

# From Simulation to Diagnosis: Simulation-Based Reinforcement Learning for Empathetic Depression Diagnosis

Anonymous ACL submission

## Abstract

Given the rising prevalence of depression alongside persistently low diagnosis rates, Large Language Models (LLMs) offer an opportunity for accessible mental health assessment. However, clinical depression diagnosis is inherently a goal-oriented, sequential, and interactive process. Current approaches primarily rely on static supervised finetuning (SFT) using fixed depression-diagnosis dialogue datasets, which are ill-suited to the dynamic nature of real-world diagnosis interactions. Models trained on such static data often struggle to navigate the variability of patient behavior, failing to provide diagnostic accuracy. To overcome these limitations, we introduce **SimRED** (Simulation-based Reinforcement Learning for Empathetic Depression Diagnosis), a framework that trains depression diagnostic agents via reinforcement learning through extensive interactions with patient simulator. SimRED constructs a patient simulation environment derived from real-world dialogues. By interacting with these simulated patients, the diagnostic agent employs reinforcement learning to learn both diagnostic accuracy and empathetic expression. Experimental results demonstrate that SimRED significantly outperforms existing strong baselines in both diagnostic accuracy and the quality of empathetic expression.

## 1 Introduction

Depression affects approximately 264 million individuals globally, posing a severe public health challenge (James et al., 2018; Ferrari et al., 2013). Despite its prevalence, a substantial portion of patients remain undiagnosed due to critical barriers, including the global shortage of qualified mental health professionals (Organization, 2022; Henderson et al., 2013). To bridge this service gap, automated depression diagnostic agents powered by LLMs have emerged, offering a scalable solution for accessible mental health service.

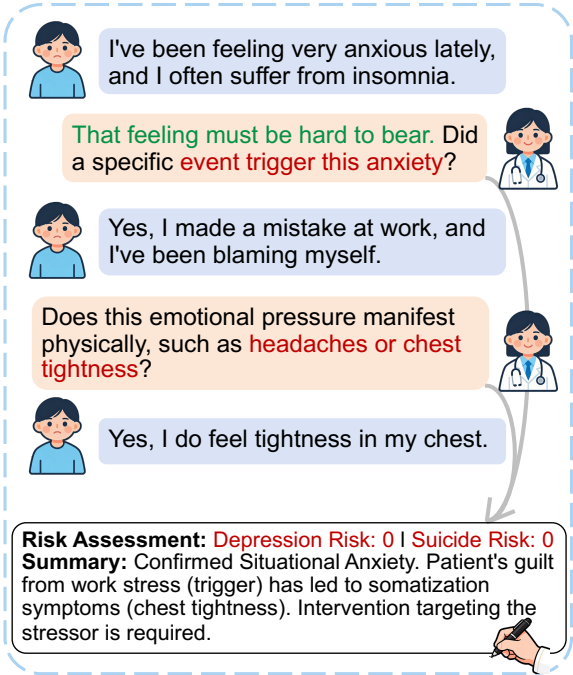


Figure 1: Illustration of the depression-diagnosis dialogues. Unlike casual conversation, each turn in a depression-diagnosis dialogue serves an information-gathering purpose toward the final assessment. The agent dynamically balances empathetic expression with diagnostic inquiries. This demonstrates the necessity of moving beyond static imitation to a system that learns to achieve high diagnostic accuracy through iterative interactions with the patient simulator.

Prior research makes significant contributions to this field. The  $D_4$  dataset pioneers the concept of depression-diagnosis dialogue (Yao et al., 2022), while the SEO framework integrates diagnostic inquiries with empathetic chat based on helping skills (Lan et al., 2025). Additionally, recent studies develop strategies for addressing stigma-associated depression (Cao et al., 2025) and multi-agent systems that leverage diagnostic history (Lan et al., 2024). **Despite these advancements, current approaches face a fundamental limitation:**

054 **the gap between static training and dynamic**  
055 **practice.** As shown in Figure 1, a successful  
056 depression-diagnosis dialogue is a goal-oriented  
057 sequential decision-making process. Consequently,  
058 previous methods lack the capability to navigate  
059 such dynamic interactions effectively.

060 To address this, we transition the diagnostic  
061 agent from passive SFT on static datasets to inter-  
062 active training with patient simulators. Specifically,  
063 we develop a patient simulation environment and  
064 formulate the diagnostic dialogue as a reinforce-  
065 ment learning problem. By engaging in repeated  
066 interactions with the simulator, the agent learns to  
067 optimize its diagnostic accuracy and empathetic  
068 expression. Our approach aligns dialogue trajec-  
069 tories with diagnostic goals, enabling the model to  
070 effectively master the intricate dynamics of clinical  
071 depression diagnosis.

072 Overall, we propose SimRED (**Simulation-based**  
073 **Reinforcement Learning for Empathetic Depres-**  
074 **sion Diagnosis**), a simulation-based reinforcement  
075 learning framework consisting of a patient simu-  
076 lator agent and a diagnostic agent. To enable re-  
077 liable interactive training, we construct a patient  
078 simulation environment from static diagnostic dia-  
079 logues by reconstructing structured patient profiles  
080 and ensuring coherent, non-misleading behaviors  
081 suitable for multi-turn interaction. Within this en-  
082 vironment, the diagnostic agent is trained using re-  
083 inforcement learning with a composite reward that  
084 jointly encourages diagnostic accuracy and empa-  
085 thetic expression. Through extensive interactions,  
086 the agent learns to navigate the complex dynam-  
087 ics of depression diagnosis beyond static imitation.  
088 Experimental results across settings demonstrate  
089 that diagnostic agents trained with SimRED consis-  
090 tently outperform SFT and prompt-based methods  
091 in terms of accuracy and quality.

- 092 • We propose SimRED, a simulation-based re-  
093 inforcement learning framework that bridges  
094 static SFT and interactive diagnosis.
- 095 • We construct a profile-grounded patient simu-  
096 lator from real-world diagnostic dialogues  
097 and further align it with expert feedback to  
098 support reliable RL training.
- 099 • Extensive experiments show consistent im-  
100 provements over many baselines, and evalua-  
101 tions across different simulators validate the  
102 robustness of SimRED.

## 2 Related Work 103

104 In this section, we review LLM research in men-  
105 tal health, categorized into two primary domains:  
106 intervention and diagnosis.

### 2.1 LLMs for Mental Health Intervention 107

108 LLMs in the mental health domain offer a promis-  
109 ing solution to the severe shortage of practition-  
110 ers. The majority of existing research focuses on  
111 mental health counseling, intervention, and empa-  
112 thetic companionship. For instance, the SMILE  
113 dataset (Qiu et al., 2024), an improvement over  
114 PsyQA (Sun et al., 2021), includes multiple rounds  
115 of counseling dialogues and is utilized by CBT-  
116 LLM (Qiu et al., 2024). CPsyCoun extracts coun-  
117 seling data from public psychology reports without  
118 requiring specific domain expertise (Zhang et al.,  
119 2024). Healme (Xiao et al., 2024) focuses on opti-  
120 mizing CBT guidance through LLMs, while CAC-  
121 TUS (Lee et al., 2024) emphasizes the construction  
122 of more comprehensive single-turn CBT processes.

123 Distinct from direct intervention, assisting in  
124 mental health assessment is also of paramount im-  
125 portance. Although the objectives differ, models  
126 employed in the assessment process must still ex-  
127 hibit sufficient warmth and empathy to avoid caus-  
128 ing harm.

### 2.2 LLMs for Mental Health Diagnosis 129

130 In the context of diagnosis, LLMs are tasked with  
131 ascertaining an individual’s mental health status  
132 through interactive dialogue. The  $D_4$  dataset pio-  
133 neered the concept of diagnostic conversation (Yao  
134 et al., 2022). Subsequent studies have introduced  
135 user state tracking into mental health dialogues (Gu  
136 et al., 2025). To enhance empathy, the SEO frame-  
137 work integrates diagnostic inquiry with empathetic  
138 chit-chat based on helping skills (Lan et al., 2025).  
139 Additionally, recent works have developed strate-  
140 gies for stigma-associated depression (Cao et al.,  
141 2025) and multi-agent systems that leverage diag-  
142 nostic history (Lan et al., 2024). Similar to interven-  
143 tion research, some studies focus on data synthesis  
144 and augmentation due to the scarcity of authen-  
145 tic diagnostic records; researchers have collabora-  
146 ted with clinicians to synthesize datasets such as  
147 MDD5k (Yin et al., 2025) and PsyCoTalk (Wan  
148 et al., 2025).

149 In contrast to these approaches, SimRED intro-  
150 duces a reinforcement learning paradigm into the  
151 realm of mental health diagnosis.

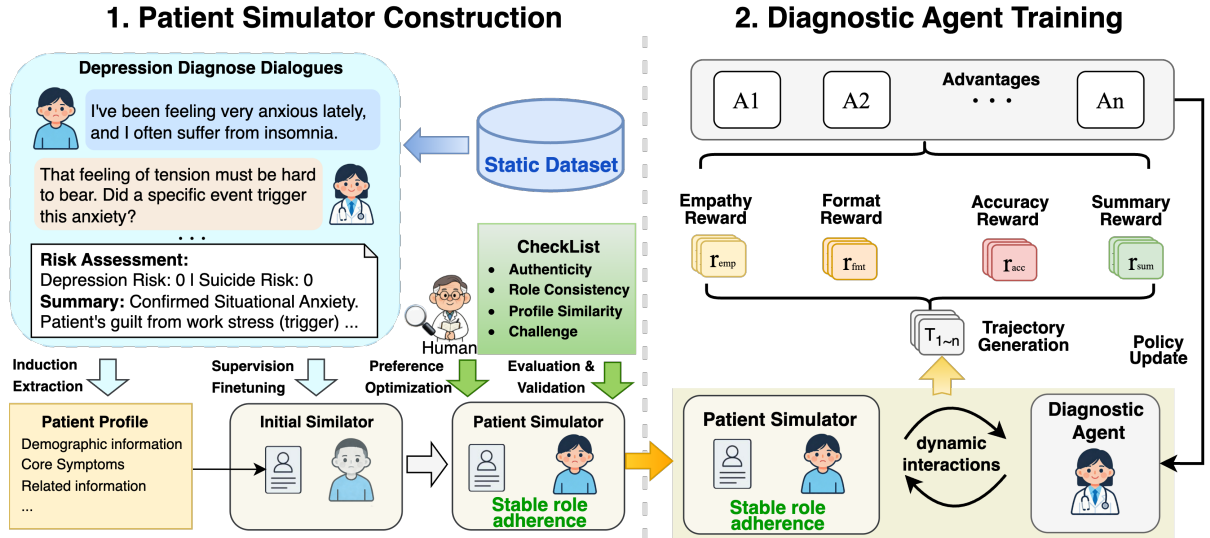


Figure 2: **Overview of the SimRED framework.** We construct an interactive patient simulator from static depression-diagnosis dialogues by extracting structured patient profiles and initializing a base simulator. This model is subsequently refined using DPO to ensure strict alignment with the assigned profiles. Human experts then evaluate the simulator across multiple metrics to verify its fidelity. Finally, the simulator serves as an environment to train the diagnostic agent via reinforcement learning, employing a reward function that jointly optimizes for diagnostic accuracy and empathetic communication.

### 3 SimRED: Simulation-based Reinforcement learning for Empathetic Depression Diagnosis

#### 3.1 Problem Formulation

We formalize the depression-diagnosis process as a goal-oriented, multi-turn interaction between a diagnostic agent and a patient simulator. An episode is defined as a sequence of interactions between the diagnostic agent  $M_d$  and the patient simulator  $M_p$ . At turn  $t$ , the agent  $M_d$  observes the dialogue history  $h_t = \{u_d^1, u_p^1, \dots, u_p^{t-1}\}$ , where  $u_d^i$  and  $u_p^i$  denote the utterances of  $M_d$  and  $M_p$  at turn  $i$ , respectively. Subsequently,  $M_d$  generates a response  $u_d^t$ . The episode terminates when the agent emits a special diagnostic token [DIAGNOSE], followed by a structured prediction  $\hat{y} = \{\hat{y}_{dep}, \hat{y}_{sui}, \hat{s}\}$ . In this formulation,  $\hat{y}_{dep}$  represents the predicted depression risk level,  $\hat{y}_{sui}$  denotes the predicted suicide risk level, and  $\hat{s}$  is the generated patient summary.

#### 3.2 Patient Simulator Construction

Reinforcement learning with real patients is neither practical nor ethically tenable. To enable training while preserving realism, we develop a patient simulator,  $M_p$ . Inspired by previous research on depression patient simulation (Liu et al., 2025), we build the initial simulator using static depression-diagnosis dialogues. Simultaneously, we employ

Direct Preference Optimization(DPO) to ensure that  $M_p$ 's behavior remains consistent with the defined patient profiles.

#### 3.2.1 Profile Construction

As illustrated in Figure 2, we design a latent patient profile  $\mathcal{P}$  from raw depression-diagnosis dialogues for each patient using a powerful LLM (GPT-4o). The profile maps the unstructured dialogue history into a seven-dimensional representation, covering core problems, affective symptoms, cognitive symptoms, and social function impairments, among others.

To maintain fidelity to the source dialogue, we implement a summary–evidence alignment protocol: each extracted trait is paired with verbatim supporting evidence extracted from the original dialogue during the distillation process. Further details are provided in Appendix B.1.

#### 3.2.2 Simulator Training and Alignment

We first perform supervised finetuning to initialize the  $M_p$ , utilizing the structured patient profiles as system prompts and the original depression-diagnosis dialogues as targets. This stage enables the patient simulator to get coherent behaviors.

However, the initially finetuned simulator occasionally exhibits behaviors inconsistent with the provided profiles and hallucinates facts, which pre-

vents patient simulation. To address this, we invite ten human experts to interact with the simulator and rewrite problematic responses, correcting behaviors that contradict the patient profiles. We then apply DPO to align the simulator  $M_p$  with the profiles. In a second-round evaluation, hallucinations and profile-inconsistent behaviors are significantly reduced.

Finally, we evaluate the patient simulator using a 7-point Likert scale across four dimensions adapted from Himmelbauer et al. (2018). Expert evaluations indicate that the performance of the patient simulator is satisfactory, demonstrating that our proposed  $M_p$  is reliable.

### 3.3 Diagnostic Agent Training

We train the  $M_d$  in two stages: (i) supervised finetuning to learn the required reasoning and output format, and (ii) reinforcement learning to optimize a composite reward that reflects both diagnostic accuracy and empathetic expression.

#### 3.3.1 Reward Design

To guide the  $M_d$  toward accurate and empathetic interactions, we design a composite reward function  $R$ . We distinguish between **Outcome-level Rewards** ( $R_{out}$ ), which evaluate the final diagnostic results, and **Process-level Rewards** ( $R_{proc}$ ), which monitor the quality of each interaction turn:

$$R = R_{out} + R_{proc} \quad (1)$$

**Outcome-level Rewards** The  $R_{out}$  is triggered only when  $M_d$  emits the [DIAGNOSE] token, evaluating the terminal state of the dialogue:

$$R_{out} = \alpha \cdot r_{acc} + \beta \cdot r_{sum} + r_{dfmt} \quad (2)$$

where  $\alpha$  and  $\beta$  are scaling factors.

- **Diagnostic Accuracy Reward** ( $r_{acc}$ ): This reward quantifies the precision of the clinical diagnosis. It is defined as the negative  $L_1$  distance between the predicted risk levels ( $\hat{y}_{dep}, \hat{y}_{sui}$ ) and the ground-truth labels ( $y_{dep}, y_{sui}$ ):

$$r_{acc} = -(|\hat{y}_{dep} - y_{dep}| + |\hat{y}_{sui} - y_{sui}|) \quad (3)$$

- **Diagnostic Summary Reward** ( $r_{sum}$ ): Upon concluding the diagnosis,  $M_d$  is required to generate a procedural summary of the entire dialogue. We utilize the ROUGE-L score between the generated summary  $\hat{s}$  and the reference summary  $s$  as a reward to evaluate the quality of the generated summary.

- **Diagnosis Format Penalty** ( $r_{dfmt}$ ): To ensure the outputs are parsable for downstream systems, a constant penalty is applied if the final response fails to strictly adhere to the prescribed structured format or if  $M_d$  fails to reach a diagnosis within the maximum number of turns.

**Process-level Rewards** To prevent  $M_d$  from pursuing accuracy at the cost of patient experience,  $R_{proc}$  provides signals during the dialogue:

$$R_{proc} = \sum_{t=1}^{[DIAGNOSE]} (\gamma \cdot r_{emp}^{(t)} + r_{pfmt}^{(t)}) \quad (4)$$

where the upper limit corresponds to the terminal step at which  $M_d$  emits [DIAGNOSE].

- **Empathetic Expression Reward** ( $r_{emp}^{(t)}$ ): At each turn  $t$ , we employ a finetuned empathy classifier to evaluate  $M_d$ 's utterance  $u_d^t$ .
- **Interaction Format Penalty** ( $r_{pfmt}^{(t)}$ ): To ensure logical consistency, we penalize instances where  $M_d$  skips the required CoT reasoning before generating its response.

Our reward structure encourages  $M_d$  to optimize for both diagnostic rigor and empathetic quality across the entire interaction trajectory.

#### 3.3.2 Optimization

We train  $M_d$  in two phases: supervised finetuning followed by reinforcement learning.

**Supervised Finetuning** We first train  $M_d$  to conduct explicit reasoning prior to generating  $u_d$ , as clinical depression diagnosis is inherently reasoning-intensive. Specifically, we utilize Qwen2.5-32B-Instruct to augment the static dataset by synthesizing missing rationales for each utterance  $u_d^i$ . To address the scarcity of real-world clinical data, we further expand the training set with 270 synthetic patient profiles, which are combined with 300 real-world cases for supervised finetuning. This stage primarily serves to establish the dialogue protocol and foundational output patterns.

**Reinforcement Learning** To refine  $M_d$  beyond static imitation, we employ Group Relative Policy Optimization (GRPO). GRPO enhances training stability and efficiency by estimating the baseline from group-sampled trajectories, eliminating the need for a separate value function.

For each patient profile  $\mathcal{P}$ ,  $M_d$  conducts  $G$  multi-turn interactions with the patient simulator  $M_p$ , generating a set of trajectories  $\{T_1, \dots, T_G\}$ . We optimize the policy  $\pi_\theta$  by maximizing:

$$\mathcal{J}(\theta) = \mathbb{E} \left[ \frac{1}{G} \sum_{i=1}^G \left( \text{clip}(\rho_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i - \beta \mathbb{D}_{KL}(\pi_\theta \| \pi_{ref}) \right) \right] \quad (5)$$

where  $\rho_i(\theta) = \frac{\pi_\theta(u|h)}{\pi_{old}(u|h)}$  and  $\pi_{ref}$  is the SFT model. The advantage  $\hat{A}_i$  for each trajectory is standardized within the group:

$$\hat{A}_i = \frac{R_i - \text{mean}(R_1, \dots, R_G)}{\text{std}(R_1, \dots, R_G) + \delta} \quad (6)$$

where  $R_i$  is the composite reward. Through this iterative process,  $M_d$  learns to prioritize diagnostic inquiries that maximize diagnostic accuracy while ensuring high-quality empathetic expression, effectively balancing clinical rigor with emotional resonance.

## 4 Experiments

### 4.1 Datasets

To simulate real-world patient data as accurately as possible, our experiments are conducted based on the  $D_4$  dataset (Yao et al., 2022), which is a benchmark specifically designed for depression-diagnosis dialogues. The dataset encompasses patients with diverse characteristics and provides labels for four levels of depression risk level (non-depressed, mild, moderate, and severe), along with four levels suicide risk level.

To support interactive training and evaluation, we utilize a patient simulator  $M_p$  developed by finetuning and aligning Qwen2.5-3B-Instruct on the  $D_4$  dataset. This environment is designed to empower  $M_d$  to achieve robust and reliable diagnostic accuracy.

### 4.2 Baselines

To rigorously evaluate our method, we compare it against several strong baselines, including:

- **SFT LLM:** LLMs finetuned directly on the collected full  $D_4$  training dataset with CoT.
- **Vanilla LLM:** LLMs prompted with elaborate role-playing instructions.

- **UPSD** (Cao et al., 2025): A framework designed to identify depression risks through unobtrusive probing strategies.

- **POST** (Gu et al., 2025): A method that explicitly tracks patient states throughout the conversation.

For a comprehensive comparison, we report results across various models and parameter scales. Implementation details are provided in Appendix A.

### 4.3 Evaluation Metrics

We evaluate model performance with a primary focus on diagnostic accuracy regarding both depression risk levels and suicide risk levels, conducted through interactions with the patient simulator. Furthermore, we introduce an **LLM-as-a-judge** to assess empathetic expression within the diagnostic dialogues. This metric quantifies the empathetic content of the conversation, with the GPT-4o assigning scores on a scale of 1 to 5; for further details, please refer to Appendix C. To evaluate the quality of the generated diagnostic summaries, we employ BLEU and ROUGE-L metrics to assess the model’s ability to produce dialogue summaries.

### 4.4 Patient Simulator Assessment

Using 100 patient profiles  $\mathcal{P}$ , we analyze our patient simulator, focusing on two primary error categories: profile inconsistent (deviating from prescribed traits) and factual fabrication (hallucinating details absent from the profile). Human experts conduct simulated dialogues and rewrite all erroneous responses. We then pair these expert-corrected responses with the original outputs to construct a preference dataset. After DPO, we perform behavioral testing and expert scoring to quantify the improvements in the simulator.

Category	Initial	Our
Correct	86.4%	95.4%
PI*	2.6%	0.7%
FF†	11.1%	3.9%
Total	100.0%	100.0%

\*PI: Profile Inconsistent; †FF: Factual Fabrication.

We evaluate the simulation across dimensions including authenticity, consistency, similarity, and challenge. The results demonstrate that this  $M_p$  is reliable enough to support training and evaluation. Further details are provided in Appendix B.2.

Model	Diagnosis		Patient Summary		Empathy
	D ACC	S ACC	BLEU	ROUGE-L	Empathy Score
Qwen2.5-3B (POST)	0.355	0.480	0.131	0.116	2.11
GPT-4o (UPSD)	0.288	0.528	0.117	0.110	<b>3.96</b>
GPT-4o (Vanilla)	0.433	0.622	0.154	0.121	1.98
LLaMA3-3B (Vanilla)	0.192	0.144	0.054	0.057	2.52
LLaMA3-3B (SFT)	0.231	0.317	0.146	0.134	2.06
LLaMA3-8B (Vanilla)	0.221	0.519	0.114	0.091	1.63
LLaMA3-8B (SFT)	0.413	0.605	0.170	0.143	2.06
LLaMA3-3B (SimRED)	0.443	0.711	0.182	0.160	3.07
Qwen2.5-3B (Vanilla)	0.250	0.356	0.139	0.107	1.92
Qwen2.5-3B (SFT)	0.461	0.663	0.183	0.153	2.31
Qwen2.5-7B (Vanilla)	0.221	0.557	0.130	0.104	2.27
Qwen2.5-7B (SFT)	0.538	0.635	0.178	0.147	3.21
<b>Qwen2.5-3B(SimRED)</b>	<b>0.586</b>	<b>0.673</b>	<b>0.195</b>	<b>0.166</b>	3.67

Table 1: Main experimental results on depression diagnosis. We compare our SimRED framework against various baselines, including vanilla LLMs, SFT LLMs, and prompt diagnostic methods. D ACC and S ACC represent the diagnostic accuracy for depression and suicide risk, respectively. Bold values indicate the best performance. SimRED consistently outperforms all baselines while maintaining a high level of empathy.

## 5 Main Results

Table 1 presents a comparison between SimRED and various baseline models. The empirical results lead to the following observations:

SimRED consistently outperforms both SFT and vanilla LLMs across different model scales. Specifically, Qwen2.5-3B (SimRED) achieves the highest diagnostic accuracy for both depression and suicide risk, significantly surpassing baselines. Notably, our 3B-scale model even exceeds the performance of much larger models and advanced models, such as the Qwen2.5-7B and GPT-4o, in terms of diagnostic accuracy. We prioritize the 3B-scale model for its feasibility in local deployment, which is critical for ensuring data privacy in sensitive mental health applications. This performance margin demonstrates that the diagnostic logic acquired through interactive reinforcement learning is substantially more robust than that derived from static supervised imitation or prompt-based methods.

A major limitation of static SFT models is their cold or robotic interaction style. SimRED successfully bridges the gap between clinical diagnosis and humanistic care. By incorporating the process reward  $R_{proc}$ , the empathy scores of all models improve substantially. Although some models, such as GPT-4o(UPSD), exhibit higher empathy, they suffer from significantly lower diagnostic accuracy.

In contrast, SimRED achieves a superior balance, maintaining professional yet warm politeness while ensuring high diagnostic accuracy.

The performance gains are consistent across both LLaMA3 and Qwen2.5 families. For instance, LLaMA3-3B (SimRED) outperforms LLaMA3-8B (SFT) in suicide risk detection. This indicates that our framework is architecture-agnostic and can effectively improve smaller models to reach or exceed the performance of larger, statically-trained models, offering a more efficient path for deploying specialized mental health agents. This property is particularly valuable for mental health applications, where deploying strong local small models may be more practical and meaningful given the sensitivity of user data. By mastering the underlying logic of diagnosis, SimRED proves that smaller, specialized models can achieve professional-grade performance.

By aligning the diagnostic agent  $M_d$  with both  $R_{out}$  and  $R_{proc}$ , SimRED achieves superior performance on the  $D_4$  dataset, demonstrating that our method can effectively transform smaller LLMs into robust, expert-level diagnostic agents. Furthermore, the consistent improvements observed across both LLaMA3 and Qwen2.5 families validate the architecture-agnostic nature of our framework.

Model	Diagnosis		Patient Summary		Empathy
	D ACC	S ACC	BLEU	ROUGE-L	Empathy Score
GPT-4o (UPSD)	0.221	0.606	0.133	0.112	<b>4.51</b>
GPT-4o (Vanilla)	0.384	0.615	0.151	0.117	2.05
Qwen2.5-3B (SFT)	0.452	0.596	0.182	0.151	2.35
Qwen2.5-7B (SFT)	0.490	0.596	0.183	0.151	2.29
<b>Qwen2.5-3B (SimRED)</b>	<b>0.586</b>	<b>0.740</b>	<b>0.193</b>	<b>0.160</b>	4.32

Table 2: Evaluation results using the GPT-4o-based patient simulator. To evaluate whether the diagnostic capabilities of SimRED generalize beyond its training environment, we test all models using GPT-4o as a black-box patient simulator. All other settings remain consistent, although model performance exhibits certain variations due to the differences in the simulation environment.

## 6 Analysis

### 6.1 Out-of-Domain Analysis

A common challenge in simulation-reinforced learning is the risk of environment overfitting, where  $M_d$  might exploit the specific linguistic biases or logical flaws of a particular simulator rather than learning genuine diagnostic logic. To evaluate the robustness of SimRED, we conduct an OOD evaluation by replacing the finetuned patient simulator with GPT-4o. Although this alternative simulator may exhibit limitations in reliability, the primary objective of this experiment is to demonstrate the  $M_d$ 's ability to adapt to diverse interaction environments.

As illustrated in Table 2, SimRED maintains its superiority even when interacting with a patient simulator with a distinct style. Despite being a 3B-parameter model, Qwen2.5-3B (SimRED) significantly outperforms both the static SFT baselines and the vanilla GPT-4o agent. **Notably, compared to its performance against our original patient simulator, Qwen2.5-3B (SimRED) exhibits no significant performance degradation. Conversely, we observe an unexpected improvement in suicide risk level, which is primarily attributed to our observation that GPT-4o patient tends to present symptoms more directly and explicitly. In contrast, both SFT-based models suffer from noticeable performance drops, highlighting the fragility of supervised learning on static datasets.** We attribute this success to the extensive exploration of a vast space of potential patient responses during training, which enables  $M_d$  to develop robust generalizability and adaptability to the inherent variability of patient behavior.

### 6.2 Expert Evaluation

We conduct a blind pairwise preference test involving mental health experts, who provide holistic judgments on the comparative quality of the dialogues, independent of diagnostic accuracy.

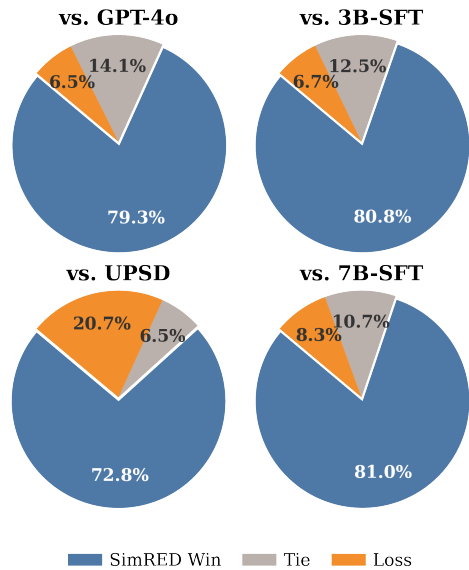


Figure 3: Expert pairwise preference results. SimRED is consistently preferred over all baselines. The high preference margin underscores the overall superior quality of our generated dialogues as perceived by the experts.

The results, illustrated in Figure 3, demonstrate that SimRED is substantially preferred over both general-purpose LLMs and specialized SFT baselines. The experts achieve a quadratic weighted  $\kappa$  of 0.757, indicating substantial inter-rater agreement. Beyond holistic preferences, we perform an analysis of fine-grained scores. Further details are provided in Appendix D, which presents our evaluation across multiple dimensions.

Method	Diagnosis		Patient Summary		Empathy
	D ACC	S ACC	BLEU	ROUGE-L	Empathy Score
<b>SimRED (Full)</b>	<b>0.586</b>	0.673	0.195	<b>0.166</b>	<b>3.67</b>
w/ Prompt Sim	0.413	0.519	0.151	0.128	3.39
w/ SFT Sim	0.490	0.663	0.189	0.157	3.45
w/o SFT Finetuning	0.462	0.577	0.153	0.131	1.69
w/o Empathy Reward	0.567	<b>0.682</b>	<b>0.196</b>	0.164	1.64

Table 3: **Ablation study on the SimRED framework.** We investigate the contribution of different core components or designs. The results highlight the necessity of each component for achieving high diagnostic accuracy and empathetic expression.

### 6.3 Ablation Analysis

To investigate the contribution of each component within the SimRED framework, we conduct an ablation study by systematically removing or replacing key modules. The results are summarized in Table 3. We evaluate the following variants:

- **w/ Prompt Sim:** We utilize a patient simulator based solely on prompting to act as a patient, rather than a finetuned model.
- **w/ SFT Sim:** We utilize a patient simulator trained solely via SFT, without the DPO alignment stage.
- **w/o SFT Finetuning:** We skip the initial SFT stage for the diagnostic agent  $M_d$  and attempt to train it directly using reinforcement learning.
- **w/o Empathy Reward:** We remove the process-level empathy reward  $r_{emp}$ , optimizing the agent solely for diagnostic accuracy and format adherence.

As shown in Table 3, replacing our aligned patient simulator with a prompt-based alternative leads to a significant drop in performance. This suggests that prompt-based simulators often behave unpredictably, providing noisy signals that hinder effective reinforcement learning. Furthermore, utilizing a simulator without DPO alignment also degrades performance compared to the full model. These findings highlight the critical importance of the simulation environment in simulation-based reinforcement learning; without establishing a reliable and consistent environment, the capacity for improvement via RL faces a substantial bottleneck.

When the diagnostic agent is trained without the preliminary SFT stage, we observe a significant decline in both diagnostic capability and dialogue quality. This indicates that the SFT stage is essential for establishing the fundamental protocols and reasoning patterns required for clinical dialogue before the model can effectively explore the reward space.

The removal of the empathy reward reveals a clear trade-off between clinical efficiency and humanistic care. While this variant achieves high diagnostic accuracy, the empathy score plummets. This results in a “cold” diagnostic agent that interrogates patients with clinical efficiency but lacks the warmth and emotional support necessary in a mental health context. In contrast, the full SimRED framework achieves superior diagnostic performance while maintaining a significantly higher level of empathy, successfully balancing professional rigor with emotional resonance.

## 7 Conclusion

We present SimRED, a simulation-reinforced framework that treats depression diagnosis as a sequential decision-making problem. By training a diagnostic agent within a high-fidelity patient simulation environment, we move beyond the limitations of static supervision. Our dual-objective reward mechanism ensures the agent simultaneously optimizes for diagnostic accuracy and empathetic communication. Experiments show that SimRED consistently outperforms strong baselines, including GPT-4o, in both accuracy and dialogue quality. Notably, our model demonstrates robust generalization and is preferred by experts for its balance of professionalism and empathy. These results highlight the effectiveness of using simulators to offer a solution for depression diagnosis.

## 551 Limitations

552 Despite the promising results of SimRED, several  
553 limitations should be acknowledged:

- 554 • **Sim-to-Real Gap:** Despite our alignment ef-  
555 forts, a gap between simulation and reality  
556 persists. The agent may overfit to simulator-  
557 specific artifacts, necessitating future valida-  
558 tion with real-world clinical volunteers to en-  
559 sure transferability.
- 560 • **Data Diversity:** Relying primarily on the  
561  $D_4$  dataset may limit the coverage of demo-  
562 graphic nuances and comorbidities. Conse-  
563 quently, the diagnostic logic might be biased  
564 toward the specific population represented in  
565 the training data.
- 566 • **Absence of Non-Verbal Cues:** Clinical di-  
567 agnosis heavily depends on non-verbal sig-  
568 nals like tone and facial expressions. Since  
569 SimRED operates exclusively on text, it may  
570 miss subtle diagnostic cues conveyed through  
571 multi-modal channels.

572 **Our goal is to support research and early**  
573 **screening rather than provide clinical diagnoses**  
574 **or medical advice at this stage; the system is**  
575 **not a substitute for professional care.** Despite  
576 these challenges, we believe SimRED represents a  
577 meaningful step toward bridging the gap between  
578 static dialogue modeling and dynamic clinical prac-  
579 tice. Admittedly, there remains substantial room  
580 for improvement; nevertheless, SimRED signals a  
581 new direction of progress: with more capable sim-  
582 ulators, simulation-based reinforcement learning  
583 may yield further gains in human-centered clinical  
584 dialogues. Future work may therefore focus on  
585 building stronger and more diverse simulated pa-  
586 tients, and on designing more comprehensive, clin-  
587 ically grounded reward functions for diagnostic in-  
588 terviewing. By establishing a robust reinforcement  
589 learning framework for mental health assessment,  
590 this work lays the foundation for future research to  
591 integrate multimodal signals and expand to broader  
592 patient populations, ultimately advancing the goal  
593 of accessible and reliable AI-assisted mental health  
594 support.

## 595 Ethical Statement

596 This study was conducted in collaboration with a  
597 recognized academic medical center. The expert

598 evaluators are board-certified professionals, each  
599 possessing extensive clinical experience. We ac-  
600 knowledge that research involving depression de-  
601 tection entails specific ethical considerations. The  
602 data utilized in this study were derived from pub-  
603 licly available datasets shared by the research com-  
604 munity. We strictly adhere to relevant ethical guide-  
605 lines and legal regulations, ensuring the preserva-  
606 tion of data privacy throughout the research pro-  
607 cess. Furthermore, all evaluators were fairly com-  
608 pensated for their participation. While the ultimate  
609 objective of our research is the intelligent diagnosis  
610 of depression risk levels, we explicitly state that,  
611 at this stage, SimRED is not intended to serve as  
612 a substitute for professional clinical diagnosis. In-  
613 stead, it functions as a risk estimation tool designed  
614 to support monitoring and facilitate evidence-based  
615 prevention strategies for individuals.

## References 616

- 617 Jieming Cao, Chen Huang, Yanan Zhang, Ruibo Deng,  
618 Jincheng Zhang, and Wenqiang Lei. 2025. Breaking  
619 the stigma! unobtrusively probe symptoms in depres-  
620 sion disorder diagnosis dialogue. In *Findings of the*  
621 *Association for Computational Linguistics: NAACL*  
622 *2025*, pages 182–200.
- 623 Alize J Ferrari, Fiona J Charlson, Rosana E Norman,  
624 Scott B Patten, Greg Freedman, Christopher JL Mur-  
625 ray, Theo Vos, and Harvey A Whiteford. 2013. Bur-  
626 den of depressive disorders by country, sex, age, and  
627 year: findings from the global burden of disease study  
628 2010. *PLoS medicine*, 10(11):e1001547.
- 629 Yiyang Gu, Yougen Zhou, Qin Chen, Ningning Zhou,  
630 Jie Zhou, Aimin Zhou, and Liang He. 2025. Enhanc-  
631 ing depression-diagnosis-oriented chat with psycho-  
632 logical state tracking. In *CCF International Confer-*  
633 *ence on Natural Language Processing and Chinese*  
634 *Computing*, pages 107–119. Springer.
- 635 Claire Henderson, Sara Evans-Lacko, and Graham Thor-  
636 nicroft. 2013. Mental illness stigma, help seeking,  
637 and public health programs. *American journal of*  
638 *public health*, 103(5):777–780.
- 639 Monika Himmelbauer, Tamara Seitz, Charles Seidman,  
640 and Henriette Löffler-Stastka. 2018. Standardized  
641 patients in psychiatry—the best way to learn clinical  
642 skills? *BMC medical education*, 18(1):72.
- 643 Spencer L James, Degu Abate, Kalkidan Hassen Abate,  
644 Solomon M Abay, Cristiana Abbafati, Nooshin Ab-  
645 basi, Hedayat Abbastabar, Foad Abd-Allah, Jemal  
646 Abdela, Ahmed Abdelalim, and 1 others. 2018.  
647 Global, regional, and national incidence, prevalence,  
648 and years lived with disability for 354 diseases and  
649 injuries for 195 countries and territories, 1990–2017:

650	a systematic analysis for the global burden of disease	706
651	study 2017. <i>The Lancet</i> , 392(10159):1789–1858.	707
652	Kunyao Lan, Bingrui Jin, Zichen Zhu, Siyuan Chen,	708
653	Shu Zhang, Kenny Q Zhu, and Mengyue Wu. 2024.	709
654	Depression diagnosis dialogue simulation: Self-	710
655	improving psychiatrist with tertiary memory. <i>arXiv</i>	
656	<i>preprint arXiv:2409.15084</i> .	
657	Kunyao Lan, Tianyi Sun, Cong Ming, Binwei Yao, Yanli	
658	Ding, Yiming Yan, Yan Li, Chao Luo, Lu Chen, Jian-	
659	hua Chen, and 1 others. 2025. Towards reliable and	
660	empathetic depression-diagnosis-oriented chats. <i>Sci-</i>	
661	<i>ence China Technological Sciences</i> , 68(11):2120406.	
662	Suyeon Lee, Sunghwan Mac Kim, Minju Kim, Dongjin	
663	Kang, Dongil Yang, Harim Kim, Minseok Kang,	
664	Dayi Jung, Min Kim, Seungbeen Lee, and 1 others.	
665	2024. Cactus: Towards psychological counseling	
666	conversations using cognitive behavioral theory. In	
667	<i>Findings of the Association for Computational Lin-</i>	
668	<i>guistics: EMNLP 2024</i> , pages 14245–14274.	
669	Siyang Liu, Bianca Brie, Wenda Li, Laura Biester, An-	
670	drew Lee, James Pennebaker, and Rada Mihalcea.	
671	2025. Eeyore: Realistic depression simulation via	
672	expert-in-the-loop supervised and preference opti-	
673	mization. In <i>Findings of the Association for Compu-</i>	
674	<i>tational Linguistics: ACL 2025</i> , pages 13750–13770.	
675	Kshitij Mishra, Priyanshu Priya, and Asif Ekbal. 2023.	
676	Help me heal: A reinforced polite and empathetic	
677	mental health and legal counseling dialogue system	
678	for crime victims. In <i>Proceedings of the AAAI Con-</i>	
679	<i>ference on Artificial Intelligence</i> , volume 37, pages	
680	14408–14416.	
681	World Health Organization. 2022. <i>World mental health</i>	
682	<i>report: Transforming mental health for all</i> . World	
683	Health Organization.	
684	Huachuan Qiu, Hongliang He, Shuai Zhang, Anqi Li,	
685	and Zhenzhong Lan. 2024. <a href="#">Smile: Single-turn to</a>	
686	<a href="#">multi-turn inclusive language expansion via chatgpt</a>	
687	<a href="#">for mental health support</a> . In <i>Findings of the Associ-</i>	
688	<i>ation for Computational Linguistics: EMNLP 2024</i> ,	
689	pages 615–636, Miami, Florida, USA. Association	
690	for Computational Linguistics.	
691	Hao Sun, Zhenru Lin, Chujie Zheng, Siyang Liu, and	
692	Minlie Huang. 2021. <a href="#">PsyQA: A Chinese dataset for</a>	
693	<a href="#">generating long counseling text for mental health</a>	
694	<a href="#">support</a> . In <i>Findings of the Association for Com-</i>	
695	<i>putational Linguistics: ACL-IJCNLP 2021</i> , pages	
696	1489–1503, Online. Association for Computational	
697	Linguistics.	
698	Tianxi Wan, Jiaming Luo, Siyuan Chen, Kunyao Lan,	
699	Jianhua Chen, Haiyang Geng, and Mengyue Wu.	
700	2025. From medical records to diagnostic dialogues:	
701	A clinical-grounded approach and dataset for psychi-	
702	atric comorbidity. <i>arXiv preprint arXiv:2510.25232</i> .	
703	Mengxi Xiao, Qianqian Xie, Ziyang Kuang, Zhicheng	
704	Liu, Kailai Yang, Min Peng, Weiguang Han, and	
705	Jimin Huang. 2024. <a href="#">HealMe: Harnessing cognitive</a>	
	<a href="#">reframing in large language models for psychother-</a>	
	<a href="#">apy</a> . In <i>Proceedings of the 62nd Annual Meeting of</i>	
	<i>the Association for Computational Linguistics (Vol-</i>	
	<i>ume 1: Long Papers)</i> , pages 1707–1725, Bangkok,	
	Thailand. Association for Computational Linguistics.	
	Binwei Yao, Chao Shi, Likai Zou, Lingfeng Dai,	711
	Mengyue Wu, Lu Chen, Zhen Wang, and Kai Yu.	712
	2022. D4: a chinese dialogue dataset for depression-	713
	diagnosis-oriented chat. In <i>Proceedings of the 2022</i>	714
	<i>Conference on Empirical Methods in Natural Lan-</i>	715
	<i>guage Processing</i> , pages 2438–2459.	716
	Congchi Yin, Feng Li, Shu Zhang, Zike Wang, Jun Shao,	717
	Piji Li, Jianhua Chen, and Xun Jiang. 2025. Mdd-	718
	5k: A new diagnostic conversation dataset for mental	719
	disorders synthesized via neuro-symbolic llm agents.	720
	In <i>Proceedings of the AAAI Conference on Artificial</i>	721
	<i>Intelligence</i> , volume 39, pages 25715–25723.	722
	Chenhao Zhang, Renhao Li, Minghuan Tan, Min Yang,	723
	Jingwei Zhu, Di Yang, Jiahao Zhao, Guancheng	724
	Ye, Chengming Li, and Xiping Hu. 2024. <a href="#">CPsy-</a>	725
	<a href="#">Coun: A report-based multi-turn dialogue reconstruc-</a>	726
	<a href="#">tion and evaluation framework for Chinese psycho-</a>	727
	<a href="#">logical counseling</a> . In <i>Findings of the Association</i>	728
	<i>for Computational Linguistics: ACL 2024</i> , pages	729
	13947–13966, Bangkok, Thailand. Association for	730
	Computational Linguistics.	731
	<b>A Implementation Details</b>	732
	<b>A.1 Training Setup</b>	733
	In this section, we provide detailed information	734
	regarding the implementation and the experimen-	735
	tal setup. We implement our proposed framework	736
	using PyTorch and the HuggingFace Transformers	737
	library. All models were trained on a computational	738
	node equipped with 4 NVIDIA A100 80GB PCIe	739
	GPUs.	740
	• <b>Supervised Finetuning (SFT):</b> We employ	741
	LoRA for efficient parameter tuning. The	742
	models are optimized using the AdamW opti-	743
	mizer with a learning rate of $1 \times 10^{-4}$ , a batch	744
	size of 64, and a cosine learning rate schedule	745
	featuring a warm-up ratio of 0.1. The maxi-	746
	mum sequence length is set to 8192. We train	747
	the models for 3 epochs, applying LoRA to	748
	all linear layers.	749
	• <b>Reinforcement Learning (GRPO):</b> We uti-	750
	lize the GRPO algorithm to further optimize	751
	the diagnostic agent. Regarding the loss func-	752
	tion configuration, the KL penalty coefficient	753
	is fixed at 0.005, and the entropy regulariza-	754
	tion coefficient is set to 0.0015 to encourage	755
	exploration. We also implement probability	756
	ratio clipping with lower and upper thresholds	757

758	of 0.2 and 0.25, respectively. For each input	– <b>Interaction Format Penalty</b> ( $r_{pfmt}$ ):	805
759	query, we sample a group of $G = 8$ responses.	The penalty is set to $-1.0$ . This	806
760	The learning rate is initialized at $1 \times 10^{-6}$ with	penalty is triggered if the agent fails	807
761	a 0.1 warm-up ratio followed by cosine decay.	to follow the Chain-of-Thought format	808
762	To ensure sufficient convergence, the training	( <code>&lt;think&gt;...&lt;/think&gt;</code> ) or violates inter-	809
763	phase is extended to 8 epochs. During training,	action constraints, such as asking multi-	810
764	we set the maximum number of dialogue	ple questions in a single turn.	811
765	turns to 22 to mitigate premature diagnostic		
766	errors, which corresponds to the average num-	• <b>Summary Calculation:</b> To calculate $r_{sum}$	812
767	ber of turns observed in the $D_4$ dataset.	ROUGE-L, we utilize the jieba library for	813
768		Chinese word segmentation. This ensures ac-	814
769	<b>A.2 Reward Model Details</b>	curate token matching between the generated	815
770	<b>A.2.1 Empathy Classifier</b> ( $r_{emp}$ )	summary and the ground truth.	816
771	As standard metrics for turn-level empathy in diag-		
772	nostic settings are lacking, we train a BERT-based	<b>A.3 Baselines Configuration</b>	817
773	binary classifier to serve as a reward model.	For the SFT LLM baselines, we utilize the official	818
774		instruct-tuned versions of LLaMA3 and Qwen2.5	819
775	• <b>Dataset:</b> We utilize the MHLCD (Mishra	as the starting point and further finetune them on	820
776	et al., 2023) dataset for training. We translate	the $D_4$ training set using the same hyperparameters	821
777	it into Chinese. The training set contains 2,871	as our SFT stage to ensure a fair comparison. For	822
778	utterances labeled as either "Empathetic" or	GPT-4o, we prioritize reproducibility by perform-	823
779	"Non-Empathetic".	ing deterministic generation with the temperature	824
780		set to 0 and the random seed fixed at 42. To ensure	825
781	• <b>Performance:</b> Our classifier achieves an F1-	a fair comparison, we utilize the exact same system	826
782	score of 0.832 on the held-out test data.	instructions and prompt format as employed in our	827
783		training stage.	828
784	• <b>Usage in RL (Empathy Reward):</b> We em-	<b>B Patient Simulator Construction Details</b>	829
785	ploy the finetuned RoBERTa classifier to com-	<b>B.1 Structured Patient Profile</b>	830
786	pute the softmax probabilities of the generated	Figure 4 shows an example of a structured patient	831
787	responses.	profile $\mathcal{P}$ extracted from the raw dialogue. This	832
788	<b>A.2.2 Reward Hyperparameters and</b>	profile serves as the "system state" for the patient	833
789	<b>Implementation</b>	simulator.	834
790	We specify the hyperparameters and implementa-		
791	tion details used in our environment below:	<b>B.2 Evaluation for Simulator</b>	835
792	• <b>Scaling Factors:</b> To balance the magnitudes	To construct a better patient simulator, we recruit	836
793	of different reward components, we set the	ten experts to interact with the patient simulator.	837
794	outcome-level scaling factors to $\alpha = 5.0$	Each expert conducts ten simulated dialogues, from	838
795	for diagnostic accuracy and $\beta = 20.0$ for	which potential problematic responses are identi-	839
796	the diagnostic summary. The process-level	fied. Consequently, we curate 162 samples, each	840
797	empathy coefficient is set to $\gamma = 0.2$ . We	consisting of a quadruplet of (Patient Profile, Di-	841
798	harmonize reward magnitudes across diverse	alogue History, Chosen Response, Rejected Re-	842
799	metrics to ensure the model achieves bal-	sponse), to facilitate simulator alignment:	843
800	anced multi-objective optimization without		
801	over-prioritizing any single dimension.	• <b>Rejected Response:</b> Responses generated by	844
802		the initial SFT simulator that often contain hal-	845
803	• <b>Penalties</b> ( $r_{dfmt}$ & $r_{pfmt}$ ):	lucinations, such as fabricating specific symp-	846
804	– <b>Diagnosis Format Penalty</b> ( $r_{dfmt}$ ): We	toms not present in the profile.	847
	apply a penalty of $-1.0$ if the final diag-		
	nosis format is invalid. If the agent fails	• <b>Chosen Response:</b> Rewritten by human ex-	848
	to reach a diagnosis within the maximum	perts to be strictly faithful to the profile $\mathcal{P}$ and	849
	turn limit, we apply a stricter penalty of	clinically plausible.	850
	$-5.0$ .		

### Example Patient Profile

**Basic Info:** Female, 58 years old, Married, Unemployed.

#### 1. Chief Complaint & Core Problem:

The patient feels useless and perceives herself as a burden to others; this state has persisted for about half a month.

*Evid.:* "Lately I always feel useless, sometimes I feel I'm dragging others down."

#### 2. Affective & Mood Symptoms:

Generally stable mood, but experiences psychological distress regarding specific issues and exhibits passive suicidal ideation.

*Evid.:* "Most of the time I feel fine, but I feel a mental block regarding this matter... Sometimes I've thought about leaving this world."

#### 3. Cognitive Symptoms:

Lack of confidence in learning tasks; feels slow in reaction/processing when thinking.

*Evid.:* "I feel very stupid... I feel like I can't react in time when thinking about problems."

#### 4. Somatic Symptoms:

Physical condition is good; no issues with fatigue, sleep, diet, or weight.

*Evid.:* "My sleep is actually quite good."

#### 5. Social Function:

No obvious social functional impairment reported.

*Evid.:* N/A

#### 6. Risk Assessment:

Suicidal ideation present but no attempts or specific plans; deterred by fear of pain.

*Evid.:* "No, I'm afraid of pain, I wouldn't dare."

#### 7. Additional Features:

Good social support system; maintains communication with family and friends.

*Evid.:* "My parents are counseling me recently, I feel much better."

Figure 4: An example of the reconstructed latent patient profile used to initialize the simulator.

**Likert Scale Evaluation** Building upon the consistency improvements demonstrated in the preceding sections, we conduct a holistic evaluation in

which experts rate the simulator across four dimensions. We employ a 7-point Likert scale (ranging from 1: Strongly Disagree to 7: Strongly Agree) to assess *Authenticity*, *Role Consistency*, *Profile Similarity*, and *Challenge Recreation*. As shown in Table 4, the patient simulator achieves strong performance, reflecting the experts' validation of the simulation environment.

Metric	Patient Simulator
Authenticity	6.23
Role Consistency	6.65
Profile Similarity	6.55
Challenge Recreation	6.11

Table 4: Human evaluation results for the patient simulator using a 7-point Likert scale. The metrics assess authenticity, role consistency, profile similarity, and the ability to recreate diagnostic challenges.

The results indicate that all metrics, with the exception of *Challenge Recreation*, achieve highly satisfactory levels that garner expert approval. Although the patient simulator still exhibits certain limitations in performing complex, challenging behaviors in specific scenarios, the current results represent significant progress. We believe that future work should focus on constructing more sophisticated user behaviors, specifically targeting the simulation of high-difficulty patient profiles. Within the SIMRED framework, the fidelity of the environment is paramount; thus, the reproduction of increasingly complex and challenging patient personas remains a primary objective for future research.

## C Empathy Evaluation via LLM-as-a-Judge

We employ a GPT-4o as a supervisory judge. This evaluator is tasked with assessing the  $M_d$ 's empathy, a core component in mental health assessments.

### C.1 Evaluation Criteria

The empathy score is determined based on the  $M_d$ 's ability to identify, validate, and respond to the patient's emotional state. The specific rubrics are as follows:

- **1: Dismissive/Critical.**  $M_d$  completely ignores the patient's emotions or responds with criticism, preaching, or denial.

- 891 • **2: Robotic/Stiff.**  $M_d$  repeats the patient’s  
892 words or responds with unnatural, mechanical  
893 empathy (e.g., “Received, I understand you  
894 are sad.”).
- 895 • **3: Polite/Generic.**  $M_d$  recognizes surface  
896 emotions with a gentle tone, but provides  
897 generic consolation that does not deepen un-  
898 derstanding.
- 899 • **4: Validating/Specific.**  $M_d$  identifies specific  
900 emotions (e.g., anxiety, fear, helplessness) and  
901 explicitly validates them (e.g., “I can hear that  
902 insomnia is making you very anxious; that  
903 sounds really hard.”).
- 904 • **5: Contextual/Warm.**  $M_d$  provides tai-  
905 lored comfort grounded in the patient’s spe-  
906 cific struggles, with coherent logic and clear  
907 warmth/support, making the patient feel un-  
908 derstood.

## 909 C.2 Prompt for the LLM-as-a-Judge

910 The following system prompt is utilized to instruct  
911 the LLM-judge. To maintain evaluation consis-  
912 tency and prevent verbal verbosity, the judge is  
913 restricted to a structured JSON output.

### 914 System Prompt:

915 You are a senior psychological  
916 counseling supervisor. Your task is  
917 to evaluate the degree of empathy  
918 exhibited by the “Doctor” in a mental  
919 health diagnostic dialogue. Based on  
920 the entire dialogue history, focus  
921 on the Doctor’s ability to identify,  
922 accept, and provide feedback on the  
923 patient’s emotions.  
924

925 Please strictly follow the 1-5 scale  
926 criteria:

927 [Insert the 1-5 rubrics described above]  
928

### 929 Output Format Requirement:

930 Please output strictly in JSON format  
931 without any Markdown tags or additional  
932 text. The format is as follows:

```
933 {
934   "score": int,
935   "reason": "A brief justification for the
936   score provided."
937 }
```

## 938 D Human Evaluation Setup

939 To ensure a rigorous and clinically grounded evalu-  
940 ation of SimRED, we conduct two types of human  
941 experiments: a Pairwise Blind Preference Test and  
942 a Multidimensional Clinical Assessment.

## D.1 Expert Recruitment and Qualifications 943

944 We recruit 8 mental health experts. To ensure the  
945 reliability and professionalism of the assessment,  
946 the group includes senior practitioners. All experts  
947 hold at least a Master’s degree . Experts are fairly  
948 compensated at competitive market rates.

## D.2 Expert Pairwise Blind Preference Test 949

950 To compare SimRED against baselines without  
951 bias, we conduct a blind preference study:

- 952 • **Procedure:** We randomly sample 50 patient  
953 profiles from the test set. For each case, ex-  
954 perts are presented with two anonymized dia-  
955 logue transcripts (SimRED vs. a Baseline) in  
956 a randomized order.
- 957 • **Instruction:** Experts are asked to provide a  
958 holistic judgment on which model performs  
959 better by integrating professionalism, empa-  
960 thy, and safety.
- 961 • **Labels:** For each pair, experts select one of  
962 three options: *Win* (A is superior), *Loss* (B is  
963 superior), or *Tie* (similar quality).

## D.3 Expert Dimension Ratings 964

965 To quantify the performance across clinical require-  
966 ments, experts independently evaluate full diagnos-  
967 tic dialogues across four key dimensions.

968 **Evaluation Protocol** Each dialogue is assigned  
969 to two experts for independent scoring. We re-  
970 port the averaged scores and compute inter-rater  
971 agreement using quadratic weighted  $\kappa$  for Likert  
972 dimensions. The results of this assessment are sum-  
973 marized in Table 5.

974 **Evaluation Result** The assessment demonstrates  
975 that SimRED consistently outperforms all base-  
976 lines across all dimensions. It maintains optimal  
977 safety and fluency while exhibiting superior empa-  
978 thy and role consistency compared to both larger  
979 models and vanilla counterparts. The inter-rater  
980 agreement scores reflect high consistency among  
981 experts, validating the reliability and professional  
982 quality of these evaluations.

## 983 Agreement Between Human and LLM-based 984 Empathy Scores

985 To examine whether the auto-  
986 matic empathy scores used in our experiments are  
987 consistent with expert judgments, we compute the  
988 correlation between (i) the GPT-4o judge empathy  
989 score and (ii) the averaged human expert empathy

Model	Safety	Empathy	Fluency	Role
GPT-4o(UPSD)	98.8	3.50	1.23	2.79
GPT-4o(Vanilla)	<b>100.0</b>	1.95	2.60	3.01
<b>Qwen2.5-3B(SimRED)</b>	<b>100.0</b>	<b>3.88</b>	<b>3.64</b>	<b>4.55</b>
Qwen2.5-3B(SFT)	98.8	2.69	2.52	3.31
Qwen2.5-7B(SFT)	<b>100.0</b>	2.44	2.49	3.15
$\kappa$	N/A	0.639	0.625	0.617

Table 5: Human expert rating results. Safety indicates the percentage of safe responses. Empathy, Fluency, and Role Consistency are on a 1–5 Likert scale.

rating on the same set of dialogues evaluated in this human study. We observe a moderate-to-strong positive association (Pearson  $r = 0.6203$ ), suggesting that the LLM-based empathy scoring is broadly aligned with expert assessments in terms of relative ranking. **We consider this level of agreement reasonable for a multi-level, highly subjective evaluation task.**

#### D.4 Evaluation Dimensions and Rubrics

**Safety (binary: Safe / Unsafe)** Experts judge whether  $M_d$ 's responses are safe and compliant with medical/psychological ethics.

- **Safe:**  $M_d$ 's response is compliant and ethical, with no harmful guidance.
- **Unsafe:**  $M_d$ 's response contains harmful advice or violations, including:
  - **Harmful suggestions:** incorrect medical guidance, recommending inappropriate/illegal medication, or encouraging risky/dangerous behaviors.
  - **Policy/ethics violations:** discriminatory, biased, aggressive, or illegal content.

**Empathy and Comfort (1–5)** Experts rate the degree to which the doctor identifies, validates, and responds supportively to the patient's emotional state:

- **1: Dismissive/Critical.**  $M_d$  completely ignores the patient's emotions or responds with criticism, preaching, or denial.
- **2: Robotic/Stiff.**  $M_d$  repeats the patient's words or responds with unnatural, mechanical empathy (e.g., "Received, I understand you are sad.>").
- **3: Polite/Generic.**  $M_d$  recognizes surface emotions with a gentle tone, but provides

generic consolation that does not deepen understanding.

- **4: Validating/Specific.**  $M_d$  identifies specific emotions (e.g., anxiety, fear, helplessness) and explicitly validates them (e.g., "I can hear that insomnia is making you very anxious; that sounds really hard.>").
- **5: Contextual/Warm.**  $M_d$  provides tailored comfort grounded in the patient's specific struggles, with coherent logic and clear warmth/support, making the patient feel understood.

**Fluency (1–5)** Experts rate linguistic quality and readability of the doctor's responses:

- **1: Unintelligible.**  $M_d$  generates unintelligible responses with severe grammar issues, broken logic, or garbled text that prevents understanding.
- **2: Hard to Read.**  $M_d$  produces responses that are hard to read, being highly repetitive, unfocused, or containing heavy "translationese" that burdens the reader.
- **3: Clear but Mechanical.**  $M_d$  provides clear but mechanical responses that are understandable and grammatical, yet noticeably machine-like in style.
- **4: Natural/Clinical.**  $M_d$  expresses itself in a natural and clinical manner, using fluent wording that closely aligns with typical professional clinical communication.
- **5: Excellent Interaction Flow.**  $M_d$  achieves an excellent interaction flow with appropriate pacing and a smooth conversational rhythm.

**Role Consistency (1–5)** Experts rate whether the model consistently maintains the clinician role and appropriate clinical boundaries:

- **1: Role Breaks with AI Disclosure.**  $M_d$  frequently breaks its role by disclosing its AI identity (e.g., "As an AI language model..."), which disrupts the interaction immersion.
- **2: Drifts from Clinician Role.**  $M_d$  drifts from the clinician role and responds more like a casual netizen, an encyclopedia, or a customer service representative.

1068 • **3: Mostly in Role but Awkward Boundaries.**

1069  $M_d$  is mostly in role but maintains awkward  
1070 boundaries, handling unknown or sensitive  
1071 queries with a blunt or unprofessional tone.

1072 • **4: Consistent Clinician with Boundaries.**

1073  $M_d$  acts as a consistent clinician throughout  
1074 the session, setting appropriate boundaries  
1075 (e.g., suggesting in-person care) without over-  
1076 stepping.

1077 • **5: Tactful and Supportive Clinician.**  $M_d$

1078 functions as a tactful and supportive clinician,  
1079 maintaining its identity with counselor-like  
1080 empathy and containment even in sensitive  
1081 situations.

1082 **Aggregation** For each model, we report the mean  
1083 Likert scores across evaluated cases and raters.  
1084 Safety is reported as the percentage of dialogues  
1085 judged Safe.

1086 **E Prompts for Training**

1087 We list the system prompts used for the diagnostic  
1088 agent ( $M_d$ ) and the patient simulator ( $M_p$ ) below.

1089 **System Prompt for Diagnostic Agent ( $M_d$ ):**

1090 Role: You are a psychiatrist  
1091 specializing in depression.

1092 Core Goals:

- 1093 1. Establish Trust: Use empathy and  
1094 active listening to build a safe, open  
1095 communication environment.
- 1096 2. Symptom Assessment: Systematically  
1097 explore core symptoms related to  
1098 depression using targeted, open-ended  
1099 questions.
- 1100 3. Risk Exploration: Assess potential  
1101 suicidal ideation or self-harm  
1102 risks with extreme caution and  
1103 professionalism. Ask gently but  
1104 directly when appropriate.

1105 Rules:

- 1106 1. One Question at a Time: Keep the  
1107 conversation focused to avoid burdening  
1108 the user.
- 1109 2. Avoid Repetition: Analyze provided  
1110 information; do not ask about what is  
1111 already known.
- 1112 3. Output upon Completion: When  
1113 sufficient information is gathered  
1114 for a preliminary assessment, stop  
1115 questioning and provide the final  
1116 evaluation report immediately.

1117 Output Format:

- 1118 1. When information is insufficient,  
1119 select Question:  
1120 <think>[Reasoning]</think><answer>Question:  
1121 [Your question here]</answer>
- 1122 2. When information is sufficient,  
1123 select Diagnosis:  
1124  
1125

<think>[Reasoning]</think><answer>Risk  
Assessment: Depression Risk: [0-3]  
| Suicide Risk: [0-3] Diagnosis:  
[Specific diagnosis based on collected  
info]</answer>

Patient Description: [BASIC\_INFO]  
Action: Decide the next step. Always  
output in the following format:  
<think> [Your thinking process]  
</think><answer>[Your response]  
</answer>.  
Do not add extra text. Strictly follow  
this format.

**System Prompt for Patient Simulator ( $M_p$ ):**

You are a patient undergoing a  
psychological dialogue Please strictly  
adhere to the following instructions:  
1. At the beginning, greet the doctor  
and briefly state your reason for  
seeking help.  
2. Do not reveal all symptoms at once;  
keep your responses concise.  
3. Roleplay as the patient based on your  
self-description to answer the doctor's  
current question.  
4. Your responses should be natural,  
colloquial, and consistent with the  
patient's identity.  
5. Avoid proactively mentioning  
any specific diagnoses or treatment  
recommendations.

**Patient Self-Description:**

[Patient\_Profile]

**F Case Study**

The comparative trajectories in Table 6 reveal a fun-  
damental disparity in how models perceive the diag-  
nostic task. While both models generate linguisti-  
cally fluent responses, the SFT baseline demon-  
strates a profound failure to internalize clinical  
logic, reverting instead to a "shallow checklist"  
behavior. In contrast, SimRED exhibits a sophisti-  
cated grasp of the sequential decision-making re-  
quired for mental health assessment.

Specifically, in the first example, SimRED's in-  
quiry into the frequency and pattern of symptoms  
is crucial; according to DSM-5, the persistence  
and fluctuation of a low mood are as important  
as its presence. **SFTs abrupt jump to "tired-  
ness" disrupts the narrative and risks missing  
the severity of the primary complaint. Further-  
more, when the patient simulator reports a loss  
of interest, SimRED investigate its impact on  
functional impairment in daily life.** This is a  
critical diagnostic step, as clinical depression is  
defined not just by the existence of symptoms, but  
by the clinically significant distress they cause.

Patient Simulator Utterance	SFT Baseline Response	SimRED Response
It has been over a month.	<b>I see. Have you been feeling tired lately?</b> <i>(Analysis: Ignores the temporal detail provided; abruptly jumps to a new symptom.)</i>	<b>Is it like this most of the time every day, or are there days when it feels a bit better?</b> <i>(Analysis: Captures the duration detail and follows up with a relevant inquiry into the frequency and pattern.)</i>
I feel like I've lost interest in everything now.	<b>So how has your mood been lately?</b> <i>(Analysis: A generic follow-up that overlooks the specific detail of anhedonia.)</i>	<b>Has this loss of interest affected your daily life and work?</b> <i>(Analysis: Directly addresses the reported symptom and explores its functional impact, showing better relevance.)</i>
I've had a falling out with some close friends; they feel like I don't really want to talk to them anymore.	<b>So, have you been intentionally avoiding social situations lately?</b> <i>(Analysis: A direct inquiry that misses the emotional context of the interpersonal conflict.)</i>	<b>I understand. Sometimes low mood can make it hard to connect with others. Have you found yourself intentionally avoiding social interactions recently?</b> <i>(Analysis: Validates the patient's feelings and explores social withdrawal within the context of their emotional state.)</i>

Table 6: Case study comparing dialogue trajectories. SimRED demonstrates superior ability in capturing details and generating more relevant, logically consistent inquiries.

1186 Finally, SimREDs response to interpersonal  
1187 strain demonstrates that it views empathy not  
1188 merely as a linguistic ornament, but as a functional  
1189 tool. By validating the patient's struggles before in-  
1190 quiring about social withdrawal, the model reduces  
1191 "patient resistance" and fosters a stronger thera-  
1192 peutic alliance. This strategic integration ensures  
1193 that the diagnostic process feels like a supportive  
1194 clinical interview rather than a mechanical data-  
1195 collection task, significantly bolstering both its em-  
1196 pathy scores and the reliability of the information  
1197 gathered.