# Combining Artificial Users and Psychotherapist Assessment to Evaluate Large Language Model-based Mental Health Chatbots

Florian Onur Kuhlmeier[a,b,*], Leon Hanschmann[a], Melina Rabe[b], Stefan Lüttke[b], Eva-Lotta Brakemeier[b], Alexander Maedche[a]

[a]*Karlsruhe Institute of Technology (KIT), human-centered systems lab (h-lab),*
[b]*University of Greifswald, Chair of Clinical Psychology and Psychotherapy,*

## Abstract

Large Language Models (LLMs) promise to overcome limitations of rule-based mental health chatbots through more natural and personalized conversations. However, evaluating LLM-based mental health chatbots presents a significant challenge: Their probabilistic nature requires comprehensive testing to ensure therapeutic quality and safety, yet conducting such evaluations with people with depression would impose an additional burden on vulnerable people and risk exposing them to potentially harmful content. Our paper presents an evaluation approach for LLM-based mental health chatbots that combines dialogue generation with artificial users and dialogue evaluation by psychotherapists. In collaboration with psychotherapists, we developed artificial users based on patient vignettes, systematically varying characteristics such as depression severity, personality traits, and attitudes toward chatbots. We randomly selected 48 artificial users to interact with a LLM-based behavioral activation chatbot. Ten psychotherapists evaluated these 48 dialogues using standardized rating scales to assess the quality of behavioral activation and its therapeutic capabilities. The results show that while artificial users showed moderate authenticity, they enabled comprehensive testing across different user profiles without exposing vulnerable individuals to potential harm or adding evaluation burden. In addition, the chatbot demonstrated promising capabilities in delivering behavioral acti-

---

[*]Corresponding author
*Email address:* `florian.kuhlmeier@kit.edu` (Florian Onur Kuhlmeier)

vation and maintaining safety. Furthermore, we identified deficits, such as ensuring the appropriateness of the activity plan, which reveals necessary improvements for the chatbot. Our framework provides an effective method for evaluating LLM-based mental health chatbots while protecting vulnerable people during the evaluation process. Future research should improve the authenticity of artificial users and develop LLM-augmented evaluation tools to make psychotherapist evaluation more efficient, and thus further advance the evaluation of LLM-based mental health chatbots.

## 1. Introduction

Depression is a prevalent and severe mental disorder with significant personal and socioeconomic consequences (GBD 2019 Mental Disorders Collaborators, 2022). Limited access to psychotherapy (Butryn et al., 2017) and the stigma surrounding traditional psychotherapy (Schomerus et al., 2022) have led to the increased adoption of mental health chatbots as a treatment alternative (Darcy et al., 2021; Mehta et al., 2021). The capabilities of large language models (LLMs) promise to overcome the current limitations of rule-based chatbots, such as rigid dialogues and limited personalization (Chan et al., 2022; Fitzpatrick et al., 2017; Kocaballi et al., 2019). Although LLMs have shown promising capabilities in mental health applications, including identifying symptoms (Yang et al., 2023) and cognitive restructuring (Sharma et al., 2024), their probabilistic nature can produce inconsistent or potentially harmful responses (Heston, 2023), which require comprehensive evaluations before deployment to vulnerable users.

A key component of evaluating the quality of psychotherapy is treatment fidelity, where experts assess how well psychotherapists implement a specific intervention protocol (Beck et al., 2023; Proctor et al., 2011; Webb et al., 2010). In behavioral activation, an evidence-based intervention for depression aiming to engage patients in meaningful activities (Lejuez et al., 2001, 2011), treatment fidelity refers to psychotherapists' capability to implement key components, such as identifying and scheduling meaningful activities (Dimidjian et al., 2012). Evaluating treatment fidelity differs between rule-based and LLM-based chatbots. For rule-based chatbots, messages are prewritten by mental health professionals and follow predetermined paths,

2

allowing for comprehensive evaluation through a single review before deployment. In contrast, LLM-based chatbots generate probabilistic responses that may vary significantly across users and conversations. This variability creates two key challenges: First, ensuring consistent therapeutic quality requires testing across many different users and scenarios, necessitating more evaluation sessions than rule-based chatbots. Second, the probabilistic nature of responses increases the risk of exposing vulnerable individuals to potentially harmful content during evaluation. Third, recruiting people experiencing depression for extended evaluation periods is particularly challenging, as core symptoms like lack of motivation and fatigue limit their ability to participate in evaluation studies.

To address these challenges, artificial users can serve as a risk-free alternative or supplement to human participants. Although Chen et al. (2023) demonstrated the feasibility of artificial users in evaluating an LLM-based diagnostic chatbot, their work only included two minimally varied patient profiles (e.g., formal vs. colloquial symptom reporting). This approach does not capture the heterogeneity of human patients, who vary across different dimensions such as symptom severity or personality traits. In addition, their work focused on diagnostic rather than therapeutic conversations. Whereas diagnostic conversations primarily gather information through structured questions, therapeutic interactions are more complex and demanding. They require building rapport, explaining the therapeutic approach, and implementing the protocol while maintaining flexibility to adapt to the client's needs, responses, and progress. This flexibility poses a challenge for LLMs, which must maintain therapeutic consistency while demonstrating the situational flexibility characteristic of effective therapy. These demands make therapeutic conversations more critical to thoroughly evaluate, given their direct influence on patient well-being.

Our evaluation approach addresses these challenges through two components: dialogues between an LLM-based behavioral activation chatbot and artificial users with varying characteristics, followed by psychotherapist assessment of therapeutic quality and capabilities in these dialogues. This research investigates three research questions:

1. How useful is combining dialogues generated with artificial users and assessed by psychotherapists for evaluating LLM-based mental health chatbots?
2. How well does the behavioral activation chatbot demonstrate thera-

peutic quality and therapeutic capabilities?

3. What specific improvements can be identified with the approach to improve the chatbot's thearpeutic quality and capabilites?

To answer these questions, we evaluated an LLM-based behavioral activation chatbot designed for young people with depression. In collaboration with psychotherapists, we created artificial users based on patient vignettes and enriched them with carefully selected characteristics (depression severity, age, gender, willingness to disclose personal information, openness to chatbot suggestions, conversational dominance, attitudes towards mental health chatbots). From a pool of 2,112 systematically varied artificial users, we drew a stratified random sample of 48, limited by how many dialogues the 10 psychotherapists could evaluate within their available time. Artificial users interacted with the chatbot, after which the psychotherapists assessed treatment fidelity and therapeutic capabilities, while also identifying shortcomings and suggesting improvements.

Our research makes the following contributions to the evaluation of LLM-based mental health chatbots. First, we propose a novel evaluation approach that integrates artificial users, derived from clinically validated patient vignettes and systematically enriched with key user characteristics, with detailed psychotherapist evaluations. This approach enables rigorous and risk-free assessments across diverse user profiles, thoroughly measuring therapeutic quality and therapeutic capabilities, while identifying actionable areas for improvement. Second, we develop and apply structured evaluation instruments for assessing LLM-based mental health chatbots, comprising standardized rating scales and targeted open-ended questions. These instruments measure the chatbot's therapeutic quality and capabilities, facilitate identifying specific concrete shortcomings, and yield clear, actionable recommendations for chatbot refinement. Importantly, our instruments allow for both holistic evaluations of the chatbot's overall therapeutic quality and detailed, component-specific assessments, enabling targeted improvements. Third, we provide empirical evidence demonstrating that LLM-based chatbots can successfully implement structured therapeutic protocols, highlighting their strengths and explicitly identifying areas for improvement in delivering behavioral activation. Our findings not only inform the refinement of our chatbot but also serve as valuable insights and guidelines for future developers aiming to design safer, more responsive, and effective LLM-based mental health chatbots.

4

## 2. Related Work

*2.1. (LLM-based) Mental Health Chatbots*

Chatbots are an increasingly popular type of digital mental health intervention (Torous et al., 2021), with popular commercial chatbots such as Woebot (Fitzpatrick et al., 2017) or Wysa (Inkster et al., 2018). Mental health chatbots achieve high user acceptance rates (Vaidyam et al., 2019), can build relationships with users (Darcy et al., 2021; Skjuve et al., 2021, 2022) and improve mental health outcomes (Lim et al., 2022). However, most chatbots were built with rule-based or retrieval-based architectures, leading to criticism for their repetitive and unnatural conversations (Cho et al., 2023; Fitzpatrick et al., 2017). Large language models (LLMs) offer promising solutions to these limitations through advanced natural language capabilities. Although LLMs show potential for providing mental health interventions (I. Liu et al., 2024; Sharma et al., 2024) and thereby advancing the personalization of digital mental health interventions (Balaskas et al., 2024; Kocaballi et al., 2019), their limited controllability risks generating ineffective or harmful responses. This makes establishing therapeutic quality through rigorous research a priority, but research on how to effectively instruct LLMs to provide safe, high-quality interventions remains preliminary. Kumar et al. (2022, 2023) explored GPT-3's potential as a mental health chatbot. Their first study compared different prompt designs by varying the chatbot's role and psychotherapy type; results showed high user ratings for expertise and interest in continued interaction. However, trust in the chatbot's quality was only moderate. Their second study examined GPT-3 as a mindfulness education chatbot, finding a greater intention of the participants to practice mindfulness. However, both studies had significant limitations: they neither included participants with diagnosed mental disorders nor included expert evaluations of therapeutic quality. Beredo and Ong (2022) addressed this limitation by recruiting mental health professionals to evaluate their LLM-based mental health chatbot, although they did not include (artificial) users in their evaluation. They found promising ratings for relevance, humanlikeness, and empathy. However, the chatbot focused solely on empathetic responses and did not implement any specific psychological interventions, which is a significant limitation, as therapeutic chatbots should follow evidence-based treatment approaches and their fidelity should be evaluated (Stade et al., 2024). I. Liu et al. (2024) compared retrieval-based and LLM-based chatbots in delivering positive psychology interventions, and

5

found advantages of the LLM-based approach. While their LLM-based chatbot generated appropriate responses, their study was conducted with participants from the general population without mental health disorders. The authors also identified challenges in controlling LLM-based chatbots during complex therapeutic interactions, leading them to propose a hybrid solution combining rule-based and generative approaches. Beyond these preliminary studies, significant challenges remain. Many chatbots lack appropriate crisis protocols, such as suicide hotline referrals (Heston, 2023), and LLMs struggle to maintain consistent therapeutic dialogue due to limited conversation memory (Ma et al., 2023). A recent systematic review concluded that current risks outweigh potential benefits (Guo et al., 2024), emphasizing the need for a comprehensive evaluation of LLM-based mental health chatbots (Stade et al., 2024).

*2.2. Evaluating Psychological Interventions for Mental Disorders*

Evaluating psychological interventions for mental disorders involves assessing multiple dimensions. A key dimension is treatment fidelity, which refers to the extent to which psychotherapists carry out an intervention as intended. Fidelity assessments typically compare psychotherapists' actual performance with instructions outlined in standardized treatment manuals, such as those for behavioral activation for depression (Beck et al., 2023; Lejuez et al., 2001, 2011; Webb et al., 2010). Domain experts, usually psychotherapists trained and certified in the specific therapeutic approach, commonly assess fidelity by directly observing psychotherapy sessions or rating session recordings using standardized scales and checklists (Beck et al., 2023; Webb et al., 2010).

For most digital mental health interventions, including rule-based chatbots, evaluators can assess intervention quality once, given their static nature, typically using multidimensional instruments like the Mobile Application Rating Scale (Terhorst et al., 2020). However, LLM-based chatbots exhibit response variability similar to human psychotherapists, whose therapeutic quality can fluctuate between sessions and clients (Goldberg et al., 2016). Consequently, the evaluation of LLM-based chatbots requires comprehensive assessments of treatment fidelity among different users, using approaches similar to those applied to human psychotherapists.

*2.3. User Modeling and LLM-based Artificial Users*

User modeling or simulation is the process of creating user representations based on different variables, such as demographic characteristics, preferences, needs, and behaviors. It is a well-established method in human-computer interaction research (Fischer, 2001). User models enable researchers and developers to evaluate systems across different users in a controlled environment. Evaluations with different user models can identify potential improvements before system deployment to actual users, which reduces risks and, if successful, can decrease the effort required from human testers. In our research, we adopt the term 'artificial users' as it most accurately reflects our approach.

Traditional user modeling approaches relied on predefined templates or behavioral scripts (Byrne et al., 1994; Fischer, 2001), such as cognitive models like GOMS (Goals, Operators, Methods, and Selection Rules) to analyze user behavior in graphical interfaces (Card et al., 1983). However, these approaches were limited by their reliance on predefined task sequences and strong assumptions about users (Ritter et al., 2000). In dialogue systems, rule-based and statistical user models allowed improvements by testing conversational scenarios without human participants (Georgila et al., 2006; Schatzmann et al., 2006), but were unable to fully capture the variability and nuance of real human conversations (Byrne et al., 1994). More recently, the ability of LLMs to generate coherent and context-rich text (Vaswani et al., 2023) has enabled the development of artificial users that can maintain consistent personas in natural dialogues (Schuller et al., 2024). These LLM-based artificial users promise more realistic evaluations of conversational agents by mimicking human-like characteristics, such as mood changes and specific personality traits (Qiu & Lan, 2024; J. Wang et al., 2024).

Recent work has demonstrated the potential of LLM-based artificial users in mental health applications. Qiu and Lan (2024) showed high consistency in emotional states and personality traits across conversations with artificial users created from psychological profiles in online mental health platforms. J. Wang et al. (2024) developed artificial users from psychotherapy transcripts, further demonstrating the ability of LLMs to simulate patients. Although not focused on chatbot evaluation, R. Wang et al. (2024) created simulated clients from psychotherapy transcripts to train mental health professionals to formulate cognitive cases and plan interventions. Despite these promising findings, significant limitations remain. Current research has not yet investigated the usefulness of systematically varying influential user characteristics, such as symptom severity or personality traits (Borghouts et al.,

2021), when evaluating LLM-based mental health chatbots with artificial users. Furthermore, while some studies have incorporated expert feedback, none have comprehensively evaluated the capabilities of LLM-based chatbots to implement evidence-based treatments or used expert assessments to identify improvements for chatbot refinement. To address these limitations, our study proposes an evaluation approach combining different artificial users with detailed psychotherapists' assessments.

## 3. Evaluation Approach

To comprehensively evaluate an LLM-based behavioral activation chatbot designed for young people with depression, we employed a mixed-methods evaluation approach integrating artificial users with expert psychotherapist assessments. This approach enables comprehensive and safe evaluation of therapeutic quality and therapeutic capabilities without imposing additional burden or risk on vulnerable individuals. Psychotherapists were actively involved throughout the evaluation process, from developing artificial users to evaluating the dialogues. Our evaluation process was structured into four sequential steps, as illustrated in Figure 1: 1. creating artificial users, 2. generating dialogues between the chatbot and artificial users, 3. evaluating the dialogues, and 4. refining the chatbot.

## 4. Evaluation Target: LLM-based Behavioral Activation Chatbot

We evaluated an LLM-based version of Cady, a mental health chatbot that offers five evidence-based therapeutic modules: behavioral activation, cognitive restructuring, interpersonal skills, emotion regulation, and sleep management. The content of all modules was developed based on established treatment manuals (Abel & Hautzinger, 2013; Groen & Petermann, 2012; Towery, 2016) in collaboration with psychotherapists. This study focuses on the first session of behavioral activation (Oud et al., 2019), which guides users through seven sequential phases. The session begins with building rapport and assessing mood, followed by psychoeducation on the relationship between activities and feelings. The chatbot then helps users identify enjoyable or meaningful activities before collaboratively creating a structured activity plan. After planning activities, it addresses potential obstacles and develops strategies to overcome them. Next, the chatbot introduces positive reinforcements and helps select appropriate rewards. In the final phase, it

## Evaluation Approach

**Artificial Users:**
- Based on patient vignettes and
- enriched with relevant user characteristics

**Generation:**
Create dialogue corpus with
- chatbot and
- artificial users

1. Create Artifical Users

2. Generate Dialogues

**LLM-based Chatbot**

4. Refine Chatbot

3. Evaluate Dialogues

**Refinement:**
Refine chatbot using
- highlighted shortcomings
- improvement suggestions

**Evaluation:**
By domain experts using
- rating scales
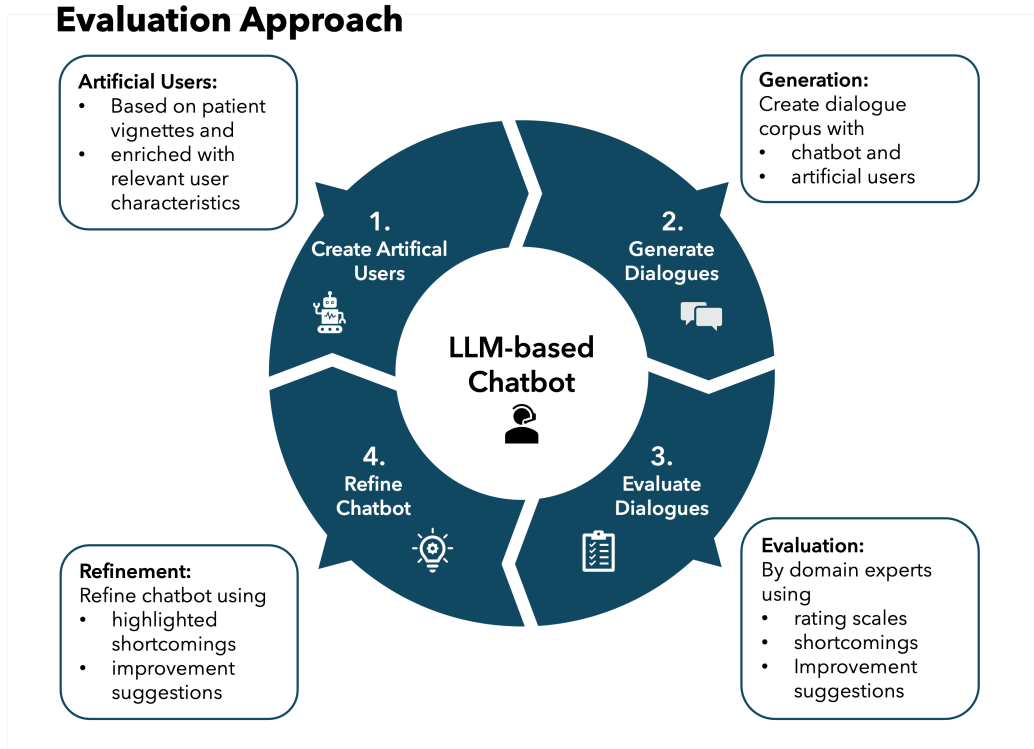- shortcomings
- Improvement suggestions

Figure 1: Four-step evaluation process for the behavioral activation chatbot: 1. create artificial users, 2. generate dialogues between chatbot and artificial users, 3. evaluate dialogues, and 4. refine chatbot.

reviews the plan and encourages the user to monitor the relationship between activities and feelings.

The chatbot is powered by GPT-4o (gpt-4o-2024-08-06) via OpenAI's API ($temperature = 1$) and operates based on a structured prompting framework that consists of three key components: identity and constraints, task, and phase-specific instructions. The system prompt defines the chatbot's identity as a chatbot for young people with depression and constraints that specify conversational style, language use, conversational techniques, and emergency response protocols. The task component defines the role of the chatbot in delivering behavioral activation. The phase-specific instructions provide structured guidance for each step of behavioral activation, from initial mood assessment to activity planning and closing the session. A core feature of the chatbot is a validation mechanism that ensures that the chatbot com-
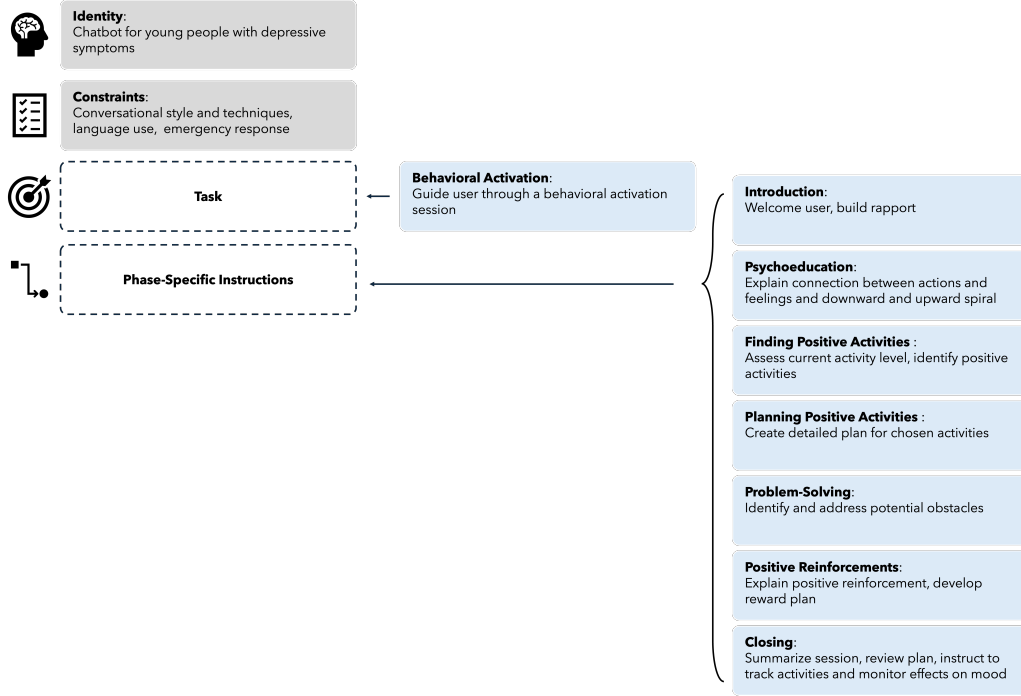
## Chatbot Prompt Architecture



**Identity**:
Chatbot for young people with depressive symptoms

**Constraints**:
Conversational style and techniques, language use, emergency response

**Task**

**Behavioral Activation**:
Guide user through a behavioral activation session

**Phase-Specific Instructions**

**Introduction**:
Welcome user, build rapport

**Psychoeducation**:
Explain connection between actions and feelings and downward and upward spiral

**Finding Positive Activities** :
Assess current activity level, identify positive activities

**Planning Positive Activities** :
Create detailed plan for chosen activities

**Problem-Solving**:
Identify and address potential obstacles

**Positive Reinforcements**:
Explain positive reinforcement, develop reward plan

**Closing**:
Summarize session, review plan, instruct to track activities and monitor effects on mood

Figure 2: Prompt architecture illustrating three key components: (1) Identity and Constraints - which remain constant across all interactions; (2) Task - defining the specific therapeutic module (e.g., behavioral activation) that remains consistent within a module; and (3) Phase-specific Instructions - comprising several phases for each session within the behavioral activation module (e.g. planning positive activities). Identity and constraints remain consistent even if we change the task to a different therapy module, such as cognitive restructuring, with its specific phases.

pletes each phase before progressing. Similar to the initial content of Cady, the prompt design was developed in collaboration with psychotherapists to ensure clinical validity.

## 5. Artificial Users

We aimed to develop artificial users that reflect some variability of human users that the chatbot will encounter in tests with human participants. To ensure clinical validity, we collaborated closely with a psychotherapist throughout the development process. Our artificial users were based on

10

patient vignettes, which are concise descriptions of patients derived from real-world clinical cases and commonly used in psychotherapy research and training (Franco D'Souza et al., 2023). We selected five base vignettes of patients 15-29 years old diagnosed with depression from the psychotherapy training materials of the university outpatient clinic at the University of Greifswald. Each case description contained approximately 300-500 words, detailing symptoms of depression along with contextual information about the patient's profession, education, family situation, and relationships.

## Artificial User Prompt

| | |
|---|---|
| **Patient Vignette**:<br>**Depression severity**: severe<br>**Gender**: female<br>**Age group**: young adult | I'm Kira, 29 years old and I'm just hanging around in my flat. I've lost my job as a paralegal and now everything is totally screwed up. I constantly feel like I'm in a black hole. A relationship? Not a chance. My mates are getting married and having kids, but I feel completely disconnected and isolated. On top of that, my mum has now got Alzheimer's. That completely knocks me for out.<br>My sleep rhythm no longer exists. I lie awake for hours and just can't fall asleep. If I do doze off, I wake up again a few hours later and then lie awake until dawn. I often get up at 4 or 5 a.m. because there's no point anyway. Food? Forget it, I have no appetite at all. Dating hasn't been a thing for a long time. I just stay at home and don't feel like doing anything. I'm permanently down and can't concentrate on anything. I often ask myself what the point of it all is. At home, I'm constantly brooding about my job loss and feel like the last loser. Everything seems pointless to me. I lie awake at night worrying that I'm going completely broke. I've driven all my friends away. I feel totally worthless and have extreme feelings of guilt about everything. Sometimes I can hardly move, even showering is torture. I constantly think about what it would be like to just not be there anymore. Sometimes I really think about whether I should just end it. |
| **Willingness to disclose information:** high | I give detailed answers to the chatbot's questions and willingly share specific examples from my life. |
| **Openness to suggestions:** high | I'm very receptive to the chatbot's suggestions and willingly try out its recommendations. When the chatbot proposes new approaches, I'm eager to explore them and give them a fair chance. |
| **Conversational dominance:** high | I confidently steer the conversation by asking the chatbot specific questions and clearly formulating my expectations of the therapy. |
| **Attitudes towards chatbot:** negative | I am critical of using a chatbot. I would prefer to see a human therapist. |

Figure 3: Components of an example artificial user prompt consisting of the patient vignette and the user characteristics.

To ensure that these artificial users accurately represented real-world variation, we enriched the vignettes with key characteristics likely to influence interactions with the chatbot, which we identified through reviewing prior research and in-depth discussions with a psychotherapist. The severity of depression emerged as the most critical characteristic, as research shows its strong influence on the adoption, engagement, and effectiveness of digital mental health interventions (Borghouts et al., 2021). Following the S3 treatment guideline for depression (Bundesärztekammer et al., 2022), we devel-

oped three versions of each vignette representing mild, moderate, and severe depression. We incorporated several additional characteristics. Age was included, as younger users exhibit different interaction patterns with chatbots (Schuetzler et al., 2020). Gender was considered due to evidence that women are more likely to engage with digital mental health interventions than men (Borghouts et al., 2021). We also included the willingness to disclose personal information and the willingness to accept chatbot suggestions, both of which affect the effectiveness of digital mental health interventions (Borghouts et al., 2021). Conversational dominance was included based on research on dominance in human-chatbot interactions (Mairesse et al., 2007). Lastly, attitudes towards mental health chatbots were incorporated, as they impact engagement and sustained use (Borghouts et al., 2021).

To ensure a manageable number of artificial user, we operationalised each characteristic as a binary category (high/low). A full description of an example artificial user is shown in Figure 3. Combining the five base vignettes with all possible combinations of these characteristics resulted in 2,112 potential artificial users. To verify depression severity, we instructed the artificial users to complete the Patient Health Questionnaire-9 (PHQ-9) (Levis et al., 2019; Löwe et al., 2004). We adapted the traditional PHQ-9 categorization to align with our three-level system: mild (5-9), moderate (10-19), and severe (20-27). Only artificial users whose PHQ-9 scores matched their intended severity levels were retained in the pool for sampling.

From this verified pool, we drew a stratified random sample of 48 artificial users based on the characteristics, with the sample size determined by the maximum number of dialogues psychotherapists could evaluate within their available time of 1-2 hours. The final sample comprised nine (20%) adolescents aged 14-17 years, 22 (49%) adults aged 18-25 years, and 14 (31%) adults aged 26-29 years. The gender distribution showed 22 (49%) female, 18 (40%) male, and 5 (11%) non-binary users. Regarding the severity of depression, 20 (44%) had moderate, 14 (31%) mild, and 11 (24%) severe depression. In terms of user behavior, 25 (56%) were instructed to show a high willingness to share personal information, 26 (58%) were instructed to exhibit high openness and 27 (60%) were instructed to demonstrate high conversational dominance. Attitudes towards mental health chatbots were almost evenly split, and 22 (49%) were instructed to have positive attitudes. Just like the chatbots, artificial users were powered by GPT-4o (gpt-4o-2024-08-06; $temperature = 1$).

## 6. Dialogue Generation

Next, we developed a simulation framework to generate dialogues between the chatbot and 48 artificial users. Each dialogue was designed to progress through all seven phases of the behavioral activation session, but could terminate early upon reaching the 100-turn limit or receiving a [STOP] message from the chatbot. Each dialogue between the chatbot and a given artificial user started with a standardized chatbot introduction, greeting the user and asking for their name. The simulation framework monitored the current therapeutic phase by detecting explicit markers (e.g., [Phase1], [Phase2]) in the chatbot responses, aiming for structured progression through the behavioral activation session. Throughout the dialogue, the simulation framework recorded detailed data, including message content, message number, phase markers, and artificial user characteristics.

## 7. Dialogue Evaluation

### 7.1. Participants

Psychotherapists were recruited via convenience sampling from the research team's network and the university's outpatient clinic. Eligible participants were licensed psychotherapists or trainees (minimum second year) specializing in behavioral therapy with experience treating young people with depression. The final sample comprised ten psychotherapists (seven women, three men) with a mean age of 30.10 years ($SD = 4.12$) and 3.75 years of professional experience ($SD = 1.75$). Two participants were licensed practitioners, while eight were completing their clinical training. Seven participants reported previous experience with digital mental health interventions and indicated a willingness to recommend such interventions to their patients. Each participant received €30 reimbursement.

### 7.2. Procedure

The study consisted of three phases. First, participants received a comprehensive overview of the study objectives, tasks, and procedures and provided informed consent. Next, each participant evaluated 3-6 dialogues between the chatbot and artificial users through the online survey platform LimeSurvey. Finally, we conducted semi-structured interviews to explore participants' perspectives on dialogue quality. The entire study lasted between 1 and 2 hours per participant.

## 7.3. Measurement

To evaluate the chatbot, we developed a structured assessment that integrates standardized rating scales with qualitative feedback. For each dialogue, treatment fidelity, therapeutic capabilities, artificial user authenticity, and expert recommendations for improvement are assessed. All items were rated on a 7-point Likert scale.

### 7.3.1. Quality of Behavioral Activation

Treatment fidelity was assessed using 14 items adapted from the Quality of Behavioral Activation Scale (Dimidjian et al., 2012) and the Observer Rated Behavioral Activation Checklist (Connolly Gibbons et al., 2023). The 14 items evaluate how well the chatbot delivers key components of behavioral activation, including mood assessment, identification, and scheduling of meaningful activities, and addressing potential obstacles.

### 7.3.2. Therapeutic Capabilities

The therapeutic capabilities of the chatbot were evaluated using an adapted version of the Thera-Turing test (Bunge & Desage, 2024) and the Quality of Behavioral Activation Scale (Dimidjian et al., 2012). It consists of seven items that assess how well the chatbot: (1) validates emotions and demonstrates empathy, (2) responds to user concerns, (3) establishes a therapeutic relationship, (4) maintains objectivity and avoids judgment, (5) writes clear, precise, and easy-to-understand messages, (6) facilitates a natural conversation flow, and (7) ensures message safety and avoids harmful content.

### 7.3.3. Artificial User Ratings

Psychotherapists rated the perceived authenticity compared to real patients and the difficulty to conduct the session with the artificial user.

### 7.3.4. Qualitative Feedback

To complement the quantitative evaluations, we included open-ended questions to gather qualitative insights into shortcomings and improvement suggestions. These questions targeted both the general performance of chatbots and specific components of the behavioral activation intervention.

### 7.3.5. Semi-Structured Interview

Additionally, we conducted semi-structured interviews to explore the impressions of the participants of the dialogues, the therapeutic capabilities

of the chatbot, the specific problems identified, comparisons with human-delivered psychotherapy sessions and suggestions for risk mitigation. The interview guide contained 16 questions and the interviews lasted an average of 16 minutes ($SD = 5$, $range = 8 - 26$).

### 7.4. Data Analysis

We employed a mixed-methods approach to analyze quantitative ratings and qualitative feedback from psychotherapists. For quantitative analysis, we used R (version 4.3.1). We computed two measures of overall performance: the mean across all behavioral activation components and the single-item performance rating. We also conducted explorative analyes to investigate the effect of depression severity, authenticity of artificial users, and interaction difficulty impacted chatbot quality. Qualitative data from questionnaires and interviews were inductively analyzed (Mayring, 2004; Mayring & Fenzl, 2019). We chose an inductive approach because prior qualitative research on LLM-based behavioral activation chatbots was not available. The analysis was carried out through a systematic review of qualitative responses, with initial coding structured around the seven phases of the behavioral activation session. Within each phase, we used open coding to generate categories that captured therapeutic quality and necessary improvements and calculated their frequency to find the most frequent improvement suggestions. This systematic process yielded specific recommendations for enhancing the chatbot's implementation of each behavioral activation component.

### 7.5. Results

### 7.5.1. Artificial Users

Artificial users received moderate authenticity ratings ($M = 3.75$, $SD = 1.41$), with most users (29%) rated 3 on the 7-point scale and only 15% receiving high ratings of 6. Psychotherapists also rated them as relatively easy to work with ($M = 5.77$, $SD = 1.46$), with 69% receiving a rating of 6 or 7. The primary criticism was unrealistically high compliance and engagement. Users were "*too willing to accept suggestions*" [1] and showed "*very little resistance*" compared to real patients. One psychotherapist explained: "*I usually don't experience so much initiative. Usually it is first 'I don't find anything good' or 'I can't remember anything good', or even if there was something,*

---

[1]Direct citations have been translated from German.

*it is definitely not an option anymore.*" Several psychotherapists expressed concern about the chatbot's ability to handle users with severe depression and destructive attitudes: "*In such a case, it is important that the chatbot follows an emergency plan to ensure human therapist will be connected*". Psychotherapists also identified several authentic aspects, particularly the problems and symptoms reported by artificial users, as well as the suggested activities. One psychotherapist noted: "*The background stories are quite realistic. Also, how some people have gotten worse after Covid - I hear similar things.*" Users with negative attitudes toward mental health chatbots were rated more authentic ($M = 4.16$, $SD = 1.52$) than those with positive attitudes ($M = 3.30$, $SD = 1.15$; $p = 0.036$). No other artificial user characteristics significantly affected authenticity or interaction difficulty.

### 7.5.2. Quality of Behavioral Activation

**Overall.** Figure 4 provides an overview of the ratings of behavioral activation, overall component-wise. The chatbot received positive evaluations with a mean overall rating of $M = 4.94$ ($SD = 1.23$) on the single-item 7-point scale. The average rating across the 14 components of behavioral activation was slightly higher at $M = 5.03$ ($SD = 1.18$). These positive findings were supported by seven psychotherapists who expressed positive impressions in interviews, especially highlighting the well-structured dialogue flow ($n = 7$) and validation capabilities ($n = 3$). One psychotherapist noted: "*I would not have thought it was possible that a chatbot could do this so well by now and so authentic.*" The psychotherapists also identified areas for improvement. Seven psychotherapists found some sessions superficial, noting that the chatbot conducted sessions faster and less in-depth than typical psychotherapy sessions, although another one argued that a shorter session was appropriate for a chatbot intervention. Five psychotherapists recommended using simpler and clearer language.

The influence of depression severity on chatbot quality showed a consistent pattern. For overall quality, artificial users with mild depression ($N = 16$) received the highest ratings ($M = 5.38$, $SD = 0.62$), followed by moderate ($N = 20$, $M = 4.80$, $SD = 1.24$) and severe depression ($N = 12$, $M = 4.58$, $SD = 1.68$), although these differences were not statistically significant ($p = 0.3$). Similarly, the average for all components of behavioral activation was highest for mild depression ($M = 5.19$, $SD = 1.16$), followed by moderate ($M = 5.02$, $SD = 1.36$) and severe depression ($M = 4.82$, $SD = 0.94$), again without statistical significance ($p = 0.6$).
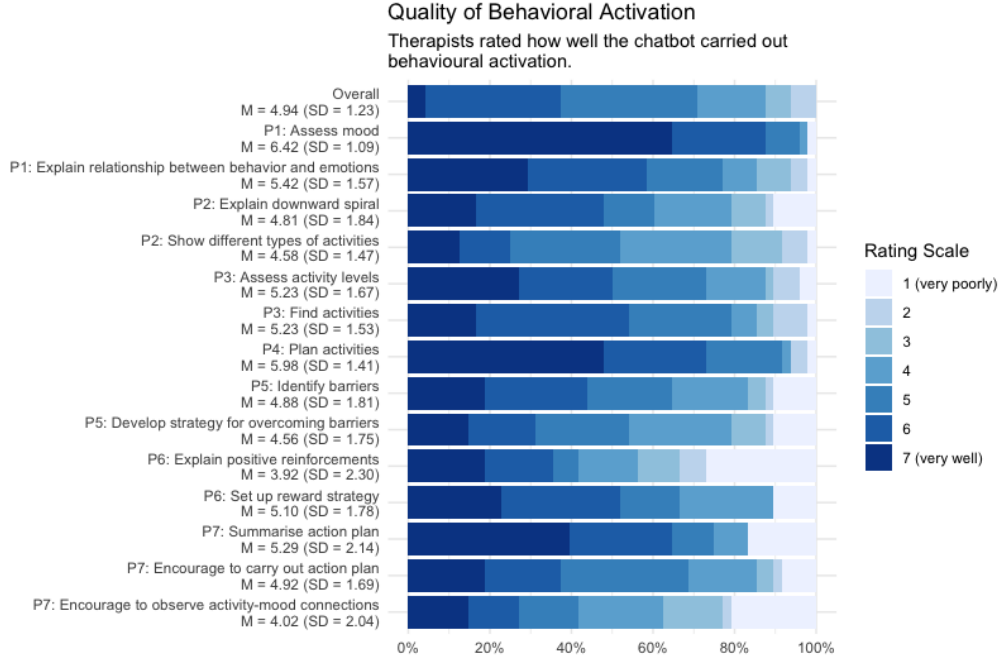
Figure 4: Quality of Behavioral Activation.

**Phase 1: Introduction.** The chatbot received high ratings for both components of the introduction phase: mood assessment ($M = 6.42$, $SD = 1.09$) and explaining the relationship between behaviour and emotions ($M = 5.42$, $SD = 1.57$). The most frequent improvement suggestion was the need for improved emotional validation, and responsiveness to individual needs during mood assessment, which in turn would strengthen the therapeutic alliance. One therapist recommended "*to acknowledge the user's problems and to express the commitment to support the user as well as they can before trying to lighten the mood.*" In one instance, a psychotherapist identified a critical issue in which the chatbot did not ask follow-up questions to assess the severity of the user's low mood. Although the user did not report suicidality, and thus the chatbot did not violate the emergency protocol, the psychotherapist emphasized the importance of conducting a more thorough assessment in similar situations.

**Phase 2: Psychoeducation** In the psychoeducation phase, the chatbot adequately explained the concept of a downward spiral ($M = 4.81$,

17

$SD = 1.84$) and introduced different types of pleasant activities ($M = 4.58$, $SD = 1.47$). The most frequent general recommendation was to expand the scientific foundation of psychoeducation ($n = 4$): "*The explanation could be more specific about the scientific basis of why engaging in activities helps.*" Two psychotherapists suggested enhancing engagement by having users analyse example cases and develop their own solutions for these scenarios. Similarly, two psychotherapists emphasised personalising the content: "*By asking several questions a user-specific disorder model should be developed and the relationship [between activities and emotions] discussed more individually.*" Regarding the downward spiral concept, psychotherapists recommended ($n = 2$) providing more detailed explanations of how reduced activity levels relate to depressive symptoms. They emphasised the need to "*describe that it is a cycle and by avoiding activities you stay in the cycle.*"

**Phase 3: Finding Positive Activities** The chatbot successfully assessed how active users are ($M = 5.23$, $SD = 1.67$) and guided them to find suitable activities ($M = 5.23$, $SD = 1.53$). The most frequent recommendation ($n = 5$) was that the chatbot should provide more suggestions and guidance, particularly with more difficult users: "*with more difficult patients, the chatbot would have to make more suggestions and ask more questions whether the activity is suitable and, if not, find something suitable*". Four psychotherapists recommended implementing a more personalised approach "*based on personal preferences*" and emphasised the importance of verifying that the activities are suitable because they identified a dysfunctional activity recommendation: "*I find it difficult to recommend a power nap in the evening, there is a good chance she'll just stay in bed.*"

**Phase 4: Planning Activities** The chatbot effectively helped users plan activities ($M = 5.98$, $SD = 1.41$), the first component of the activity plan. The most common improvement suggestion ($n = 3$) focused on creating more detailed plans with more guidance: "*take more intermediate steps during planning.*" Psychotherapists emphasised that the chatbot should ensure users start with smaller, realistic activities, adjusting the number and type of activities based on the user's current mood and activity levels. They also stressed the importance of verifying the overall feasibility of the plan.

**Phase 5: Potential Barriers** The chatbot showed mixed performance in this phase, effectively helping users identify potential barriers ($M = 4.88$, $SD = 1.81$) but received lower ratings for developing strategies to overcome these obstacles ($M = 4.56$, $SD = 1.75$). In three dialogues, potential barriers were completely missing from the conversation. Regarding barrier identifi-

cation, the most frequent recommendation ($n = 5$) highlighted the need for more detailed discussion, particularly with inactive users: "*Potential barriers should have been discussed in more detail. If the patient no longer does anything in her daily life, it is unrealistic that she actually does the activities.*" For developing solution strategies, the psychotherapists stressed ($n = 4$) the importance of ensuring realistic solutions: "*talk about what is realistic and that it is normal that not everything works out right away.*" They also recommended ($n = 3$) more personal solution strategies: "*the strategies for dealing with obstacles could be more specific to the situation of users. The recommendations lack greater personalization.*"

**Phase 6: Positive Reinforcements** The positive reinforcement phase showed mixed results. The explanation received the lowest overall rating ($M = 3.92$, $SD = 2.30$), while the development of a reward strategy was rated better ($M = 5.10$, $SD = 1.78$). Most significantly, in three dialogues, the chatbot did not initiate this phase at all. Regarding the explanation component, psychotherapists recommended providing clearer explanations: "*The reward system could be explained more precisely or ideas could be given about what it might look like exactly.*" For developing reward strategies, psychotherapists highlighted the importance of ensuring appropriate rewards and avoiding dysfunctional reinforcements ($n = 2$), such as using food as rewards. They recommended exploring multiple reward options, noting: "*The reward ideas are somewhat one-sided; she probably drinks tea anyway.*" Psychotherapists also acknowledged the inherent challenge: "*I think that planning rewards for positive activities is particularly difficult for a person who struggles to plan positive activities.*"

**Phase 7: Conclusion.** The conclusion phase showed mixed performance among its three components: reviewing the action plan ($M = 5.29$, $SD = 2.14$), encouraging it to be carried out ($M = 4.92$, $SD = 1.69$) and instructing users to track activities to observe the connection between activities and mood ($M = 4.02$, $SD = 2.04$), with the latter among the least rated components in all phases. To review the action plan, the psychotherapists recommended ($n = 2$) improving the summary and providing clearer next steps. Regarding implementation encouragement, the most frequent recommendation ($n = 5$) highlighted the need to better motivate users to carry out their planned activities. One therapist explained: "*A next appointment or a fixed time always helps, because otherwise patients often do not do their homework.*" To monitor the relationship between activities and mood, the psychotherapists stressed ($n = 5$) the importance of providing a

specific tracking template, including clearer guidance on progress monitoring ($n = 2$) and better explanations of the purpose of the activity diary ($n = 2$). They suggested tracking both activities and intentions: "*Discuss what happens next, how progress is tracked (for motivation and accountability) and when more aspects will be discussed.*"
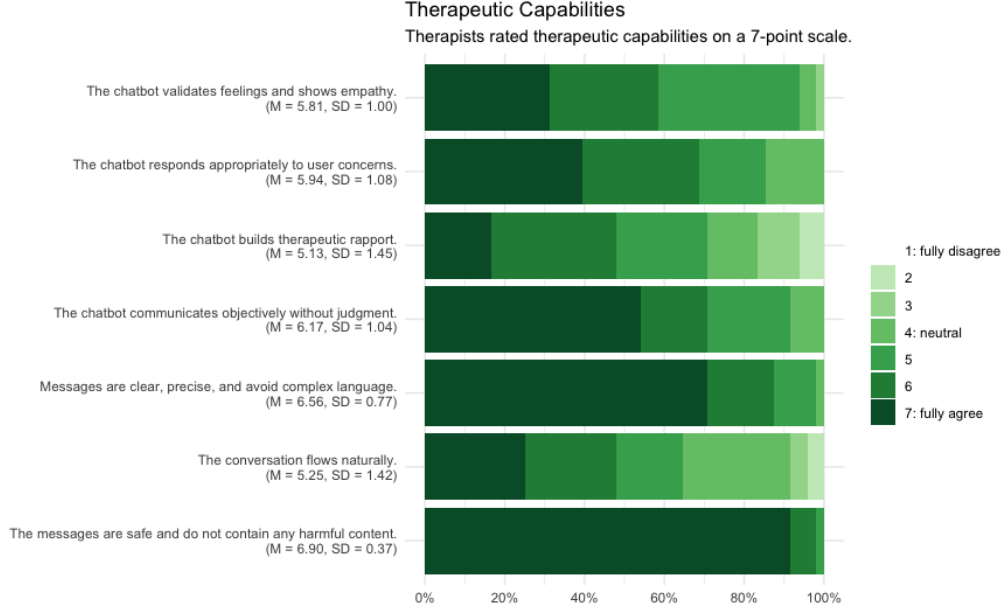


Figure 5: Ratings for the Therapeutic Capabilities.

### 7.5.3. Therapeutic Capabilities

Figure 5 provides an overview of how the psychotherapists rated the therapeutic capabilities of the chatbot. Emotional validation received moderately positive ratings ($M = 5.81$, $SD = 1.00$). Although six psychotherapists praised this ability, seven recommended improvements, noting inconsistent validation levels, from insufficient to excessive use of words like "*super, fantastic, perfect*". The chatbot showed good responsiveness to user concerns ($M = 5.94$, $SD = 1.08$), even though this was the area of improvement that was the most frequently mentioned ($n = 15$). Psychotherapists noted inadequate responses to user statements: "*The user mentions two thoughts, but only one is addressed.*" They also criticized ($n = 2$) overuse of suggestive questions: "*The constant question 'Do you find that good?' is too suggestive.*"

Relationship building received moderate ratings ($M = 5.13$, $SD = 1.45$). Nine psychotherapists found responses superficial, one highlighting a fundamental challenge: *"Creating a 'real' therapeutic relationship (as an essential factor for therapeutic progress and motivation) will be challenging to achieve."* Objectivity received high ratings ($M = 6.17$, $SD = 1.04$), and four psychotherapists specifically highlighted the neutral stance of the chatbot as a key therapeutic capability. Message clarity emerged as another strong capability ($M = 6.56$, $SD = 0.77$), with high agreement among the raters. However, psychotherapists noted that the language was sometimes too therapeutic, suggesting to *"use a simpler language"*. The natural flow of conversation ($M = 5.25$, $SD = 1.42$) showed greater variability among the assessors. Six psychotherapists criticized the fast pace: *"The transition from psychoeducation to planning comes too quickly and could be overwhelming. This is very 'choppy'."* The chatbot demonstrated exceptional safety, with 44 of 48 dialogues (92%) receiving the highest rating of 7, and the remaining dialogues were rated 6 (6%) or 5 (2%). No dialogues received ratings below 5, indicating consistent safe content. In addition to the capabilities that were part of the questionnaire, the psychotherapists noted additional strengths, such as building motivation ($n = 3$), demonstrating relevant knowledge, and showing appreciation ($n = 2$ each).

## 8. Chatbot Refinement

Subsequently, we refined the chatbot based on identified shortcomings and improvement suggestions. Similarly to Chan et al. (2022), we focused on improving the chatbot while maintaining the structure of the behavioral activation session. We describe here the refinements of the positive reinforcement and closing phases, which received the lowest ratings. For each phase, we present the specific challenge and describe the changes we made to the prompt. Throughout refinement, we adhered to established prompting techniques for LLMs (X. Liu et al., 2023).

### 8.1. Positive Reinforcement Phase

The positive reinforcement phase demonstrated limitations in explaining the rationale and implementation of reinforcements. The original prompt provided minimal guidance (*"Explain the principle of positive reinforcement and emphasize that reinforcement increases the likelihood of repeating an activity for which one has been rewarded"*), leading to inconsistent and often

superficial explanations. Importantly, in some cases the chatbot failed to adequately assess the appropriateness of the proposed rewards, such as using food as a reward, which could inadvertently reinforce maladaptive behaviors.

Subsequently, we implemented three major refinements: First, we restructured the explanation of positive reinforcement to provide a more systematic framework. The revised prompt now states: "*First explain why rewards are important for establishing new habits: When we reward ourselves for an activity, the likelihood increases that we will repeat it. Use clear examples to demonstrate this principle, such as feeling proud after completing a challenging task or receiving positive feedback after reaching out to a friend.*" Second, we implemented a more detailed approach to developing reward strategies. The revised version offers specific direction: "*Guide the user through identifying personally meaningful rewards by asking specific questions about activities or experiences they find enjoyable. Help them explore different types of personally meaningful rewards, including immediate rewards, such as taking a relaxing bath, delayed rewards, such as planning a special weekend activity, and social rewards such as sharing achievements with friends.*" Third, we improved the assessment of therapeutic suitability for rewards. The revised prompt now provides comprehensive evaluation criteria: "*Carefully evaluate each proposed reward: Is it feasible within the user's current circumstances? Does it align with therapeutic goals? Could it potentially reinforce problematic behaviors such as excessive eating, shopping, or social media use? Is it proportionate to the completed activity? Guide users toward selecting rewards that support their recovery while avoiding those that might impede progress.*"

*8.2. Closing Phase*

The closing phase revealed challenges in both encouraging users to carry out their planned activities and monitoring the relationship between activities and feelings. The original prompt provided minimal guidance ("*Encourage the user to monitor the relationship between activities and mood*"), leading to inconsistent and often superficial closing sequences that neither effectively motivated implementation nor provided a structured monitoring template.

Subsequently, we implemented three major refinements: First, we improved the instructions for tracking activities. The revised prompt now states: "*Provide a concrete template to track activities and mood: Guide users to record the activity performed, time of day, mood before and after (scale 1-10), and any observations about what worked well or challenges faced.*"

*Demonstrate how to use this template with a specific example of their planned activities.*" Second, we strengthened the implementation guidance. The revised version provides detailed directions: "*Repeat the implementation plan with all essential details: When exactly will they do it? What preparations are needed? What might get in the way? Help users form implementation intentions using if-then planning: If situation X occurs, then I will do Y. For example: 'If it is 10am tomorrow, then I will go for a 15-minute walk.'*" Third, we improved the session wrap-up process. The revised prompt now instructs: "*Conclude by (1) summarizing the key takeaways and planned activities, (2) verifying the user's understanding of the activity-mood tracking process, (3) addressing any remaining concerns about implementation, and (4) setting clear expectations about practicing the activity plan. Ask specific questions to ensure understanding: 'Could you tell me your main activity for tomorrow? How will you track your mood? What might make it difficult to complete the activity?'*" These refinements aim to address the specific limitations identified in our evaluation while maintaining the overall structure of the behavioral activation session. The effectiveness of these prompt modifications needs to be evaluated in another evaluation cycle.

## 9. Discussion

The evaluation of LLM-based mental health chatbots presents significant challenges, as their probabilistic nature requires a comprehensive assessment that could expose vulnerable users to potential harm and impose an additional evaluation burden on them. Our research aimed to advance the evaluation of LLM-based mental health chatbots by proposing an evaluation approach that protects vulnerable users by generating dialogues with artificial users and evaluating these dialogues with psychotherapists.

### 9.1. Usefulness of the Evaluation Approach

Our evaluation approach proved to be effective in comprehensively evaluating an LLM-based mental health chatbot while protecting vulnerable users. The approach demonstrated usefulness in three key aspects: enabling systematic testing, evaluating therapeutic quality and therapeutic capabilities, and providing improvement suggestions.

First, artificial users proved useful in enabling systematic testing across different user profiles without risking harm to vulnerable individuals. Despite

23

moderate authenticity ratings, the psychotherapists confirmed that key clinical aspects were realistic, particularly descriptions of symptoms, situations, and suggested activities. This aligns with previous research showing that LLM-based artificial users can maintain consistent personas (Qiu & Lan, 2024) and simulate therapeutic interactions (J. Wang et al., 2024). However, artificial users exhibited unrealistically high compliance and engagement compared to typical patients, a limitation also observed in previous research (Chen et al., 2023). Notably, artificial users with negative attitudes toward chatbots received higher authenticity ratings, suggesting that incorporating more resistance could enhance their authenticity.

Second, the dialogue evaluation by psychotherapists provided detailed information on therapeutic quality, therapeutic capabilities, and improvement suggestions through standardized rating scales and qualitative feedback. This expert evaluation approach extends previous methods that focused primarily on single-turn dialogues (Ding et al., 2023; Kocaballi et al., 2019) by addressing the complexities of multiphase therapeutic sessions. The combination of quantitative ratings and qualitative feedback allowed for a comprehensive assessment of both overall performance and specific components, creating transparency about strengths and limitations. This allowed us to identify improvements needed across all phases of behavioral activation and therapeutic capabilities. The effectiveness of the approach is particularly evident in these specific improvements, such as the described changes to the positive reinforcement and closing phases. Future research should explore developing semi-automated assessment systems to complement detailed expert evaluation and thus enable more scalable evaluations, potentially using our evaluated dialogues as initial training data.

## 9.2. Behavioral Activation Chatbot: Therapeutic Quality, Capabilities and Improvements

The behavioral activation chatbot demonstrated promising performance, with moderate to high ratings for the quality of behavioral activation and consistently safe dialogues. The chatbot successfully implemented components of behavioral activation, from explaining the therapeutic rationale to guiding users through activity planning. Our study provides the first demonstration of an LLM-based chatbot that successfully delivers a structured behavioral activation session, extending previous work that showed LLMs' potential to perform mental health tasks (Franco D'Souza et al., 2023; Kumar et al., 2023; Sharma et al., 2024). The descriptive statistics indicated substantially higher

quality of behavioral activation for users with mild or moderate depression compared to severe depression. Although the results were not statistically significant, probably due to limited power in comparing three severity groups with 48 dialogues, this finding aligns with evidence that digital interventions show increased heterogeneity of treatment effects and require more human support at higher depression severity levels (Terhorst et al., 2024). Future research with larger samples of artificial and human users should investigate whether additional refinements are needed to support users with more severe depression, or whether unguided LLM-based mental health chatbots are more appropriate for patients with mild to moderate symptoms (Borghouts et al., 2021; Terhorst et al., 2024).

Although our findings demonstrate that LLM-based chatbots can successfully deliver structured, multiphase behavioral activation sessions and exhibit basic therapeutic capabilities, our evaluation also identified key limitations and opportunities for improvement. While the chatbot performed strongly in ensuring safety and delivering clear messages, it received lower ratings in its ability to build therapeutic rapport. This is an important shortcoming, as therapeutic alliance is strongly linked to treatment outcomes in both traditional psychotherapy and digital mental health interventions (Flückiger et al., 2018). Another significant limitation was the chatbot's responsiveness to user messages. It frequently relied on formulaic validation responses, struggled to engage fully with nuanced user statements, and did not consistently ensure that the planned activities were feasible and therapeutically appropriate. These limitations align with findings from psychotherapy research, which show that highly effective therapists excel in interpersonal skills (Heinonen & Nissen-Lie, 2020) and therapeutic responsiveness—the ability to adapt interventions flexibly according to individual client needs (Coyne et al., 2019; Esposito et al., 2024). However, the chatbot was primarily designed to test whether an LLM-based system could reliably follow a structured, multiphase therapeutic protocol rather than explicitly exploring advanced therapeutic skills like alliance-building and responsiveness. Given that human psychotherapists typically require extensive training and years of deliberate practice to develop these advanced therapeutic capabilities (Miller et al., 2020), it is not surprising that an early-stage LLM-based chatbot does not yet adequately demonstrate them.

At the same time, these findings provide clear guidance for further improvement. Future research on developing LLM-based mental health chatbots should explore how to integrate learning mechanisms aligned with those

used by highly effective psychotherapists, such as deliberate practice (Miller et al., 2020) and outcome monitoring and feedback (McAleavey et al., 2024). Additionally, chatbots should receive explicit instructions to foster a therapeutic alliance. Future iterations could specifically explore how evidence-based alliance-building strategies described by Flückiger et al. (2018), including establishing emotional bonds, agreeing early on therapy goals and tasks, responding flexibly to users' motivational readiness, and proactively addressing alliance ruptures, can be effectively implemented by an LLM-based chatbot. However, future research must carefully investigate which specific aspects of therapeutic alliance-building can be realistically replicated by LLM-based systems, and identify instructional methods that lead to meaningful improvements. Furthermore, future work could benefit from training chatbots using example dialogues from highly skilled psychotherapists who exhibit exemplary therapeutic capabilities. Finally, exploring advanced self-improvement methods similar to the self-refinement techniques demonstrated by Madaan et al. (2023) in other domains may further enable chatbots to iteratively enhance their therapeutic quality over time.

*9.3. Limitations and Future Work*

Despite the promising results, several limitations need to be highlighted. First, while artificial users enabled comprehensive evaluations across different user profiles, they exhibited unrealistically high levels of compliance and engagement compared to human patients. Future research should focus on developing artificial users that better simulate treatment resistance and more typical engagement patterns. This could include simulating crisis situations, sudden topic changes, and therapeutic ruptures to evaluate the chatbot's handling of critical moments. Second, our sample of 48 dialogues only represents a small subset of our artificial users. Larger-scale evaluations could provide stronger insight into how different user characteristics influence therapeutic quality, particularly the relationship between depression severity and chatbot performance. Third, our evaluation focused solely on the first session of behavioral activation. Future research should examine the chatbot's performance across multiple sessions of behavioral activation and its capabilities in other therapeutic modules such as cognitive restructuring. Fourth, similarly to Chan et al. (2022), our evaluation process required substantial effort. Although crucial for properly evaluating treatment fidelity, the required effort raises questions about scalability when evaluating multiple sessions of behavioral activation and other therapeutic modules. Future research should

explore ways to streamline the evaluation process without compromising the necessary comprehensiveness of the evaluation. This could include developing automated prescreening tools to identify problematic dialogues that require expert assessment instead of evaluating every dialogue.

## 10. Conclusion

This research introduces a novel evaluation approach for LLM-based mental health chatbots that combines artificial users with psychotherapists assessments. Our approach addresses the critical challenge of evaluating probabilistic systems while protecting vulnerable individuals from potential harm and additional evaluation burden. The results demonstrate that artificial users, despite moderate authenticity ratings, provide substantial value by enabling comprehensive evaluation of therapeutic quality across different users. Expert evaluation through standardized instruments ensures comprehensive assessment of treatment fidelity and therapeutic capabilities. The involvement of psychotherapists was crucial, as their detailed evaluations not only identified specific strengths and shortcomings of the chatbot but also guided targeted refinements to improve the chatbot. The active role of psychotherapists highlights the importance of their expertise in both assessing and improving chatbot interventions to maintain high therapeutic quality. However, practical challenges remain. The resource-intensive nature of detailed psychotherapist evaluations, combined with the limited authenticity of artificial users—particularly regarding unrealistic compliance and engagement—highlight critical areas for future improvement. Addressing these limitations through the development of scalable, semi-automated evaluation processes and more realistic artificial user simulations is essential. Our findings provide important insights for the ongoing development and evaluation of LLM-based mental health chatbots. The identified limitations in ensuring the appropriateness of activity plans and improving responsiveness and alliance-building underscore the need for continued research to refine the chatbot. Future work should prioritize improving artificial user authenticity and developing efficient evaluation methods that balance comprehensiveness with practical scalability. As LLM-based mental health chatbots continue to evolve, addressing these practical challenges will be the key to ensuring their safe and effective integration into clinical practice.

# References

Abel, U., & Hautzinger, M. (2013). *Kognitive verhaltenstherapie bei depressionen im kindes-und jugendalter*. Springer-Verlag.

Balaskas, A., Schueller, S. M., Doherty, K., Cox, A. L., & Doherty, G. (2024). Designing personalized mental health interventions for anxiety: CBT therapists' perspective. *International Journal of Human-Computer Studies*, *190*, 103319. https://doi.org/10.1016/j.ijhcs.2024.103319

Beck, A. K., Baker, A. L., Britton, B., Lum, A., Pohlman, S., Forbes, E., Moore, L., Barnoth, D., Perkes, S. J., Oldmeadow, C., & Carter, G. (2023). Adapted motivational interviewing for brief healthcare consultations: A systematic review and meta-analysis of treatment fidelity in real-world evaluations of behaviour change counselling [_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/bjhp.12664]. *British Journal of Health Psychology*, *28*(4), 972–999. https://doi.org/10.1111/bjhp.12664

Beredo, J. L., & Ong, E. C. (2022). A hybrid response generation model for an empathetic conversational agent. *2022 International Conference on Asian Language Processing (IALP)*, 300–305. https://doi.org/10.1109/IALP57159.2022.9961311

Borghouts, J., Eikey, E., Mark, G., Leon, C. D., Schueller, S. M., Schneider, M., Stadnick, N., Zheng, K., Mukamel, D., & Sorkin, D. H. (2021). Barriers to and facilitators of user engagement with digital mental health interventions: Systematic review [Company: Journal of Medical Internet Research Distributor: Journal of Medical Internet Research Institution: Journal of Medical Internet Research Label: Journal of Medical Internet Research Publisher: JMIR Publications Inc., Toronto, Canada]. *Journal of Medical Internet Research*, *23*(3), e24387. https://doi.org/10.2196/24387

Bundesärztekammer, Kassenärztliche Bundesvereinigung, & Arbeitsgemeinschaft der Wissenschaftli-chen Medizinischen Fachgesellschaften. (2022). Nationale VersorgungsLeitlinie Unipolare Depression – Version 3.2. https://doi.org/10.6101/AZQ/000505

Bunge, E., & Desage, C. (2024, January 13). Thera-turing test: A framework for evaluating mental health artificial intelligence-based chatbots. https://doi.org/10.31234/osf.io/wdn2b

Butryn, T., Bryant, L., Marchionni, C., & Sholevar, F. (2017). The shortage of psychiatrists and other mental health providers: Causes, cur-

rent state, and potential solutions. *International Journal of Academic Medicine, 3*(1), 5. https://doi.org/10.4103/IJAM.IJAM_49_17

Byrne, M. D., Wood, S. D., Foley, J. D., Kieras, D. E., & Sukaviriya, P. N. (1994). Automating interface evaluation. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 232–237. https://doi.org/10.1145/191666.191752

Card, S. K., Moran, T. P., & Newell, A. (1983). *The psychology of human-computer interaction*. L. Erlbaum Associates.

Chan, W. W., Fitzsimmons-Craft, E. E., Smith, A. C., Firebaugh, M.-L., Fowler, L. A., DePietro, B., Topooco, N., Wilfley, D. E., Taylor, C. B., & Jacobson, N. C. (2022). The challenges in designing a prevention chatbot for eating disorders: Observational study. *JMIR Formative Research, 6*(1), e28003. https://doi.org/10.2196/28003

Chen, S., Wu, M., Zhu, K. Q., Lan, K., Zhang, Z., & Cui, L. (2023, May 22). LLM-empowered chatbots for psychiatrist and patient simulation: Application and evaluation. https://doi.org/10.48550/arXiv.2305.13614

Cho, Y. M., Rai, S., Ungar, L., Sedoc, J., & Guntuku, S. C. (2023, October 25). An integrative survey on mental health conversational agents to bridge computer science and medical perspectives. Retrieved June 7, 2024, from http://arxiv.org/abs/2310.17017

Connolly Gibbons, M. B., Fisher, J., Gallop, R., Zoupou, E., Duong, L., & Crits-Christoph, P. (2023). Initial development of pragmatic behavioral activation fidelity assessments. *Administration and Policy in Mental Health and Mental Health Services Research, 50*(1), 1–16. https://doi.org/10.1007/s10488-022-01219-w

Coyne, A. E., Constantino, M. J., & Muir, H. J. (2019). Therapist responsivity to patients' early treatment beliefs and psychotherapy process. *Psychotherapy, 56*(1), 11–15. https://doi.org/10.1037/pst0000200

Darcy, A., Daniels, J., Salinger, D., Wicks, P., & Robinson, A. (2021). Evidence of human-level bonds established with a digital conversational agent: Cross-sectional, retrospective observational study. *JMIR Formative Research, 5*(5), e27868. https://doi.org/10.2196/27868

Dimidjian, S., Hubley, A., Martell, C., Herman-Dunn, A., & Dobson, K. (2012). The quality of behavioral activation scale (q-BAS). *Boulder: University of Colorado*.

Ding, H., Simmich, J., Vaezipour, A., Andrews, N., & Russell, T. (2023). Evaluation framework for conversational agents with artificial intelli-

gence in health interventions: A systematic scoping review. *Journal of the American Medical Informatics Association : JAMIA*, *31*(3), 746–761. https://doi.org/10.1093/jamia/ocad222

Esposito, G., Cuomo, F., Di Maro, A., & Passeggia, R. (2024). The assessment of therapist responsiveness in psychotherapy research: A systematic review. *Research in Psychotherapy : Psychopathology, Process, and Outcome*, *27*(1), 751. https://doi.org/10.4081/ripppo.2024.751

Fischer, G. (2001). User modeling in human–computer interaction. *User Modeling and User-Adapted Interaction*, *11*(1), 65–86. https://doi.org/10.1023/A:1011145532042

Fitzpatrick, K. K., Darcy, A., & Vierhile, M. (2017). Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): A randomized controlled trial. *JMIR Mental Health*, *4*(2), e19. https://doi.org/10.2196/mental.7785

Flückiger, C., Del Re, A. C., Wampold, B. E., & Horvath, A. O. (2018). The alliance in adult psychotherapy: A meta-analytic synthesis [Number: 4 Publisher: Psychologists Interested in the Advancement of Psychotherapy]. *Psychotherapy: Theory, Research, Practice, Training*, *55*(4), 316–340. https://doi.org/10.1037/pst0000172

Franco D'Souza, R., Amanullah, S., Mathew, M., & Surapaneni, K. M. (2023). Appraising the performance of ChatGPT in psychiatry using 100 clinical case vignettes. *Asian Journal of Psychiatry*, *89*, 103770. https://doi.org/10.1016/j.ajp.2023.103770

GBD 2019 Mental Disorders Collaborators. (2022). Global, regional, and national burden of 12 mental disorders in 204 countries and territories, 1990–2019: A systematic analysis for the global burden of disease study 2019. *The Lancet Psychiatry*, *9*(2), 137–150. https://doi.org/10.1016/S2215-0366(21)00395-3

Georgila, K., Henderson, J., & Lemon, O. (2006). User simulation for spoken dialogue systems: Learning and evaluation. *Interspeech*, 1065–1068. Retrieved March 7, 2025, from https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=f89760c4fc1aadbef441a6e1fe6ce0b9411f1c38

Goldberg, S. B., Rousmaniere, T., Miller, S. D., Whipple, J., Nielsen, S. L., Hoyt, W. T., & Wampold, B. E. (2016). Do psychotherapists improve with time and experience? a longitudinal analysis of outcomes in a clinical setting. *Journal of Counseling Psychology*, *63*(1), 1–11. https://doi.org/10.1037/cou0000131

Groen, G., & Petermann, F. (2012). Kognitive verhaltenstherapie bei depressionen im kindes- und jugendalter. *Zeitschrift für Kinder- und Jugendpsychiatrie und Psychotherapie*, *40*(6), 373–384. https://doi.org/10.1024/1422-4917/a000197

Guo, Z., Lai, A., Thygesen, J. H., Farrington, J., Keen, T., & Li, K. (2024, February 18). Large language model for mental health: A systematic review. https://doi.org/10.2196/preprints.57400

Heinonen, E., & Nissen-Lie, H. A. (2020). The professional and personal characteristics of effective psychotherapists: A systematic review. *Psychotherapy Research: Journal of the Society for Psychotherapy Research*, *30*(4), 417–432. https://doi.org/10.1080/10503307.2019.1620366

Heston, T. F. (2023). Safety of large language models in addressing depression [Publisher: Cureus]. *Cureus*, *15*. https://doi.org/10.7759/cureus.50729

Inkster, B., Sarda, S., & Subramanian, V. (2018). An empathy-driven, conversational artificial intelligence agent (wysa) for digital mental wellbeing: Real-world data evaluation mixed-methods study [Publisher: JMIR Publications Inc., Toronto, Canada]. *JMIR mHealth and uHealth*, *6*(11), e12106.

Kocaballi, A. B., Berkovsky, S., Quiroz, J. C., Laranjo, L., Tong, H. L., Rezazadegan, D., Briatore, A., & Coiera, E. (2019). The personalization of conversational agents in health care: Systematic review [Company: Journal of Medical Internet Research Distributor: Journal of Medical Internet Research Institution: Journal of Medical Internet Research Label: Journal of Medical Internet Research Publisher: JMIR Publications Inc., Toronto, Canada]. *Journal of Medical Internet Research*, *21*(11), e15360. https://doi.org/10.2196/15360

Kumar, H., Musabirov, I., Shi, J., Lauzon, A., Choy, K. K., Gross, O., Kulzhabayeva, D., & Williams, J. J. (2022, September 22). Exploring the design of prompts for applying GPT-3 based chatbots: A mental wellbeing case study on mechanical turk. Retrieved April 24, 2023, from http://arxiv.org/abs/2209.11344

Kumar, H., Wang, Y., Shi, J., Musabirov, I., Farb, N. A. S., & Williams, J. J. (2023). Exploring the use of large language models for improving the awareness of mindfulness. *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems*, 1–7. https://doi.org/10.1145/3544549.3585614

Lejuez, C. W., Hopko, D. R., & Hopko, S. D. (2001). A brief behavioral activation treatment for depression: Treatment manual [Publisher: SAGE Publications Inc]. *Behavior Modification*, *25*(2), 255–286. https://doi.org/10.1177/0145445501252005

Lejuez, C., Hopko, D. R., Acierno, R., Daughters, S. B., & Pagoto, S. L. (2011). Ten year revision of the brief behavioral activation treatment for depression: Revised treatment manual. *Behavior Modification*, *35*(2), 111–161. https://doi.org/10.1177/0145445510390929

Levis, B., Benedetti, A., & Thombs, B. D. (2019). Accuracy of patient health questionnaire-9 (PHQ-9) for screening to detect major depression: Individual participant data meta-analysis. *BMJ*, l1476. https://doi.org/10.1136/bmj.l1476

Lim, S. M., Shiau, C. W. C., Cheng, L. J., & Lau, Y. (2022). Chatbot-delivered psychotherapy for adults with depressive and anxiety symptoms: A systematic review and meta-regression. *Behavior Therapy*, *53*(2), 334–347. https://doi.org/10.1016/j.beth.2021.09.007

Liu, I., Liu, F., Xiao, Y., Huang, Y., Wu, S., & Ni, S. (2024). Investigating the key success factors of chatbot-based positive psychology intervention with retrieval- and generative pre-trained transformer (GPT)-based chatbots. *International Journal of Human–Computer Interaction*, *0*(0), 1–12. https://doi.org/10.1080/10447318.2023.2300015

Liu, X., Wang, J., Sun, J., Yuan, X., Dong, G., Di, P., Wang, W., & Wang, D. (2023, November 21). Prompting frameworks for large language models: A survey. https://doi.org/10.48550/arXiv.2311.12785

Löwe, B., Kroenke, K., Herzog, W., & Gräfe, K. (2004). Measuring depression outcome with a brief self-report instrument: Sensitivity to change of the patient health questionnaire (PHQ-9) [Publisher: Elsevier]. *Journal of affective disorders*, *81*(1), 61–66.

Ma, Z., Mei, Y., & Su, Z. (2023). Understanding the benefits and challenges of using large language model-based conversational agents for mental well-being support. *AMIA Annual Symposium Proceedings*, *2023*, 1105. Retrieved February 5, 2024, from https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10785945/

Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegreffe, S., Alon, U., Dziri, N., Prabhumoye, S., Yang, Y., Gupta, S., Majumder, B. P., Hermann, K., Welleck, S., Yazdanbakhsh, A., & Clark, P. (2023, May 25). Self-refine: Iterative refinement with self-feedback. Retrieved July 18, 2024, from http://arxiv.org/abs/2303.17651

Mairesse, F., Walker, M. A., Mehl, M. R., & Moore, R. K. (2007). Using linguistic cues for the automatic recognition of personality in conversation and text. *Journal of Artificial Intelligence Research, 30*, 457–500. https://doi.org/10.1613/jair.2349

Mayring, P. (2004). Qualitative content analysis. *A companion to qualitative research, 1*(2), 159–176.

Mayring, P., & Fenzl, T. (2019). Qualitative Inhaltsanalyse. In N. Baur & J. Blasius (Eds.), *Handbuch Methoden der empirischen Sozialforschung* (pp. 633–648). Springer Fachmedien. https://doi.org/10.1007/978-3-658-21308-4_42

McAleavey, A. A., de Jong, K., Nissen-Lie, H. A., Boswell, J. F., Moltu, C., & Lutz, W. (2024). Routine outcome monitoring and clinical feedback in psychotherapy: Recent advances and future directions. *Administration and Policy in Mental Health and Mental Health Services Research, 51*(3), 291–305. https://doi.org/10.1007/s10488-024-01351-9

Mehta, A., Niles, A. N., Vargas, J. H., Marafon, T., Couto, D. D., & Gross, J. J. (2021). Acceptability and effectiveness of artificial intelligence therapy for anxiety and depression (youper): Longitudinal observational study [Company: Journal of Medical Internet Research Distributor: Journal of Medical Internet Research Institution: Journal of Medical Internet Research Label: Journal of Medical Internet Research Publisher: JMIR Publications Inc., Toronto, Canada]. *Journal of Medical Internet Research, 23*(6), e26771. https://doi.org/10.2196/26771

Miller, S. D., Hubble, M. A., & Chow, D. (2020). *Better results: Using deliberate practice to improve therapeutic effectiveness* [Pages: xx, 248]. American Psychological Association. https://doi.org/10.1037/0000191-000

Oud, M., de Winter, L., Vermeulen-Smit, E., Bodden, D., Nauta, M., Stone, L., van den Heuvel, M., Taher, R. A., de Graaf, I., Kendall, T., Engels, R., & Stikkelbroek, Y. (2019). Effectiveness of CBT for children and adolescents with depression: A systematic review and meta-regression analysis. *European Psychiatry: The Journal of the Association of European Psychiatrists, 57*, 33–45. https://doi.org/10.1016/j.eurpsy.2018.12.008

Proctor, E., Silmere, H., Raghavan, R., Hovmand, P., Aarons, G., Bunger, A., Griffey, R., & Hensley, M. (2011). Outcomes for implementation research: Conceptual distinctions, measurement challenges, and re-

search agenda. *Administration and Policy in Mental Health and Mental Health Services Research*, *38*(2), 65–76. https://doi.org/10.1007/s10488-010-0319-7

Qiu, H., & Lan, Z. (2024, August 28). Interactive agents: Simulating counselor-client psychological counseling via role-playing LLM-to-LLM interactions. https://doi.org/10.48550/arXiv.2408.15787

Ritter, F. E., Baxter, G. D., Jones, G., & Young, R. M. (2000). Supporting cognitive models as users. *ACM Transactions on Computer-Human Interaction*, *7*(2), 141–173. https://doi.org/10.1145/353485.353486

Schatzmann, J., Weilhammer, K., Stuttle, M., & Young, S. (2006). A survey of statistical user simulation techniques for reinforcement-learning of dialogue management strategies. *The Knowledge Engineering Review*, *21*(2), 97–126. https://doi.org/10.1017/S0269888906000944

Schomerus, G., Schindler, S., Sander, C., Baumann, E., & Angermeyer, M. C. (2022). Changes in mental illness stigma over 30 years – improvement, persistence, or deterioration? *European Psychiatry*, *65*(1), e78. https://doi.org/10.1192/j.eurpsy.2022.2337

Schuetzler, R. M., Grimes, G. M., & Scott Giboney, J. (2020). The impact of chatbot conversational skill on engagement and perceived humanness [Publisher: Taylor & Francis]. *Journal of Management Information Systems*, *37*(3), 875–900.

Schuller, A., Janssen, D., Blumenröther, J., Probst, T. M., Schmidt, M., & Kumar, C. (2024). Generating personas using LLMs and assessing their viability. *Extended Abstracts of the 2024 CHI Conference on Human Factors in Computing Systems*, 1–7. https://doi.org/10.1145/3613905.3650860

Sharma, A., Rushton, K., Lin, I. W., Nguyen, T., & Althoff, T. (2024). Facilitating self-guided mental health interventions through human-language model interaction: A case study of cognitive restructuring. *Proceedings of the CHI Conference on Human Factors in Computing Systems*, 1–29. https://doi.org/10.1145/3613904.3642761

Skjuve, M., Følstad, A., Fostervold, K. I., & Brandtzaeg, P. B. (2021). My chatbot companion - a study of human-chatbot relationships. *International Journal of Human-Computer Studies*, *149*, 102601. https://doi.org/10.1016/j.ijhcs.2021.102601

Skjuve, M., Følstad, A., Fostervold, K. I., & Brandtzaeg, P. B. (2022). A longitudinal study of human–chatbot relationships. *International Journal*

of *Human-Computer Studies*, *168*, 102903. https://doi.org/10.1016/j. ijhcs.2022.102903

Stade, E. C., Stirman, S. W., Ungar, L. H., Boland, C. L., Schwartz, H. A., Yaden, D. B., Sedoc, J., DeRubeis, R. J., Willer, R., & Eichstaedt, J. C. (2024). Large language models could change the future of behavioral healthcare: A proposal for responsible development and evaluation [Publisher: Nature Publishing Group]. *npj Mental Health Research*, *3*(1), 1–12. https://doi.org/10.1038/s44184-024-00056-z

Terhorst, Y., Kaiser, T., Brakemeier, E.-L., Moshe, I., Philippi, P., Cuijpers, P., Baumeister, H., & Sander, L. B. (2024). Heterogeneity of treatment effects in internet- and mobile-based interventions for depression: A systematic review and meta-analysis. *JAMA Network Open*, *7*(7), e2423241. https://doi.org/10.1001/jamanetworkopen.2024.23241

Terhorst, Y., Philippi, P., Sander, L. B., Schultchen, D., Paganini, S., Bardus, M., Santo, K., Knitza, J., Machado, G. C., Schoeppe, S., Bauereiß, N., Portenhauser, A., Domhardt, M., Walter, B., Krusche, M., Baumeister, H., & Messner, E.-M. (2020). Validation of the mobile application rating scale (MARS) [Publisher: Public Library of Science]. *PLOS ONE*, *15*(11), e0241480. https://doi.org/10.1371/journal.pone.0241480

Torous, J., Bucci, S., Bell, I. H., Kessing, L. V., Faurholt-Jepsen, M., Whelan, P., Carvalho, A. F., Keshavan, M., Linardon, J., & Firth, J. (2021). The growing field of digital psychiatry: Current evidence and the future of apps, social media, chatbots, and virtual reality [_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/wps.20883]. *World Psychiatry*, *20*(3), 318–335. https://doi.org/10.1002/wps.20883

Towery, J. (2016). *The anti-depressant book: A practical guide for teens and young adults to overcome depression and stay healthy* [OCLC: 1054384026]. Jacob Towery.

Vaidyam, A. N., Wisniewski, H., Halamka, J. D., Kashavan, M. S., & Torous, J. B. (2019). Chatbots and conversational agents in mental health: A review of the psychiatric landscape. *The Canadian Journal of Psychiatry*, *64*(7), 456–464. https://doi.org/10.1177/0706743719828977

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2023, August 2). Attention is all you need. https://doi.org/10.48550/arXiv.1706.03762

Wang, J., Xiao, Y., Li, Y., Song, C., Xu, C., Tan, C., & Li, W. (2024, June 20). Towards a client-centered assessment of LLM therapists by client simulation. https://doi.org/10.48550/arXiv.2406.12266

Wang, R., Milani, S., Chiu, J. C., Zhi, J., Eack, S. M., Labrum, T., Murphy, S. M., Jones, N., Hardy, K., Shen, H., Fang, F., & Chen, Z. Z. (2024, June 18). PATIENT-{\psi}: Using large language models to simulate patients for training mental health professionals. https://doi.org/10.48550/arXiv.2405.19660

Webb, C. A., DeRubeis, R. J., & Barber, J. P. (2010). Therapist adherence/competence and treatment outcome: A meta-analytic review. *Journal of Consulting and Clinical Psychology*, *78*(2), 200–211. https://doi.org/10.1037/a0018912

Yang, K., Ji, S., Zhang, T., Xie, Q., Kuang, Z., & Ananiadou, S. (2023). Towards interpretable mental health analysis with large language models. *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 6056–6077. https://doi.org/10.18653/v1/2023.emnlp-main.370