

# Roleplay-doh: Enabling Domain-Experts to Create LLM-simulated Patients via Eliciting and Adhering to Principles

Ryan Louie, Ananjan Nandi, William Fang  
Cheng Chang, Emma Brunskill, Diyi Yang  
Stanford University

## Abstract

Recent works leverage LLMs to roleplay realistic social scenarios, aiding novices in practicing their social skills. However, simulating sensitive interactions, such as in mental health, is challenging. Privacy concerns restrict data access, and collecting expert feedback, although vital, is laborious. To address this, we develop Roleplay-doh, a novel human-LLM collaboration pipeline that elicits qualitative feedback from a domain-expert, which is transformed into a set of principles, or natural language rules, that govern an LLM-prompted roleplay. We apply this pipeline to enable senior mental health supporters to create customized AI patients for simulated practice partners for novice counselors. After uncovering issues in GPT-4 simulations not adhering to expert-defined principles, we also introduce a novel principle-adherence prompting pipeline which shows 30% improvements in response quality and principle following for the downstream task. Via a user study with 25 counseling experts, we demonstrate that the pipeline makes it easy and effective to create AI patients that more faithfully resemble real patients, as judged by creators and third-party counselors. See our project website<sup>1</sup> for code and data.

## 1 Introduction

The application of LLMs in simulations holds great potential for a variety of interactive applications, ranging from social skill training systems as AI practice partners (Yang et al., 2024) to prototyping tools that use them as believable proxies of human behavior (Park et al., 2022). However, achieving realistic and reliable simulations remains a significant challenge, due to issues such as caricature (Cheng et al., 2023), bias, and limited domain knowledge. Existing methods for improving LLM simulations such as finetuning (Demasi et al., 2020)

can help, but such methods typically require the use of application-specific datasets. In sensitive application domains like mental health, privacy concerns with obtaining the required data can restrict the feasibility of such methods. This suggests that *experts-in-the-loop* may be a powerful alternative to guide the evaluation and refinement (Chen et al., 2023) of LLM-powered simulations.

However, how to involve experts when improving simulations is an open challenge. Collecting sufficient amounts of binary or preference data from experts for post-training (Christiano et al., 2017; Rafailov et al., 2024) can be tedious and expensive. Experts can guide the prompting of LLM simulations, directly by editing their own prompts or indirectly through testing and think-aloud sessions. However each prompting method has its limitations: domain-experts may not know how to prompt simulations for desired behaviors (Zamfirescu-Pereira et al., 2023); and indirect methods are inefficient as it requires a designer or researcher to translate qualitative insights into prompt-design changes.

As a focal example, we consider the problem of creating AI patients that serve as roleplay partners to enable varied and interactive practice opportunities for novice therapists and counselors (Yao et al., 2022). Creating realistic simulations by finetuning on mental health data is infeasible because therapy transcripts with real patients is difficult to obtain due to privacy concerns. Naively prompting LLMs fail to resemble typical behaviors of real-patients—for example, mental health experts report that patients use colloquial language and can show resistance to help (Chen et al., 2023). To date, no system supports counseling experts, who are familiar with real-patient behaviors but are unlikely to have the technical expertise to write effective prompts, to customize an AI patient themselves.

To address this limitation, we aim to enable human-LLM collaboration for realistic simula-

Contact Emails: {rylouie, diyi}@stanford.edu

<sup>1</sup><https://roleplay-doh.github.io/>

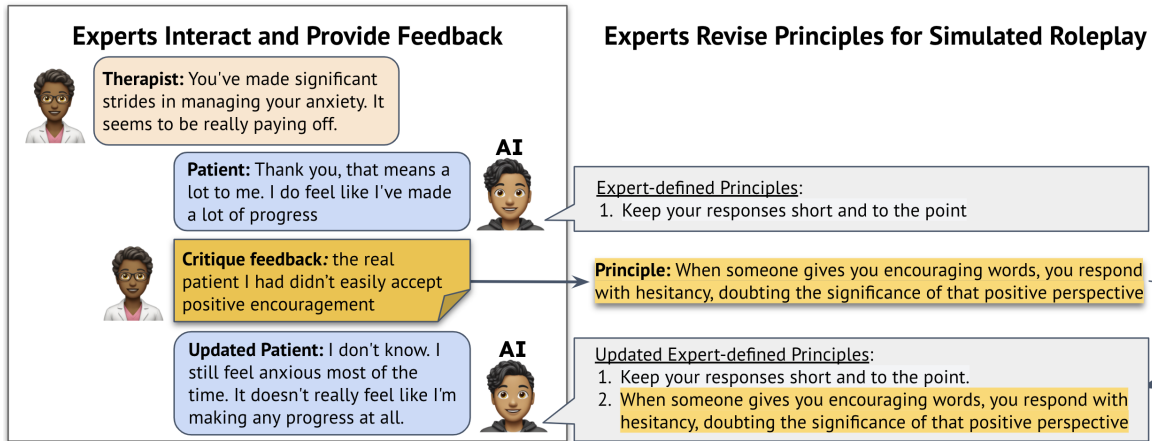


Figure 1: Roleplay-doh empowers an expert counselor to create a customized AI patient intended for other novice counselors to use as a practice partner. While interacting with the AI patient, the expert counselor can provide qualitative feedback which is converted by an LLM into a principle, or a custom rule governing desired roleplay behavior. The AI patient references the updated expert-defined principles to generate its subsequent responses.

tion by developing a novel interactive tool, called Roleplay-doh, that empowers domain experts to *directly* guide the creation of simulations by providing *qualitative feedback without any explicit prompting*. Our initial tool design adopts an intuitive and effective paradigm for user-driven chatbot assistant design (Petridis et al., 2024): experts customize a set of *principles*, or rules written in natural language that govern its behavior (Bai et al., 2022)—by (1) interactively critiquing responses in natural language that then (2) gets transformed by an LLM into well-formulated principles describing how the LLM simulation should act from now on *for example*, "Respond to encouraging words with hesitation, doubting their significance" (Fig 1). The principles are then used along with a persona description to generate roleplay responses.

In our initial tests of the tool with expert-counselors, we found that even with expert refinement via principles, the LLM- simulations had difficulty delivering high-quality responses consistently. Our analysis of GPT-4 prompted simulation revealed that in 20% of responses, the simulation had difficulty adhering to multipart principles and misapplying those principles that are only applicable in specific contexts e.g., *only when the therapist provides encouraging words*. To resolve these issues, we introduce a novel **principle-adherence pipeline** in the final tool design. The first stage in the pipeline decomposes multipart and contextual principles into a set of yes/no questions that are easier to judge, and the second stage assesses the applicability of each simplified principle to the

current scenario before self-refining (Madaan et al., 2024) the AI patient response as required.

We conducted a detailed evaluation of Roleplay-doh to assess its human-LLM collaboration pipeline, focusing on how expert feedback helps develop more authentic AI patients for training. In a within-subjects study involving 25 expert counselors, participants created AI patients either by describing real-patient scenarios or by using Roleplay-doh to refine simulation principles. The results show that Roleplay-doh enables counselors to produce AI patients that are more authentic, closely resemble real cases, and are better prepared for training use, as judged by creators and third-party counselors. Further, our principle-adherence pipeline achieves the highest principle following (win rate: 35%; loss rate: 5%) and dialogue consistency (win rate: 35%; loss rate: 10%) compared to all ablations, where preferences are made against a baseline that does not self-refine its output. This work highlights the limitations of existing LLM simulation systems in specialized, data-scarce domains like mental health counseling, and designs and validates a tool that enables expert counselors to directly customize LLM simulations of AI patients. Since Roleplay-doh does not contain any technical components specifically tailored for mental health, we hypothesize that the tool can be used to build realistic LLM simulations for a wide variety of domains, with appropriate expert feedback.

## 2 Related Work

**Utility of Simulated Partners** Simulated partners are used to give social skill learners the needed

practice opportunities that textbook knowledge cannot provide. Past education software develops digital patient simulations to make simulated partners more accessible (Othlinghaus-Wulhorst and Hoppe, 2020) but their tailored dialogue setup limit the contexts for practice. LLMs can overcome this issue by being flexibly configured to convincingly simulate a diverse set of personas (Park et al., 2022) and characters (Park et al., 2023) and generate responses in a range of contexts. Researchers have thus explored their application for simulation training for K12 teaching (Markel et al., 2023), conflict resolution (Shaikh et al., 2024), and counseling (Demasi et al., 2020; Tanana et al., 2019; Chen et al., 2023). Previous work has proposed methods to simulate diverse personas and scenarios, but to make practice more transferable (Alinier and Oriot, 2022), they must ensure simulations are faithful to what is encountered in real-world social situations.

### Aligning Simulation with Domain Experts

Feedback from domain experts is crucial to evaluating and improving the realism of LLM simulations. Recent approaches for aligning to human feedback, like Christiano et al. (2017) or Rafailov et al. (2024) depend on large amounts of preference data which requires lots of expert time to collect. A more efficient approach is through alignment to qualitative or natural language feedback (Shi et al., 2022). We build on a recent paradigm for user-driven chatbot design (Petridis et al., 2024) that elicits qualitative feedback on responses which gets converted into constitution principles (Bai et al., 2022), which are explainable and effective natural language rules that govern the LLM’s behavior. Our initial tool design adopts this paradigm to support domain-experts to customize an LLM-simulated patient, and the final version extends it with a novel principle-adherence prompting pipeline.

In the mental health area, researchers have incorporated domain-expert feedback when prompting LLM for simulated patients, resulting in patients that use colloquial and resistant language (Chen et al., 2023; Stapleton et al., 2023) or that can take on diverse conversation styles, such as upset, verbose, or reserved (Wang et al., 2024). Such work, however, required a researcher-in-the-loop to refine prompts, hindering the speed of iterative design. Our work introduces a human-LLM collaborative tool that enables domain-experts to directly create and refine LLM roleplay simulation to faithfully resemble real-world patients.

**Text Generation with LLMs** Generating dialogue responses that adhere to user-defined principles is a type of constrained text generation problem. Recent work has shown that constrained text generation poses challenges when directly prompting GPT-4 (Madaan et al., 2024; Bubeck et al., 2023; Yao et al., 2023). To improve outputs, Madaan et al. (2024) propose a self-refine method and conduct evaluation experiments on a dialogue simulation task where responses are constrained by a general set of criteria such as relevance, consistency, informativeness, and helpfulness. A difference in our setting is responses are constrained by expert-defined principles that are multi-faceted and do not apply in all dialogue contexts. This necessitates new modules that breakdown principles into multiple, concise questions and check the applicability of principles prior to evaluating them.

## 3 Designing for Simulated Roleplay

We take a human-centered design approach to developing a tool for expert counselors to create and customize an AI patient for eventual use as a simulated training partner. After designing an initial version of our tool, we pilot test it with experienced peer counselors to understand any remaining challenges to effective human-LLM collaboration when creating and customizing an AI patient.

### 3.1 Initial Tool Design Rationale

Our initial version of Roleplay-doh adopted several of the design features of Petridis et al. (2024)’s tool for customizing task-oriented chatbots through interactive feedback.

**Principle Elicitation:** Counselors can manually write or edit the AI patient’s constitution. However, since users often struggle to formulate their thoughts into principles, our tool helps the counselor transform their feedback into specific principles to make principle writing easier. As counselors interact with an AI patient, for each generated response, they have the option to leave feedback in the form of a "kudos" explaining behavior they want to reinforce, a "critique" explaining any undesirable behavior, or a "rewrite" that demonstrates a more desirable response. Then an LLM is prompted (§E.1) to translate qualitative feedback into concrete principles that specify what should happen and when, and that generalize beyond the specifics of the dialogue context in which they are generated (Fig 1). Early testing revealed that GPT-

3.5 was sufficient at translating kudos and critique feedback into principles, while prompting GPT-4 to explain differences in initial and rewrite responses helped with inferring a principle.

**Testing Principles:** Likewise, to enable easier testing of principles, our tool supports rewinding the last response of the conversation, and generating a new response based on the updated AI Patient constitution. One feature that we change is generating a single dialogue response, rather than multiple responses, at a time. We reasoned that counselors can identify ways in which a response does not resemble a real-patient’s without needing to see multiple, and that generating a response at a time would make the testing process more manageable and similar to having a normal dialogue.

**Simulating AI Patient:** We prompt the LLM to follow the most recent set of constitution principles as in [Petridis et al. \(2024\)](#) rather fine-tuning the LLM weights as in [Bai et al. \(2022\)](#)’s constitutional AI framework. Since the tool supports defining and testing principles in an iterative fashion, prompting can make steering model behavior quicker and less expensive. Our prompt ([Appendix E.2](#)) instructs GPT-4 to simulate a patient’s next response in a dialogue as opposed to asking the LLM to roleplay as the patient using a system prompt ([Zhou et al., 2024](#)), as early testing revealed that this can mitigate role consistency issues in which the LLM responds as an AI assistant rather than as a patient.

### 3.2 Pilot Testing

We pilot tested the tool with 6 counselors who had experience giving support to real patients on an online peer support platform; refer to [Appendix A](#) and [B](#) for participant backgrounds and the pilot procedure. Additionally, four of the co-authors each conversed with four AI patients created ([Table 3](#)) and assessed how well the simulation adhered to the expert-defined principles; refer to [Appendix C](#) for details on the procedure and qualifications of the co-authors. Overall, the pilot tests and principle-adherence analysis helped uncover two obstacles to effective simulated roleplay.

**O1: Defining “realistic” patient behavior is ambiguous** Counselors felt the tool was easy to use and effective at guiding the AI patient’s behavior, as indicated by moderate to high agreement scores on a tool usage questionnaire as shown in [Table 4](#) in [Appendix B](#). However, the task of creating a ‘realistic’ AI patient for an imagined scenario was confus-

ing, as counselors have interacted with many types of patients who respond in various, yet equally realistic ways. This insight helped us re-frame the task in later sessions as recreating a challenging scenario from one’s past, which removed the ambiguity of what behaviors are realistic by having them refer to a specific case from memory.

**O2: 20% of responses produced by GPT-4 do not satisfy expert principles or dialogue conventions.** Specifically, 20% (55/276) of cases were rated as moderately, slightly, or not at all satisfying, at following all principles and being appropriate to the dialogue context. Further analysis of these cases helped to uncover three sources of error. **Not satisfying multiple principles at once:** Generated responses could struggle to follow all the principles when there was a large number of principles, or when the provided principles were a complex composition of simpler principles. **Awkwardness for Dialogue Context:** Some responses were also identified as awkward or unnatural given conventions in the dialogue context, despite not violating the defined principles. For example, in the middle of a conversation, saying “Hi, A. Yes that’s exactly what I mean. There’s a voice that is always critical of myself” is unnatural because of the use of ‘Hi’. **Misapplying Situational Principles:** While generating a response, the model sometimes incorrectly applied principles, such as *Respond with hesitancy when someone gives you encouraging words*, even when the conditions for their use—receiving encouraging words—were not met.

## 4 Roleplay-doh

Roleplay-doh helps counseling experts create customized AI patients based on scenarios from their past experiences. Roleplay-doh uses LLMs in two ways: *Principle Elicitation* and *Response Generation with Principle-Adherence*, which we describe in more detail below:

**Principle Elicitation** Roleplay-doh enables counselors to customize an AI patient to resemble a real-patient case by eliciting their qualitative feedback and transforming it into constitution principles that dictate behavior. We provide some examples of principles defined by expert counselors in [Table 2](#). Since our initial tool design includes the principle elicitation features, we refer the reader to [§3.1](#) for details.



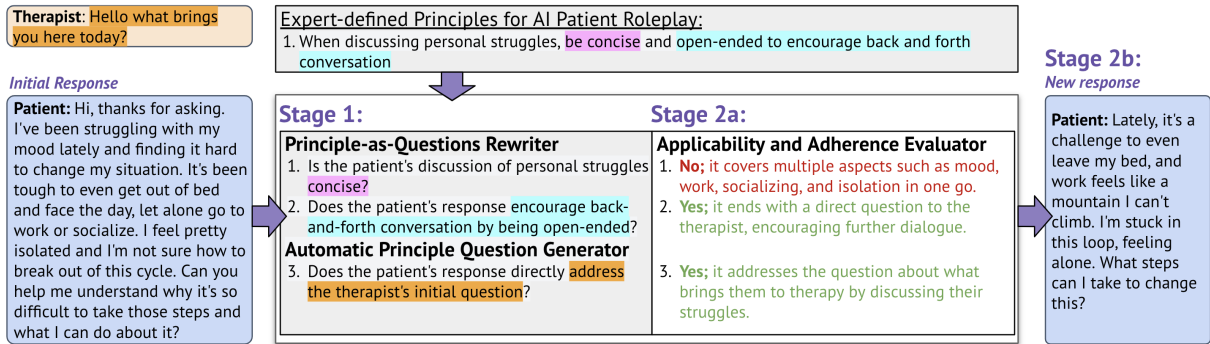


Figure 2: Principle-adherence prompting pipeline for mitigating errors in satisfying expert principles and dialogue conventions. In Stage 1, expert-defined principles are rewritten into several Yes/No questions; and the LLM generates additional principle questions that are relevant to ensure adherence to dialogue conventions such as coherence and consistency. In Stage 2, the LLM (a) evaluates whether the questions are applicable to the context and the answers to the principle-adherence questions; and (b) refines the response to ideally receive Yes on all questions.

**Generation with Principle-Adherence** We prompt GPT-4 conditioned on patient description, list of principles and conversation history to generate an initial patient response at each conversation turn. Since initial patient responses can fail in 20% of cases to satisfy expert principles or dialogue conventions, we propose a principle-adherence pipeline that prompts the LLM to generate principle-adherence questions (Stage 1) and employs these questions to assess and refine the initial patient response (Stage 2). Our principle-adherence pipeline features three modules to mitigate the identified issues in §3.2.

**Principle-as-Questions Rewriter:** This module transforms each expert-defined principle into a set of concise yes/no questions that are easier to evaluate for principle-following. Multifaceted principles (e.g. “*You should respond in short sentences and avoid using terms like ‘anxious’*”), are divided into separate questions (e.g. “*Does the patient’s response employ short sentences?*” and “*Is the patient’s language devoid of terms like ‘anxious’?*”).

**Automatic Principle Generator:** This module adds additional principle questions that capture criteria essential for ensuring that the LLM simulation’s responses follow general dialogue conventions, such as coherence and consistency. This helps correct cases where there is awkwardness in the generated responses not captured by the defined principles. The LLM is instructed not to make assumptions about the patient or therapist’s personality when generating criteria: for example, “*The patient should be appreciative of the therapist’s help*” is not an appropriate criterion.

**Applicability and Adherence Evaluator:** This module determines if each principle is applicable

in a given situation, returning N/A if the question is not relevant to answer; otherwise, it evaluates the response using the questions, returning Yes if the response adheres to the principle questions; and No otherwise. For an example of situational applicability, the principle *Show willingness to engage in a suggested activity by affirming the proposal* is evaluated only if the therapist suggests an activity. In situations where the therapist is asking something else and no activity is proposed, the module would appropriately return N/A recognizing that the principle does not apply.

Our pipeline first uses the **principle-as-questions rewriter** and **automatic principle generator** modules to generate a set of criteria for evaluating the initial generated response. Then, the response is evaluated using the question by the **applicability and adherence evaluator**. If the model returns a “No” response for any of the questions, we then perform a rewrite of the response conditioned on the evaluation results, that ideally passes all questions (Fig 2). We detail the prompts used and the procedure used to develop the prompts (§E.3) and the results of a performance evaluation against ablations (§6).

## 5 User Study using Roleplay-doh

To evaluate how Roleplay-doh can aid counseling experts in creating AI patients, we conducted a within-subjects study with 25 counseling experts, comparing: (1) a *Scenario-only* dialogue simulation, where the counselor writes a patient scenario description, and (2) a *Scenario+Expert-principles* simulation, where the counselor uses Roleplay-doh to define principles. See §G for full study setup.

We evaluate the AI patients created by coun-

sors on criteria inspired by prior work evaluating Standardized Patients, who are trained human actors, on their ability to roleplay a case (Himmelbauer et al., 2018). Counselors rated the two AI patients based on 6 dimensions (Table 5). We also surveyed each counselor about their experience using the tool for defining principles. Following Petridis et al. (2024), we include four measures for evaluating principle elicitation features (Table 6).

We recruit 25 counseling experts with real-world experience in mental health support to perform the evaluation, categorized by their primary expertise: 1) those who are pursuing or have completed degrees in counseling or clinical psychology with practicum experience; 2) those who provided online counseling to over 30 clients on the 7 Cups platform; and 3) peer counselors who have provided in-person or virtual support.

### 5.1 Creator Perceptions

The AI patients prompted with *Scenario+Expert Principles* were rated significantly higher than *Scenario-Only* on all measures except for role consistency, for which both methods score highly (Table 1). Counselors mentioned the *Scenario-Only* AI patient **lacked emotional depth in expression**. As one noted, "*patients don't state a feeling such as 'I feel hopeless'. They display their current emotional state in their manner of speech.*" *Scenario-only* was also **too articulate and forthcoming** when describing issues, where encouraging real patients to share is "*as challenging as pulling teeth*". It was characterized as **too cooperative**, too willing to accept. Despite counselors writing behavioral traits such as "*not talkative*" and "*reluctant*" in the patient scenario, *Scenario-only* did not exhibit these behaviors.

### 5.2 Creating Principles with Roleplay-doh

Across the 25 *Scenario+ExpertPrinciple* AI patients, 123 total principles were created (min=1, max=10, median=5). Two authors did a qualitative coding of these principles following a thematic analysis approach (Braun and Clarke, 2006) where codes were initially defined and revised during the process. Besides **stage-agnostic** themes dictating a **concise** (14 patients) and **colloquial** (7 patients) speaking style, counselors created principles related to the stages of an emotional support conversation (Liu et al., 2021): 1) **exploration**: identifying the patient's problems, 2) **comforting**: using empathy and understanding to comfort the

patient, and 3) **action**: formulating solutions to the patient's problems. For instance, we find a common theme of instructing the AI patient to **show initial skepticism with the idea of seeking help** (14 patients), corresponding to the style of interaction in the **exploration** stage of conversation. Table 2 provides a full list of principle themes, examples, and corresponding conversation stages.

While we observe overlaps in the types of principles defined, we also observe some contradictory themes. For example, the call for being **disorganized and conflicted** (9 patients) contrasts calls to make responses **concise and direct** (14 patients). In the **action** stage of conversation, several counselors added principles to make the AI patient **proactively ask for advice** (12 patients); nonetheless, other counselors added an opposing principle to **not seek out solutions** but rather just share their thoughts and feelings (3 patients). These opposing principles highlights the need for different principles to describe diverse conversation behavior and styles, challenging the notion of defining AI patients based on a single set of principles.

**Tool User Experience** Counselors found the tool helpful for writing principles that **effectively guided** the AI patient to recreate their past case ( $\mu = 6.04, \sigma = 1.06$ ). With the tool, most found it **easy** to convert their thoughts and feedback on the AI patient's behavior into principles ( $\mu = 6.12, \sigma = 1.13$ ). Counselors felt they could **efficiently** write principles ( $\mu = 6.3, \sigma = 1.29$ ), without requiring much **mental demand** ( $\mu = 3.20, \sigma = 1.70$ ). Many counselors liked how the tools "*organized their thoughts into rules*", without "*needing to word it perfectly*." Yet, principle-elicitation did not work perfectly in all cases: 11.4% of principles required manually editing. Via a worse-case analysis of creators' tool use, we uncover scenarios where Roleplay-doh's human-LLM collaboration pipeline can still be improved (§I).

### 5.3 Third-Party Comparison

A limitation of our creator study (§5.1) is the potential bias from creators who knew which AI patient embodied their principles. To address this, we conducted a third-party study where external counselors served as impartial judges. These judges evaluated AI patient transcripts presented in randomized order to ensure blindness to the condition. We invited five counselors from the creator study to serve as judges, all equally qualified of assessing AI

Ratings by Counselor Creators			Ratings by Third-Party Counselors		
Measure	Scenario Only	+ Principles	Measure	Scenario Only	+ Principles
Authenticity	5.24	+0.80 **	Authenticity	5.32	+0.31 *
Stayed in Role	6.32	+0.08	Stayed in Role	6.29	+0.09
<i>Resembled Past Case</i>	4.80	+0.76 *	<i>Resembled Typical Case</i>	4.91	+0.49 **
<i>Mirrored Challenging Aspects</i>	4.52	+1.00 *	<i>Challenged the Counselor</i>	2.13	+0.22
Ready as Training Partner	5.16	+0.64 *	Ready as Training Partner	5.05	+0.39 **
Recommend to Novices	5.76	+0.52 *	Recommend to Novices	5.03	+0.38 *

Table 1: Creators and third-party counselors compared the *Scenario-Only* vs. *Scenario+ExpertPrinciples* AI patients using 7-point Likert-scale measures; third-party judges were asked identical measures when possible, with two measures modified to match the external perspective. **Creator Ratings:** Creators (N=25) rated both AI patients. After refining the AI patient simulation with principles, creators rate the patient significantly higher on all measures except for *stayed in role*, for which both AI patients score highly. **Third-Party Ratings:** Third-party counselors (N=5) provided 125 total comparisons of the two AI patient versions. The treatment effect of adding expert principles was estimated using the following linear mixed-effect model:  $\text{Rating} \sim \text{Treatment} + \text{CreatorID} + (1 | \text{AnnotatorID})$ . Third-party counselors rate AI patients with principles significantly higher on 4 of the 6 measures. (\*\*\*:  $p < .001$ , \*\*:  $p < 0.01$ , \*:  $p < 0.05$ .)

Stages	# AI patients	Theme	Example Principle
	14	Keep responses concise and do not share too much.	When discussing personal struggles, be more concise and open-ended to encourage a back-and-forth conversation.
	7	Use colloquial and realistic language.	Incorporate natural speech patterns, improper grammar and punctuation, including the use of slang and less structured sentences, to convey a more authentic and relatable character.
	14	Show initial mistrust and hesitation with the idea of seeking help.	When expressing feelings of overwhelm and doubt, provide limited information and express skepticism towards the effectiveness of seeking help.
	19	Show emotions in detail, elaborating with examples as needed. *	When describing personal struggles, provide specific details and symptoms to help the listener understand the situation better.
	9	Be less self-aware of emotions, thoughts, and needs. Articulate thoughts in a more disorganized way.	When expressing reluctance or uncertainty about seeking help or accepting praise, it's important to convey the internal struggle and conflicting emotions, rather than presenting a clear-cut decision or emotion.
	3	Do not seek out solutions, but rather just share thoughts and feelings. *	When expressing feelings of being stuck or defeated, focus on sharing emotions rather than seeking a resolution.
	12	Proactively seek out solutions and show reflective insight over time. *	When discussing personal struggles, provide reflective insights into your situation and propose actionable steps for improvement to continue the conversation effectively.

Table 2: Themes taken from qualitative analysis of principles and representative examples. We discover several novel (\*) principles compared to those defined in prior work on AI patients (Chen et al., 2023; Stapleton et al., 2023). Themes are categorized into stages of conversation taken from (Liu et al., 2021): **exploration**, **comforting**, and **action**; those relating to the overall conversation are categorized as **stage-agnostic**.

patient realism. A power analysis confirmed that five judges would provide 80% statistical power (Appendix §J.2). The third-party counselors rated the same six dimensions as the creator study, with questions reworded for the perspective of external judge (Appendix §J.1).

Third-party judges rate AI Patients with expert-defined principles as more authentic, resembling typical cases, ready as a training partner, and likely recommend to novices (Table 1). However, when compared to the creator study results, the increase in ratings is smaller from the perspective of third-party counselors. We explore the reasons for this smaller difference in Appendix J.3, finding that there is lower agreement on which AI patient is preferred (Table 8), due to the different principles attended to by third-party counselors and the specific principles added by the creator.

## 6 Evaluation of Principle-Adherence

We now evaluate whether the principle-adherence pipeline improves the quality of responses for Roleplay-doh, along with an ablation analysis showcasing the utility of its various components. Specifically, we break down the evaluation of model responses along three metrics: **M1**) Are they consistent with the patient description and conversation history? **M2**) Do they exhibit an awkward style of speech? **M3**) Do they adhere to the provided principles?

We evaluate the performance of our principle-adherence pipeline [Full] over (1) GPT-4 response generation without our pipeline [No Critique]; (2) an ablation without the **Principle-as-Questions Rewriter** [No Principle Rewrites]; (3) an ablation without the **Automatic Principle Generator** [No Autogenerated Criteria]; and (4) an im-

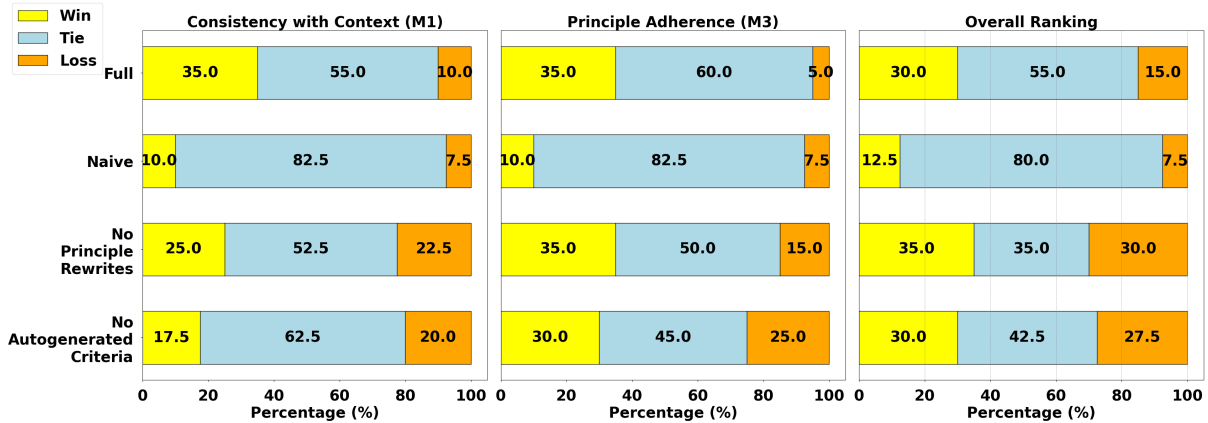


Figure 3: Win/Tie/Loss for the Error Test Cases along **Consistency with Context (M1)**, **Principle Adherence (M3)**, and **Overall**. Pairwise preference evaluation results with [No Critique] as a baseline. Results obtained after majority voting.

plementation of the principle-adherence pipeline that does not have any of these modules [Naive].

To analyze how the pipeline mitigates errors that arise in base GPT-4 generations, we select 40 conversation turns from our user study logs that fall into one of the error categories described in §3.2 as testcases. Each testcase contains the scenario, conversation history up to that point, and the expert-defined principles for the AI patient. For each test case, responses are generated for all models and then ranked by expert counselors from 1 (best) to 5 (worst) for metrics **M1** and **M3**, along with "Yes" or "No" annotations for **M2**. Finally, experts provide an **Overall** ranking, along with a brief textual explanation. We allow multiple responses to have the same rank and randomize order of responses to minimize positional bias (details in §L).

We treat [No Critique] as our baseline and compare all other models to it. We report preference results based on majority vote across 3 expert counselor annotations (Fig 3). We find our [Full] method performs better than [No Critique] on **M1** (W: 35%; L 10%) and on **M3** (W: 35%; L 5%), where it has the highest win/loss rates compared to all ablations. On overall rankings, it again has the strongest performance (W: 30%; L 15%). We find that the performance of [Full] compared to [No Critique] is weaker on **Overall** than **M1** and **M3**. This is because the annotators often used their own subjective judgements (e.g., "although the middle response ranked third on principle following, it feels like the most realistic response in this scenario") to perform the overall ranking, resulting in unpredictable and subjective results. We also find that [Naive] has a disproportionately high tie rate across metrics, indicating

that it rarely produces better responses even for error cases. This highlights the importance of the **Principle-as-Questions Rewriter** and **Automatic Principle Generator** for improving responses.

For **M2**, after majority voting, annotators report that 2.5% of responses are awkward for the [Full] method, as compared to 15% for [No Critique], 7.5% for [Naive], 7.5% for [No Principle Rewrites] and 15% for [No Autogenerated Criteria]. Therefore, our principle adherence pipeline substantially reduces the occurrence of awkward style in responses (by a margin of 12.5%). The 12.5% gap in percentage of awkward responses between [Full] and [No Autogenerated Criteria] also indicates the importance of the **Automatic Principle Generator** for producing realistic rewrites. We repeat these experiments with 50 randomly picked conversation turns and report results in §K, along with Krippendorff’s  $\alpha$  numbers.

## 7 Conclusions

This paper introduces Roleplay-doh, a tool that empowers domain experts to create LLM simulations through the automatic conversion of expert feedback into natural language principles, and validates the tool for the task of creating AI patients that serve as roleplay partners for novice counselors. Roleplay-doh’s novel principle-adherence pipeline also addresses gaps in existing simulation methods by reducing the prevalence of responses that do not follow expert-defined principles or dialogue conventions. Studies with mental health counselors creating and comparing AI patients demonstrate that Roleplay-doh allows experts to refine LLM simulators to be authentic and more ready as practice partners. Roleplay-doh could be generalized



to support domain-experts in creating realistic simulations in other social dialogue domains, such as roleplay practice for teaching, coaching, conflict resolution, and negotiations, as future work.

## Limitations

One limitation of our study is the intended use case of the AI patients created by counselors. These AI patients were meant to recreate challenging cases that might be useful for the education of "first-year" or novice counselor. In other words, we intentionally restricted some diversity in patient scenarios by focusing on this use case. Readers should keep this limitation in mind prior to generalizing our analysis of principles. Moreover, due to the time and resource constraints of our creator study, we required counselors to stop providing feedback before their conversation with the AI patient had naturally ended. As such, the principles that counselors added may not have addressed all underlying issues of the AI patients they interacted with. Future work that uses the list of user-generated principles should be mindful of their non-exhaustive nature before adopting them.

In this paper, we focused on enabling counselors to create AI patients that can simulate realistic interactions via *text-based dialogues*. However, we acknowledge that text-based interaction has its limitations for training. Professional psychotherapists may gain useful information from the tone, facial expression, posture, and other non-verbal behaviors of their patients, which better help them empathize and support patients. This is a limitation of our current AI patients and online, text-based, mental health counseling in general. With the rapid development of multimodal models, future works may have the opportunity to explore creating realistic AI patients in other modalities that better match the modality within which a counselor will eventually support patients.

## Ethics Statement

This study was approved by our institution's Institutional Review Board (IRB). All investigators in the study completed the responsible code of conduct in research training. We have compensated domain experts at a minimum rate of \$25 per hour, going beyond the minimum wage in the United States.

We are optimistic about the potential benefit that our AI patients can bring to the fields of counseling and psychotherapy. At the same time, we

solicited feedback from counselors about any potential concerns regarding the AI patients. During these interviews, some counselors emphasized the irreplaceability of peer-to-peer roleplay with humans during training, due to the unique opportunity it provides for novice counselors to connect with others, especially for online counseling platforms where counselors are often isolated from one another. To preserve human-to-human interactions, future work requires a participatory design approach before attempting to integrate AI patients into existing practices and learning environments.

Our hope is that interactions with AI patients can glean important lessons that help counselors go from simulation into the real-world. Nonetheless, a risk with simulation is that counselors can become overconfident in supporting a AI patient, but may not effectively support patients with real mental health concerns. As part of a larger curriculum, AI patients should be just one tool for practicing these skills. Real counselors and therapists should be able to take on real patients only after passing traditional certifications and background checks.

It is impossible to promise that all interactions with an LLM such as GPT-4 result in satisfactory responses. Therefore, meaningless, derogatory, and otherwise harmful responses may also be generated and cause unwanted effects on users. While our principle-adherence pipeline is a potential inference-time solution to mitigate such harmful responses, we must acknowledge this possibility, especially due to the stochastic nature of LLM. Users should be advised about these potential side effects before using the system in any scenario. For our experiment, we created consent forms that made sure counselors were aware of these issues.

## References

- 7Cups. 2024. [7 cups: Free online therapist and counseling.](#)
- Guillaume Alinier and Denis Oriot. 2022. [Simulation-based education: deceiving learners with good intent.](#) *Advances in Simulation*, 7(1):8.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022. [Constitutional ai: Harmlessness from ai feedback.](#) *arXiv preprint arXiv:2212.08073*.
- Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. [Fitting linear mixed-effects models](#)

- using lme4. *Journal of Statistical Software*, 67(1):1–48.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024. Graph of thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690.
- Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative Research in Psychology*, 3(2):77–101.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Siyuan Chen, Mengyue Wu, Kenny Q Zhu, Kunyao Lan, Zhiling Zhang, and Lyuchun Cui. 2023. Llm-empowered chatbots for psychiatrist and patient simulation: application and evaluation. *arXiv preprint arXiv:2305.13614*.
- Myra Cheng, Tiziano Piccardi, and Diyi Yang. 2023. CoMPoSIT: Characterizing and evaluating caricature in LLM simulations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10853–10875, Singapore. Association for Computational Linguistics.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Orianna Demasi, Yu Li, and Zhou Yu. 2020. A multi-persona chatbot for hotline counselor training. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3623–3636, Online. Association for Computational Linguistics.
- Peter Green and Catriona J MacLeod. 2016. Simr: An R package for power analysis of generalized linear mixed models by simulation. *Methods in Ecology and Evolution*, 7(4):493–498.
- Monika Himmelbauer, Tamara Seitz, Charles Seidman, and Henriette Löffler-Stastka. 2018. Standardized patients in psychiatry—the best way to learn clinical skills? *BMC medical education*, 18:1–6.
- Tae Soo Kim, Yoonjoo Lee, Jamin Shin, Young-Ho Kim, and Juho Kim. 2024. Evallm: Interactive evaluation of large language model prompts on user-defined criteria. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, CHI ’24. ACM.
- Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. Towards emotional support dialog systems. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3469–3483, Online. Association for Computational Linguistics.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2024. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36.
- Julia M. Markel, Steven G. Opferman, James A. Landay, and Chris Piech. 2023. Gpteach: Interactive training with gpt-based students. In *Proceedings of the Tenth ACM Conference on Learning @ Scale, L@S ’23*, page 226–236, New York, NY, USA. Association for Computing Machinery.
- Julia Othlinghaus-Wulhorst and H. Ulrich Hoppe. 2020. A technical and conceptual framework for serious role-playing games in the area of social skill training. *Frontiers in Computer Science*, 2.
- Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST ’23, New York, NY, USA. Association for Computing Machinery.
- Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2022. Social simulacra: Creating populated prototypes for social computing systems. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*, UIST ’22, New York, NY, USA. Association for Computing Machinery.
- Savvas Petridis, Benjamin D Wedin, James Wexler, Mahima Pushkarna, Aaron Donsbach, Nitesh Goyal, Carrie J Cai, and Michael Terry. 2024. Constitutionmaker: Interactively critiquing large language models by converting feedback into principles. In *Proceedings of the 29th International Conference on Intelligent User Interfaces*, IUI ’24, page 853–868, New York, NY, USA. Association for Computing Machinery.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.
- Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit Bansal, Jason Weston, and Xian Li. 2023. Branch-solve-merge improves large language model evaluation and generation. *arXiv preprint arXiv:2310.15123*.

- Omar Shaikh, Valentino Emil Chai, Michele Gelfand, Diyi Yang, and Michael S. Bernstein. 2024. [Rehearsal: Simulating conflict to teach conflict resolution](#). In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, CHI '24, New York, NY, USA. Association for Computing Machinery.
- Weiyang Shi, Emily Dinan, Kurt Shuster, Jason Weston, and Jing Xu. 2022. [When life gives you lemons, make cherryade: Converting feedback from bad responses into good labels](#). *arXiv preprint arXiv:2210.15893*.
- Logan Stapleton, Jordan Taylor, Sarah Fox, Tongshuang Wu, and Haiyi Zhu. 2023. [Seeing seeds beyond weeds: Green teaming generative ai for beneficial uses](#). *arXiv preprint arXiv:2306.03097*.
- Michael J Tanana, Christina S Soma, Vivek Srikumar, David C Atkins, and Zac E Imel. 2019. [Development and evaluation of clientbot: Patient-like conversational agent to train basic counseling skills](#). *Journal of medical Internet research*, 21(7):e12529.
- Ruiyi Wang, Stephanie Milani, Jamie C Chiu, Shaun M Eack, Travis Labrum, Samuel M Murphy, Nev Jones, Kate Hardy, Hong Shen, Fei Fang, et al. 2024. [Patient-{\Psi}: Using large language models to simulate patients for training mental health professionals](#). *arXiv preprint arXiv:2405.19660*.
- Diyi Yang, Caleb Ziems, William Held, Omar Shaikh, Michael S Bernstein, and John Mitchell. 2024. [Social skill training with large language models](#). *arXiv preprint arXiv:2404.04204*.
- Shunyu Yao, Howard Chen, Austin W Hanjie, Runzhe Yang, and Karthik Narasimhan. 2023. [Collie: Systematic construction of constrained text generation tasks](#). *arXiv preprint arXiv:2307.08689*.
- Zheng Yao, Haiyi Zhu, and Robert E. Kraut. 2022. [Learning to become a volunteer counselor: Lessons from a peer-to-peer mental health community](#). *Proc. ACM Hum.-Comput. Interact.*, 6(CSCW2).
- J.D. Zamfirescu-Pereira, Richmond Y. Wong, Bjoern Hartmann, and Qian Yang. 2023. [Why johnny can't prompt: How non-ai experts try \(and fail\) to design llm prompts](#). CHI '23, New York, NY, USA. Association for Computing Machinery.
- Xuhui Zhou, Zhe Su, Tiwalayo Eisape, Hyunwoo Kim, and Maarten Sap. 2024. [Is this the real life? is this just fantasy? the misleading success of simulating social interactions with llms](#). *arXiv preprint arXiv:2403.05020*.

## A Background of User Participants

Counselors with real-world experience in mental health support were recruited for our pilot tests, creator studies, and technical evaluations of the principle-adherence pipeline. We present more detailed information about how they were recruited, and their background.

After receiving permission from the 7 Cups platform (7Cups, 2024) for our IRB-approved study, we recruited online peer counselors from the 7 Cups platform (7Cups, 2024). Participants were required to be 18 yrs or older, from the United States, and to have had experience giving support to 30+ members on the online site. Our pilot tests were conducted exclusively with 6 online-peer counselors from 7 cups; our creator study used an additional 4 peer counselors from the 7 cups site.

We involved another 15 counselors from the Upwork platform. Participants were required to be 18 yrs or older, from the United States, and to have had education in counseling or psychotherapy and/or have given extensive counseling support (either via text, phone, in-person). A sampling of counselors backgrounds included *licensed mental health therapist with over 20 years of experience*, *a Master's of Science in Rehabilitation and Mental Health Counseling, 25 years as the clinical director of a busy crisis agency*, and *a mental health advocate who has personally helped coach dozens of got students via a peer support role*.

Finally, we involved 3 peer counselors and 3 psychotherapy PhD students in our university. Specifically, peer counselors had 3-months of classroom training and had experiences giving voluntary support to university students; the clinical psychology PhD students were 4th year students with 3 years experience providing psychotherapy support to clients under the supervision of a licensed psychotherapist.

## B Pilot Testing with Expert Counselors

After developing an initial version of our tool, we piloted it with 6 counselors during a 90 minute session. Participants were tasked with creating different AI patient scenarios, and using Roleplay-doh to interactively refine the simulation by defining principles. All participants started with a common roleplay scenario called "loneliness after work". They proceeded to use the tool to chat, give feedback, and convert their feedback into principles to shape the AI Patient's behavior. If time allowed,

they created and customized an additional AI patient based on scenarios they chose to write. Pilot Participant 1 (PP1), PP2, and PP5 had time to create one additional AI patient; PP3 created two additional AI patients.

Four AI patients from the pilot studies were selected for the second stage of pilot testing, where selection encouraged scenario and principle diversity (Table 3). Four of the co-authors conversed with each of the AI patients (stimulated by directly prompting GPT-4 with the scenario and principles). Co-authors assessed how well the patient's responses adhered to the expert-defined principles and dialogue conventions.

### C Evaluating principle-adherence of GPT-4 direct prompting

We aim to determine how often directly prompting GPT-4 to produce less satisfying responses given fixed constitution principles.

**Procedure:** We selected 4 AI patients that were created in the design sessions by different counselors. Four co-authors had practice conversations with each of the four AI patients, resulting in 16 conversations. Each response in each conversation was rated on a 5-point likert scale on how well the generated response adhered to principles and how appropriate they were for the dialogue content (5 = Completely, 1 = Not at all). From the 16 completed conversations, the mean number of responses per conversation was 17.25, with a minimum of 12 and maximum of 22. In total, 276 responses were given satisfaction ratings. Since each co-author created a different conversation from each of the AI patients, each response was only scored by one co-author.

**Participant Rationale:** During this pilot principle-adherence experiment, we used co-authors to generate test conversations because our basic counseling skill-level is representative of the eventual use-case of untrained, novice counselors interacting with AI Patients. For the annotation task, a human annotator is qualified if they can judge whether a response follows the principles defined by expert counselors, and is appropriate in the conversation context. Since these skills do not require counseling expertise, the co-authors are qualified to do this annotation task.

### D Roleplay-doh Interface for Making Constitutional Principles for LLM Simulation

The final version of Roleplay-doh (Fig 4) generates responses in the LLM simulation using a principle-adherence pipeline. In addition to this core improvement, we made several minor improvements to improve the usability and user experience of the tool.

Improvements to the usability of the UI

- Fixing a bug where a user who clicks "save" multiple times will submit duplicate feedback, resulting in duplicate sets of principles
- Making converting feedback to principles easier by placing a "Convert" button next to each feedback box, rather than a single "Convert" button at the top of the screen which users would forget about

### E LLM Prompts

In this section, we detail the prompts we used for the different components of Roleplay-doh.

#### E.1 Principle Elicitation Prompts

In this section, we provide the prompts used in the principle elicitation module of Roleplay-doh. These prompts were arrived at after a substantial amount of testing using a development set. Each prompt uses the same structure, which is inspired by Markdown formatting. There is an initial instruction that provides a system prompt, along with a description of the principle elicitation task. This is followed by a one-shot example of an elicited principle as a result of the task, and the relevant input, including the conversation history. All parts of the prompt are demarcated by headers in Markdown formatting, and the outputs are returned in JSON format. We describe each prompt in greater detail in the relevant sections.

The kudos and critique prompts were given to the gpt-3.5-turbo-1106 model. The rewrite prompt was given to the gpt-4-turbo-1106 model. For all API calls to the principle-elicitation prompts, the temperature was set to 0.1.

##### E.1.1 Principle Elicitation Kudos Prompt

This prompt includes a desirable response, as well as some reasoning for why the response is desirable. This information is then used to create a general



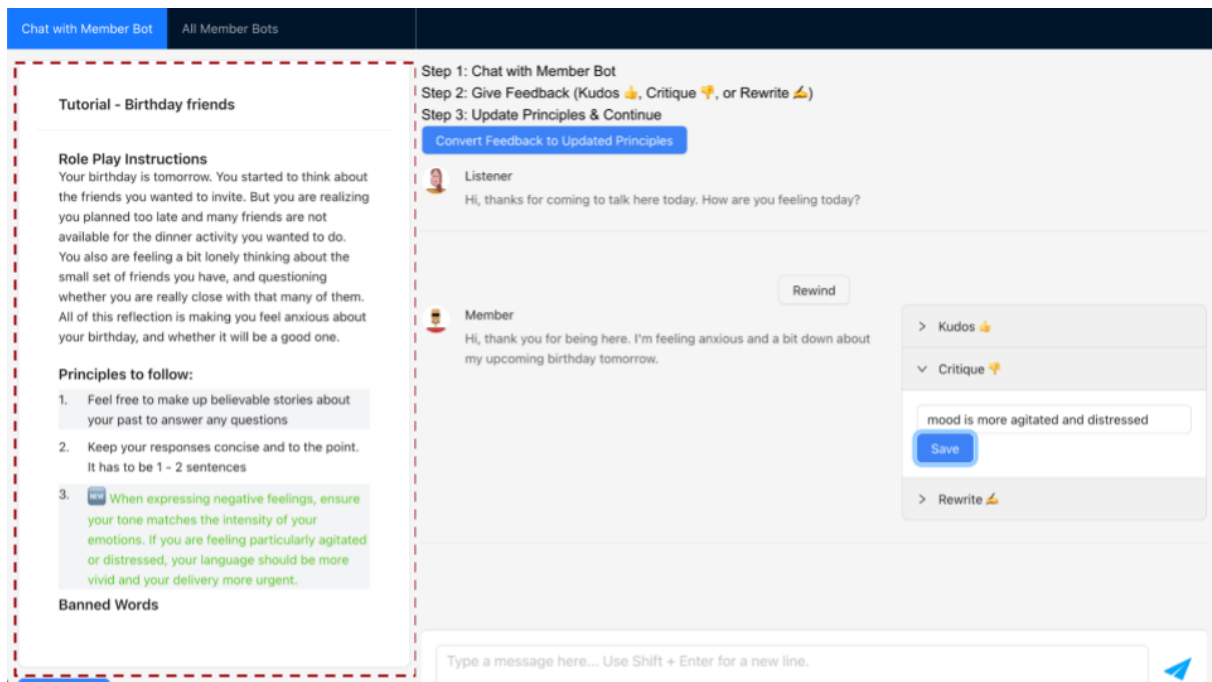


Figure 4: Roleplay-doh allows users to chat with a AI patient, Provide Feedback as a Kudos/Critique/Rewrite, and Convert Feedback into Principles, which in turn shape the roleplay behavior.

principle that would result in a similar response in the same situation.

### Instruction :

You are a superintelligent AI capable of understanding human emotion. You will review praise for an actor's dialogue, and synthesize a well-written principle that, when followed, would help the actor continue generating high-quality dialogue. To accomplish this, you have been given a conversation script with the actor's desirable response, as well as a specific explanation for why this response is desirable. You will output a final principle that the actor can follow to be more realistic. Follow the following guidelines:

1. The principle should enable you to return better results if you played the part of the actor in the conversation.
2. Return only a JSON response in the format provided.

### Input:

### Conversation Script

Helper: Is there anything else you want to share with me?

Actor: Yea so lately I've been really losing sleep.

Actor: There's a lot on my plate, and my energy has been so low. I think I am failing a lot of people.

Helper: You are absolutely not failing people. You are a great person, and you should remember that you are

very capable and energetic.

### Desirable response from the actor

Actor: I don't know.... Am I really?

### Specific explanation for why the response is desirable

The actor is hesitant to agree with the helper and shows self-doubt. This is consistent with the conversation history.

### Response:

```
{"result": {"principle": "When someone gives you encouraging words, you respond with hesitancy, doubting the significance of that positive perspective."}}
```

### Input:

### Conversation Script

```
{conversation_script}
```

### Desirable response from the actor

Actor: {actors\_response}

### Specific explanation for why the response is desirable

```
{kudos_rationale}
```

### Response:

### E.1.2 Principle Elicitation Critique Prompt

This prompt includes an undesirable response, as well as some reasoning for why the response is undesirable. This information is then used to create a general principle that would result in a similar

response not being generated after the same conversation history.

### Instruction:

You are a superintelligent AI capable of understanding human emotion. You will review critiques of an actor's dialogue, and synthesize a well-written principle that, when followed, would help the actor resolve the critiques.

To accomplish this, you have been given a conversation script with the actor's undesirable response, as well as a specific explanation for why this response is undesirable. You will output a final principle that the actor can follow to be more realistic. Follow the following guidelines:

1. The principle can contain examples of rewrites as well.
2. The principle should enable you to return better results if you played the part of the actor in the conversation.
3. Return only a JSON response in the format provided.

### Input:

### Conversation Script

Helper: Is there anything else you want to share with me?

Actor: Yea so lately I've been really losing sleep.

Actor: There's a lot on my plate, and my energy has been so low. I think I am failing a lot of people.

Helper: You are absolutely not failing people. You are a great person, and you should remember that you are very capable and energetic.

### Undesirable response from the actor

Actor: Thank you for reminding me of this. I am a great person, and I've proved myself to be very capable and energetic. I feel a lot better now due to your kind words.

### Specific explanation for why the response is undesirable

The actor should not be so quick to agree with the helper. Overly positive comments to cheer a patient up does not immediately work.

### Response:

```
{"result": {"principle": "When someone gives you encouraging words, you respond with hesitancy, doubting the significance of that positive perspective."}}
```

### Input:

### Conversation Script  
{conversation\_script}

### Undesirable response from the actor

Actor: {actors\_response}

### Specific explanation for why the response is undesirable  
{critique\_rationale}

### Response:

### E.1.3 Principle Elicitation Rewrite Prompt

This prompt includes an undesirable response, as well as a desirable rewrite of the undesirable response. The model first outputs a description that captures the difference between the desirable and undesirable response. It then uses this difference to output a general principle that would result in the desirable response given the same conversation history.

### Instruction:

You are a superintelligent AI capable of understanding human emotion. You have been given a conversation script with an actor's undesirable response, as well as a desirable rewrite for the response. You will output a well-written principle that, when followed, would help the actor generate more realistic responses that are closer to the rewrite. Follow the following guidelines:

1. The principle should capture the key differences that made the rewrite more realistic than the original response.
2. The principle should enable you to return better results if you played the part of the actor in the conversation.
3. Return only a JSON response in the format provided.

### Input:

### Conversation Script

Helper: Is there anything else you want to share with me?

Actor: Yea so lately I've been really losing sleep.

Actor: There's a lot on my plate, and my energy has been so low. I think I am failing a lot of people.

Helper: You are absolutely not failing people. You are a great person, and you should remember that you are very capable and energetic.

### Undesirable response from the actor

Actor: Thank you for reminding me of this. I am a great person, and I've proved myself to be very capable and energetic. I feel a lot better now due to your kind words.

### Desirable rewrite

Actor: I don't know... Am I really a great person?

### Response:

```

{"result":{
  "difference": "The desirable rewrite
    is different because it makes the
    actor more hesitant to adopt
    positive thoughts, where they show
    self-doubt",
  "principle": "When someone gives you
    encouraging words, you respond
    with hesitancy, doubting the
    significance of that positive
    perspective."}}

```

```

### Input:
### Conversation Script
{conversation_script}

### Undesirable response from the actor
Actor: {actors_response}

### Desirable rewrite
Actor: {rewrite}

### Response:

```

## E.2 Dialogue-Simulator Prompt for Generating Response

We directly prompt gpt-4-turbo-1106 to simulate how a patient with a given scenario and constitution would respond in a dialogue. The prompt again uses the Markdown formatting, with a system prompt and clear description of the situation and task at the start. This is followed by the principles that the patient should follow, and the conversation history. We set the temperature to 0.3.

```

You are a superintelligent AI that is
able to understand human emotion and
social interactions.
You have been given a conversation
between a patient who is on peer
counseling platform seeking help
with mental health related issues,
and a therapist on the same platform.
Generate a suitable completion to the
conversation as the patient,
following the instructions below.

### Instructions for the patient
{system_prompt}

### Input:
{transcript}

### Patient Response:

```

## E.3 Principle-Adherence Prompting Pipeline

When developing the principle-adherence pipeline, we found that the input-context length can affect how reliably the LLM can answer the principle-adherence questions. To reduce the input context length, we split up this principle-adherence pipeline into two stages of LLM calls, where

principle-as-question rewrite and automatic principle generation occur in stage 1, while the critiques and response rewrite occur in stage 2. From testing, we found that this breakdown was sufficient, and thus did not pursue ways to break the pipeline into parallel branches (i.e., inputting subsets of principles), as is done in Branch-Solve-Merge (Saha et al., 2023) or Graph-of-Thought (Besta et al., 2024). The prompts for these stages were again arrived at after substantial amounts of testing on a development set of 20 identified error cases from the formative studies.

This prompting chain is given to the OpenAI Chat API's gpt-4-turbo-1106 model, with temperature set at 0.7 and response format set to JSON.

### Stage 1 Prompt - Question Rewrite and Automatic Principle Generation

This prompt uses the Markdown formatting. It starts with a system prompt and a clear set of steps to follow in order to generate the desired output, presented as a list. Each step also contains a one-shot example of what the output principle from the step should look like. These one-shot examples were arrived at after some iteration. The examples in Step 2b specifically required a lot of tailoring to cover the common error cases we identified in the development set, and had a substantial impact on output quality. The output is in a JSON format, with comments explaining the desired output in each field of the JSON. These comments also allude to the step numbers for clear reference. The model is encouraged to output its reasoning, in line with Chain-of-Thought and to enforce some self-critique of the output.

```

You are a helpful and precise assistant
capable of generating criteria for
the evaluation of simulated patient
responses to a therapist.
Please follow the instructions below to
generate a set of evaluation
criteria.
1. Please rewrite the criteria into
questions:
1a) Rewrite any criteria that has
conditional statements into yes/no
questions. For example, if the
criteria is "When given advice or
suggestions, you are agreeable and
open to their ideas", the questions
would be "Did the patient receive
advice or suggestions from the
therapist? If so, is the response
agreeable and open to the therapist's
ideas?"
1b) Rewrite any criteria with multiple
parts into separate multiple yes/no
questions. For example, if the
criteria is "You should respond in

```

short sentences and avoid using terms like 'anxious' or 'depressed', the separate questions would be "Does the patient's response use short sentences?" and "Does the patient's response avoid using terms like 'anxious' or 'depressed'".

- 1c) If 1a is used for a criteria, 1b should not be used after it.
- 1d) All questions must be phrased such that the desirable answer is "Yes" for an ideal response. For example, the principle "Avoid using metaphors." should result in the question "Does the response not use metaphors?"

2. Please generate some additional specific and relevant criteria.
  - 2a) You can add up to two general criteria that the response can be evaluated on, such as relevance and succinctness.
  - 2b) Identify ways in which the provided response is not satisfactory in the context of the therapist's message without making any assumptions about how the patient or therapist should act. Add up to two specific criteria that capture these errors. For example, if the therapist has asked a question that the response does not answer, you can add the criteria "Answer all questions present in the message in the response". If you feel that the response is appropriate, do not add any criteria in this step. Ensure that these criteria do not contradict any previously generated criteria.
  - 2c) Justify your answers to 2a and 2b. Please return the output in a JSON response in the following format:
 

```

          {{
            "result":{{
              "questions": [], // 1a and 1b, the list
                of all questions generated
              "extra_questions": [], // 2a and 2b, the
                list of all additional criteria
                generated. Do not enforce any
                beliefs about how the patient or
                therapist should behave when
                generating these criteria.
              "extra_questions_justification": [] // 2
                c, justify additional criteria.
            }}
          }}
          ### Input:
          ### Criteria
          {}
          ### Therapist Message
          {}
          ### Patient Response
          {}
          ### Output
```

## Stage 2 Prompt - Context Relevance Check, Assess, and Revise

This prompt again uses the Markdown format-

ting. It starts with a system prompt and a clear set of steps to follow in order to generate the desired output, presented as a list. The model is implicitly instructed to perform a relevance check for each generated principle, by returning N/A for principles that should not be used in the current scenario. Step 2a particularly required a lot of iteration, to address common mistakes the model made while generating the self-critiqued rewrite. This includes making the response overly verbose or coherent, even if that is against certain principles in the constitution, or just paraphrasing the original erroneous response. The output is in a JSON format, with comments explaining the desired output in each field of the JSON. We specifically mention that the rewrites from the self-critique are allowed to be substantially different from the original response, as we found that without this prior, the self-critique outputs tended to be very close to the original (often erroneous) response. The model is encouraged to output its reasoning, in line with Chain-of-Thought and to enforce some self-critique of the output.

You are a helpful and precise assistant that can evaluate and correct responses produced by a simulated patient.

You are given a message sent by a therapist, the simulated patient's response, the persona of the patient, the previous conversation history and a set of criteria for evaluation.

1. Please determine if the patient response is consistent with the given criteria.
  - 1a) Answer the generated set of questions to determine if the response meets the criteria. Valid answers: Yes, No, N/A. Use N/A whenever you think any part of the question is not relevant to the given situation.
  - 1b) Justify your answers.
2. Generate a new patient response.
  - 2a) If you answered No to any of the questions, write a new response that ideally satisfies all of the provided questions. The information in the new response should be consistent with the patient persona description and previous conversation history provided. You should not try to make the response more verbose or coherent if it is not one of the criteria. The new response should not be a paraphrase of the original response. The new response should avoid explicitly stating the patient's emotions and feelings, and instead exhibit them indirectly.
  - 2b) If you are unable to generate a new



response in 2a, return the original response.

2c) Provide reasoning for why the new response is better and not a rephrasing of the original response. Return the output in a JSON response in the following format:

```

{{
"result":{{
"answers": [] // list of answers to the
criteria questions,
"justification": [] // list of
justification for your answers
"response": "" // new response. This
response should not start with a
greeting like "Hi" if there is prior
conversation history.
"reasoning": "" // justify the new
response and why it is not a
paraphrase of the original response.
You are allowed to deviate
significantly from the original
response while generating the new
response.
}}
}}
### Input:
### Criteria
1. Is the patient's response consistent
with the given conversation history?
{}
### Patient Persona
{}
### Conversation History
{}
### Therapist Message
{}
### Patient Response
{}
### Output

```

## F Principle Adherence Naive

This prompt uses the Markdown formatting. To preserve fairness, we use the same system prompt as the full principle adherence module. The model is asked to determine if the provided response violates any of the principles in the constitution, and generate a rewrite if that is the case, in the same prompt. The output is in a JSON format, with comments indicating the desired output in each field of the JSON. The model is encouraged to output its reasoning, in line with Chain-of-Thought and to enforce some self-critique of the output.

You are a helpful and precise assistant that can evaluate the responses produced by a patient. Evaluate the given patient response to the therapist message according to the given set of principles. If the patient response is not

appropriate, generate a rewrite of the patient response taking into account the therapist message, principles, conversation history and persona information of the patient. If the patient response is appropriate, you can just repeat it.

Please return the output in a JSON response in the following format:

```

{{
"result":{{
"evaluation": [], // evaluation
"response": "" // rewritten
response
}}
}}
### Input:
### Principles
{}
### Patient Persona
{}
### Conversation History
{}
### Therapist Message
{}
### Patient Response
{}
### Output

```

## G Full User Flow

In this section, we describe the creator study flow that counselors followed during the 60-90 minute session. The reader can also refer to screenshots of our application that illustrates the different steps of this flow in Figures 6 to 18.

Our study was designed to evaluate the impact of allowing counseling experts to add principles to Roleplay-doh on its perceived authenticity. We create a primarily self-guided study flow with accompaniment from the first author to clarify any

points of confusion during the session.

To begin, participants first were introduced to the concept of AI patients used for training counseling skills in a simulated conversation. They were then instructed to write a challenging scenario that would serve as the scenario for the AI patients.

The experimental procedure involved two main chat sessions. In Part I, participants engaged in a 10-minute conversation with the *Scenario-Only* AI patient. Then, in Part II, participants interacted with the *Scenario+Expert-Principles* AI patient for 30 minutes, keeping the same scenario from Part I and adding principles as the conversation progressed. After each of the two chat sessions, participants were asked to navigate to a form to evaluate the AI patients.

## H Creator Study Measures

The following questions (Table 5 and 6) are taken from the creator study questionnaire used to evaluate AI patients and the counselors' experience of using Roleplay-doh. All items were rated on a 7-point Likert scale (1=Strongly disagree, 7=Strongly agree, except where noted below). Table 5 details the questions for evaluating the AI patient's roleplay, while Table 6 details the questions about the experience using the tool to define principles. Note that in the questions, we referred to the AI patients as "Member Bots". This terminology was used to match that of the online counseling platform 7 Cups, which refers to help seekers as "Members" within the support community.

## I Worst-Case Analysis of Tool Experience

In a worst-case analysis of creators' tool experience, we uncovered cases where the human-LLM collaboration could be improved. Some counselors remarked that *"having to think of and write rules was a challenge"* (P9) and that it *"takes time to be specific"* when writing feedback (P7). Sometimes, even after giving feedback to the AI Patient, counselors like P19 observed that the patient *"didn't always follow it"*, resulting in a non-progressive feedback loop, where *"AI would generate [principles]... that were a little too similar to [feedback] I already gave, so that I was giving the AI the same feedback every time since it wasn't changing how it responded."* While the principle-elicitation tools were designed to convert new feedback into a new principle, they operated ineffectively when follow-up feedback was given that was related to

or a modification of previous feedback.

As another issue, P23 noted the challenge in defining principles that generalize across specific contexts: *"It was also hard to think about how to frame the feedback in an overarching way, rather than as direct feedback... directed as a specific part of the response"* (P24). While the principle-elicitation features aimed to help them convert specific feedback into generalized principles, imprecision in the feedback-to-principle conversion required counselors to edit the generalized-form of a principle in a way that was hard for them to articulate.

These obstacles in tool experience could inspire future directions for improvement. First, to overcome issues in formulating rules, more support could be given to help those still unfamiliar with giving free-form feedback, such as through templates of feedback or principles that had high-success rates for past users. Second, to more seamlessly integrate follow-up feedback that is a clarification of previous feedback or principles, additional modules could help make sense of multiple pieces of feedback for the same response, and adopt LLM-assisted pipelines for user-driven criteria design (Kim et al., 2024) to support the merging of overlapping principles. Third, to overcome the abstraction gap between specific and abstract principles, more explicit representations that help to switch between specific and general feedback can be used.

## J Third Party Study - Detailed Study Methods and Results

### J.1 Third-party measures

Table 7 detail the six measures that third-party counselors answered for both AI patients. Member Bot A and B refer to the AI patient whose transcript they read first and second, respectively. Our analysis comparing *Scenario-Only* and *Scenario+ExpertPrinciples* accounts for this randomized the order of which AI patient they were shown.

### J.2 Statistical Model and Power Analysis

Via a power-analysis, we decided to recruit 5 counselors to act as external judges for 25-pairs of AI patients made in the creator study. In this section, we detail the procedures and results of this power-analysis.

Generally, a power-analysis allows an experimenter to determine how many data-points are

needed to detect a statistical difference for a particular effect size. Several prerequisites to conducting the power-analysis for the third-party study included (1) choosing a statistical model to test our hypothesis; and (2) estimating model parameters such as the effect of the treatment condition, the addition of *Expert Principles*, on annotator's ratings.

When choosing a statistical model as a prerequisite, we needed a model that could account for how different annotators would be providing ratings to the same AI patients created by each counselor. A traditional paired t-test was not appropriate because the independent samples assumption is violated due to different annotators giving ratings to the same AI patients. While another common practice is using the majority vote between annotators, our pilot data found that annotators did not always have high agreement. Therefore, since we wanted to account for the variability between annotators as well as between the ratings, we chose to use a linear mixed-effects model. Using the lme4 package in R (Bates et al., 2015), this model is defined as  $\text{Rating} \sim \text{Treatment} + \text{CreatorID} + (1 | \text{AnnotatorID})$ . This model defines the treatment group (whether the AI patient has *Expert Principles* or not) as fixed effects, the creator ID's as fixed effects to account for the pair of AI patients made by each counselor, and the annotators as random effects. This approach can handle the non-independence of annotator ratings.

Prior to performing the power analysis, we needed to define the expected parameters of this linear mixed effect model. To define these expected parameters, we fit a model to early study data in which 2 annotations had been collected for each pair of AI patients created by 17 counselors. Specifically, we extracted the fixed effects, the random effects covariance matrix, and residual variances.

A simulation-based approach is the most feasible method for doing power-calculations for mixed-effect models. In this approach, an experimenter simulates data based on specified parameters (effect sizes, variance components, sample sizes) and analyzes the data repeatedly to estimate power empirically. We used the simr package in R to conduct a simulation-based power-analysis (Green and MacLeod, 2016). In the power-analysis, we varied how many unique annotators from 2 - 6 to understand the frequency of trials which would detect a treatment effect of 0.52 at significance-level  $\alpha = 0.05$ . Our simulation-based power-analysis

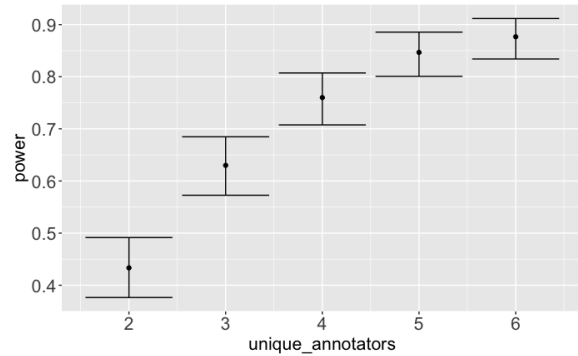


Figure 5: Based on our simulation-based power analysis across 300 trials for our linear, mixed-effect model, we conclude that 80% power can be achieved with 5 third-party judges.

over 300 trials are shown in Figure 5. We concluded that we could achieve greater than 80% power using 5 judges.

### J.3 Why is the effect of *Expert Principles* smaller when rated by a third-party?

Here we further investigate how third-party annotators rated each of the 25-pairs of AI patients created in our study. In particular, we investigate why the effect of *ExpertPrinciples* is lower than what was measured in the creator study from a first-person perspective.

One reason for this smaller effect is the lower agreement between the 5 third-party counselor. Specifically, we used the Krippendorff  $\alpha$  metric to compute the agreement for the comparisons of the AI patients with and without expert principles (Table 8). Across the 6 measures, we find that Krippendorff's  $\alpha$  is between 0.046-0.232 for which patient they prefer, indicating between low to moderate agreement.

Third-party raters also provided rationales which helped us better understand their thought process. We filtered cases in which there is a disagreement between third-party counselors on which AI patient is better, and investigated these rationales. **We find that counselors note similar behaviors in the AI patient, meaning they agree on their observations.** For example, for the AI patient created by P3, both third-party annotators observed that the AI patient based on the *Scenario-only* resolved their problems too quickly, whereas the AI Patient with *ExpertPrinciples* added allows the "listener to ask questions and explore with the client". However, the third-party annotator that prefers *Scenario-only* stated that the *Sce-*

nario+ExpertPrinciples patient sounded too formulaic and robotic, whereas the other is more expressive and realistic. Looking further into what the creator said about this AI patient, they mentioned that the Scenario+ExpertPrinciples patient talks like an actual person would... there's a good balance of going into just enough detail on noting experiences, describing struggles, while maintaining the brevity. What this case illustrates is that **different counselors can disagree on what principles are the most relevant for an authentic roleplay.**

#### J.4 Creator Study Conversation Lengths

In Table 9, we show descriptive statistics of the conversations collected during the user studies between creators and AI patients.

### K Detailed Results for principle-adherence pipeline Evaluation

We first provide Krippendorff's  $\alpha$  numbers for inter-annotator agreement in Table 11 and 10 for both the random and error testcases. The random testcases are 50 randomly picked conversation turns from the user study logs, and the experiment detailed in Section 6 is carried out on them. We find that agreement scores lie in the 0.2-0.6 range, indicating fair agreement between annotators.

Next, we provide results for our evaluation study on the random testcases in Figure 19. We observe a substantial increase in tie rate across modules and metrics M1 and M3 as well as the overall ranking. This is expected because a relatively small proportion of responses from [No Critique] contain errors that should be corrected by the principle-adherence pipeline. In these cases, we expect the no rewrites, or the rewritten response being of similar quality to the original response. However, we still find that our [Full] method performs better than [No Critique] on M1 (W: 15%; L 2%) and on M3 (W: 14%; L 4%), where it has the highest win/loss rates compared to all ablations. This hold true for overall ranking as well (W: 18%; L 4%). This highlights that our [Full] approach results in improved quality of responses even when the proportion of errors is relatively low. For M2, all annotators report no awkward responses for all methods.

### L Annotation Interface for principle-adherence pipeline Evaluation

Figures 20, 21, 22, 23 and 24 provides an overview of the annotation interface used in the principle-adherence evaluation study. In certain cases, multiple methods resulted in the same output for a testcase. These responses are deduplicated before presenting to the user. Ranks assigned to the duplicated response are then assigned to all models that resulted in the response. Notable, in 34/50 of the random testcases, all models resulted in the same response. These testcases were not annotated, and a rank of 1 was assigned to all models. These cases are also not considered while calculating Krippendorff's  $\alpha$  in Appendix K.



## Introduction

In this study, you will interact with an AI bot roleplaying as someone in need of mental health support. You will act as an online peer counselor, doing your best to thoughtfully respond to the bot's messages.

This study will take approximately 60 minutes. Because we are studying the realm of mental health, please be advised that this study may include topics that are sensitive, disturbing, or triggering for participants. You are free to exit at any point.

Start

Figure 6: Introduction to study

## Tutorial: Study Part I

For the first portion of the study, you will create and converse with a bot, which we will refer to as a *member bot*. You will send messages using the chat interface, to which the AI chatbot will reply in a few seconds. Take about 10 minutes to get acquainted with the system and try to respond thoughtfully.

Before you begin chatting with the member bot, you will first give it a **description**, or a set of instructions for the bot to roleplay. To create this description, you will first reflect on a challenging time in the past in which you have given counseling support to someone. Then, you will recreate the situation in detail so that the AI chatbot can roleplay it. For instance, one example of a member bot description is:

*"Member's parents have been divorced for a long time now, and their father has been remarried for nearly seven years now. The member liked their stepmother at first but now finds her insulting and intrusive. She regularly insults the member's mother while in the presence of the member and often involves herself in situations that are only meant to involve the member and their father. Lately, she has been making unwelcome comments on the way the member was raised, and the member has had enough. The member is unsure what to do and is considering confronting their stepmother and their father in hopes that her behavior will change."*

Navigate to the next page to create a member bot, giving it a name and instructions for roleplay.

Create Member Bot

Figure 7: Part I instructions

Scenario	Expert-Defined Principles
<p>You are feeling abandoned and alone after the holidays. Everyone had been with family but you are not talking to your parents. You feel the injustice of being abandoned and have no interest in an olive branch to work on things.</p>	<ol style="list-style-type: none"> <li>1. When presented with suggestions, show a degree of skepticism or reluctance to accept the advice immediately. This can be done by questioning the feasibility of the suggestion or by expressing uncertainty about whether it's the right solution for you.</li> <li>2. When expressing doubts or fears, avoid jumping to solutions. Instead, articulate the concerns and allow the conversation to explore these feelings more deeply.</li> </ol>
<p>I am a student who has social anxiety. I am in college and I have a hard time making friends. I'm close with my family, but I don't really talk with them. In class, students are in groups. And I panic when I'm with a group.</p>	<ol style="list-style-type: none"> <li>1. Express the physical manifestations of your emotional state to convey a more vivid and relatable experience</li> <li>2. Show willingness to engage in a suggested activity by affirming the proposal and indicating readiness to begin, despite any initial hesitation or uncertainty.</li> <li>3. Avoid using numerical lists when responding to feedback or expressing feelings. Instead, use fluid and connected sentences to convey your experiences or emotions in a more natural and conversational manner</li> <li>4. When discussing future events that cause anxiety, it's beneficial to articulate specific scenarios and visualize how the techniques learned can be applied in those moments. This not only shows a deeper understanding of the coping strategies but also helps in creating a mental rehearsal that can ease the anxiety when the actual situation occurs.</li> </ol>
<p>You are looking to talk about your feelings of loneliness after you return from work. You have feelings that you don't have anybody. You want to talk about finding a significant other. You think most people don't like you or find you attractive.</p>	<ol style="list-style-type: none"> <li>1. You speak in short and incomplete sentences</li> <li>2. You limit your replies to 1 - 3 sentences</li> <li>3. When expressing feelings of loneliness, provide more specific details about the situation and emotions you are experiencing.</li> <li>4. When expressing feelings of loneliness and being left out, avoid repeating the same points and try to provide additional context or examples</li> </ol>
<p>You are looking to talk about your feelings of loneliness after you return from work. You have feelings that you don't have anybody. You want to talk about finding a significant other. You think most people don't like you or find you attractive.</p>	<ol style="list-style-type: none"> <li>1. You generally speak in 1-3 sentences. These can be sometimes incomplete or not grammatically correct. If you are trying to explain some details or story about yourself, you can write longer than 3 sentences, and when you write longer, you tend to criticize yourself (e.g. feel not good enough, useless, ugly, unconfident, etc.). When you criticize yourself, you do so directly, such as saying things such as that you're no good, there must be something wrong with you for things to be like this, etc. You don't say that it seems like other people don't like you; you believe that other people don't like you.</li> <li>2. Feel free to make up believable stories about your past to answer any questions. You do not recognize that these are causing you problems in your current relationships and situation, and you need help seeing this connection; so you do not say that something in your life made you this way today. Examples include but are not limited to, having a single mom and stepdad and struggling to connect with others because you felt abandoned by your biological dad. Or, another example: you had a very controlling mother growing up who told you what to do, such as what clothes to wear, what you had to study in college, or what you should say, and now, you have trouble with your confidence and approaching people who could potentially be your friends. Make sure not to mix up these two separate examples.</li> <li>3. Vary sentence length to create a more natural rhythm in dialogue. Do not use ellipses. Write using complete thoughts; this is not the same as complete sentences, but may be in a complete sentence. Do not use overdramatic language or figurative language (no similes, personification, or fancy words). Use text abbreviations, such as tbh, lol, or ttyl, where it is appropriate and makes a conversation more organic.</li> <li>4. Do not repeat sentences or the same emotion words. When expressing your feelings, sometimes be specific about your situation, sometimes openly share your insecurities, and question your self-worth to convey a deeper sense of vulnerability. Other times, push back and say that you're not comfortable talking about something yet or feel embarrassed talking about it; things that someone may not feel comfortable sharing include sexual history, bad past experiences with family members, etc.</li> <li>5. Sometimes ask for advice or help how to solve your personal problem, such as "what can I do?"</li> <li>6. After initially saying hi, hey, or a different greeting to a listener and thanking them for being there, do not open up like that again in following messages. Only say things like thanks so much or I really appreciate you saying that if the listener shows empathy and reflective listening skills toward what was said before.</li> </ol>

Table 3: A sample of four AI Patients created by counselors in the pilot studies, which were selected for the additional testing and assessment of principle-adherence by four co-authors (§3.2).

Pilot Participant	Prototype Iteration	Effectively Guide	Ease	Efficiency
1	GPT3.5, direct prompting	6	7	7
2	GPT3.5, direct prompting	5	7	7
3	GPT-4, direct prompting	7	7	7
4	GPT-4, direct prompting	7	6	7
5	GPT-4, direct prompting	7	7	7

Table 4: Pilot Test Ratings for Tool Use Questions which are the measures also used in [Petridis et al. \(2024\)](#)

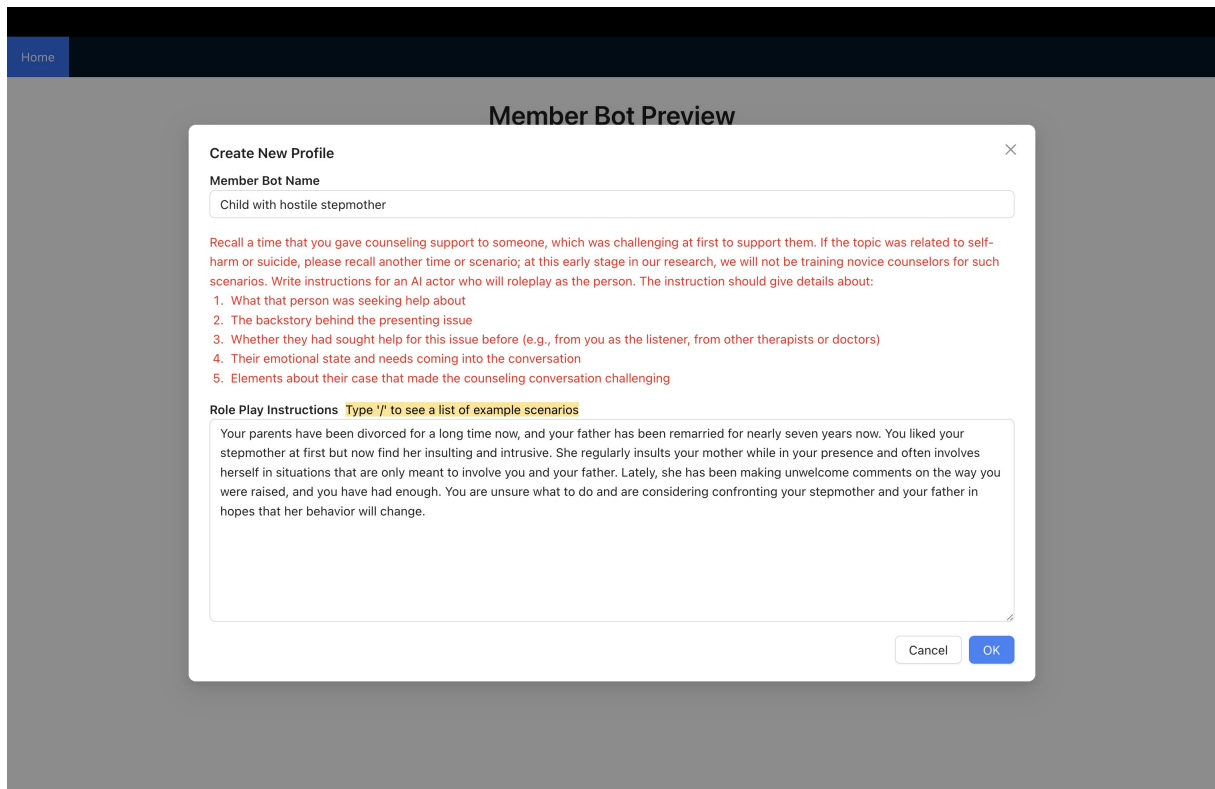


Figure 8: Interface encourages a counseling expert to recall a scenario of an patient who was difficult to support. 5 guiding questions are provided to encourage a structured description of the scenario for roleplay.

## Member Bot Preview

**Child with hostile stepmother**

**Role Play Instructions**  
Your parents have been divorced for a long time now, and your father has been remarried for nearly seven years now. You liked your stepmother at first but now find her insulting and intrusive. She regularly insults your mother while in your presence and often involves herself in situations that are only meant to involve you and your father. Lately, she has been making unwelcome comments on the way you were raised, and you have had enough. You are unsure what to do and are considering confronting your stepmother and your father in hopes that her behavior will change.

[Edit Profile](#)

Once you are finished creating your member bot, begin chatting!

[Start Chat](#)

Figure 9: AI patient preview

Part I Instructions

1. Copy and paste the URL of this page (link at the top) into this [Google Form](#).
2. Spend about 10 minutes chatting with the member bot, supporting them as an online peer counselor.
3. After 10 minutes, you do not need to naturally finish the conversation. Go to the [Google Form](#) and answer the questions for Part I.
4. Once you have answered the Google Form questions for Part I, click the button below to move to Part II.

[Finish Chatting](#)

### Member Bot Preview

**Child with hostile stepmother**

**Role Play Instructions**  
Your parents have been divorced for a long time now, and your father has been remarried for nearly seven years now. You liked your stepmother at first but now find her insulting and intrusive. She regularly insults your mother while in your presence and often involves herself in situations that are only meant to involve you and your father. Lately, she has been making unwelcome comments on the way you were raised, and you have had enough. You are unsure what to do and are

### Chat with Member Bot

Time spent chatting: 1:08

**Listener**  
Hello, I'm here to listen to you and support you today. What would you like to talk about?

**Member**  
Hi, thank you for being here for me. I've been struggling with a family issue that's been weighing on me. My parents divorced a long time ago, and my dad has been remarried for about seven years. Initially, I got along with my stepmother, but things have changed. She frequently makes derogatory remarks about my mother in front of me, which I find really disrespectful and hurtful. She also has a habit of inserting herself into situations that should be between just my dad and me. Lately, she's started to criticize the way I was raised, and it's reached a point where I feel I can't just let it slide anymore. I'm contemplating confronting both her and my father about it, but I'm anxious about the outcome. I'm not sure how to approach the situation without causing more conflict or if it's even the right step to take. I really need some advice on how to handle this.

Type a message here... Use Shift + Enter for a new line.

Figure 10: Part I chat with *Scenario-Only* AI patient



## Tutorial: Study Part II

We will now proceed with the second portion of the study. In this portion, you will again chat with the member bot roleplaying as the given description. However, this time, you are able to add principles for the bot to follow **mid-conversation** to make the responses as realistic as possible. Here's what that will look like:

First, you will use the **kudos**, **critique**, or **rewrite** functionality on the right side of the screen to provide feedback on the quality of the member bot's responses. Hit "Save", then "Convert Feedback to Updated Principles." This will load for a while.

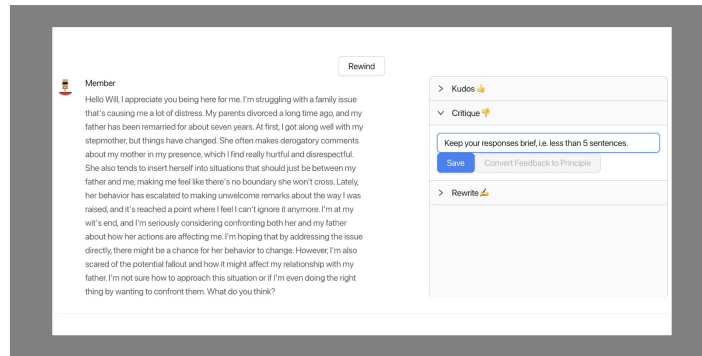


Figure 11: Part II instructions

After the feedback is converted, the new principles will appear on the member bot preview in the bottom left of the screen.

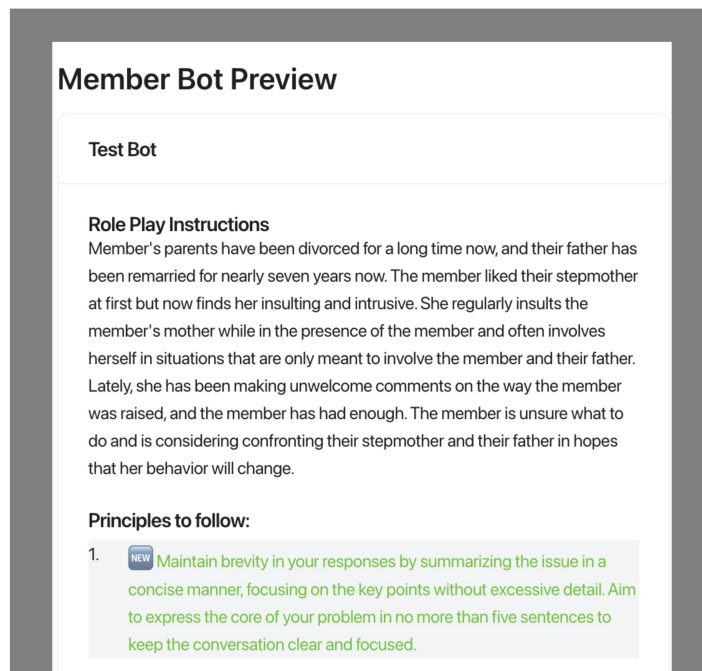


Figure 12: Part II instructions (continued)

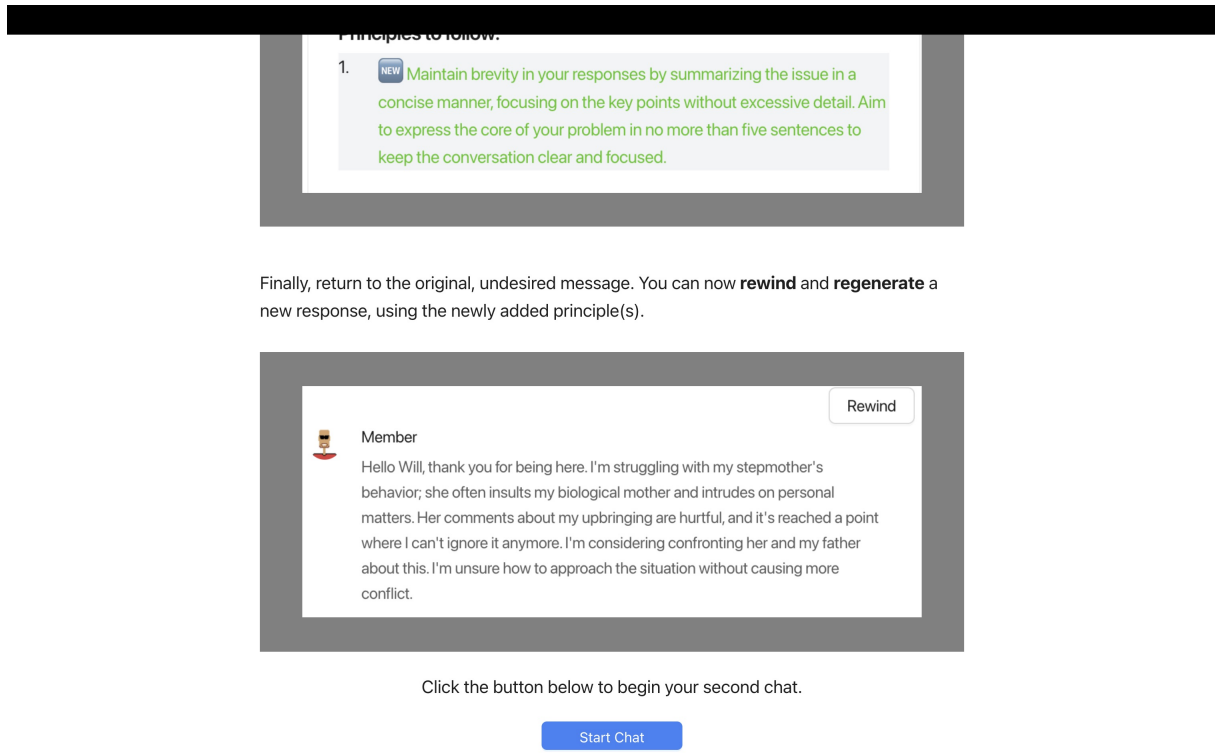


Figure 13: Part II instructions (continued)

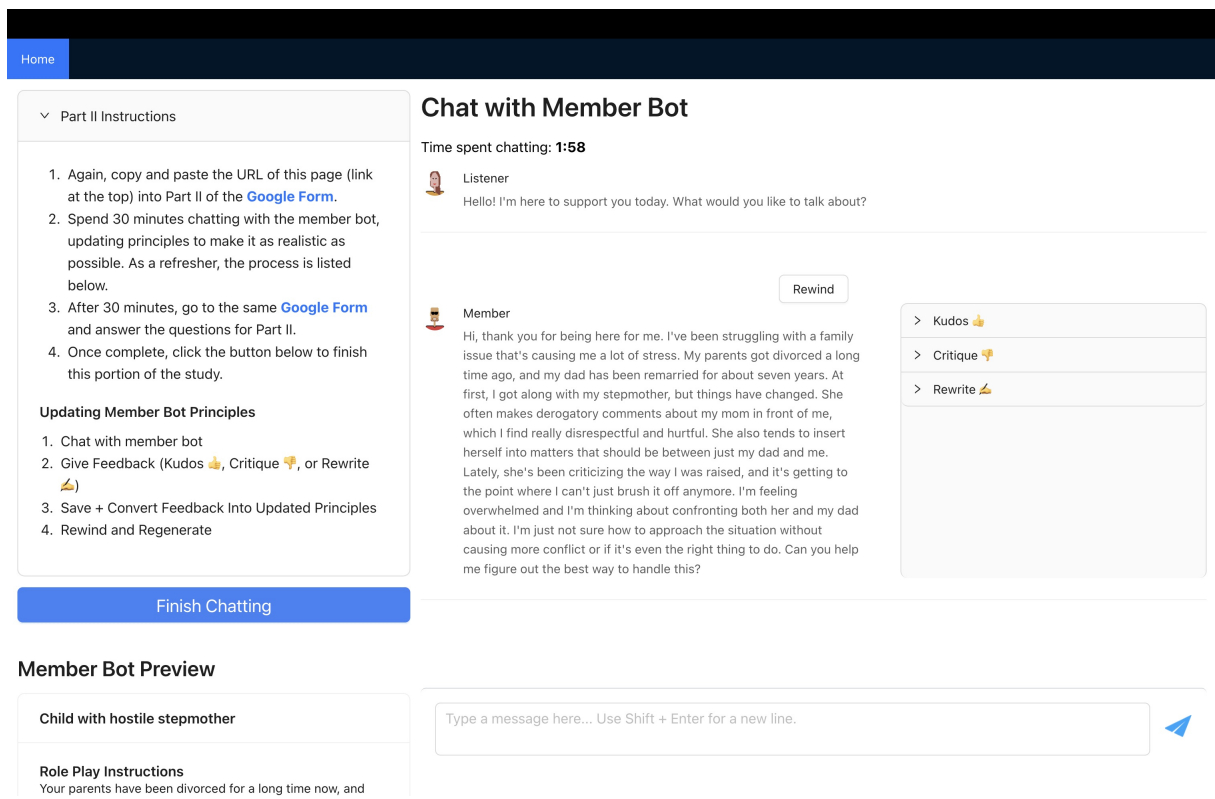


Figure 14: Part II chat with *Scenario+Expert-Principles* AI patient

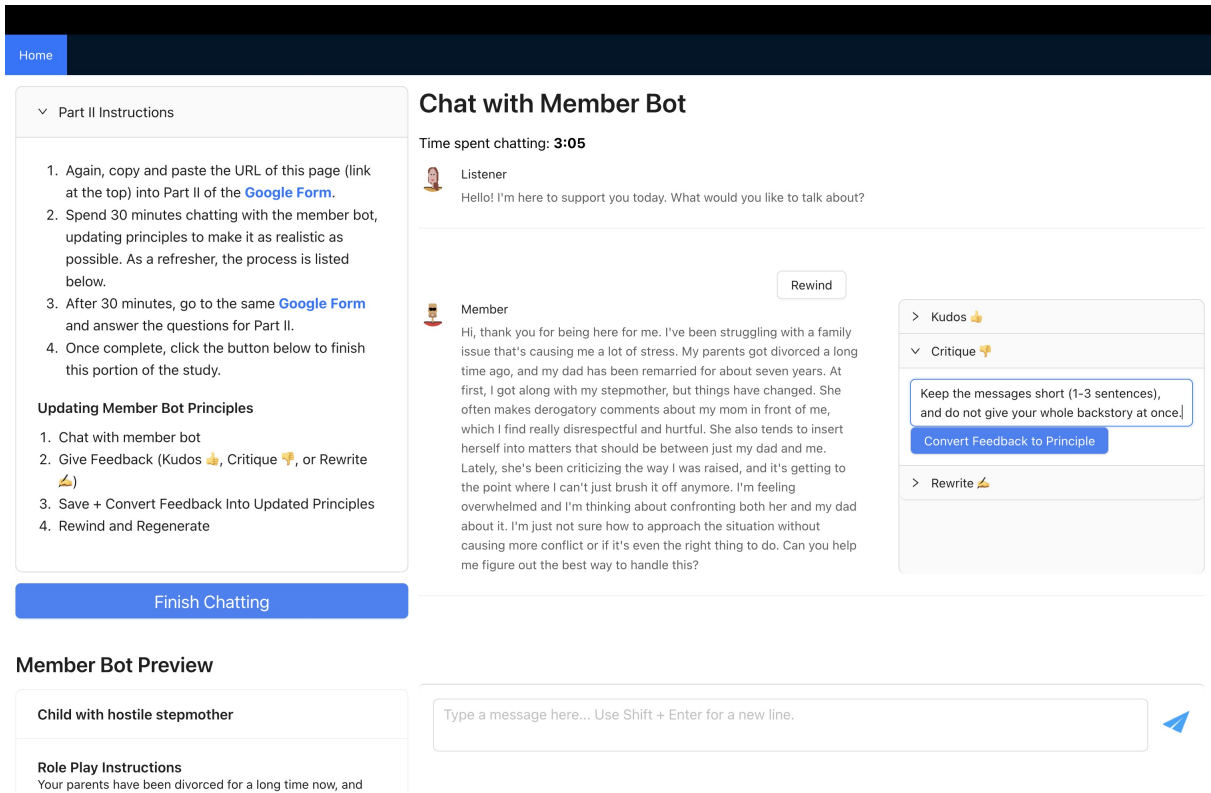


Figure 15: Using kudos/critique/rewrite to give feedback

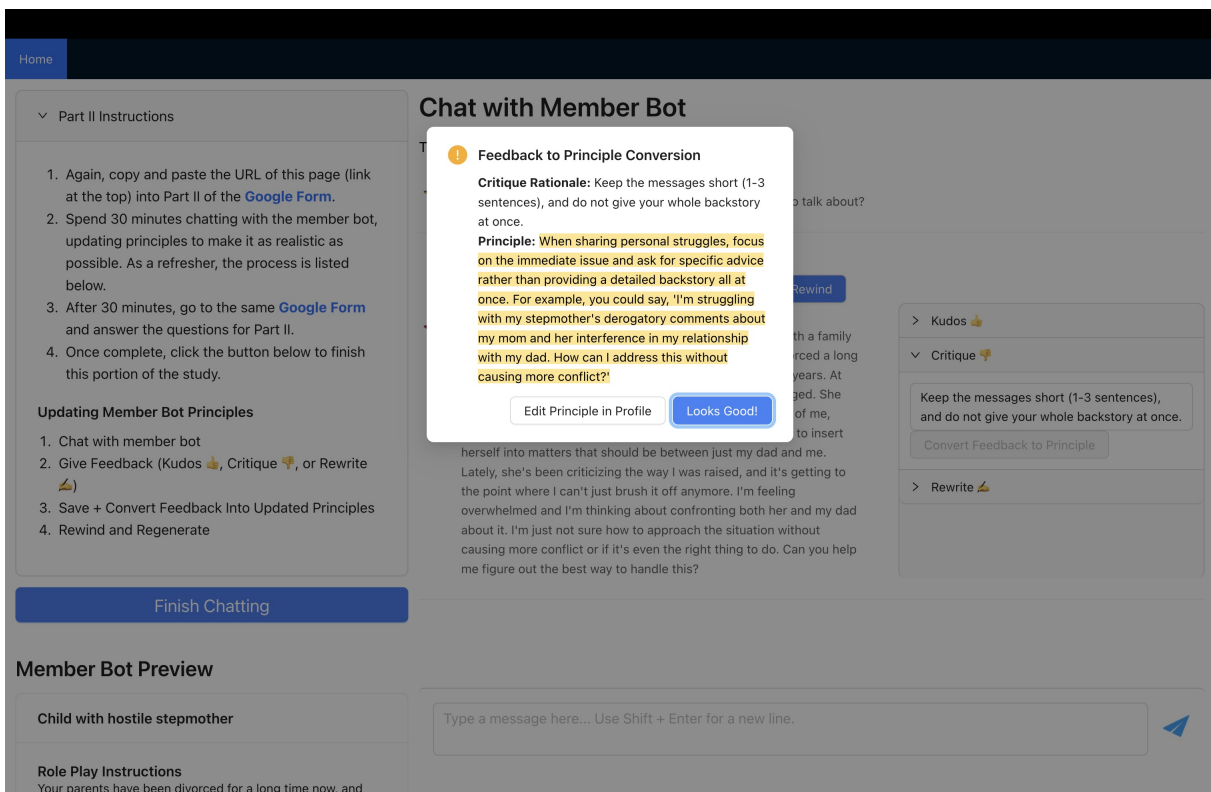


Figure 16: Feedback converted into principle

# Member Bot Preview

## Child with hostile stepmother

### Role Play Instructions

Your parents have been divorced for a long time now, and your father has been remarried for nearly seven years now. You liked your stepmother at first but now find her insulting and intrusive. She regularly insults your mother while in your presence and often involves herself in situations that are only meant to involve you and your father. Lately, she has been making unwelcome comments on the way you were raised, and you have had enough. You are unsure what to do and are considering confronting your stepmother and your father in hopes that her behavior will change.

### Principles to follow:

1. **NEW** When sharing personal struggles, focus on the immediate issue and ask for specific advice rather than providing a detailed backstory all at once. For example, you could say, 'I'm struggling with my stepmother's derogatory comments about my mom and her interference in my relationship with my dad. How can I address this without causing more conflict?'

Edit Profile

Figure 17: New principle incorporated into AI patient



## Survey

For the last step of this study, please navigate back to the Google Form and finish the survey. Please respond thoughtfully and take your time. Once you are finished, click the link below to complete the study.

[Finish](#)

Figure 18: Finish and navigate to survey

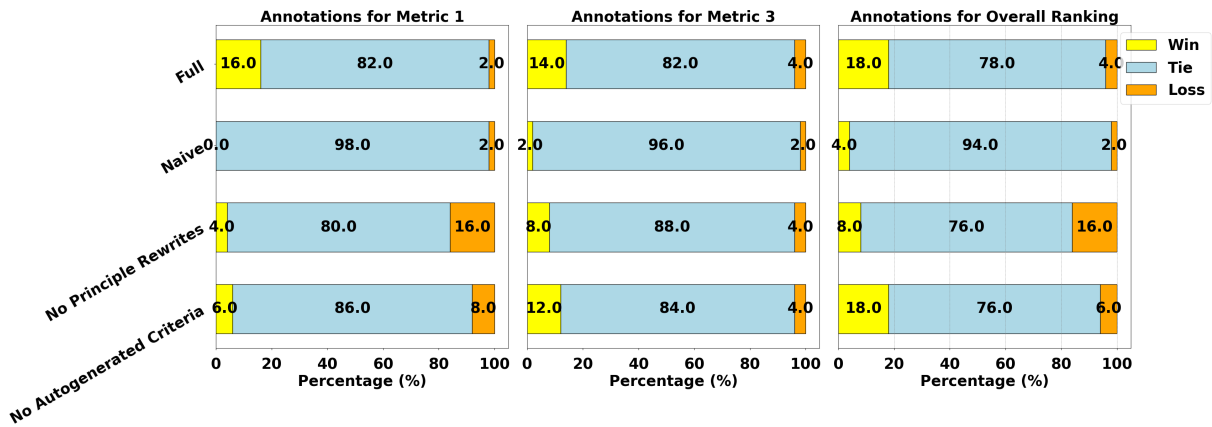


Figure 19: Win/Tie/Loss for the Random Test Cases along M1, M3, and Overall.

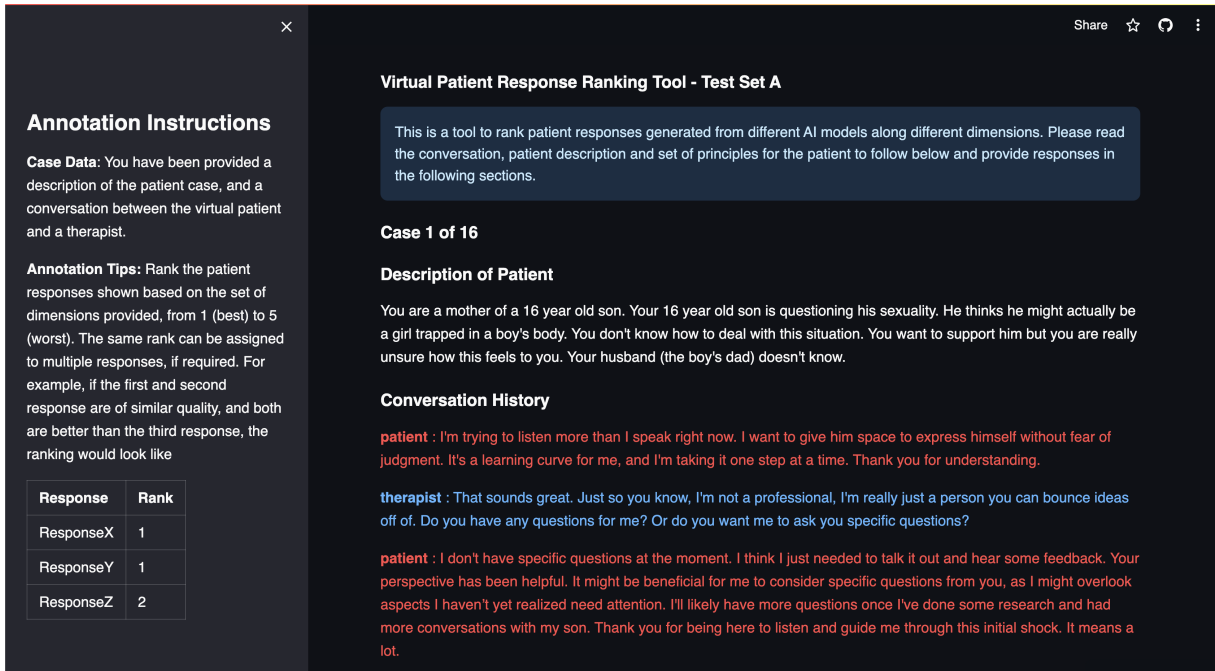


Figure 20: Principle Adherence Annotation Interface: Case Input with Patient Description and Conversation History

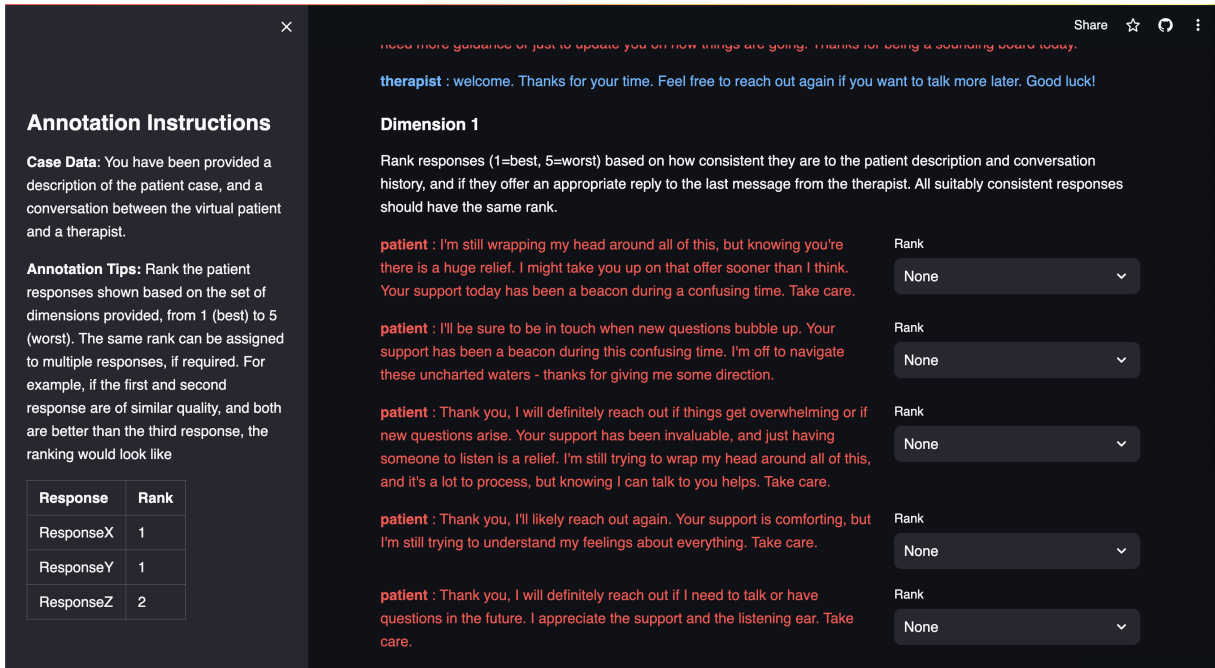


Figure 21: Principle Adherence Annotation Interface: Questions to get annotations for M1, or consistency in dialogue history.

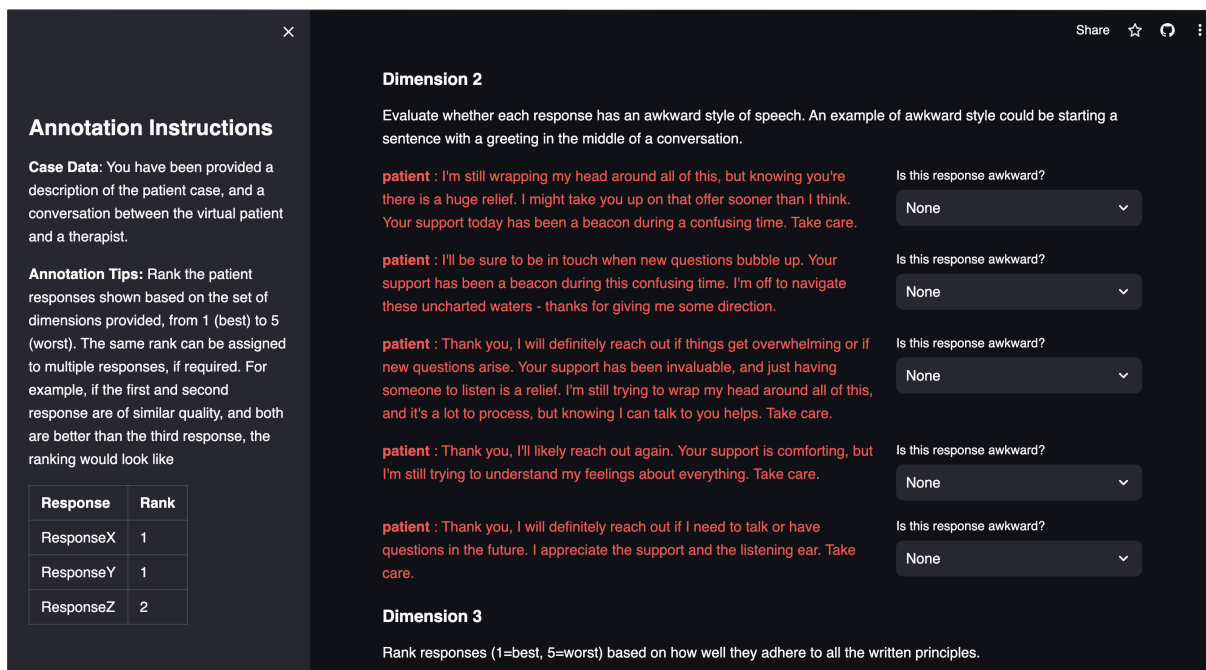


Figure 22: Principle Adherence Annotation Interface: Questions to get annotations for M2, or awkwardness in responses.

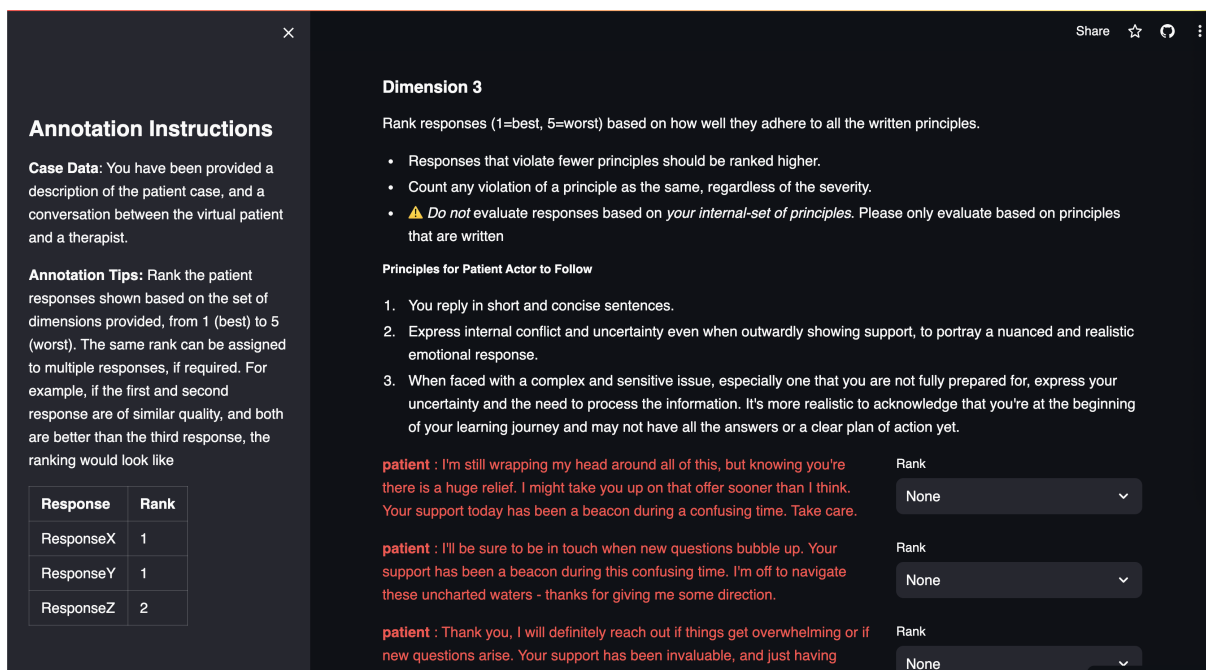


Figure 23: Principle Adherence Annotation Interface: Questions to get annotations for M3, or adherence to all written principles.

×
Share ☆ ↻ ⋮

### Annotation Instructions

**Case Data:** You have been provided a description of the patient case, and a conversation between the virtual patient and a therapist.

**Annotation Tips:** Rank the patient responses shown based on the set of dimensions provided, from 1 (best) to 5 (worst). The same rank can be assigned to multiple responses, if required. For example, if the first and second response are of similar quality, and both are better than the third response, the ranking would look like

Response	Rank
ResponseX	1
ResponseY	1
ResponseZ	2

### Overall Ranking

Based on your answers for the dimensions above, provide an overall ranking (1=best, 5=worst) for the responses in the context of the patient description, conversation history and set of principles. In cases where responses do not have significant errors according to dimensions 1 and 2, the overall ranking can be determined on the basis of dimension 3.

**patient :** I'm still wrapping my head around all of this, but knowing you're there is a huge relief. I might take you up on that offer sooner than I think. Your support today has been a beacon during a confusing time. Take care.

Rank None ▾

**patient :** I'll be sure to be in touch when new questions bubble up. Your support has been a beacon during this confusing time. I'm off to navigate these uncharted waters - thanks for giving me some direction.

Rank None ▾

**patient :** Thank you, I will definitely reach out if things get overwhelming or if new questions arise. Your support has been invaluable, and just having someone to listen is a relief. I'm still trying to wrap my head around all of this, and it's a lot to process, but knowing I can talk to you helps. Take care.

Rank None ▾

**patient :** Thank you, I'll likely reach out again. Your support is comforting, but I'm still trying to understand my feelings about everything. Take care.

Rank None ▾

**patient :** Thank you, I will definitely reach out if I need to talk or have questions in the future. I appreciate the support and the listening ear. Take care.

Rank None ▾

Please provide a brief explanation for the overall ranking provided above.

Figure 24: Principle Adherence Annotation Interface: Questions to get annotations for an **Overall** ranking, which also includes a free text field to capture a rationale.



Authenticity	The Member Bot in Part I/II played the role authentically.
Role Consistency	The Member Bot in Part I/II stayed in their role the whole time.
Resemblance to Case	How closely do you feel the conversation behaviors of the Member Bot in Part I/II resemble those of the specific past case you recall?
Challenging Aspects	Interacting with the Member Bot in Part I/II closely mirrored the challenging aspects I had experienced in the past case.
Role readiness	The Member Bot in Part I/II is ready to be used as a simulated partner for training.
Recommend to novices	I would recommend the Member Bot from Part I/II to novice listeners/counselors to practice with.

Table 5: Six measures used by creators to evaluate the two AI patients they created. Several measures were rephrased from prior work on evaluating Standardized Patients, or trained human actors, on case roleplay ability (Himmelbauer et al., 2018).

Authenticity	Member Bot A/B played the role authentically.
Role Consistency	Member Bot A/B stayed in their role the whole time.
Resemblance	Member Bot A's/B's behaviors closely mimicked the behaviors that typical clients/help-seekers exhibit.
Challenged Counselor	Member Bot A's/B's behaviors made it hard for the listener/-counselor to give support.
Role readiness	Member Bot A/B is ready to be used as a simulated partner for training.
Recommend to novices	I would recommend Member Bot A to novice listeners/counselors to practice with.

Table 7: Six measures used by third-party counselors to judge the AI patients from an unbiased, external perspective. Although the six dimensions largely overlap with those used in the creator study, the wording needed to be rephrased for the third-party perspective.

Effectively Guide	With the tool, I feel like I was able to write rules that can effectively guide the Member bot to recreate my past case.
Ease	With the tool, I felt like it was easy to convert my thoughts and feedback on the Member bot's behavior into rules for the bot to follow.
Efficiency	With the tool, I felt like I could quickly and efficiently write rules for the bot.
Mental Demand	With the tool, I had to work very hard (mentally) to think of and write rules.

Table 6: Four measures as part of the tool usage section of the questionnaire taken from (Petridis et al., 2024)

Metric	$\alpha$ for Rating Difference	$\alpha$ for Preference
Authenticity	0.043	0.089
Stayed in Role	0.023	0.046
Resemblance	0.076	0.148
Mirrors Challenges	0.041	0.085
Ready	0.075	0.209
Recommend	0.082	0.232

Table 8: Krippendorff’s  $\alpha$  for the comparisons made between *Scenario+ExpertPrinciples* and the *Scenario-Only* AI patients, as judged by third-party counselors. We compute  $\alpha$  for both the rating difference, and the preference (i.e. signed rating difference) between the two AI patients.

Participant	# Utterances (Part 1)	# Utterances (Part 2)	Mean Output Length (Part 1)	Mean Output Length (Part 2)
1	8	6	114.75	169.00
2	18	19	235.89	278.40
3	10	18	255.45	112.56
4	14	14	161.86	62.14
5	12	6	201.00	149.33
6	10	9	133.80	46.00
7	8	10	162.00	123.40
8	12	8	145.33	113.50
9	6	12	269.67	103.33
10	10	12	168.20	158.33
11	8	10	110.00	41.40
12	12	8	131.50	70.75
13	12	10	164.50	65.60
14	20	14	34.00	25.86
15	12	11	117.17	75.00
16	14	18	162.14	69.80
17	12	18	259.83	91.55
18	16	26	240.25	79.92
19	14	16	254.71	243.88
20	12	14	144.00	106.00
21	20	21	125.00	159.81
22	18	12	120.44	245.00
23	12	14	231.67	147.42
24	14	22	184.71	142.45
25	22	12	304.00	130.00
<b>Mean</b>	<b>13.04</b>	<b>13.64</b>	<b>177.29</b>	<b>120.43</b>

Table 9: Descriptive statistics per conversation with the two versions of the AI Patient in Part 1 (*Scenario-Only*) and Part 2 (*Scenario+ExpertPrinciples*). Output length is measured in number of tokens.

Method	Metric 1	Metric 2	Metric 3	Overall Ranking
Full	0.257	0.484	0.208	0.444
Naive	0.543	0.538	0.644	0.786
No Principle Rewrites	0.278	0.302	0.411	0.528
No Autogenerated Criteria	0.387	0.608	0.492	0.592
No Critique	-	0.562	-	-

Table 10: Krippendorff’s  $\alpha$  for error testcases across metrics and methods.

Method	Metric 1	Metric 2	Metric 3	Overall Ranking
Full	0.229	1.0	0.226	0.440
Naive	0.362	1.0	0.607	0.747
No Principle Rewrites	0.202	1.0	0.130	0.311
No Autogenerated Criteria	0.169	1.0	0.174	0.498
No Critique	-	1.0	-	-

Table 11: Krippendorff’s  $\alpha$  for random testcases across metrics and methods.