self-talk, and disorganized speech, which could be associated with postpartum mood disorder. However, considering the current symptoms, the possibility of **bipolar disorder** must be excluded. The patient experienced emotional agitation postpartum but did not seek medical attention, and her symptoms improved on their own after one month, suggesting a mild postpartum mood disorder. However, the recent appearance of emotional agitation, disorganized speech, self-talk, denial of some parts of her history, and lack of insight are more consistent with the characteristics of bipolar disorder. Therefore, the treatment should comprehensively consider the patient's condition and select appropriate medications. **Recommended Medications:** Considering the patient's condition, the following medications are recommended for treatment: **Valproate:** To control mood swings and prevent manic episodes. **Lamotrigine:** To stabilize mood and reduce emotional fluctuations. **Paroxetine:** To improve depressive symptoms and enhance emotional stability. (The patient is currently in a manic episode, and the use of antidepressants may exacerbate the manic symptoms.) **Lorazepam:** To relieve anxiety and tension, and improve sleep. **Aripiprazole:** To control manic symptoms and improve emotional stability. **Monitoring and Follow-up:** It is important that these medications are used under medical guidance. The drug types and dosages should be adjusted based on the patient's symptoms and drug responses. Additionally, regular monitoring of biochemical indicators, thyroid function, and other relevant health parameters is necessary to ensure the patient's physical health. Furthermore, the patient should undergo psychological therapy to help her adjust her mindset, improve her mood, and enhance her quality of life. | |

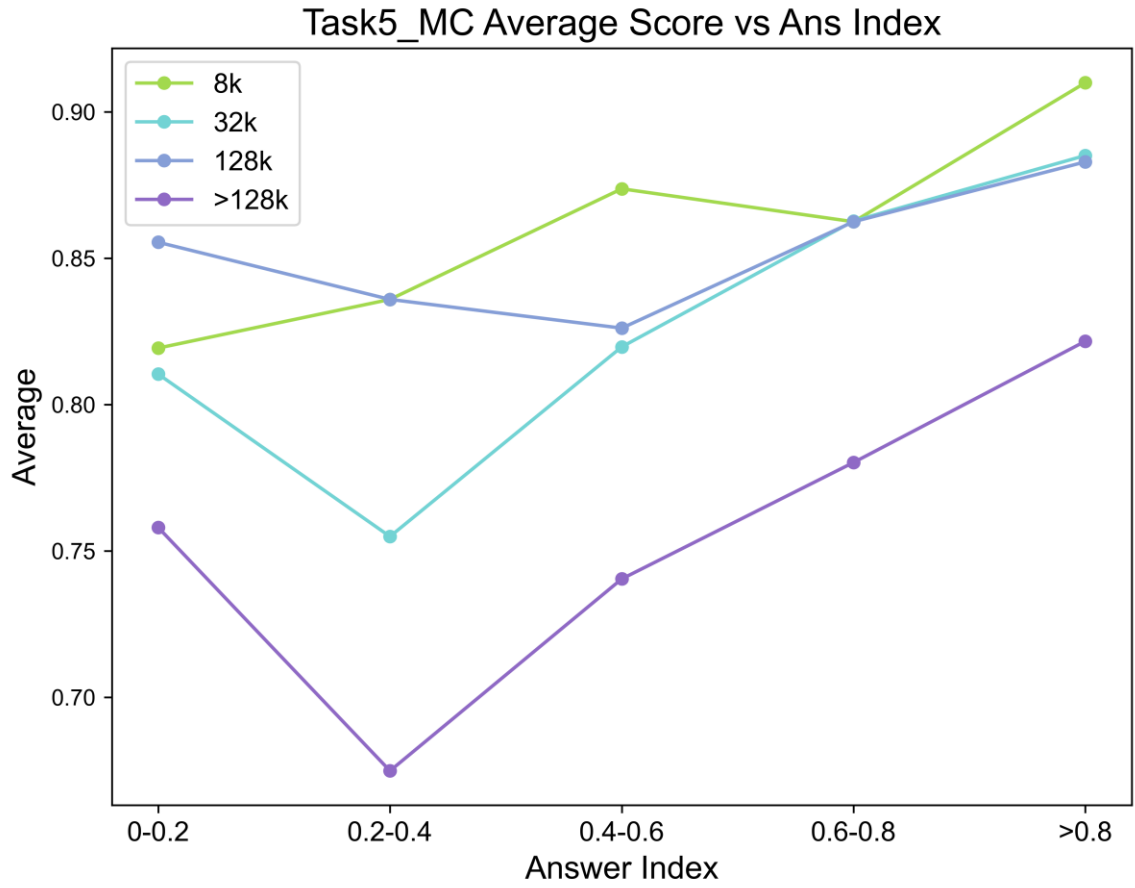**Extended Data Fig. 7 | Error analysis of Task 4.**

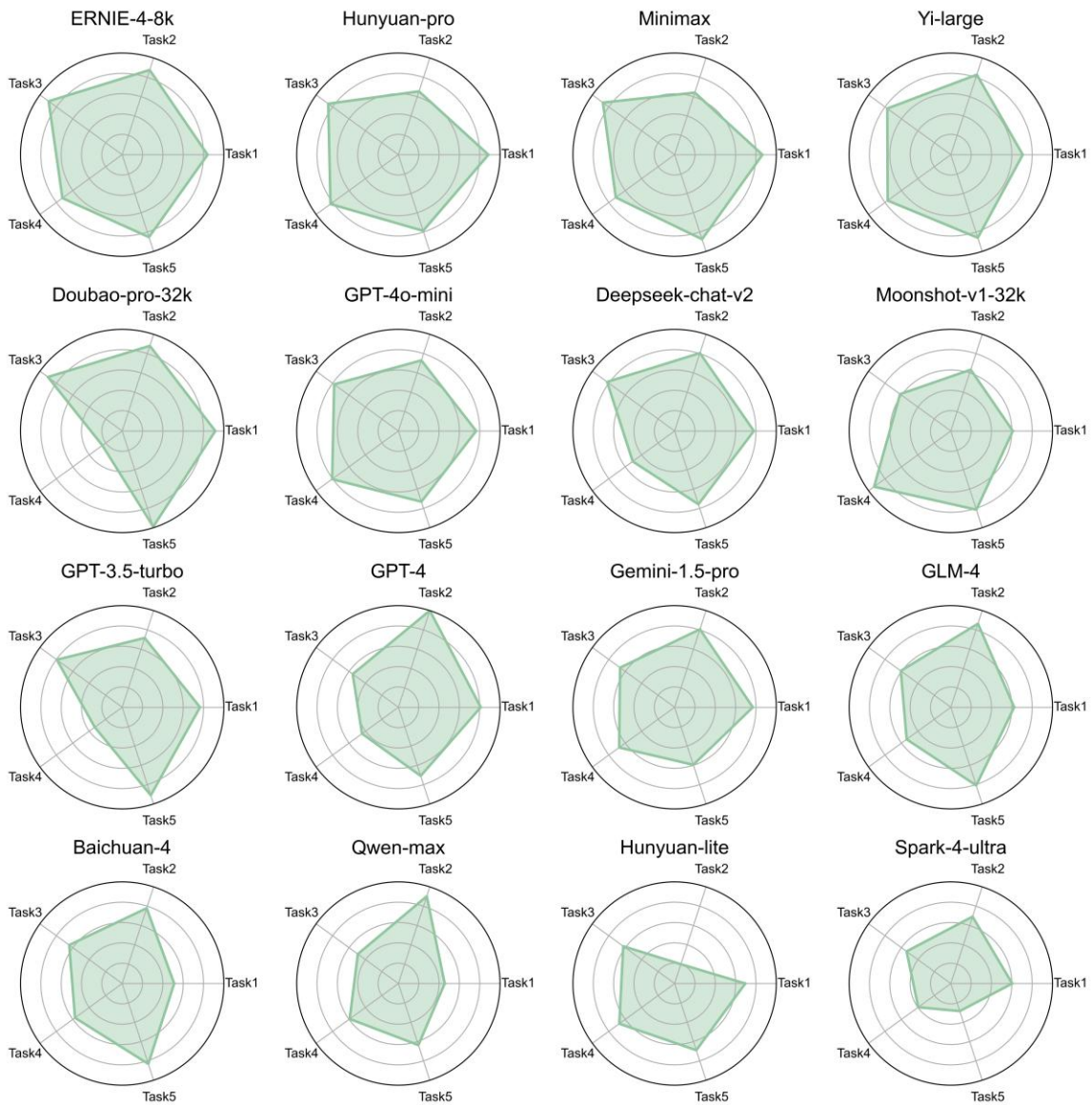| Error type | Example | Green: Correct answer<br>Red: Wrong answer |
|---|---|---|
| **Variation information judgement error** | **[Question]**<br>Which measure was adjusted in the patient's treatment to improve their mental state?<br>A. Increase sertraline to 100mg QD<br>B. Adjust MECT treatment to three times a week<br>C. Increase nifedipine extended-release tablets to 20mg in the morning<br>D. Keep the dose of mirtazapine tablets unchanged<br>**[Answer]**<br>A (B) | |
| **Summary information judgement error** | **[Question]**<br>According to the medical records, what were the main findings in the patient's mental examination?<br>A. The patient has obvious hallucinations<br>B. The patient denies having psychiatric symptoms<br>C. The patient has anxiety<br>D. The patient has regained insight<br>**[Answer]**<br>C (B) | |
| **Comparative information judgement error** | **[Question]**<br>During the patient's hospitalization, which medication had the highest dosage?<br>A. Sertraline 150mg QD<br>B. Sertraline 125mg QD<br>C. Sertraline 100mg QD<br>D. Lithium carbonate 0.25g BID<br>**[Answer]**<br>A (D) | |
| **Locational information judgement error** | **[Question]**<br>On which day did the patient begin taking Quetiapine fumarate, Lorazepam, and Escitalopram oxalate?<br>A. Day 13<br>B. Day 16<br>C. Day 19<br>D. Day 22<br>**[Answer]**<br>C (B) | |
| **Knowledge deficiency caused error** | **[Question]**<br>Which of the following laboratory test results is below the reference value, with the nature of the abnormality undetermined?<br>A. Complement C3 0.73g/L<br>B. High-density lipoprotein 1.15mmol/L<br>C. Activated partial thromboplastin time 34.70 seconds<br>D. Prolactin 31.12ng/ml<br>**[Answer]**<br>A (B) | |

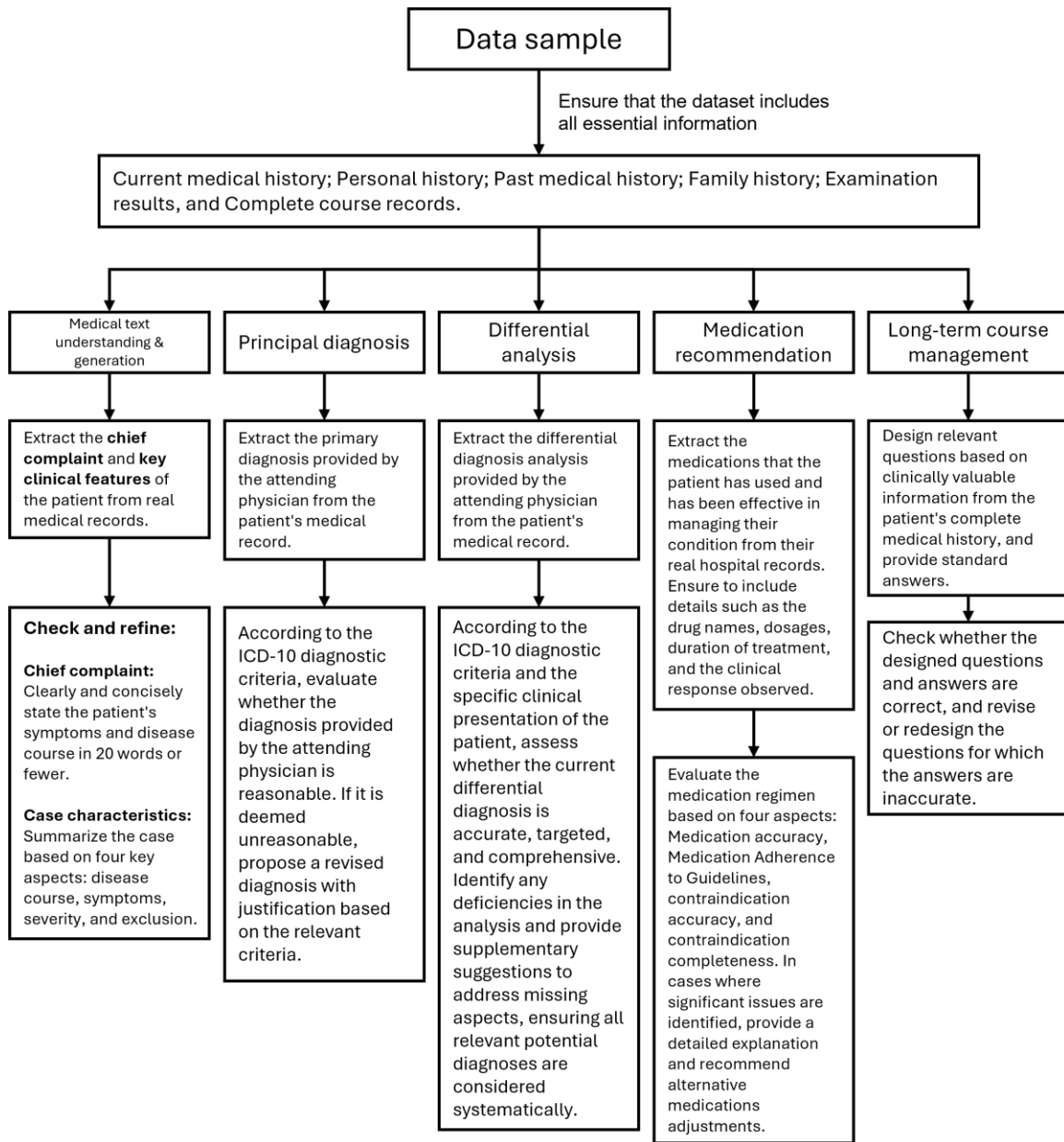**Extended Data Fig. 8 | Error analysis of Task 5.**

**Extended Data Fig. 9 | The statistics regarding model errors and quantitative metric scores in Task 1 (A,B) and Task 3 (C,D,E).**

**Extended Data Fig. 10 | The impact of answer index on the model performance on task 5(multi-choice).**

**Extended Data Fig. 11 | The normalized quantitative results of evaluated LLMs across five tasks (0-shot).**

**Extended Data Fig. 12 | The guidelines for data verification and annotation.**