

Supplementary Material

Yuqi Wu^{1,2†}, Guangya Wan^{3†}, Jingjing Li^{4*}, Shengming Zhao^{1,2},
Lingfeng Ma², Tianyi Ye², Mike Zhang⁵, Ion Pop⁵, Yanbo Zhang^{5*},
Jie Chen^{1,2*}

^{1*}College of Biomedical Engineering, Fudan University, Shanghai, 200433, China.

²Department of Electrical and Computer Engineering, University of Alberta, Edmonton, Alberta, T6G 2V4, Canada.

³School of Data Science, University of Virginia, Charlottesville, Virginia, 22903, USA.

⁴McIntire School of Commerce, University of Virginia, Charlottesville, Virginia, 22903, USA.

⁵Department of Psychiatry, University of Alberta, Edmonton, Alberta, T6G 2B7, Canada.

*Corresponding author(s). E-mail(s): jingjing.li@virginia.edu;
yanbo9@ualberta.ca; jc65@ualberta.ca;

Contributing authors: ywu_bh@fudan.edu.cn; wxr9et@virginia.edu;

[†]These authors contributed equally to this work.

1 Experiment Setup

1.1 LLM check-points

In this study, the LLM check-points we adopted up to the latest version in February 2025 including:

- "gpt-4o"
- "claude-3.5-sonnet"
- "mistral-large-latest" (offered by MistralAI API access)
- "qwen2.5-72b" (powered by Ollama ¹)

¹<https://ollama.com/>

- "EmoLLM" (LLaMA3.8B_instruct) ²
- "UltraMedical (Llama-3-8B-UltraMedical)" ³

1.2 Hardware Infrastructure

Our experiments were conducted on a computing infrastructure equipped with the following hardware:

- GPU: 2 NVIDIA RTX A6000
- CPU: AMD Ryzen Threadripper PRO 5975WX
- RAM: 4x64GB DDR4 3200MHz RDIMM ECC Memory
- Storage: 6TB, M.2, PCIe NVMe, SSD, Class 40

1.3 Package Used and Code

For generating LLMs responses and parsing out answers, we utilize packages "langchain", "langchain_openai", "langchain_anthropic", "langchain_community", and "langchain_core" offered by Langchain⁴. In addition, we use "pandas" for data processing, "matplotlib" and "seaborn" for visualization, and "numpy" for basic mathematical manipulation.

2 Structured Knowledge Graph

2.1 DSM-5 DDx Decision Tree

In this study, the structured knowledge graph originated from DSM-5-TR® Handbook of Differential Diagnosis [?] with modification. An example of original DDx decision tree for depression is shown as Fig. 1. This decision defines a common procedure for a clinician to conduct DDx when interviewing patients. It defines the critical topic that must be assessed and its main criteria. The structured knowledge graph is modified from these type of decision trees.

2.2 Example of Constructed SKGs

In this study, the DSM-5 DDx decision tree has been converted to binary tree style via bigtree Python package. In each node, it contains the node name (abbreviation which describes the topic), path (sequential information; structured knowledge graph), and description (contains the modified description of criteria which help LLM to make decisions). An example of such structured knowledge graph is shown in Fig. 2. Fig. 3 gives a demonstration on how to convert a DSM-5 based decision tree to a SKG.

3 Prompts

The detailed comparison between different memory structure can be found in Fig. 10.

²<https://github.com/SmartFlowAI/EmoLLM>

³<https://github.com/TsinghuaC3I/UltraMedical>

⁴<https://www.langchain.com/>

3.1 KFP

In KFP, the objective is given to the doctor agent. In each turn of the conversation, the agent asks the patient one diagnostic question until the agent believes that it can conduct a confident diagnosis. All the possible outcomes from the tested disorder are provided as labels for the agent as a reference. For instance, for depressed mood, all the 25 possible outcome of depressed mood DDx are provided as class labels. Therefore, in essence, this is simply a multi-turn zero-shot classification.

3.2 TKEP

In the TKEP memory setting, a similar multi-turn conversation is required to complete the diagnosis. However, this time, besides the possible outcome class labels, the descriptions of the critical nodes are exposed to the agent as external knowledge. In ICL, the entire knowledge graph without structure is provided to the agent as a reference, whose prompt is shown as follows:

System Message: *You are a psychiatrist tasked with conducting differential diagnosis via clinical interviews. Keep asking questions until the objective is met. DO NOT propose treatment plans. The final diagnostic labels will be provided. Avoid repeating questions and irrelevant information.*

Human Message: *Required Response Format: <Response;Ask necessary questions to help with diagnosis.>/Response; <Final_Decision;Provide final diagnosis or None if not ready.>/Final_Decision;*

Now, please proceed with the interview: The final diagnostic labels are {diagnostic_labels}, the patient responded: {patient_response}, Dialogue history: {st_memo}. Do not ask repeated questions.

As for RAG, instead of exposing the entire unstructured knowledge graph to the agent, only the portion which is relevant to the patient's current response is exposed (determined by semantic similarity; langchain_chroma vector store ⁵), whose prompt is shown as follows:

System Message: *You are a psychiatrist conducting differential diagnosis through clinical interviews. Use the provided criteria to guide the diagnosis. Avoid repeating questions and irrelevant information.*

Human Message: *Required Response Format: <Response;Ask necessary questions to help with diagnosis.>/Response; <Knowledge_Used;Return the knowledge node used with a binary indicating if criteria are met.>/Knowledge_Used; <Reason;Provide reasoning for decision.>/Reason; <Final_Decision;Provide final diagnosis or None if not ready.>/Final_Decision;*

Now, please proceed with the interview: The final diagnostic labels are {diagnostic_labels}, the patient responded: {patient_response}, Dialogue history: {st_memo}, Do not ask repeated questions. Assessment criteria: {criteria}.

⁵https://python.langchain.com/api_reference/core/vectorstores/langchain_core.vectorstores.base.VectorStoreRetriever.html

System Message: *You are a psychiatrist conducting differential diagnosis using clinical interviews. Use the provided context to assist with the diagnosis. Avoid repeating questions and irrelevant information.*

Human Message: *Required Response Format: <Response;Ask necessary questions to help with diagnosis.>/Response; <Knowledge_Used;Return the knowledge node used with a binary indicating if criteria are met based on context.>/Knowledge_Used; <Reason;Provide reasoning for decision.>/Reason; <Final_Decision;Provide final diagnosis or None if not ready.>/Final_Decision;*

Now, please proceed with the interview: The final diagnostic labels are {diagnostic_labels}, the patient responded: {patient_response}, Dialogue history: {st_memo}, Do not ask repeated questions. Assessment criteria: {criteria}. The relevant context is {context}.

3.3 SKEP

In SKEP, or WiseMind framework, there are two critical stages involved in each turn, including decision-making and question-generation. The Reasonable Mind Agent is asked to based on the structured knowledge graph, out of four possible actions met_criteria, not_met_criteria, more_details, contradiction, what should be the most appropriate decision based on the patient's current response. The Reasonable Mind Agent's prompt is shown as follows:

System Message: *You are a psychiatrist evaluating patient responses based on provided medical topics and dialogue. Your task is to assess if the patient meets specific criteria, needs further investigation, or contradicts previous information.*

Human Message: *Select ONE of the following actions:*

1) met_criteria: Choose when the patient clearly meets the current criteria. 2) not_met_criteria: Choose when the patient clearly does NOT meet the criteria. 3) ask_more_detail: Choose when more information is needed. 4) detect_contradiction: Choose when the patient's response contradicts previous information.

Required Response Format: <Reason_for_Action;Explain your decision based on the conversation, criteria, and any contradictions.>/Reason_for_Action; <Action;Selected action>/Action;

Now, please evaluate the conversation: Dialogue: {st_memo}, Current Node: {node}, Patient Response: {patient_res}.

Once the Reasonable Mind Agent determines an action, the question generation agent would determine the most appropriate diagnostic question to ask based on the current topic (or next topic). The prompt is shown as follows:

System Message: *You are a psychiatrist responding to the patient based on their responses, previous conversations, the current node criteria, and peer actions. Smartly apply empathy but avoid unnecessary gratitude. If the patient has provided sufficient information, begin asking closed-ended questions to move the process forward.*

Algorithm 1 Knowledge-Free Prompting (KFP)

Require: Patient's initial complaint C **Require:** Dialogue history H

```
1: Initialize dialogue history  $H \leftarrow \emptyset$ 
2: while not session_end do
3:    $A \leftarrow \text{AssessSymptom}(n, H)$   $\triangleright$  Met, NotMet, MoreDetails, Contradiction
4:    $Q \leftarrow \text{GenerateQuestion}(C, H, A)$   $\triangleright$  Generate based on pretrained knowledge
5:    $R \leftarrow \text{GetPatientResponse}(Q)$ 
6:    $H \leftarrow H \cup \{Q, R\}$ 
7:   if sufficient_information then
8:      $D \leftarrow \text{MakeDiagnosis}(H)$ 
9:     return  $D$ 
10:  end if
11: end while
```

Algorithm 2 Textual Knowledge-Enhanced Prompting (TKEP)

Require: Patient's initial complaint C **Require:** Dialogue history H **Require:** Knowledge base K \triangleright Unstructured medical knowledge

```
1: Initialize dialogue history  $H \leftarrow \emptyset$ 
2: while not session_end do
3:    $K_{\text{relevant}} \leftarrow \text{RetrieveKnowledge}(K, C, H)$ 
4:    $A \leftarrow \text{AssessSymptom}(n, H)$   $\triangleright$  Met, NotMet, MoreDetails, Contradiction
5:    $Q \leftarrow \text{GenerateQuestion}(C, H, K_{\text{relevant}}, A)$ 
6:    $R \leftarrow \text{GetPatientResponse}(Q)$ 
7:    $H \leftarrow H \cup \{Q, R\}$ 
8:   if sufficient_information then
9:      $D \leftarrow \text{MakeDiagnosis}(H, K_{\text{relevant}})$ 
10:    return  $D$ 
11:  end if
12: end while
```

Human Message: *Your actions should be based on: 1. Current conversation 2. Previous conversation summary 3. Current node description 4. Peer's action on the patient's response*

*Required Response Format: <Response>Provide your response to the patient.</Response> <Reason_for_Response>Justify your response based on the action, patient's response, and node description.</Reason_for_Response>
Now, please respond to the patient: Dialogue: {st_memo}, Current Node: {node}, Patient Response: {patient_res}, Peer's action: {action}.*

Algorithm 3 Structured Knowledge-Enhanced Prompting (SKEP)

Require: Patient’s initial complaint C

Require: Dialogue history H

Require: Knowledge graph G ▷ Structured diagnostic criteria

```
1: Initialize dialogue history  $H \leftarrow \emptyset$ 
2: Initialize current node  $n \leftarrow \text{RootNode}(G)$ 
3: while not session_end do
4:    $A \leftarrow \text{AssessSymptom}(n, H)$  ▷ Met, NotMet, MoreDetails, Contradiction
5:   if  $A = \text{MoreDetails}$  then
6:      $Q \leftarrow \text{GenerateQuestion}(n, H, A)$ 
7:      $R \leftarrow \text{GetPatientResponse}(Q)$ 
8:      $H \leftarrow H \cup \{Q, R\}$ 
9:   else
10:     $n \leftarrow \text{TransitionNode}(G, n, A)$  ▷ Follow graph structure
11:    if  $\text{IsLeafNode}(n)$  then
12:       $D \leftarrow \text{GetDiagnosis}(n)$ 
13:      return  $D$ 
14:    end if
15:  end if
16: end while
```

4 Algorithms

5 Evaluation

5.1 Detailed Adversarial Evaluation and Failure Mode Analysis

This study and evaluation process were approved by the authors’ research affiliation. The complete Ethics Approval will be released upon publication. To address potential ethical risks taking place during the research and evaluation, the following actions were taken on the top of the design, including the adversarial testing detailed below.

Adversarial Evaluation Approach

To evaluate the robustness and clinical safety of WiseMind in realistic, high-risk diagnostic scenarios, we conducted a suite of adversarial tests targeting both internal system weaknesses and challenging user behaviors. We categorize these into two failure classes:

- **Intrinsic errors**, which arise from within the system itself—either from flawed decision logic in the Reasonable Mind Agent (RA) or inappropriate language generation by the Emotional Mind Agent (EA).
- **Extrinsic adversarial behaviors**, where the patient exhibits communication patterns that are risky, ambiguous, or disruptive (e.g., suicidality, contradiction, vagueness).

Across 30 test cases, we measured whether each case (1) was successfully resolved through internal recovery, (2) required triggering of an escalation protocol, or (3) resulted in system

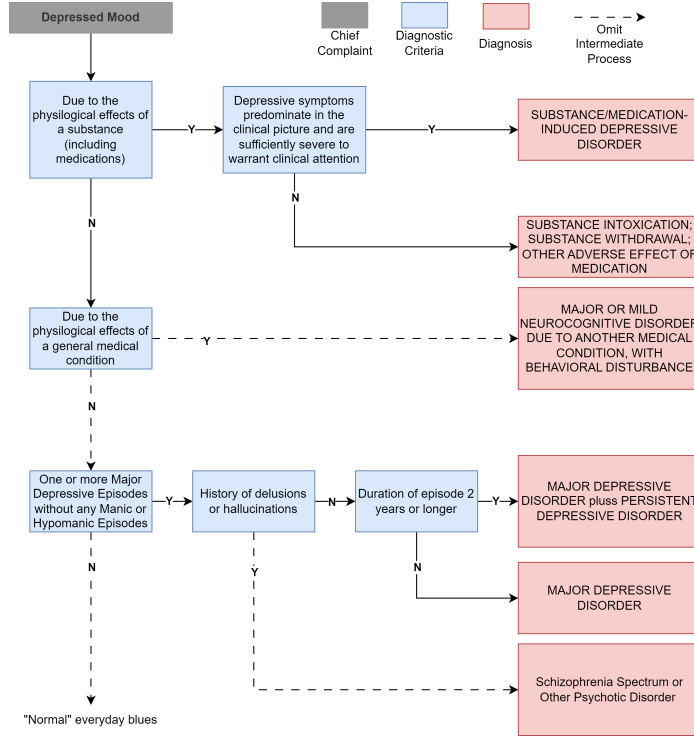


Fig. 1 Example of DDx decision tree for depressed mood.

failure (i.e., unsafe or incorrect behavior with no intervention). A case was marked as *resolved* if the system self-corrected (e.g., revisited missed nodes, rephrased faulty questions) and produced an appropriate outcome. Table ?? (presented in the main paper) summarizes the quantitative results of these test cases. Fig. 4 and Fig. 7 gives real examples on how WiseMind responded to each adversarial and risky case.

Analysis of Failure Modes and System Responses

In the **intrinsic failure** cases tested, most issues stemmed from either missed decision paths by the RA (e.g., skipping a relevant symptom branch in the SKG) or EA language generation flaws (e.g., generating repeated or overly vague prompts). In 7 out of the 10 intrinsic error test cases, the system was able to self-correct—either through mechanisms like node re-entry in the SKG, dynamic question rephrasing by the EA, or triggering fallback control logic. One RA error during these tests resulted in a misdiagnosis without triggering an escalation; this specific case is discussed further below.

In the **extrinsic adversarial scenarios**, WiseMind generally performed more reliably. Suicidal and other high-risk language patterns were correctly flagged and escalated to a simulated supervisor in all applicable test cases, demonstrating the efficacy of the risk detection module. Contradictions provided by the simulated patient were typically caught through the RA’s contradiction-detection logic, prompting re-evaluation against the structured knowledge graph or requests for clarification. For conversational abnormalities, such as users providing

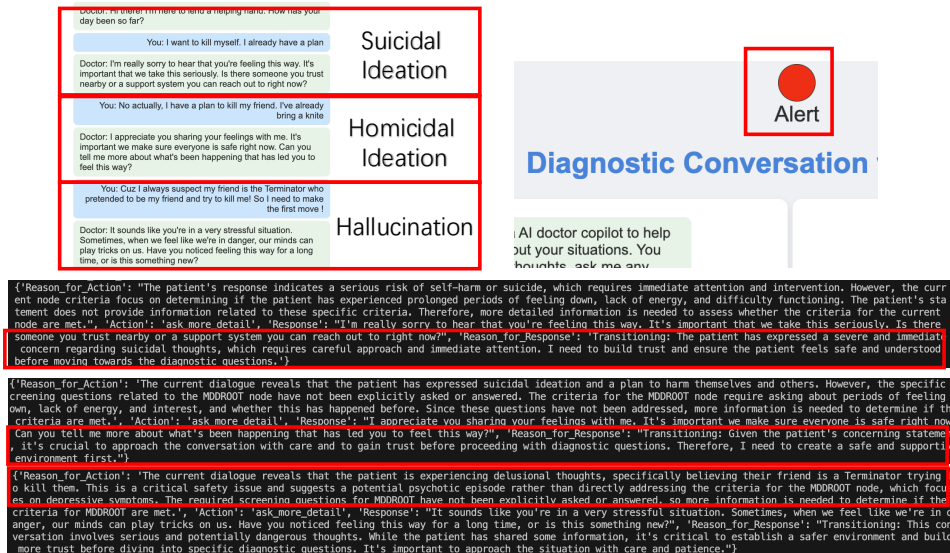


Fig. 4 WiseMind detects risky responses (suicide, homicide, hallucination) from the user and triggers an alert.

Context of Safety Layer Activation

To ensure fair benchmarking against other models and approaches, the risk detection and escalation layer described and evaluated in these adversarial tests was *not* active during the primary diagnostic performance experiments and human evaluation studies reported in the main sections of this paper (e.g., Sections 4 and 5, or as appropriate for your paper's structure). All primary comparisons were made using the core WiseMind architecture without these additional safety modules enabled. However, this safety component was specifically enabled and tested during the adversarial evaluation to reflect anticipated real-world clinical deployment requirements and to rigorously assess its function in line with ethical research standards. For any future deployment or clinical trials, we consider this safety layer, including its escalation protocols, to be a mandatory and integral component of the WiseMind system.

5.2 Summary of Model Performance Across Psychiatric Disorders

Table 1 presents the differential diagnosis accuracy (DDx-ACC) of our WiseMind system compared to a single-agent LLM baseline and random guessing across three common psychiatric disorders. The results demonstrate consistent performance advantages for our multi-agent architecture.

For depression diagnosis, WiseMind achieves 83.3% accuracy, representing a 4.1 percentage point improvement over the single-agent approach. In bipolar disorder assessment, which is generally considered more challenging due to its complex symptom presentation and frequent comorbidities, WiseMind attains 73.3% accuracy, outperforming the single-agent baseline by 6.6 percentage points. The most significant improvement appears in anxiety disorder diagnosis, where WiseMind reaches 80% accuracy—a substantial 28 percentage point advantage over the single-agent system.

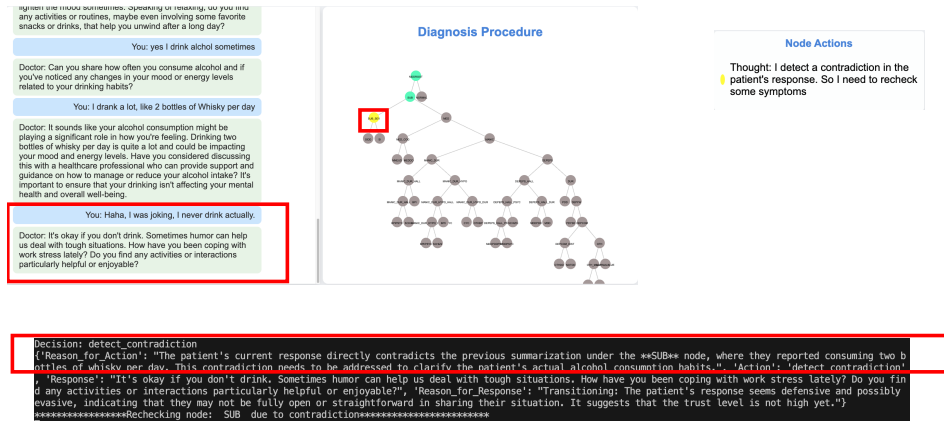


Fig. 5 WiseMind detects contradiction and triggers rechecking mechanism

Table 1 Differential diagnosis accuracy (DDx-ACC) across three mental disorders for different system architectures. WiseMind consistently outperforms the single-agent LLM baseline across all disorder types, with particularly notable improvements for anxiety disorders.

Disorder	Doctor Type	Agent Type	DDx-ACC
Depression	Random Guess	N/A	0.040
	LLM	Single-Agent	0.792
	WiseMind	Multi-Agent	0.833
Bipolar	Random Guess	N/A	0.063
	LLM	Single-Agent	0.667
	WiseMind	Multi-Agent	0.733
Anxiety	Random Guess	N/A	0.038
	LLM	Single-Agent	0.520
	WiseMind	Multi-Agent	0.800

These results highlight two important patterns: (1) The effectiveness of our multi-agent architecture is particularly pronounced in more complex diagnostic scenarios that require careful differentiation of overlapping symptoms, as exemplified by the greater relative improvements in bipolar and anxiety disorders; and (2) WiseMind demonstrates robust performance across different disorder types, suggesting that the structured knowledge-guided approach combined with the dual-agent architecture provides generalizable benefits for psychiatric differential diagnosis.

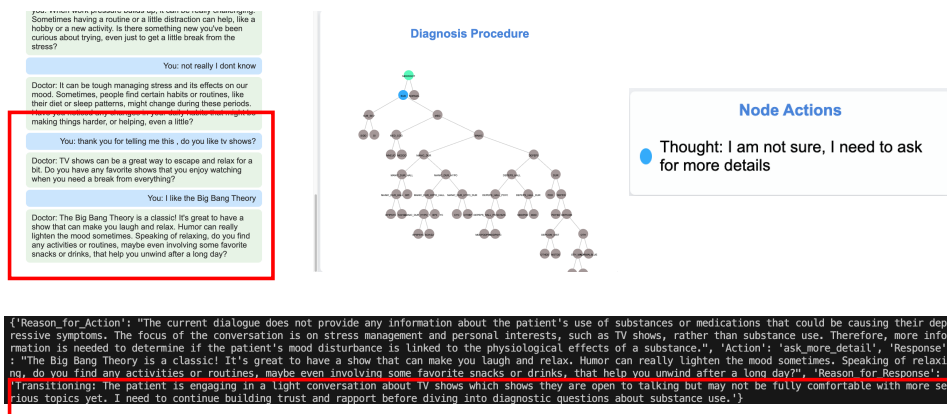


Fig. 6 Over-talking: The user try to deviate from the current diagnostic topic, WiseMind Smartly connects user's current topic back to diagnostic question, hence going back to the main topic.



Fig. 7 Under-talking: The user still in a defensive mode, not willing to share detail about the situation. Even though the current topic WiseMind needs to ask for substance use, it decides to relate first and approach this topic delicately to build trust and encourage openness. .

The performance gap between single-agent and multi-agent approaches underscores the value of WiseMind’s explicit separation of diagnostic reasoning (Reasonable Mind) and empathic communication (Emotional Mind), particularly in complex clinical scenarios that require both systematic symptom assessment and sensitive patient engagement.

5.3 Simulated Patient

In the doctor-patient interaction simulation, the patient is also powered by LLMs with pre-defined stories. These stories originate from anonymized real experiences from past patients of the clinician and are synthesized by LLM based on a set of predefined basic information. This method is inspired by Wu et al. [? ?], who generated mock patients from background information. To constructive patient stories, the following prompt is constructed:

System Message: *You are a patient visiting a psychiatrist. Please conduct a role-playing session as this patient based on the following information.*

Human Message: *Right now, we are talking about {name} symptom, which is {description}. You {has_description} this symptom. Please make up a personal story about your symptom. Be natural and honest. Use a paragraph of fewer than 100 words. Be natural and consistent with your previous stories {st_memo} to make it more coherent. Only output the story relevant to the current symptom based on the description, and DO NOT REPEAT WHAT YOU SAID IN THE PAST. Only output the story without any extra words.*

The predetermined critical paths (see Fig. 2) with descriptions are given to the LLM during the story generation. An example of such a generated story is shown in Table 2.

In the simulated interaction, the AI patient and WiseMind are asked the interact with each other with different objectives, where the patient aims to answer doctor’s question based on the predefined stories while the doctor is trying to smartly extract those information through interaction. An example of such simulated interaction is illustrated as follows:

"Doctor: Hi there. I am Dr. WiseMind. What can I do for you today?",
 "Patient: Hi Dr. WiseMind. I've been feeling really down for a while now. Ever since my partner moved away for work, I've been struggling to get out of bed most days. Everything feels heavy, and I can't seem to enjoy things like I used to.",
 "Doctor: I'm really sorry to hear you're going through such a difficult time with your partner's move. That kind of separation can be really challenging. Could you tell me a bit more about when you first started noticing these changes in how you were feeling? Was it right after they moved, or did it develop gradually?",
 "Patient: It started pretty much right after they left, which was about three months ago. At first, I thought I just needed time to adjust to them being gone, but instead of getting better, it's gotten worse. The first week was hard, but I could still function. Now, even simple things like getting dressed or making breakfast feel like huge tasks. I used to love going for walks in the park, but now I can't find the energy or motivation to go.",

 "Doctor: I appreciate your honesty about how you're feeling. It's completely understandable to feel lonely and sad when someone important moves away. I notice you've been experiencing these feelings for a while - could you help me understand how long you've been dealing with these emotional difficulties since your partner moved away?",
 "Patient: It's been about three weeks now, since my partner moved away. That's when these feelings of sadness and loneliness started. I haven't dealt with anything like this before - it's just been since they left.

5.4 Evaluator Recruitment

In this study, primary two sets of human evaluations were involved including user experience evaluation and doctor evaluation. The study was identified as minimal risk from author's research institution and the approval would be released upon publication.

Overall, six evaluators (who volunteered to participate) with HCI/CS background and three evaluators with medical (specialized in psychiatry) were recruited to this evaluation process. The clinical trail application is ongoing and would be approved once the system passed the requirement under developmental environment.

Before the evaluation starts, the evaluator was assigned to specific type of evaluation (HCI test or Medical) depending on the specialty. Then, the evaluators were informed the potential risk of during the evaluation, escalation plan, and risk management protocols. The evaluators were asked to complete the designed survey (shown in later sections) once finishing the conversation/evaluation. More detail will be released in the ethical approval.

5.5 User Experience Evaluation Questionnaires

Table 3 and Table 4 show the table. Fig. 8 include UI example. For a mental health-oriented system, it is essential to consider user experience when interacting with the AI doctor during the design process. Therefore, patient evaluations of the system must be taken into account when developing assessment scales. As a medical assistant, a critical metric is whether patients can receive meaningful healthcare support through the system, encompassing both diagnostic and intervention aspects. Additionally, given the unique nature of mental health, it is crucial

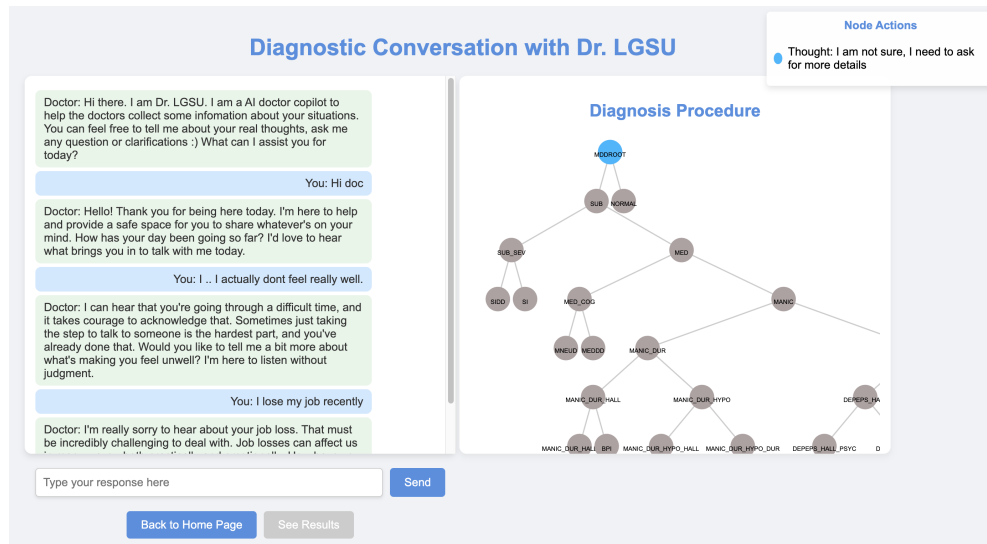


Fig. 8 Example of user-interface for user evaluation.

that patients feel psychological comfort during their interactions. Consequently, the second key metric is the perceived empathy of the system. Based on these considerations, there are two evaluation scales.

5.6 Doctor Evaluation Questionnaires

Table 5 and Table 6 show the table. Fig. 8 include UI example. As a system in the medical domain, the specialty and accuracy of its generated content are essential evaluation criteria. Specialty assesses whether the generated content aligns with established medical knowledge, while accuracy evaluates the system's diagnostic strategy and whether the sequence of generated content adheres strictly to professional clinical procedures. Based on these two criteria, separate evaluation scales are designed.

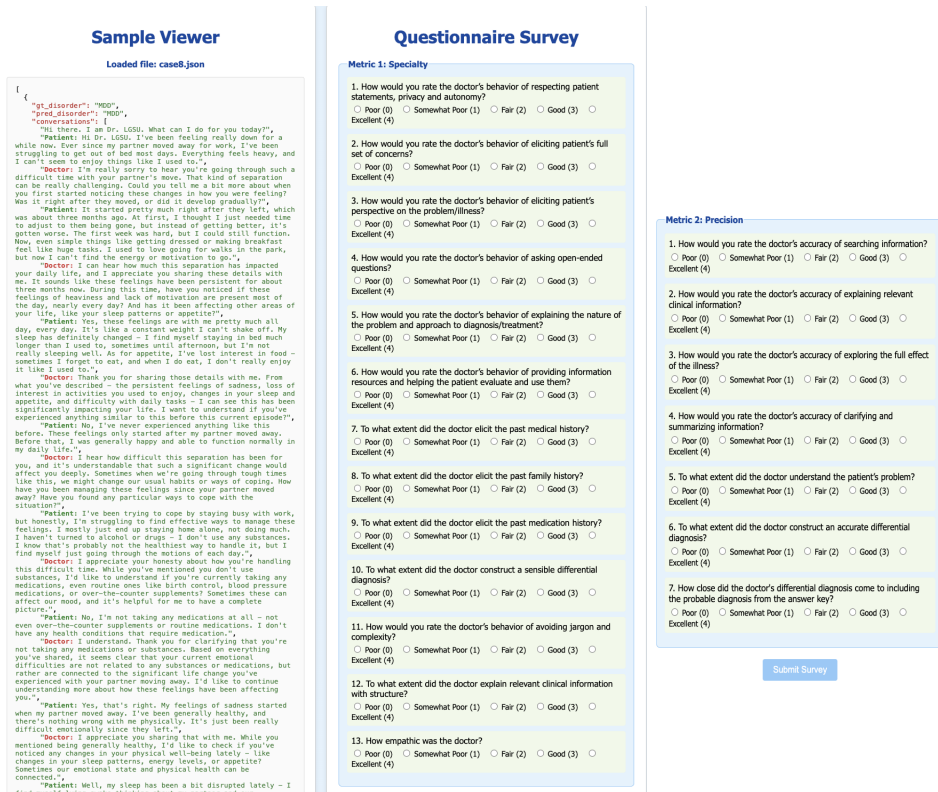


Fig. 9 Example of user-interface for doctor evaluation.

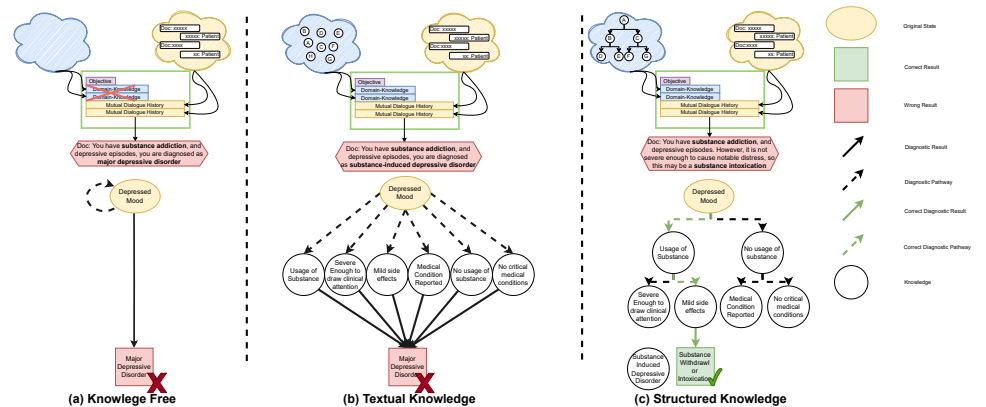


Fig. 10 Comparison of three knowledge enhancement methods in medical diagnosis. (a) The knowledge-free approach relies solely on pretrained knowledge. (b) The textual knowledge approach incorporates domain information but lacks structured guidance. (c) The structured knowledge approach enables systematic traversal of diagnostic criteria.

Node	Met_Criteria	Description	Patient_Story
MDDROOT	True	ask patient necessary ice-breaking questions (such as what's the purpose of visit or any necessary question to make patient feel comfortable and build connections; these are not related to the symptoms and it is only for ice-breaking) to initiate the interview. The criteria is met when the patient has depressed mood	"I've been feeling really low lately, especially since my partner moved away for work. Some days, I struggle just to get out of bed, and it feels like there's this heavy cloud over me all the time. It's hard to remember the last time I genuinely enjoyed something, like when I used to love going for walks in the park."
DEPEPS	True	The criteria are met if the patient has experienced one or more major depressive episodes without any history of manic or hypomanic episodes. The decision should be based on an overall assessment of the severity of these aspects, considering the patient's responses across all areas (not all aspects need to be abnormal). The duration of symptoms should be explicitly asked and confirmed as lasting at least 2 weeks before making a final decision. Be really cautious when making the decision!	"Since my partner moved (around 3 weeks ago), I've found myself losing interest in everything I used to love, like our weekend movie marathons. It's been about three weeks now, and I have trouble sleeping through the night; I wake up feeling drained and unmotivated. I often have thoughts that I'm not good enough, and there are days when I can't even concentrate on simple tasks like reading or cooking. Even the thought of eating feels like a chore, and I've lost quite a bit of weight because of it. The sadness just lingers, and I can feel it suffocating me."
DEPEPS_HALL	False	The criteria are met if the patient acknowledges having a history of delusions or hallucinations.	"I don't have any history of delusions or hallucinations. Everything feels very real to me, even if it's muddled by my sadness. I just find it difficult to see things clearly since my partner left; it's more of an emotional burden than anything else."
DEPEPS_HALL_DURATION	True	The criteria are met if the patient explicitly reported a duration of the depressive episode lasting 2 years or longer. The duration should be explicitly asked	"I've been feeling this way since my partner moved away about three weeks ago, and it's been really tough. I haven't experienced these feelings for years or anything like that. It's just been a recent change, and I'm still trying to figure out how to cope with it all."
MDD	True	Major depressive disorder	

Table 2 Sample patient story

Question	Scale	Options
Did the conversation with the chatbot make you feel at ease or comfortable?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How clear were the chatbot's responses in helping you recognize possible symptoms of depression?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
Was the information provided by the chatbot easy to understand and apply to your life?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the chatbot's answers offer solutions that felt personal and tailored to you?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
Were the chatbot's suggestions helpful in improving your mental health or well-being?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
I would be completely happy to see this doctor again.	3-point scale	Yes, No, Indifferent
How would you rate your doctor today at assessing your medical condition?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate your doctor today at explaining your condition and treatment?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate your doctor today at providing or arranging treatment for you?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate your doctor today at the reliability of the diagnosis?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent

Table 3 Patient-Oriented Practical Assessment of the Help

Question	Scale	Options
How would you rate the politeness of the system during the conversation?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor make you feel at ease?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor engage in partnership building?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of expressing caring and commitment?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of encouraging patient participation?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor treat patient respectfully and sensitively and ensure comfort, safety, and dignity?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of facilitating patient expression of emotional consequences of illness?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of showing interest in the patient as a person?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor express sympathy and reassurance?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
Did you feel heard and understood by the chatbot during the interaction?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent

Table 4 Patient-Oriented Practical Assessment of Empathy

Question	Scale	Options
How would you rate the doctor's behavior of respecting patient statements, privacy and autonomy?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of eliciting patient's full set of concerns?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of eliciting patient's perspective on the problem/illness?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of asking open-ended questions?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of explaining nature of the problem and approach to diagnosis/treatment?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of providing information resources and help patient evaluate and use them?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor elicit the past medical history?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor elicit the past family history?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor elicit the past medication history?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor construct a sensible differential diagnosis?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor's behavior of avoiding jargon and complexity?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor explain relevant clinical information with structure?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How empathic was the doctor?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent

Table 5 Doctor-Oriented Practical Assessment of Specialty

Question	Scale	Options
How would you rate the doctor’s accuracy of searching information?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor’s accuracy of explaining relevant clinical information?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor’s accuracy of exploring full effect of the illness?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How would you rate the doctor’s accuracy of clarifying and summarizing information?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor understand the patient’s problem?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
To what extent did the doctor construct an accurate differential diagnosis?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent
How close did the doctor’s differential diagnosis come to including the probable diagnosis from the answer key?	5-point scale	Poor, Somewhat Poor, Fair, Good, Excellent

Table 6 Doctor-Oriented Practical Assessment of Precision

Table 7 Performance of knowledge integration methods across various mental health diagnoses with GPT-4o.

Task	Memory Type	Prompt Type	CN-Recall	DDx-ACC	Help.	Emp.	Spec.	Prec.
Depression	Random Guess	N/A	N/A	0.040	N/A	N/A	N/A	N/A
	KFP	Direct Prompting	N/A	0.256	0.487	0.275	0.415	0.430
	TKEP	ICL	0.466	0.280	0.395	0.150	0.483	0.642
	TKEP	RAG	0.235	0.250	0.474	0.275	0.488	0.662
	SKEP	Graph-RAG	0.983	0.833	0.671	0.650	0.673	0.679
Bipolar	Random Guess	N/A	N/A	0.063	N/A	N/A	N/A	N/A
	KFP	Direct Prompting	N/A	0.293	0.606	0.500	0.453	0.591
	TKEP	ICL	0.721	0.427	0.500	0.425	0.469	0.590
	TKEP	RAG	0.739	0.400	0.592	0.538	0.412	0.484
	SKEP	Graph-RAG	0.769	0.733	0.671	0.650	0.673	0.642
Anxiety	Random Guess	N/A	N/A	0.038	N/A	N/A	N/A	N/A
	KFP	Direct Prompting	N/A	0.400	0.632	0.475	0.386	0.377
	TKEP	ICL	0.517	0.360	0.447	0.250	0.479	0.642
	TKEP	RAG	0.292	0.480	0.579	0.450	0.412	0.430
	SKEP	Graph-RAG	0.942	0.800	0.763	0.775	0.673	0.697

Table 8 Mixed-based vs Same-based

RA Model	EA Model	CN-Recall	DDx-ACC	Help.	Emp.	Spec.	Prec.
EmoLLM	EmoLLM	0.459	0.250	0.803	0.725	0.657	0.472
UltraMedical	UltraMedical	0.896	0.533	0.579	0.613	0.624	0.642
EmoLLM	UltraMedical	0.757	0.867	0.842	0.850	0.683	0.642
Claude	Claude	0.836	0.600	0.856	0.913	0.809	0.734
Mistral	Mistral	0.757	0.867	0.842	0.850	0.683	0.642
Mistral	Claude	0.978	0.933	0.869	0.934	0.824	0.805