

Personalizing Mental Health Chatbots for Young People

Zur Erlangung des akademischen Grades
eines Doktors der Wirtschaftswissenschaften

(Dr. rer. pol.)

von der KIT-Fakultät für Wirtschaftswissenschaften
des Karlsruher Instituts für Technologie (KIT)

genehmigte

DISSERTATION

von

Florian Onur Kuhlmeier, M. Sc.

Tag der mündlichen Prüfung

25.06.2025

Referent:

Prof. Dr. Alexander Mädche

Korreferentin:

Prof. Dr. Eva-Lotta Brakemeier

Karlsruhe, Mai 2026

Abstract

Depression is a prevalent mental disorder among young people, yet many affected individuals do not receive adequate support because of stigma, attitudinal barriers, and long waiting times for treatment. Mental health chatbots offer anonymous, easily accessible, and conversational support, but existing systems often fail to address young users' specific needs, the effects of different content personalization strategies remain poorly understood, and the rise of large language models (LLMs) introduces new design and evaluation challenges. This dissertation investigates how personalized mental health chatbots for young people can be designed, instantiated, and evaluated to improve user engagement and effectiveness across three complementary levels: youth-sensitive design, personalized content selection, and the personalization of therapeutic content during ongoing conversations. The first part examines youth-specific needs through a qualitative study of young people with depression and derives design recommendations for chatbots to treat depression, highlighting the importance of age-sensitive content and personalization. The second part develops and evaluates TheraBot in two field experiments that compare standardized content with chatbot-led, user-led, and hybrid content personalization. The results show that content personalization can improve therapeutic outcomes, but only hybrid personalization consistently outperformed standardized content; moreover, the effects of personalization depended on individual characteristics such as depression severity, stated preferences, and personality traits, indicating that the personalization strategy itself should be adapted to the user. The third part investigates how therapeutic content can be personalized during ongoing conversations through an LLM-based behavioral activation chatbot and how such a system can be iteratively refined. Using a structured prompt architecture, artificial users, and clinical expert fidelity assessment, the study shows that the chatbot can deliver behavioral activation with robust protocol adherence in most sessions while also revealing shortcomings in clinical judgment, rapport building, and handling skeptical or information-seeking user steering. Overall, the dissertation provides evidence-based design recommendations, empirically validated knowledge on content personalization, and a reusable evaluation workflow for LLM-based systems, thereby advancing research and practice on personalized mental health chatbots for young people.

Contents

List of Figures	ix
List of Tables	xii
List of Abbreviations	xiii
1 Introduction	1
1.1 Motivation	1
1.2 Research Questions and Thesis Structure	3
2 Foundations	7
2.1 Mental Health Chatbots	7
2.2 Personalization	8
2.2.1 Personalization of Information Systems	9
2.2.2 Personalization of Psychotherapy	9
2.2.3 Personalization of Mental Health Chatbots	10
3 Part I: How to Design a Chatbot to Treat Depression among Youth? Insights from a Qualitative Study	13
3.1 Introduction	13
3.2 Methods	15
3.2.1 Study Design	15
3.2.2 Participants	15
3.2.3 Procedure	16
3.2.4 Material	16
3.2.5 Data Analysis	20
3.3 Results	20
3.3.1 Participant Characteristics	20
3.3.2 Research Question 1: Problems with Depression	21

3.3.3	Research Question 2: Coping Strategies	22
3.3.4	Research Question 3: Attitudes and Expectations toward Chatbots to Treat Depression	23
3.3.5	Research Question 4: Chatbot Design Preferences	24
3.4	Discussion	27
3.4.1	Primary Findings	27
3.4.2	Limitations	32
3.4.3	Conclusions and Future Work	33
4	Part II: User-Led, Chatbot-Led, or Hybrid? Design and Effect of Content Personalization Approaches in Mental Health Chatbots for Young People	35
4.1	Introduction	35
4.2	Background	38
4.2.1	Digital Mental Health Interventions	38
4.2.2	Personalization	38
4.2.3	Personalization of Digital Mental Health Interventions	41
4.3	Designing Content Personalization in Mental Health Chatbots	42
4.3.1	Design Process	42
4.3.2	Design Artifact: TheraBot	43
4.4	Study 1: Investigating the Effect of Content Personalization	49
4.4.1	Experimental Design	49
4.4.2	Participants	50
4.4.3	Procedure	50
4.4.4	Measurement	52
4.4.5	Data Analysis and Results	52
4.5	Study 2: Examining the Effects of Chatbot-led, User-led, and Hybrid Content Personalization	54
4.5.1	Experimental Design	55
4.5.2	Participants	56
4.5.3	Procedure	56

4.5.4	Measurement	56
4.5.5	Data Analysis	58
4.6	Discussion	66
4.6.1	Implications for Research	68
4.6.2	Implications for Practice	70
4.6.3	Limitations and Future Research	71
5	Part III: Designing an LLM-Based Behavioral Activation Chatbot for Young People with Depression: Insights from Artificial Users and Clinical Experts	73
5.1	Introduction	73
5.2	Background	75
5.2.1	Mental Health Chatbots	75
5.2.2	Behavioral Activation	76
5.2.3	Evaluating LLM-Based Mental Health Chatbots	77
5.3	Designing an LLM-based Behavioral Activation Chatbot	78
5.3.1	Intervention Development	79
5.3.2	Prompt Architecture	79
5.4	Evaluation Study	81
5.4.1	Generating Behavioral Activation Sessions with Artificial Users	81
5.4.2	Clinical Expert Assessment	85
5.4.3	Results	88
5.5	Discussion	95
5.5.1	LLMs Can Adhere to Instructions but Struggle with Clinical Reasoning	96
5.5.2	Tension between Ensuring Protocol Adherence and Allowing for Personalized Intervention Delivery	97
5.5.3	Formulaic Validation versus Genuine Responsiveness: Limitations in Relational Capabilities	98
5.5.4	Usefulness of the Evaluation Approach	99
5.5.5	From Evaluation Insights to Targeted Prompt Refinement	101
5.5.6	Limitations and Future Work	104

5.6	Conclusion	105
6	Discussion	107
6.1	Research Contributions	109
6.2	Practical Contributions	113
6.3	Limitations and Future Research	116
6.4	Outlook	117
7	Conclusion	121
	Bibliography	148
	Appendix	153
A	Supplementary Material for Part I	153
A.1	Interview Guide	153
A.2	Category System	157
B	Supplementary Material for Part II	167
B.1	Review of Commercial Mental Health Chatbots	167
B.2	Implicit Personalization	172
B.3	Study 1 and 2	173
B.4	Robustness Checks and Additional Analyses	178
C	Supplementary Material for Part III	181
C.1	Quality of Behavioral Activation	181
C.2	Therapeutic Capabilities	183
C.3	Effects of Artificial User Characteristics	184
C.4	Prompt of the Behavioral Activation Chatbot	190
C.5	Artificial Users	212
C.5.1	Overview and Variation Matrix	212
C.5.2	Base Personas	212

C.5.3 Variation Dimensions 216

List of Figures

3.1	Screenshot of an Example Conversation between Cady and a Participant	19
3.2	Results of the Chatbot Capabilities Questionnaire	25
4.1	Design Science Research Approach based on Peffers et al. (2007)	43
4.2	TheraBot’s System Architecture	47
4.3	Screenshots of TheraBot’s Personalization Features	48
4.4	Changes in depression severity by personalization approach.	53
4.5	Number of completed modules by personalization approach.	54
4.6	Changes in Depression Severity for each Experimental Condition	61
4.7	Number of Completed Modules across the three personalization approaches	63
4.8	Effect of Pre-Study PHQ-8 score on Number of Completed Modules across the three personalization approaches	64
5.1	Quality of Behavioral Activation	88
5.2	Therapeutic Capabilities	93
C.1	Quality of Behavioral Activation Heatmap	181
C.2	Therapeutic Capabilities Heatmap	183

List of Tables

1.1	Dissertation Structure and Publications	5
3.1	Items of the Chatbot Capabilities Questionnaire	17
3.2	Excerpt from the Interview Guide	18
3.3	Participant Characteristics	21
3.4	Key Design Recommendations	28
4.1	Sample Characteristics of Study 1	51
4.2	Results of Linear Mixed Effects Models for Hybrid vs. Standardized Content in Study 1	53
4.3	Experimental Conditions	55
4.4	Sample Characteristics of Study 2	57
4.5	Key Variables Across Experimental Groups in Study 2	59
4.6	Results of Linear Mixed Effects Models for Personalized vs. Standardized Content in Study 2 and from both Studies Combined	60
4.7	Results from Pairwise Contrasts between the four Experimental Conditions in Study 2 and the Combined Data	62
4.8	Final Design Requirements	70
5.1	Prompt Architecture Components and Functions	80
5.2	Artificial User Characteristics and Rationale for Selection	82
5.3	Example of an Artificial User Profile	83
5.4	Participant Characteristics	85
6.1	Research Contributions of this Dissertation	109
6.2	Practical Contributions of this Dissertation	114
A.1	Interview Guide for Designing Chatbots to Treat Depression in Youth	153
A.2	Problems with Depression	157
A.3	Adaptive Coping Strategies	159

A.4	Attitudes and Expectations	161
A.5	Design Preferences	163
B.1	Review of Commercial Mental Health Chatbots.	168
B.2	Items used for Implicit (Chatbot-led) Personalization	172
B.3	Measurement Items for Study 1 and 2.	173
B.4	Sample Characteristics for Study 1 and 2.	175
B.5	Results of Linear Mixed Effects Models for Study 1 ($N = 74$).	176
B.6	Results of Linear Mixed Effects Models for Study 2 ($N = 160$).	176
B.7	Results of Linear Mixed Effects Models for Combined Data ($N = 234$).	177
B.8	Results from Pairwise Contrasts between the Personalization Approaches in Study 2 and the Combined Data.	177
C.1	Q-BAS component adequacy rates across sessions	182
C.2	Effects of Artificial User Characteristics on Single-Item Holistic Rating	184
C.3	Effects of Artificial User Characteristics on Q-BAS Average	185
C.4	Effects of Artificial User Characteristics on Safety Ratings	185
C.5	Effects of Artificial User Characteristics on Natural Conversation Flow	186
C.6	Effects of Artificial User Characteristics on Message Clarity	186
C.7	Effects of Artificial User Characteristics on Objectivity	187
C.8	Effects of Artificial User Characteristics on Therapeutic Rapport	187
C.9	Effects of Artificial User Characteristics on Response to User Concerns	188
C.10	Effects of Artificial User Characteristics on Validation and Empathy	188
C.11	Effects of Artificial User Characteristics on Authenticity	189
C.12	Effects of Artificial User Characteristics on Difficulty	189
C.13	Overview of variation dimensions for artificial users.	212

List of Abbreviations

BADS-SF	Behavioral Activation for Depression Scale—Short Form
CBT	Cognitive Behavioral Therapy
CES-D	Center for Epidemiologic Studies Depression Scale
CFQ	Cognitive Fusion Questionnaire
DMHIs	Digital Mental Health Interventions
DSM-5	Diagnostic and Statistical Manual of Mental Disorders, Fifth Edition
DSR	Design Science Research
ERSQ	Emotion Regulation Skills Questionnaire
IS	Information Systems
K-SADS-PL	Schedule for Affective Disorders and Schizophrenia for School-Age Children, Present and Lifetime Version
LLMs	Large Language Models
NLP	Natural Language Processing
PHQ	Patient Health Questionnaire
Q-BAS	Quality of Behavioral Activation Scale
SCARED	Screen for Child Anxiety Related Disorders

1 Introduction¹

1.1 Motivation

Depression is a common mental disorder among young people between 15 and 29 years (European Union, 2023) that significantly reduces their quality of life and is a major risk factor for disability and suicide (GBD 2019 Mental Disorders Collaborators, 2022; Thapar et al., 2012). Despite the availability of effective treatments, many young people do not seek professional help due to attitudinal barriers, such as stigma, a preference to solve their problems on their own, or skepticism of traditional psychotherapy (Ebert et al., 2019; Gulliver et al., 2010). Even among those who seek treatment, structural barriers, including waiting periods of several months, often impede access to care (Thapar et al., 2012).

These challenges have spurred the development of digital mental health interventions (DMHIs), among which chatbots (e.g. Woebot, Wysa, Youper) have become particularly popular (Darcy et al., 2021; Mehta et al., 2021). Mental health chatbots offer anonymous, easily accessible, and conversational support (Darcy et al., 2021; Mehta et al., 2021). Although recent studies have shown that chatbots can alleviate depressive symptoms (Fitzpatrick et al., 2017; Lim et al., 2022; Vaidyam et al., 2019) and are well accepted by young people (Mehta et al., 2021), user engagement and effectiveness remain limited (Lim et al., 2022; Limpanopparat et al., 2024).

This thesis identifies and addresses three critical issues that underlie these limitations. First, mental health chatbots have been designed and tested primarily with adult populations (Abd-Alrazaq et al., 2020; Bendig et al., 2019; Fitzpatrick et al., 2017; Inkster et al., 2018; Vaidyam et al., 2020), a significant shortcoming for several reasons. One, depression manifests differently in youth because of ongoing developmental changes in biological, psychological and social systems (Cicchetti & Toth, 2009), with distinct symptom patterns, particularly during puberty (Rice et al., 2019). Two, young people interact with smartphones and chatbots differently than adults (Andone et al., 2016; Huffman, 2014). Three, youth consistently report that existing DMHIs do not address their specific needs, concerns, and communication styles adequately (Agapie et al., 2022). Four, studies that have examined mental health chatbots for youth have focused predominantly on preventing mental disorders instead of treating depression (Grové,

¹This chapter is based on the following studies, which are published, available as preprints, or under review: (Kuhlmeier et al., 2022a; Kuhlmeier et al., 2022b; Kuhlmeier et al., 2025a; Kuhlmeier et al., 2025b; Kuhlmeier et al., 2026).

2021; Høiland et al., 2020).

Second, the design and effects of content personalization approaches in mental health chatbots are poorly understood, which contributes to limited user engagement and effectiveness (Kocaballi et al., 2019). In traditional psychotherapy, personalization is considered crucial for successful treatment (Cohen et al., 2021). Psychotherapists personalize treatment by matching patient characteristics and preferences with therapeutic options (Cohen et al., 2021; Stumpp & Sauer-Zavala, 2022), such as selecting the most relevant cognitive behavioral therapy (CBT) modules based on clinical expertise and patient needs. Although commercial mental health chatbots attempt to replicate this by selecting relevant modules for users, no prior study had evaluated content personalization in mental health chatbots (Kocaballi et al., 2019). Evidence from other types of DMHIs was also discouraging: a recent systematic review found that five of six randomized trials reported no benefit of personalized over standardized content (Schaeuffele et al., 2025). Importantly, these studies examined only purely chatbot-led or user-led personalization, and the relative effects of chatbot-led, user-led, and hybrid approaches had not been empirically compared (Kankanhalli et al., 2021; Kocaballi et al., 2019; Schaeuffele et al., 2025). Understanding which strategy works best, whether a single approach suffices for all users, and how individual characteristics and stated preferences, such as depression severity and personality traits, interact with personalization approaches is necessary to increase user engagement with and the effectiveness of such mental health chatbots.

Accordingly, this dissertation addresses personalization at two complementary levels: Part II focuses on content selection, that is, choosing the most relevant therapy modules, whereas Part III focuses on personalizing therapeutic content during ongoing conversations.

Third, the recent emergence of LLM-based systems opens new possibilities for overcoming the limitations of earlier rule-based mental health chatbots, which relied on predetermined dialogue paths and often felt rigid or repetitive (Fitzpatrick et al., 2017). However, designing LLM-based mental health chatbots remains challenging because their probabilistic nature can produce inconsistent or potentially harmful responses (Heston, 2023). Existing work has often emphasized ad-hoc metrics, user experience, or downstream outcomes without first establishing whether such systems deliver structured interventions with adequate therapeutic quality (Y. Hua et al., 2025). To support design, approaches are needed that can expose chatbots to diverse interaction patterns, identify concrete shortcomings, and guide iterative refinement without exposing vulnerable users to untested systems. Combining artificial users with clinical expert

fidelity assessment offers a promising way to meet this need.

1.2 Research Questions and Thesis Structure

In this dissertation, I investigate the design of personalized mental health chatbots to improve user engagement and effectiveness by addressing these three issues. I organize the dissertation around the following overarching research question:

RQ: How can personalized mental health chatbots for young people be designed to improve user engagement and effectiveness?

In this dissertation, I focus on young people as the overarching target group. Within this broader population, I focus more specifically on youth in **Part I** to identify age-sensitive needs and derive initial design requirements. In the subsequent parts, I build on these insights while addressing young people more broadly.

In **Part I**, I present a qualitative study in which we investigate how to design a chatbot to treat depression in youth and derive specific design recommendations. The study addresses the following research questions:

- RQ1a: What problems do youth with symptoms of depression face?
- RQ1b: What adaptive coping strategies do they apply?
- RQ1c: What attitudes and expectations do they have for chatbots designed to treat depression?
- RQ1d: What are their design preferences?

In **Part II**, I present research on the design and effects of content personalization approaches in mental health chatbots, specifically focusing on the selection of therapy modules to form a personalized treatment plan. Following the design science research (DSR) paradigm, we develop and evaluate a chatbot for young people with depression that uses hybrid content personalization and then compare the effects of three distinct personalization approaches in two field experiments. In this work, we address the following research questions:

- RQ2a: How can content personalization in a mental health chatbot be designed for young people with depression?

- RQ2b: Does personalized content improve user engagement and effectiveness compared to standardized content?
- RQ2c: What are the relative effects of chatbot-led, user-led, and hybrid content personalization approaches on user engagement and effectiveness?

By addressing these questions, we provide recommendations for designing effective content personalization in mental health chatbots and help explain inconsistencies in previous research. The studies show that hybrid personalization—where the chatbot recommends modules but users can adapt the selection—is the most effective single strategy, while the optimal approach for individual users depends on their symptom severity, stated preferences, and personality characteristics.

In **Part III**, I present a study in which we investigate how to design and iteratively refine an LLM-based behavioral activation chatbot for young people with depression using artificial users and clinical expert fidelity assessment. The study addresses the following research questions:

- RQ3a: How can an LLM-based behavioral activation chatbot be designed to maintain clinical fidelity while delivering a structured psychological intervention?
- RQ3b: What design limitations and opportunities for refinement can be identified through artificial users and clinical expert fidelity assessment?

This fidelity assessment showed that the chatbot generally followed the behavioral activation protocol, but systematic weaknesses emerged in clinical judgment, such as assessing whether selected activities and rewards were feasible and therapeutically appropriate and maintaining therapeutic direction when users steered the conversation off course. These findings informed concrete prompt-level refinements and broader architectural considerations. Thus, this part provides an empirically grounded starting point for designing LLM-based mental health chatbots for young people and illustrates how artificial users and clinical expert fidelity assessment can support iterative refinement.

Together, these three parts allow me to address the overarching research question of this dissertation and contribute theoretically and practically to the design of more effective personalized mental health chatbots for young people. The dissertation is structured as follows (see Table 1.1):

Ch.	Chapter and Publication
1	Motivation, Research Questions, Thesis Structure
2	Foundations
3	<p>Part 1: How to Design a Chatbot to Treat Depression among Youth? Insights from a Qualitative Study. (RQ1) Kuhlmeier, F. O., Bauch, L., Gnewuch, U., & Lüttke, S. (2025). Designing Chatbots to Treat Depression in Youth: Qualitative Study. <i>JMIR Human Factors</i>, 12(1), e66632. https://doi.org/10.2196/66632</p>
4	<p>Part 2: User-Led, Chatbot-Led, or Hybrid? Design and Effect of Content Personalization Approaches in Mental Health Chatbots for Young People. (RQ2) Kuhlmeier, F. O., Gnewuch, U., Scheu, S., Luettke, S., & Maedche, A. (2025). User-Led, Chatbot-Led, or Hybrid? Design and Effect of Content Personalization approaches in Mental Health Chatbots for Young People. Under review at the <i>Journal of the Association for Information Systems (JAIS)</i>.</p>
5	<p>Part 3: Designing an LLM-Based Behavioral Activation Chatbot for Young People with Depression: Insights from Artificial Users and Clinical Experts. (RQ3) Kuhlmeier, F. O., Hanschmann, L., Rabe, M., Luettke, S., Brakemeier, E.-L., & Maedche, A. (2026). A Large Language Model-Based Behavioral Activation Chatbot for Young People with Depression: Mixed-Methods Evaluation Using Artificial Users and Clinical Experts. <i>JMIR Preprints</i>. https://doi.org/10.2196/preprints.94781. Under review at <i>JMIR Mental Health</i>.</p>
6	Discussion
7	Conclusion

Table 1.1: Dissertation Structure and Publications

After motivating the research, articulating the research questions (Chapter 1), and presenting relevant foundations (Chapter 2), I use the three principal parts of this dissertation to address the identified research gaps and answer the research questions. In Chapter 3, I present a qualitative study in which we investigate how to design chatbots for youth with depression. In Chapter 4, I present two field experiments in which we explore the design and effects of different content personalization approaches in mental health chatbots. In Chapter 5, I present a study in which we investigate how to design and iteratively refine an LLM-based behavioral activation chatbot for young people using artificial users and clinical expert fidelity assessment. In Chapter 6, I discuss the findings and their theoretical and practical contributions, examine overarching limitations, and explore future research directions for personalized mental health chatbots. In Chapter 7, I conclude the dissertation.

2 Foundations²

2.1 Mental Health Chatbots

DMHIs are software applications that provide mental health support through technologies such as computers, tablets, and mobile phones (Schueller & Torous, 2020). Chatbots have become a particularly popular type of DMHIs by providing support through conversational interfaces (Mehta et al., 2021; Vaidyam et al., 2019). They provide therapeutic content from evidence-based psychotherapy approaches, such as CBT (Martinengo et al., 2021), through structured dialogues, including psychoeducation about mental health conditions, behavioral activation, cognitive restructuring, problem solving techniques and relaxation and mindfulness exercises (Lin et al., 2023). Mental health chatbots also monitor users' mental health, for example, by assessing mood or mental health symptoms and their progress with therapeutic exercises. Their goal is to provide support through empathetic responses, non-judgmental listening, crisis resources when needed, and ongoing encouragement and motivation (Martinengo et al., 2022) as well as checking in regularly and using gamification elements (Lin et al., 2023; Martinengo et al., 2022). Research has shown that mental health chatbots' conversational nature improves user engagement (Perski et al., 2019) and supports the formation of a working alliance (Darcy et al., 2021), an essential component of effective psychotherapy (Cameron et al., 2018). It also indicates that chatbots can reduce depression severity (Lim et al., 2022; Linardon et al., 2022). However, maintaining users' engagement and increasing effectiveness remains a challenge (Limpanopparat et al., 2024).

A key limitation is that research on mental health chatbots to date has primarily targeted adult populations (Abd-Alrazaq et al., 2020; Bendig et al., 2019). These findings cannot be generalized straightforwardly to young people, because developmental, psychological, and social characteristics differ from those of adults, particularly during youth (Cicchetti & Toth, 2009). Depressive symptoms and interaction patterns with digital technologies also differ from those of adults (Andone et al., 2016; Huffman, 2014; Rice et al., 2019). Furthermore, young people have expressed that DMHIs often do not adequately address their specific concerns (Agapie et al., 2022). Recent exploratory studies have begun to address some of these needs, such as

²This chapter is based on the following studies, which are published, available as preprints, or under review: (Kuhlmeier et al., 2022a; Kuhlmeier et al., 2022b; Kuhlmeier et al., 2025a; Kuhlmeier et al., 2025b; Kuhlmeier et al., 2026).

preferences for peer support, engaging personalities, and humor (Grové, 2021; Høiland et al., 2020), but they focus primarily on prevention rather than the treatment of depression.

Another key limitation is that content personalization in mental health chatbots has largely been unexplored. First, content personalization approaches in mental health chatbots remain poorly understood, as no prior study had evaluated their effects (Kankanhalli et al., 2021; Kocaballi et al., 2019). Evidence from other types of DMHIs was also discouraging: a recent systematic review found that five of six randomized trials reported no benefit of personalized over standardized content (Schaeuffele et al., 2025). Importantly, these studies examined only purely chatbot-led or user-led personalization, and the relative effects of different content personalization approaches had not been empirically compared (Kankanhalli et al., 2021; Kocaballi et al., 2019; Schaeuffele et al., 2025). Commercial mental health chatbots nevertheless employ diverse personalization features, leaving limited evidence-based guidance for designing content personalization. Second, the recent emergence of LLM-based systems opens new possibilities for overcoming the limitations of earlier rule-based and retrieval-based mental health chatbots, which relied on predetermined dialogue paths and often felt rigid or repetitive (Cho et al., 2023; Fitzpatrick et al., 2017). LLMs promise more flexible, natural, and highly personalized conversations. However, designing LLM-based mental health chatbots remains challenging because their probabilistic nature can produce inconsistent or potentially harmful responses (Heston, 2023). Existing work has often emphasized ad-hoc metrics, user experience, or downstream outcomes without first establishing whether such systems deliver structured interventions with adequate therapeutic quality (Heinz et al., 2025; Y. Hua et al., 2025; Thieme et al., 2020). To support design, methods are needed that can expose chatbots to diverse interaction patterns, identify concrete shortcomings, and guide iterative refinement without exposing vulnerable users to untested systems.

2.2 Personalization

The Cambridge Dictionary defines personalization as "the act of making something suitable for the needs of a particular person" (Cambridge Dictionary, 2023). In IS and related disciplines, personalization refers to adapting systems to user requirements to increase their effectiveness in helping users achieve their goals (Fan & Poole, 2006). However, there is little consensus on the way to conceptualize personalization, as the term is often used interchangeably with "cus-

tomization", "adaptation", or "individualization" (Fan & Poole, 2006). This thesis adopts Fan and Poole's definition: "a process that changes the functionality, interface, information access and content, or distinctiveness of a system to increase its personal relevance to an individual" (Fan & Poole, 2006, p. 183).

2.2.1 Personalization of Information Systems

Personalization of IS can be described according to three key dimensions (Benlian, 2015; Fan & Poole, 2006): the element of the system that is personalized (content, interface, functionality, channel), the target of personalization (individuals or categories of individuals), and who performs the personalization (implicit or explicit). One can broadly distinguish between content personalization and design personalization (Benlian, 2015; Fan & Poole, 2006). In explicit (user-led) personalization (Adomavicius & Tuzhilin, 2005; Fan & Poole, 2006), users specify their preferences, whereas in implicit (system-led) personalization, the system determines content based upon user data (Ebrahimi et al., 2022). These two approaches can be combined into a hybrid approach (Fan & Poole, 2006). IS research has shown that personalization can enhance user outcomes across multiple domains. In e-commerce, personalization can reduce information overload and improve attention and engagement (Tam & Ho, 2005; Xiao & Benbasat, 2007). In healthcare, it has shown benefits for personalized treatment planning and content recommendations (W. Chen et al., 2021; Zhou et al., 2023). However, the effects of personalization depend on system design, context, and user characteristics (Ebrahimi et al., 2022; T. Li & Unger, 2012; Liang et al., 2006; Thirumalai & Sinha, 2013). These differences suggest that optimal personalization approaches cannot be transferred straightforwardly across domains and usage contexts.

2.2.2 Personalization of Psychotherapy

In the context of psychotherapy, personalization is important because patients differ in their symptoms, preferences, and responsiveness to different therapy approaches (Cohen et al., 2021). Therefore, it is important to develop and implement personalized therapy approaches to achieve optimal treatment outcomes. The overarching aim of personalized therapy is to match patients with their specific symptom profiles and characteristics to available treatment options (Cohen et al., 2021). However, this process can become complex because of multiple therapeutic di-

mensions and practical implementations (Stumpp & Sauer-Zavala, 2022). Cohen et al. (2021) distinguish personalization along three dimensions: (1) level of intervention, (2) structure, and (3) time. The level of intervention refers to the level at which treatment is personalized, such as selecting the overall psychotherapy modality (e.g. psychodynamic psychotherapy or CBT), selecting therapy modules (e.g. behavioral activation or cognitive restructuring), or personalizing the content of a therapy module. Structure refers to the method of personalization, such as clinician judgment, patient preferences, decision aids, or statistical models. Time specifies whether personalization is performed before, during, or after treatment. More recently, Stumpp and Sauer-Zavala (2022) reviewed six personalized therapy approaches, such as selecting therapy modules or choosing and sequencing therapeutic skills, each with multiple implementation options, thereby highlighting the complexity of personalized psychotherapy. Although there is promising evidence that personalization can improve treatment outcomes in face-to-face psychotherapy (Nye et al., 2023), determining how it should be implemented remains a complex and unresolved challenge, particularly when psychotherapy is delivered through software applications such as mental health chatbots.

2.2.3 Personalization of Mental Health Chatbots

Research on the personalization of mental health chatbots has focused either on formative work collecting patients' views on personalized content (Abd-Alrazaq et al., 2021) or on personalizing chatbot design, such as the avatar, name, or personality (Ahmad et al., 2022; Six et al., 2022), but not the therapy content itself. Prior to this dissertation, no research had investigated content selection in mental health chatbots Kocaballi et al. (2019). Following Cohen et al. (2021), content personalization in mental health chatbots can be conceptualized at two complementary levels. The first level involves *content selection*: choosing the therapeutic modules most relevant to the user's symptom profile to form a personalized treatment plan. The second level involves *personalizing the content itself* during an ongoing conversation based on users' messages, such as providing personalized explanations or tailored feedback on therapy exercises.

For content selection, the available evidence came only from other types of DMHIs and was largely discouraging. A recent systematic review found that five of six randomized trials showed no benefit of personalized over standardized content (Schaeuffele et al., 2025). For

example, Johansson et al. (2012) and Berger et al. (2014) found no significant differences between personalized and standardized conditions, whereas Chaturvedi et al. (2023) reported greater engagement and Levin et al. (2019) demonstrated increased effectiveness for personalized content selection.

Against this background, three key gaps remained. First, empirical evidence on the impact of content personalization in mental health chatbots was absent (Kankanhalli et al., 2021; Kocaballi et al., 2019). Importantly, all prior DMHI studies had tested either purely system-led or purely user-led personalization, and no study had directly compared multiple personalization approaches under controlled conditions (Kankanhalli et al., 2021; Schaeuffele et al., 2025). Second, little was known about how different personalization approaches interact with user characteristics and stated preferences to influence engagement and health outcomes (Diederich et al., 2022; Kankanhalli et al., 2021). Third, theory-grounded prescriptive design knowledge for implementing content personalization in mental health chatbots was lacking. Although commercial mental health chatbots (e.g. Woebot, Wysa, Youper) employ diverse personalization features and have shown overall effectiveness (Fitzpatrick et al., 2017; Inkster et al., 2018; Mehta et al., 2021), their underlying personalization approaches and effects remain largely opaque (Kocaballi et al., 2019).

With respect to the personalization of therapeutic content itself during ongoing sessions, traditional rule-based chatbots are constrained by predefined static content, which limits their ability to provide open-ended personalized dialogues (Cho et al., 2023; Fitzpatrick et al., 2017). In contrast, LLMs offer a promising solution because they can generate more flexible, context-dependent messages that respond to users' input and enable higher personalization (Kocaballi et al., 2019). However, their probabilistic nature can produce inconsistent or potentially harmful responses (Heston, 2023). Prior work on LLM-based mental health chatbots has explored prompt engineering, fine-tuning, and hybrid architectures, but has often relied on user ratings, ad-hoc metrics, or downstream outcomes rather than validated fidelity instruments (Heinz et al., 2025; Y. Hua et al., 2025; Stade et al., 2024; Thieme et al., 2020). Consequently, it remains unclear how such systems should be designed to deliver structured psychological interventions with adequate therapeutic quality, and methods are needed to identify concrete shortcomings and support iterative refinement without exposing vulnerable users to untested systems (Heston, 2023).

In summary, while personalizing mental health chatbots holds considerable promise, rigorous

design and empirical research are essential to address the current limitations and gaps in our understanding to increase user engagement with and effectiveness of chatbots to ultimately provide better care to young people with depression.

3 Part I: How to Design a Chatbot to Treat Depression among Youth? Insights from a Qualitative Study

3.1 Introduction

Depression is a prevalent mental disorder in youth with significant personal and socioeconomic consequences that requires professional care (Clayborne et al., 2019; Thapar et al., 2012). Despite the availability of effective treatments, such as CBT (Oud et al., 2019), accessing them remains challenging. Even when professional services are free and accessible (Eisenberg et al., 2007), many youth avoid seeking professional support due to perceived stigma, a preference to solve problems on their own, and fear of psychotherapeutic settings (Domhan et al., 2023; Ebert et al., 2019; Gulliver et al., 2010; Radez et al., 2020). Interestingly, such attitudinal barriers may be more important reasons to not seek help than structural barriers, such as limited treatment resources and long waiting periods before starting psychotherapy (Domhan et al., 2023). DMHIs are a promising and effective way to overcome these barriers because they provide anonymous and self-empowered access to effective professional care (Wu et al., 2023). However, to leverage the full potential of DMHIs in the treatment of depressive symptoms in youth, two main limitations must be overcome: low adherence (Leech et al., 2021) and difficulty in establishing a therapeutic alliance, which is a crucial factor for effective psychotherapy (Cameron et al., 2018).

Chatbots, software systems that interact with users using natural language (Følstad & Brandtzæg, 2017; Vaidyam et al., 2019), have shown the potential to address these limitations. They are well accepted, feasible, and have shown promising effectiveness in strengthening mental health (Abd-Alrazaq et al., 2020; Bendig et al., 2019; Fitzpatrick et al., 2017; Inkster et al., 2018). Additionally, incorporating chatbots into DMHIs improves user engagement (Perski et al., 2019) and mental health outcomes (Linardon et al., 2024). Notably, users seem to develop a therapeutic alliance with chatbots (Bae Brandtzæg et al., 2021; Beatty et al., 2022; Darcy et al., 2021) partly because of social cues such as empathetic messages and humor (Feine et al., 2019; Fitzpatrick et al., 2017). Despite these encouraging results, most studies on mental health

chatbots, including those targeting depression symptoms, have focused on adult populations (Abd-Alrazaq et al., 2020; Bendig et al., 2019; Fitzpatrick et al., 2017; Inkster et al., 2018; Vaidyam et al., 2020). This is a key shortcoming because the results from adult populations cannot be generalized to youth. Youth face significant developmental changes in their biological, psychological, and social systems (Cicchetti & Toth, 2009), and depression symptoms differ from those in adulthood, especially at the onset of puberty (Rice et al., 2019). Moreover, youth interact with smartphones and chatbots differently than adults (Andone et al., 2016; Huffman, 2014), and have expressed that existing DMHIs often fail to adequately address their specific concerns (Agapie et al., 2022).

To address this research gap, recent studies have explored the design of chatbots for youth mental health. Høiland et al. (2020) involved youth in designing a chatbot for high school health services aimed at preventing mental disorders. Through focus groups, they identified four key support needs: (1) receiving information about mental health, (2) relating to someone beyond their immediate network, (3) receiving support for self-help, and (4) being referred to mental health services. Similarly, Grové (2021) developed a preventive mental health chatbot with high school students. Participants suggested topics such as school, family, friends, sexuality, and identity, as well as resources for adaptive coping strategies, mindfulness, and distractions. They expressed a preference for chatbots with inspiring, charismatic, and fun personalities using emojis, humor, and GIFs. While these two studies provide valuable insights, they focused on the prevention rather than the treatment of mental disorders and addressed a broad spectrum of mental health issues rather than specifically targeting depression. Thus, there is a critical need for research on how to design chatbots that focus on specific problems of youth with depressive symptoms to achieve sufficient engagement and optimal treatment outcomes.

Our study aims to address this research gap by investigating the following research questions:

1. What problems do youth with depressive symptoms face?
2. What adaptive coping strategies do they apply?
3. What attitudes and expectations do they have for chatbots designed to treat depression among youth?
4. What are their design preferences?

By addressing these questions, we aim to provide a comprehensive foundation for developing

chatbot-based DMHIs tailored to the unique needs and preferences of youth with depression, which will facilitate the development of engaging and effective DMHIs for this vulnerable population.

3.2 Methods

3.2.1 Study Design

We conducted a qualitative study to examine how to design a chatbot to treat depression in youth. The study included a questionnaire, a semistructured interview, and a concurrent think-aloud session with a chatbot prototype. We chose interviews because they allowed us to gather rich, detailed data on participants' experiences with depression, coping strategies, attitudes and expectations, and chatbot design preferences. The Institutional Review Board of the Medical Faculty of the University of Tübingen approved the study (project ID: 595/2021B01). The participants and caregivers provided written informed consent prior to screening for eligibility and enrollment.

3.2.2 Participants

Participants were eligible if they were between 14 and 17 years of age, owned a smartphone, and met diagnostic criteria for a current or remitted depressive episode. Participants with suicidal ideation or psychotic symptoms were excluded. Participants were recruited through a resident child and adolescent psychiatrist and the University's newsletter between June and August 2021. The target sample size was calculated based on the goal of data saturation (Guest et al., 2006), considering age group (14 - 15 and 16 - 17 years), depression status (remitted and acutely depressed) and gender (women and men), assuming homogeneity within each subgroup. Given that data saturation can be achieved at six cases per homogeneous group (Guest et al., 2006), the study aimed to recruit 48 participants. The sample was recruited through convenience sampling but the sampling was constrained by the predetermined recruitment period and available recruitment channels. The final sample consisted of 14 participants (12 women and 2 nonbinary).

3.2.3 Procedure

After participants and caregivers provided informed consent, we assessed their eligibility to participate in the study. Next, enrolled participants answered a questionnaire on sociodemographic information and mental health on a laptop. We conducted semistructured interviews to explore participants' problems with depression, their adaptive coping strategies, attitudes and expectations toward chatbots for depression, and their design preferences. Finally, the participants interacted with a prototype chatbot using the concurrent think-aloud method. After the study, participants received a reimbursement of 30€ (equivalent to US \$33) for taking part in the study. The study lasted between one and two hours.

3.2.4 Material

Eligibility Interview

We conducted a semistructured interview to evaluate whether interested candidates were eligible to participate in the study. The interview guide included questions regarding age, smartphone ownership, suicidal ideation, symptoms of depression and symptoms of psychotic disorders. The questions on suicidal ideation, symptoms of depression and symptoms of psychotic disorders were based on the DSM-5 criteria (American Psychiatric Association, 2013) and K-SADS-PL (Kaufman et al., 1997).

Questionnaire

Participants completed a questionnaire on demographic characteristics (age, gender, and level of education), history of mental disorders, and previous experience with psychotherapy. We assessed the prevalence of current symptoms of depression using the 8-item Patient Health Questionnaire (PHQ-8) (Martin et al., 2006) as well as current symptoms of anxiety using the Screen for Child Anxiety Related Disorders (Weitkamp et al., 2010). To complement the insights from the semistructured interview on the design preferences, participants answered a questionnaire on potential chatbot capabilities. This questionnaire is based on CBT manuals and literature on the content of DMHIs for depression (Auerbach et al., 2016; Huguet et al., 2016; Towery, 2016) and includes 13 capabilities that a chatbot to treat depression in youth could implement. For each capability, we asked the participants to indicate the extent to which

the chatbot should support them on a scale from 1 ("strongly disagree") to 5 ("strongly agree"). The capabilities are presented in Table 3.1. The participants completed all questionnaires on a laptop using SoSciSurvey.

Table 3.1: Items of the Chatbot Capabilities Questionnaire

The chatbot should support me with...

- ... how I can become more physically active or do sports.
- ... how to sleep better.
- ... how I can change negative/self-critical thinking.
- ... how I can do more activities that are important to me or that I have enjoyed in the past.
- ... tracking my mood.
- ... learning more about my depression.
- ... improving my social skills.
- ... reminding me to take my medication.
- ... connecting with an expert (e.g., psychotherapist) if I feel very bad.
- ... connecting with other people who have similar problems.
- ... how I can use my breath to make me feel better.
- ... writing a journal with things that concern me or for which I am grateful for.
- ... writing about events from my life.

Study Interview

We conducted semistructured interviews to investigate the problems participants faced due to depression, their adaptive coping strategies, their attitudes and expectations toward chatbots for depression, and their chatbot design preferences. Table 3.2 provides an excerpt of the interview guide (see Table A.1 in the Appendix for the full interview guide). We asked several questions about each topic and prepared sub-questions to follow up on specific details or offer suggestions based on the participants' initial responses.

Chatbot Prototype: Cady

Cady is the prototype of a chatbot for the treatment of depressive symptoms in youth that we developed for this and subsequent studies. Cady guides participants through a behavioral activation exercise, which is a key component of CBT to treat depression in youth (Oud et al., 2019). Behavioral activation aims to encourage patients to engage in pleasant activities to

Table 3.2: Excerpt from the Interview Guide

Topic	Question
Problems	Think of a period where you felt down. How did it look like?
Adaptive Coping Strategies	What have you tried in the past to feel better?
Attitudes and Expectations	Do you know what a chatbot is? What do you think it would be like for you to use a chatbot that is there to help you with one of the problems you have mentioned? For example, what could be the advantages and disadvantages?
Design Preferences	Imagine you are thinking about using one of these chatbots. How should it need to be designed, so that you would download and use it? How do you imagine an ideal conversation with the chatbot? For example, what topics would you like to discuss or what its personality should be like? If it was accessible as a mobile app, what should the app look like? (What) Would you like to personalize? How often would you like to use such a chatbot (per week)? How long should every session last?

overcome positive reinforcement deficits (Abel & Hautzinger, 2013; Lejuez et al., 2011). The behavioral activation exercise consists of the following sections: (1) introduction; (2) mood check on a scale from 1 (lowest) to 5 (highest) with adaptive empathetic responses; (3) psychoeducation on the relationship between behavior, thoughts and feelings; (4) find pleasant activities; (5) plan pleasant activities; (6) advice on how to overcome potential barriers; and (7) feedback and goodbye. The conversation was designed by licensed psychologists based on established CBT manuals for youth (Abel & Hautzinger, 2013; Groen & Petermann, 2015) with a focus on a positive, activating, and age-appropriate language style. Cady also uses emojis in its messages and sends GIFs.

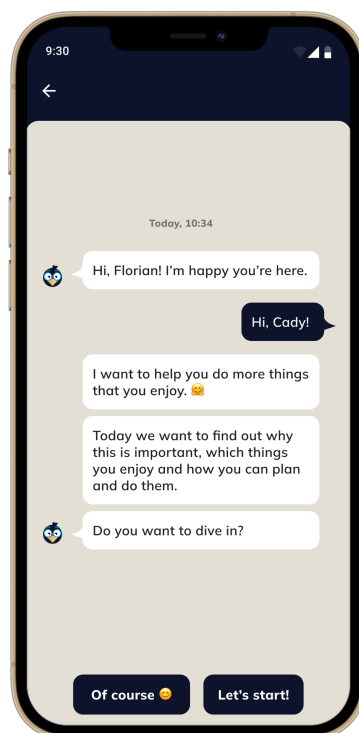


Figure 3.1: Screenshot of an Example Conversation between Cady and a Participant

Figure 3.1 shows a screenshot of an example conversation between Cady and a participant. We named the chatbot Cady in accordance with the title of our research project and did not specify its age, gender, or other demographic characteristics to prevent specific demographic characteristics from influencing the results. We developed Cady using the prototyping software Botsociety (Botsociety, n.d.). The chatbot was built with a rule-based architecture, where each user input triggers a predetermined response pathway following a decision tree structure. To ensure appropriate responses and maintain therapeutic quality, users interacted with the chatbot primarily by selecting from pre-defined options presented as buttons. We allowed free text input in specific situations, for example, to identify and schedule pleasant activities. However, these free-text responses did not alter the chatbot's conversation pathway.

Think-Aloud

We asked the participants to interact with Cady using the concurrent think-aloud method (Jaspers et al., 2004), in which users were instructed to express their thoughts and feelings loudly while interacting with the prototype. At the beginning of the think-aloud session, the interviewer introduced the chatbot prototype Cady and explained how the think-aloud method works. The participants interacted with the prototype using a laptop that we provided. While the partic-

Participants interacted with the prototype, we recorded the laptop screen and their comments. The interviewer stayed in the room because all participants felt it was easier to share their thoughts and feelings when the interviewer was in the room. The interviewer asked for further explanation if the statements were unclear or nonspecific.

3.2.5 Data Analysis

The interviews and think-aloud sessions were recorded using an audio recorder (Zoom Q2n) and transcribed verbatim. Two coders analyzed the semistructured interviews and concurrent think-aloud sessions using qualitative content analyses by Mayring and Fenzl (Mayring & Fenzl, 2019) and Mayring (Mayring, 2015) with QCAmap (Fenzl & Mayring, 2017). Similarly to Zehetmair et al. (2020), we chose an inductive coding approach to achieve the most unbiased and thorough description of the data (Mayring, 2015; Mayring & Fenzl, 2019), which we deemed important for the exploratory nature of our study. Both coders coded the material independently, discussed their coding results, and jointly developed a category system. When disagreements between the coders could not be resolved, a third coder (FOK) was consulted. We chose this process because intercoder comparison is not feasible with inductive category development (Mayring, 2015; Mayring & Fenzl, 2019). In an inductive approach, categories are created bottom-up from the material rather than selected before data analysis, as in deductive approaches (Mayring, 2015; Mayring & Fenzl, 2019). We report the frequencies of the categories and codes to increase transparency and demonstrate the saliency of the codes in the data. The questionnaires were analyzed using R (version 4.1.3). For PHQ-8, SCARED and the chatbot capabilities questionnaire, we calculated descriptive statistics (mean and standard deviation). For the chatbot capabilities questionnaire, we also generated frequency distributions to visualize response patterns.

3.3 Results

3.3.1 Participant Characteristics

In total, 14 youth aged 14–17 years ($M = 16.1$, $SD = 1.14$) participated in the study. Furthermore, 12 participants identified as women and two as nonbinary. All participants reported a current or remitted depressive disorder. In addition, 10 participants had a PHQ-8 score of 10 or

higher, indicating acute depressive symptoms (Kroenke et al., 2009). Twelve participants had a SCARED score of 25 or higher, indicating current symptoms of anxiety (Caporino et al., 2017). Nine participants showed symptoms of both depression and anxiety. Nine (64%) participants were currently receiving psychotherapy, while eight (57%) had received psychotherapy in the past. See Table 3.3 for the full sample characteristics.

Table 3.3: Participant Characteristics

Characteristic	Category	Value
Age	M (SD)	16.1 years (1.14)
Gender	Women	12 (86%)
	Nonbinary	2 (14%)
PHQ-8 (sum score)	M (SD)	13.57 (5.58)
	\geq Cut-off (10)	10 (71%)
SCARED (sum score)	M (SD)	40.86 (17.55)
	\geq Cut-off (25)	12 (86%)
Psychotherapy Experience	No experience	2 (14%)
	Current (\leq 6 months)	9 (64%)
	Remitted ($>$ 6 months)	8 (57%)
	Both current and remitted	6 (43%)
Diagnosis of Depression	Current (\leq 6 months)	8 (57%)
	Remitted ($>$ 6 months)	10 (71%)
	Recurrent (min. 2 episodes)	5 (36%)

3.3.2 Research Question 1: Problems with Depression

All 14 participants ($n = 14$, 100%) experienced problems categorized as depressive symptoms, including lack of motivation and energy, depressed mood, and self-doubt. Furthermore, 2 participants explained their main problems: *"I just don't have the strength to do anything"* and *"I'm just tired of everyday life, or of all these everyday activities like brushing my teeth"*. Some participants (9/14, 64%) reported comorbidities, such as panic attacks, anxiety, and excessive alcohol consumption, while four participants (28%) reported physical problems, such as tension, headaches, or stomach pain. All 14 participants (100%) also reported interpersonal problems, particularly withdrawal from social relationships, stress caused by parents, and problems with friends. Furthermore, 11 participants (78%) expressed concerns about school or their future. They shared negative experiences at school and experienced pressure to plan their future, which was reflected by statements like: *"When it comes to school, I have a very, very*

big fear of the future" and *"Success is a big topic right now, because of the end of school year grades"*.

Six participants (42%) reported experiencing stigma associated with depression. They mentioned that adults stigmatized or trivialized their problems or feared negative reactions when sharing their problems. One participant was concerned that disclosing that they are seeing a psychotherapist would lead others to perceive them as *"sick in the head"*. Most of the participants reported difficulties with mental health care. Twelve individuals (85%) reported experiencing treatment barriers, including attitudinal and structural barriers. Attitudinal barriers included fear of being judged, self-assessment that the problems are not severe enough to deserve support, and negative experiences during psychotherapy ($n = 9$, 64%). Participants reported ineffective techniques, difficulty with disclosing personal information, pressure to perform, trust violations, and condescending treatment from therapists. They also reported structural barriers, such as waiting lists or being considered not severe enough by healthcare providers or parents to qualify for professional support.

3.3.3 Research Question 2: Coping Strategies

All participants ($n = 14$, 100%) identified social support as an adaptive coping strategy. The majority received social support from friends, family members, partners, teachers, or people they interacted with online. In total, 13 participants (93%) engaged in activities to distract themselves or to have a positive experience such as spending time outdoors or consuming media. A total of 9 participants (64%) reported receiving professional treatment and nine (64%) used cognitive strategies, including positive self-talk, self-reflection, and reducing rumination. Six participants (42%) implemented a daily structure that included forming habits, getting up early, scheduling positive activities, and setting goals. Seven (50%) prioritized their needs as a coping mechanism, such as intentionally allocating time to self-care or avoiding stressful social situations. Three individuals (21%) engaged in mindfulness practice, such as breathing techniques and meditation. Two individuals (14%) sought online information on their depressive symptoms, interrelated problems, and additional information.

3.3.4 Research Question 3: Attitudes and Expectations toward Chatbots to Treat Depression

Most participants had positive attitudes toward and expectations of chatbots to treat depression. In total, 12 participants (85%) stated that they would be less anxious to use a chatbot than to see a human therapist. They pointed out that using a chatbot would be a suitable option to discuss sensitive topics that they would not share with others, primarily because they would not fear negative reactions. In addition, texting was considered less intimidating than speaking to a human therapist. Furthermore, 11 participants pointed out the unlimited capacity and flexibility of chatbots (79%): *"I think it is also an advantage that you can really chat with it at any time, because in therapy you just have one appointment per week. If you feel bad in the evening or at night, or something like that, you can still text the chatbot."* Furthermore, the participants indicated that a chatbot is more accessible and requires less effort than seeing a human therapist. A total of 8 participants (57%) expressed confidence in the effectiveness of chatbots for the treatment of depression. They believed that such a chatbot would be capable of addressing a wide range of issues and particularly welcomed the idea of having a helpful everyday chatbot. All participants demonstrated a great interest in using a chatbot to treat depression either because they expected it to be effective or because they were curious about using it. One participant stated that a daily usable chatbot could alleviate feelings of loneliness. Several others noted that a chatbot's personal and human-like nature would increase the motivation to use it, particularly compared to other less interactive DMHIs.

However, participants also reported several concerns about the use of chatbots to treat depression. Some participants ($n = 9$, 64%) were skeptical about the intelligence and natural language capabilities of the chatbot. They expressed the concern that the chatbot would not be able to solve individual, diverse or unusual problems effectively and feared being disappointed if the chatbot was unable to do so. Participants were particularly concerned about inappropriate answers to intimate and emotional topics or inappropriate advice for their problems. Nine participants (64%) were concerned that a conversation with the chatbot would not feel like a conversation with a human therapist, due to a potential lack of emotional intelligence or overly robotic or analytical responses. However, chatbots that appeared too human were believed to be uncanny. Two participants (14%) were skeptical about sufficient data security and privacy. Finally, two participants (14%) worried that their symptoms of depression, such as little moti-

vation or forgetfulness, would result in low engagement with the chatbot, thus preventing them from effectively using it. They also pointed out that the lack of social pressure when using a chatbot, compared to seeing a human therapist, could contribute to low engagement, which might not be solvable by the chatbot.

3.3.5 Research Question 4: Chatbot Design Preferences

Category 1: Dialogue Topics and Therapeutic Content

All participants ($n = 14, 100\%$) shared preferences regarding the dialogue topics that the chatbot should be able to cover. These suggestions include comprehensive assessments of depressive symptoms and interrelated problems, psychoeducation, therapeutic exercises that address specific issues, reminders to take care of basic needs, support to regulate emotions, deal with intrusive thoughts, and be distracted when needed. Furthermore, they emphasized the importance of discussing their daily lives and, more specifically, sharing current problems and receiving suggestions on how to solve them. One participant highlighted the importance of the chatbot explicitly asking the user what type of support they require, such as emotional support (i.e., listening and validation) or solution-oriented support: "*Getting advice on how to solve a problem isn't always helpful, even if it's well-intentioned. [...] It would be important for the chatbot to ask or understand if you want advice or if you only want to share your feelings.*" Others elaborated on how they imagined talking about daily life and receiving advice: "*I would just talk about the things that depress me at the time, to which I don't know the answer*", and "*I would probably just chat about everyday situations that were unpleasant to me or something like that.*" These preferences were complemented by the chatbot capabilities questionnaire (see Figure 3.2).

The responses indicate the importance of three key components of CBT: (1) cognitive restructuring ("*changing negative thinking*"; $M = 4.9, SD = 0.4$), (2) behavioral activation ("*pursuing activities that are important to me or have brought me joy in the past*"; $M = 4.5, SD = 0.7$) and (3) psychoeducation ("*learning about depression*"; $M = 4.5, SD = 1$). Furthermore, participants expressed a preference for "*improving social skills*" ($M = 4.4, SD = 0.9$). On the other hand, therapeutic writing, represented by "*writing about events in my life*" ($M = 4, SD = 1.3$) and "*writing a diary about things that bother me or that I'm grateful for*" ($M = 3.7, SD = 1.5$), was rated lower than most items. Finally, "*connecting with other people who face*

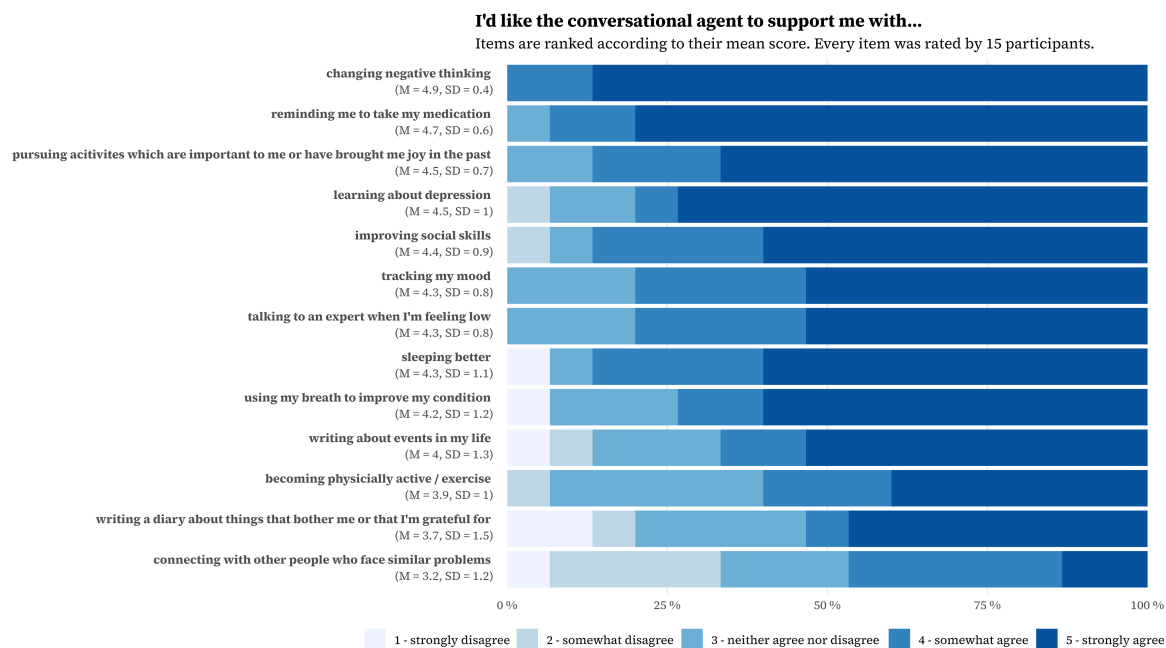


Figure 3.2: Results of the Chatbot Capabilities Questionnaire

similar problems" was ranked the lowest ($M = 3.2$, $SD = 1.2$).

Category 2: Personality, Interaction Style, and Social Role

The participants made various suggestions for the chatbot personality ($n = 14$, 100%), such as understanding, caring, friendly, empathic, encouraging, humorous, interested, neutral, and nonjudgmental. During the think-aloud sessions, most participants characterized Cady's personality as one of its main strengths. One participant said: *"I like that you communicate in a friendly way, like an internet friend."* Two participants were pleased that Cady was interested and asked them personal questions. For example, one participant approved that Cady had asked if she had experienced becoming less active herself: *"I just think it's good that she is asking if I know something like that or haven't experienced it yet"* In terms of the interaction style and language used by the chatbot, participants suggested that it should be appropriate for the user's age without complicated terms. This issue was also raised during think-aloud sessions. One participant asked for more age-specific language use, such as short sentences: *"So, this was a therapeutic chat, but if it is supposed to be like a friend, I think it is better if it writes like people my age, e.g. using short sentences."* Another participant perceived the language as too therapeutic: *"It was relatively authentic, but you also notice that it is not written for my age"* The human-likeness of the chatbot was also a significant factor, as participants found it

important to react to feelings and to understand irony and humor. However, participants also expressed concerns about the chatbot being too human-like, which could lead to irritation or fear. During the think-aloud sessions, one participant appreciated that Cady disclosed personal information such as enjoying conversations with nice people, which resembled a human-like trait.

Regarding the social role, the participants ($n = 14, 100\%$) had divergent views: some wanted the chatbot to resemble a friend, while others preferred it to assume the role of a therapist or a combination of both. One participant explained their preference for a therapist-like role: "*I actually think it's better if it were like a therapist [...], because with friends you don't necessarily want to talk about everything.*" One participant proposed using situation-specific roles: "*A therapist if you're doing therapeutic exercises and a friend if you just want to talk.*".

Category 3: Personalization

Participants frequently discussed the importance of personalization ($n = 13, 92\%$). They reported various aspects that should be personalized, including the chatbot's personality, gender, avatar, message style, dialogue topics, therapeutic content, the mobile app's color theme, notifications, and user profile picture and name. However, the participants had divergent opinions on who should control the personalization of the chatbot. While some suggested that the user should be in control, others believed that the chatbot itself should control personalization. Others, in turn, proposed a mixed approach in which both the user and the chatbot share control. Some participants preferred the personalization to occur only once at the beginning, whereas others preferred dynamic personalization, where personalization is performed continuously. In addition to these explicit preferences, personalization also emerged from divergent design preferences. All participants expressed a preference for an attractive user interface, but they had differing opinions on the specific design elements, with some suggesting a bright and colorful design, whereas others preferred a plain and minimalist design with black and white colors. Most participants preferred the chatbot as a standalone application, but some preferred integration into a popular instant messaging app such as WhatsApp. Regarding the interaction modalities, most of the participants preferred a chatbot over a voice agent. In terms of text input, the participants believed that a mix of predefined responses and free text input was optimal, as it combines the minimal effort of clicking predefined answers with the flexibility of unrestricted text input. Two participants (14%) expressed concerns about personalization, stat-

ing that personalized content could lead to avoidance behavior, in which users avoid topics that they find difficult or uncomfortable. They also believed that personalization of the chatbot's appearance, such as its gender or avatar, could reduce its seriousness as an application to treat depression.

Table 3.4 summarizes key chatbot design recommendations derived from the results of all four research questions.

3.4 Discussion

3.4.1 Primary Findings

This study explored the problems youth face with depression, their adaptive coping strategies, and their attitudes and design preferences for chatbots to treat depression. Our findings indicate that youth experience diverse problems beyond core depressive symptoms, employ a range of coping strategies, hold predominantly positive attitudes towards chatbots for depression, and express various, sometimes contradictory, design preferences for chatbots to treat depression. In the following section, we discuss the results according to the four research questions guiding this study: (1) the problems youth with depression face, (2) their adaptive coping strategies, (3) their attitudes and expectations towards chatbots to treat depression, and (4) their design preferences. As presented in Table 3.4, our findings across all research questions offer specific guidance for designing chatbots tailored specifically to youth experiencing depression.

Problems of Youth with Depression

First, participants reported various problems with depression beyond the symptoms described in the DSM-5 (American Psychiatric Association, 2013), including mental health and somatic comorbidities, interpersonal issues and concerns about school and the future. These findings align with prior research showing that youth with depression face an increased risk for anxiety disorders (Thapar et al., 2012), substance abuse (Thapar et al., 2012) and physical health issues (Agnafors et al., 2019), as well as social withdrawal, lack of friendship (Mullarkey et al., 2019), diminished academic achievement, school dislike, and pessimism (Mullarkey et al., 2019). To address these problems, engaging and effective chatbots should (1) repeatedly assess core depressive symptoms and related problems to inform therapeutic decisions and strengthen

Table 3.4: Key Design Recommendations

Topic	#	Design Recommendation
Assessment	1	Incorporate repeated multidimensional assessment, covering depressive symptoms and related problems, to monitor and understand the key issues a user faces.
	2	Ensure assessments are comprehensive yet concise to minimize the effort for users.
Scope and Limitations	3	Communicate the chatbot's scope and limitations clearly, e.g., when it cannot properly handle a topic. Redirect users to external resources (e.g., helplines, professional counseling) in such cases to ensure users receive appropriate support.
Dialogue Content	4	Curate a diverse content database covering a wide range of topics that are relevant to youth with depression.
	5	Integrate evidence-based therapeutic techniques alongside user-preferred coping strategies (e.g. behavioral activation, cognitive restructuring, social support).
	6	Enable the chatbot to engage in conversations about users' daily life, including discussions on urgent problems.
	7	Support follow-up discussions on user-reported problems, track progress, and encourage reflection on how successfully they carried out the learned techniques.
Personality, Social Role and Language Style	8	Design a chatbot with an understanding, caring, friendly, empathic, encouraging, humorous, interested, and non-judgmental personality.
	9	Ensure the chatbot uses age-appropriate language without sounding like technology.
	10	Balance human-likeness (e.g. recognizing feelings, irony, humor) with clear boundaries to avoid unrealistic or overly human-like interactions.
Personalization	11	Personalize the content selection and presentation.
	12	Personalize the chatbot's persona (e.g. personality, social role, gender, avatar, message style).
	13	Explore the impact of persona personalization (e.g. switching between formal and casual depending on user feedback, interaction style preferences and situation).
	14	Explore the impact of who controls personalization (the user, the chatbot or both) and when / how often personalization is performed (initial vs ongoing).
	15	Ensure personalization enhances user experience while maintaining the chatbot's therapeutic quality and seriousness.

the therapeutic relationship by understanding the users' individual challenges, and (2) curate a diverse content database to address the wide range of problems, aligning with S. H. Li et al. (2022). Second, participants reported significant barriers to receiving professional treatment, including poor mental health literacy (e.g., awareness of symptom severity), attitudinal (e.g., perceived stigma), and structural barriers (e.g., waiting times). These findings align with research showing that youth struggle to recognize symptoms (Radez et al., 2022), hesitate to disclose personal information to professionals due to fear of judgement or not being taken seriously (Radez et al., 2020), and question the effectiveness of professional treatment (Gulliver et al., 2010; Radez et al., 2020), a concern supported by evidence that 60% of youth with depression do not benefit from psychological treatments (Cuijpers et al., 2023a). To address these barriers, chatbots should (1) encourage emotional disclosure by providing a non-judgmental conversational environment, something DMHIs cannot fully achieve due to their static, non-conversational design, and (2) clearly communicate their scope and limitations to set realistic user expectations and direct to external support when needed.

Adaptive Coping Strategies

Participants reported various adaptive coping strategies, including seeking social support, establishing positive activities, using cognitive strategies and receiving professional treatment. Many also emphasized the importance of focusing on their needs and establishing a daily structure. Most of these strategies fall into the CBT strategies behavioral activation, cognitive restructuring and problem solving (Wenzel, 2017), which are effective and first line treatments for depression in youth (Luxton & Kyriakopoulos, 2022; Oud et al., 2019). The prevalence of CBT-based coping strategies may be due to the prior psychotherapy experience of most participants. Similarly, among youth with non-suicidal self-injury, those with experience in dialectical behavioral therapy have suggested incorporating it into DMHIs designed for them (Čuš et al., 2021), showing how prior therapy experience influences preferences for DMHIs. The use of CBT-aligned adaptive coping strategies is important, as research shows that youth with depression use adaptive cognitive strategies less frequently than healthy controls (Mihailescu et al., 2023) and that applying these strategies is associated with fewer depressive symptoms (Schäfer et al., 2017). Chatbots for youth with depression should incorporate these strategies to create engaging and effective content that aligns with users' existing coping strategies and effective therapeutic techniques. For example, the chatbot could facilitate social support by helping

users identify supportive individuals and suggesting personalized ways to reach out, such as a video call with a friend who is good at cheering them up or a message to a family member who provides good advice. Chatbots can even draft customizable messages to support users who struggle to reach out.

Attitudes and Expectations towards Chatbots to treat Depression

Participants predominantly held positive attitudes towards chatbots to treat depression. Many participants reported feeling less anxious about using a chatbot than seeing a human psychotherapist, especially when discussing sensitive topics. This finding extends research on post-traumatic stress disorder, showing that digital agents lead to greater disclosure of sensitive or stigmatized information (Lucas et al., 2017; Pickard et al., 2016), and supports evidence that youth value the privacy and anonymity offered by DMHIs (Kenny et al., 2016). Given these findings, chatbots could help users become comfortable with sharing sensitive information and practicing therapeutic conversations, with potential benefits for subsequent sessions with psychotherapists. However, some participants were skeptical about the chatbots' ability to address individual problems and provide appropriate advice on sensitive topics, which is supported by evidence that commercial chatbots for adults with depression often fail to match user inputs, understand messages, or respond appropriately (Martinengo et al., 2022). Although privacy and data security are frequently cited as primary risks of DMHIs for young people (Wies et al., 2021), few participants raised these concerns. Instead, they viewed chatbots as privacy-enhancing compared to human therapists. Nevertheless, robust privacy and data security remain essential, and it is unclear whether participants were genuinely unconcerned or simply assumed chatbots would have strong protections in place.

Design Preferences

Our study revealed important insights into design preferences for chatbots for youth with depression. First, participants expressed diverse preferences for dialogue topics, with three areas standing out: (1) chatting about daily life, (2) discussing urgent problems and receiving advice, and (3) working through therapy exercises that address specific problems. Interestingly, while youth with nonsuicidal self-injury (Čuš et al., 2021), emotional problems (Ludlow et al., 2023) and from the general population (Kenny et al., 2016) included being connected to others facing

similar challenges as a key feature of their DMHIs, our participants ranked such social connection as their lowest priority, highlighting how design preferences may vary across mental health conditions. In summary, a chatbot to treat depression in youth should reflect these preferences. However, available applications have not implemented chats about daily life, and relied on rule-based approaches for problem-solving (Fitzpatrick et al., 2017; Inkster et al., 2018), likely due to insufficient conversational capabilities. Structured therapy techniques, such as behavioral activation or cognitive restructuring are standard components of DMHIs (Huguet et al., 2016) and chatbots (Ahmed et al., 2023). However, their current implementation remains predominantly static (Denecke et al., 2022), without personalized advice or feedback during therapeutic exercises. Large language models (LLMs) such as OpenAI's GPT or Anthropic's Claude promise to overcome current limitations of DMHIs and chatbots. LLMs enable chatbots to better address our participants' design preferences, such as natural conversations about daily life and advice on urgent problems. Additionally, LLMs can enhance static implementations of therapeutic exercises by offering personalized guidance and feedback. In behavioral activation, LLMs can personalize the explanation of the relationship between behavior and feelings and evaluate the proposed activity plan for feasibility and therapeutic appropriateness. In cognitive restructuring, Sharma et al. (2024) demonstrated that LLMs can assist users in identifying thinking traps and generating reframed thoughts. Although these enhanced capabilities have the potential to improve understanding and skill development, key challenges remain. LLMs can generate false or harmful messages (Birkun & Gautam, 2023), posing risks to vulnerable users. Future research needs to explore how to capitalize on their advanced conversational abilities while ensuring therapeutic quality and safety (Stade et al., 2024). Second, our findings show that engaging and effective chatbots for youth with depression need personalization beyond the conversational capabilities of LLMs. Drawing on the framework by Cohen et al (Cohen et al., 2021) for psychotherapy personalization, effective chatbots require personalization in content selection and interaction style. Personalized content selection ensures that chatbots address the diverse problems and dialogue topics our participants reported, aligning with findings from S. H. Li et al. (2022) and Ludlow et al. (2023). A personalized interaction style accommodates the preferences for different social roles and language use. The benefits of personalizing the chatbot persona have been shown in prior research (Ahmad et al., 2022; Nißen et al., 2022). LLMs enable the personalization of interaction styles via prompt instructions, reducing the need to craft different responses for each style manually (Jiang et al., 2024).

Despite the clear need for personalization, participants disagreed on who should control it (Fan & Poole, 2006), preventing a clear design recommendation. While some preferred user-led personalization, aligning with Kenny et al. (2016), others favored chatbot-led personalization or a hybrid approach. Given this ambiguity and the lack of empirical evidence (Kocaballi et al., 2019), future research is needed to guide the design of chatbot personalization.

3.4.2 Limitations

Our study has two main limitations. First, our sample size was limited, and the participants were predominantly women. Although depression is more prevalent in women (Avenevoli et al., 2015; Z. Hua et al., 2024; Seedat et al., 2009), the absence of men participants limits the generalizability of our findings across genders. The gender imbalance resulted from convenience sampling through a single youth psychiatrist. One man was informed about the study but declined to participate. While our study included two nonbinary participants, an underrepresented demographic in research, recruiting men would have required targeted efforts beyond our study's constraints. Additionally, our sample consisted entirely of participants who had actively sought mental health support most of whom had prior psychotherapy experience. As a result, perspectives from those who avoid professional help or resist treatment may be underrepresented. Due to these sampling limitations, we likely did not achieve full data saturation. Although qualitative research can reach saturation with as few as twelve interviews (Guest et al., 2006), our sample composition suggests that some perspectives may have been missed. However, our findings on youth depression and related problems replicated findings from studies with much larger sample sizes (Chevance et al., 2020; Mullarkey et al., 2019), indicating that our results are comprehensive. Future studies should recruit a more diverse sample, specifically including men, individuals reluctant to seek professional treatment, and those without prior psychotherapy experience, to validate and extend our findings.

Second, our study captured attitudes, preferences, and hypothetical usage scenarios rather than actual usage behavior. While participants interacted with a prototype during the think-aloud sessions, this brief, controlled interaction may not reflect real-world interactions. The gap between stated preferences and actual behavior is well documented (De Corte et al., 2021), and our study design does not allow us to determine whether the design preferences of the participants would translate into sustained engagement or effective use. Future research should

follow the next steps of an iterative human-centered design process (Harte et al., 2017; Yardley et al., 2015) with functional prototypes to assess whether the identified design preferences lead to actual user engagement and therapeutic benefits. Despite these limitations, our study provides in-depth findings on the design of chatbots to treat depression in youth, highlighting the value of qualitative research in the iterative development of DMHIs (Harte et al., 2017; Yardley et al., 2015).

3.4.3 Conclusions and Future Work

Our study provides valuable insights into the problems and coping strategies of youth with depression, and their attitudes, expectations and design preferences for a chatbot to treat depression in youth. We found complex user needs, predominantly positive attitudes towards chatbots, and various design preferences, including the need for diverse dialogue topics and personalization. Our findings led to concrete design recommendations that lay a crucial foundation for developing engaging and effective chatbots to treat depression in youth. Despite these contributions, several research directions remain. First, future studies should validate and extend these findings with larger, more diverse samples to ensure broader representation of youth with depression. Second, examining actual usage patterns and long-term engagement with functional chatbot prototypes will help assess the effectiveness of our design recommendations. Third, investigating the responsible integration of LLMs is important, including the development of robust safeguards and evaluating dialogue quality and therapeutic outcomes. Finally, researchers must design and evaluate effective personalization features, particularly regarding whether users, the chatbot or both should control personalization, while balancing implementation efforts and impact. By addressing these research directions, we can further improve chatbots to treat depression in youth and ultimately contribute to more accessible, engaging, and effective mental health support for this vulnerable population.

4 Part II: User-Led, Chatbot-Led, or Hybrid? Design and Effect of Content Personalization Approaches in Mental Health Chatbots for Young People

4.1 Introduction

Depression is a common mental disorder with serious individual and socioeconomic consequences (Reinert et al., 2022). Young people between 15 and 29 years old (European Union, 2023), in particular, are often reluctant to seek treatment due to stigma, the desire to solve problems on their own, or lack of interest in face-to-face psychotherapy (Ebert et al., 2019; Gulliver et al., 2010). Even when seeking treatment, they face long waiting periods before starting psychotherapy (Thapar et al., 2012). As a result, DMHIs, such as chatbots (e.g., Woebot, Wysa, Youper), have become popular as easily accessible conversational alternatives to treat depression (Darcy et al., 2021; Mehta et al., 2021).

Although recent studies have demonstrated the potential of mental health chatbots (Lim et al., 2022; Linardon et al., 2024), improving user engagement and effectiveness remains a challenge (Jabir et al., 2024; Lim et al., 2022; Limpanopparat et al., 2024). Personalization, in the context of psychotherapy defined as matching patient characteristics and preferences with available psychotherapeutic treatment options (Cohen et al., 2021; Stumpp & Sauer-Zavala, 2022), has been proposed as a promising approach to improve user engagement and effectiveness. In routine face-to-face psychotherapy, the psychotherapist chooses the modules most relevant to the individual patient, such as behavioral activation or cognitive restructuring in CBT. Similarly, many commercial mental health chatbots provide features to select the most relevant modules for the user, thereby personalizing their content.

Considerable IS research has shown the benefits of personalization in business contexts (e.g., product recommendation agents) (Benlian, 2015; Ebrahimi et al., 2022; Komiak & Benbasat, 2006). In contrast, no research has investigated content personalization in mental health chatbots, and research on the personalization of DMHIs has produced conflicting results. Specif-

ically, a recent systematic review found that 5 out of 6 randomized trials showed no benefit of personalization compared to standardized content (Schaeuffele et al., 2025). This finding is problematic because transferring personalization findings from e-commerce to mental health faces three key challenges. First, the stakes in mental health contexts are considerably higher. Inadequate recommendations can lead users to discontinue using the chatbot, delay improvements, or even worsen symptoms (Baird et al., 2025; Kankanhalli et al., 2021). Without human oversight, personalization should not lead users to circumvent modules that are especially helpful for their symptom profiles or select modules that are ineffective for them. Second, users with depression often suffer from cognitive impairments such as decision fatigue, making the trade-off between cognitive effort and decision quality more difficult than in e-commerce (Rock et al., 2014). Third, mental health interventions require repeated engagement and sustained practice of techniques, such as completing CBT homework, practicing therapeutic skills, and applying these techniques in daily life, unlike one-off purchase decisions. Furthermore, previous research has not examined the effects of different personalization approaches, such as implicit (chatbot leads personalization), explicit (user leads personalization), and hybrid (both chatbot and user are involved in the personalization process) personalization (Kocaballi et al., 2019). Thus, it remains unclear why previous personalization approaches have not outperformed standardized content delivery and whether finding a balance between system-led and user-led personalization could lead to better outcomes.

Our research addresses these gaps by investigating the design and effects of content personalization in chatbots for young people with depressive symptoms. Following the DSR paradigm (Peppers et al., 2007), we developed a chatbot called "TheraBot". Grounded in self-determination theory, the effort-accuracy framework, and cognitive load theory, TheraBot provides users with hybrid content personalization by selecting the most relevant therapy modules, but also allowing users to adapt the module selection to their own needs and preferences. Using TheraBot, we conducted two field experiments ($N_1 = 74$, $N_2 = 160$) over two weeks to investigate the effects of personalized content (vs. standardized) and the effects of different personalization approaches on user engagement and depression severity. Our results suggest that TheraBot led to a greater reduction in depression severity than standardized content, contrasting with the null results found in previous studies (Schaeuffele et al., 2025). In addition, hybrid content personalization was more promising than chatbot-led or user-led personalization. However, we also found that initial depression severity affected user engagement differently across person-

alization approaches. Users with more severe symptoms were more engaged in the user-led personalization condition, whereas in the chatbot-led condition, higher initial severity was associated with lower engagement. Finally, we explored the mechanisms underlying TheraBot's effectiveness and synthesized our findings into design requirements, highlighting the need to adapt the personalization approach itself to different users rather than applying the same approach to everyone.

Our research contributes to the emerging IS literature on digital mental health interventions (Angst et al., 2024; Peng et al., 2024; Sjöström et al., 2022) by advancing our understanding of the design and impact of content personalization in mental health chatbots. Responding to calls for systematic, healthcare domain-specific research (Baird et al., 2025; Kankanhalli et al., 2021), we provide a nuanced view of how different content personalization approaches affect users of mental health chatbots. Our findings show that hybrid personalization is more effective than standardized content delivery and most preferred by users, while neither chatbot-led nor user-led personalization show benefits over standardized content delivery. This finding may help explain why previous studies evaluating system-led or user-led personalization approaches alone did not show improved mental health outcomes (Schaeuffele et al., 2025). However, our results also reveal important nuances: participants with higher depression severity in the user-led condition tend to complete more modules relative to those with lower severity, whereas this trend is reversed in the chatbot-led condition, challenging the assumption that a single personalization approach suffices for all users. Therefore, our findings emphasize the importance of personalizing the personalization approach as a promising direction for optimizing user engagement and effectiveness across a diverse user base. This insight bridges our first contribution with our second: we offer novel findings about how different personalization approaches (chatbot-led, user-led, and hybrid) affect mental health outcomes. Earlier studies have focused only on system-led approaches that do not involve users. Our finding that hybrid personalization seems more promising than the other two strategies advances the literature on DMHIs and helps clarify why system-led personalization approaches have sometimes not outperformed standardized content (Berger et al., 2014; Johansson et al., 2012). Third, our research offers theory-grounded prescriptive design knowledge for content personalization in mental health chatbots. Existing commercial mental health chatbots act as "black boxes" (Kocaballi et al., 2019), and although design knowledge for chatbots exists, it typically focuses on social roles or interface design rather than therapy content (Ahmad et al., 2022; Diederich et al., 2022; Paul

et al., 2024). We address this gap by deriving design requirements from kernel theories (self-determination theory, effort-accuracy framework, cognitive load theory) and validating them through our artifact TheraBot. Finally, we provide actionable design prescriptions for developers of mental health chatbots to implement effective content personalization approaches, while also highlighting the need for more nuanced strategies that consider individual user characteristics and preferences.

4.2 Background

Our research is positioned at the intersection of two research fields: digital mental health and personalization. In the following, we briefly summarize previous research in each of these streams before reviewing studies at their intersection, which are most closely related to our work. We then highlight the research gaps that our work addresses.

4.2.1 Digital Mental Health Interventions

DMHIs are software applications that use technologies such as computers, tablets, or mobile phones to provide mental health care (Schueller & Torous, 2020). They are typically based on established psychotherapy approaches, such as CBT, and their use has been shown to reduce the severity of depression (Torous et al., 2021). Mental health chatbots are an increasingly popular type of DMHIs that mimic human conversations (Vaidyam et al., 2019), with popular commercial examples such as Woebot, Wysa, and Youper. Mental health chatbots have been perceived to be more engaging than traditional DMHIs (Perski et al., 2019) and shown to establish a working alliance with users (Darcy et al., 2021), which is important to achieve success in traditional psychotherapy (Cameron et al., 2018). Although there is initial evidence that chatbots can alleviate symptoms of depression (Lim et al., 2022; Linardon et al., 2022), increasing the user engagement and effectiveness of mental health chatbots remains a challenge (Limpanopparat et al., 2024).

4.2.2 Personalization

The Cambridge Dictionary (Cambridge Dictionary, 2023) defines personalization as "the act of making something suitable for the needs of a particular person". In IS and related disciplines,

various studies have examined personalization. However, there is little agreement on how to conceptualize it, and the term is often used interchangeably with "customization", "adaptation", or "individualization" (Fan & Poole, 2006). For our research, we adopt a broad definition of personalization as "a process that changes the functionality, interface, information access and content, or distinctiveness of a system to increase its personal relevance to an individual" (Fan & Poole, 2006, p. 183). One can broadly distinguish between two basic elements of an IS that can be personalized: content (i.e., the information itself) and design (i.e., how and when the content is presented in the user interface) (Benlian, 2015; Fan & Poole, 2006). One key dimension of personalization is who controls it. Two strategies can be distinguished: explicit and implicit personalization (Adomavicius & Tuzhilin, 2005; Fan & Poole, 2006). In explicit (user-led) personalization, also called tailoring or customization, users specify their preferences or make choices regarding the content and design of an information system (Germonprez et al., 2007; Miah et al., 2019). In implicit (system-led) personalization, the system collects user data, such as interactions, browsing history, and feedback, to make informed decisions about the content that best fits the user (Ebrahimi et al., 2022). These two strategies can also be combined into a hybrid approach in which both the system and the user are involved in personalization (Fan & Poole, 2006). IS research has demonstrated that personalization can enhance user outcomes across multiple domains. In e-commerce, system-led personalization improves attention and engagement (Tam & Ho, 2005; Xiao & Benbasat, 2007). In healthcare, system-led personalization has shown to outperform standardized methods, for example, by personalizing treatment plans in breast cancer care (W. Chen et al., 2021) and by providing personalized content recommendations in digital weight-loss interventions (Zhou et al., 2023). In digital news and media, personalized news article lists increase satisfaction (Liang et al., 2006), while culturally-adapted user interfaces outperform standardized designs (Reinecke & Bernstein, 2013). However, system design, user characteristics and context shape the effects of personalization. System design shapes how personalization mechanisms influence adoption: interactive recommendation agents that explicitly elicit user preferences achieve adoption primarily through recommendation quality, whereas non-interactive systems that implicitly infer preferences from behavioral data depend predominantly on user trust (Ebrahimi et al., 2022). Contextual factors moderate personalization's business impact: Personalized product recommendation systems demonstrate differential effectiveness based on retailer characteristics. Large retailers with diverse inventories benefit substantially, while price-focused competi-

tors see minimal benefits (Thirumalai & Sinha, 2013). Conversely, checkout personalization universally enhances customer loyalty regardless of business model. Domain-specific privacy considerations further shape effectiveness, with users accepting greater disclosure risks for enhanced personalization quality in financial and news services compared to other domains (T. Li & Unger, 2012). Finally, user characteristics moderate outcomes as personalized news particularly satisfies those with informational or social goals (Liang et al., 2006). These differences suggest that optimal personalization approaches must account for heterogeneous user populations and cannot easily be transferred between domains and usage contexts. In the context of psychotherapy, personalization is important because patients differ in their symptoms, preferences, and responsiveness to different therapy approaches (Cohen et al., 2021). Therefore, it is important to develop and implement personalized therapy approaches to achieve optimal treatment outcomes. The overarching aim of personalized therapy is to match patients with their specific symptom profiles and characteristics to available treatment options. However, this process can become complex because of multiple therapeutic dimensions and practical implementations. For instance, in their organizational framework for personalized psychotherapy, Cohen et al. (2021) provide multiple options for personalization along three dimensions: (1) level of intervention, (2) structure, and (3) time. First, the level of intervention categorizes the level at which psychotherapy is personalized, such as selecting the overall psychotherapy modality (e.g., psychodynamic psychotherapy or CBT), selecting therapy modules (e.g., choosing from CBT modules, such as behavioral activation or cognitive restructuring), or personalizing the content of the therapy module, e.g. by changing the interaction style. Second, structure classifies the personalization methods, such as relying on clinician's judgement, patient preferences, decision aids, or complex statistical models. And third, time specifies whether personalization is performed before, during, or after (e.g., relapse prevention) treatment. More recently, Stumpp and Sauer-Zavala (2022) reviewed six specific personalized therapy approaches, such as the selection of therapy modules or the choice and order of therapeutic skills, each with multiple implementation options, thus showing the complexity of personalized psychotherapy. In summary, although the need to adopt and investigate the personalization of psychotherapy is clear, and there is promising evidence that it can improve treatment outcomes in face-to-face psychotherapy (Nye et al., 2023), determining how personalization should be implemented remains a complex and unresolved challenge, particularly when psychotherapy is delivered through software applications, such as mental health chatbots.

4.2.3 Personalization of Digital Mental Health Interventions

Studies on the personalization of mental health chatbots have focused on formative research collecting the opinions of patients on personalized content (Abd-Alrazaq et al., 2021) or on the personalization of chatbot design (Ahmad et al., 2022; Six et al., 2022), but not the therapy content. So far, the effect of content personalization has only been investigated using other types of DMHIs than chatbots. A recent systematic review summarized these findings and concluded that content personalization generally yields no benefit in digital mental health, with 5 out of 6 randomized trials showing no improvement over standardized content (Schaeuffele et al., 2025). For example, studies by Johansson et al. (2012) and Berger et al. (2014) found no significant difference between personalized and standardized interventions.

Despite considerable research, three critical gaps persist in our understanding of content personalization in mental health chatbots. First, empirical evidence on the impact of content personalization in mental health chatbots remains absent (Kankanhalli et al., 2021; Kocaballi et al., 2019). For content personalization in digital mental health interventions more broadly, a recent systematic review revealed that five of six randomized trials found no benefits for user engagement and effectiveness compared to standardized content delivery (Schaeuffele et al., 2025). All studies used either purely system-led or user-led personalization, which may explain why their personalization efforts have been insufficient. Moreover, no studies have compared multiple personalization approaches under controlled conditions (Kankanhalli et al., 2021; Schaeuffele et al., 2025). Second, the field lacks understanding of how different personalization approaches interact with user characteristics and stated preferences to influence engagement and health outcomes (Diederich et al., 2022; Kankanhalli et al., 2021). Current implementations assume uniform effects across users, overlooking potential differential impacts based on user characteristics or stated preferences, which could have contributed to the null findings. Third, existing literature lacks theory-grounded prescriptive design knowledge for implementing content personalization in mental health chatbots. While commercial mental health chatbots employ various personalization features and their overall effectiveness has been demonstrated (Fitzpatrick et al., 2017; Inkster et al., 2018; Mehta et al., 2021), the underlying personalization approaches, their theoretical foundations and empirical effects remain unclear. Previous design research has addressed chatbot social roles, anthropomorphic features, and interface personalization (Ahmad et al., 2022; Paul et al., 2024), yet content personalization

lacks design guidance. To address these gaps, we followed the DSR paradigm to derive theory-grounded design requirements, implement and evaluate the effect of multiple personalization approaches on user engagement and effectiveness, and examine their differential effects across user characteristics and symptom severity in controlled field experiments.

4.3 Designing Content Personalization in Mental Health Chatbots

We combined the DSR paradigm (Hevner et al., 2004) to explore the design of content personalization in mental health chatbots with a behavioral research approach to empirically investigate its effects on the symptoms of depression in young people. The DSR paradigm is well-suited to guide our research because it allows us to develop an innovative solution for a real-world problem. In addition, by empirically investigating the effects of our artifact on real users, we can contribute to the understanding of the effects and mechanisms of content personalization (Baskerville et al., 2018). In this section, we outline our design process and describe our design artifact 'TheraBot', which subsequently served as the foundation for our two studies on the effects of content personalization.

4.3.1 Design Process

We adopted the DSR process model by Peffers et al. (2007) to guide our interdisciplinary research project and its design process (see Figure 4.1). Following a problem-centered initiation, we conducted our research through three consecutive design cycles that iteratively explored the design of content personalization in mental health chatbots. In the first design cycle, we began with an extensive exploration of the problem space by reviewing the literature on the personalization of DMHIs and mental health chatbots, as well as commercial mental health chatbots. In addition, we conducted interviews with five mental health experts and 14 young people (avg. 16 years old) with depression. Based on our findings and kernel theories (self-determination theory, effort-accuracy framework, and cognitive load theory), we developed interactive prototypes and gathered feedback. In the second design cycle, we incorporated participants' feedback and developed a fully functional chatbot (TheraBot) with hybrid content personalization capabilities. This cycle involved deploying TheraBot as a mobile application and conducting a

field experiment (Study 1) with 74 participants to evaluate the effects of hybrid content personalization compared to standardized content delivery. The divergent preferences regarding who controls personalization, identified in the first cycle, shaped our understanding of the design space. While the hybrid approach aimed to reconcile user-led and chatbot-led personalization, the divergent opinions required us to compare these approaches directly. Therefore, the third design cycle focused on comparing the effects of chatbot-led, user-led, and hybrid content personalization approaches on user engagement and depression severity in a second field experiment (Study 2) with 160 participants.

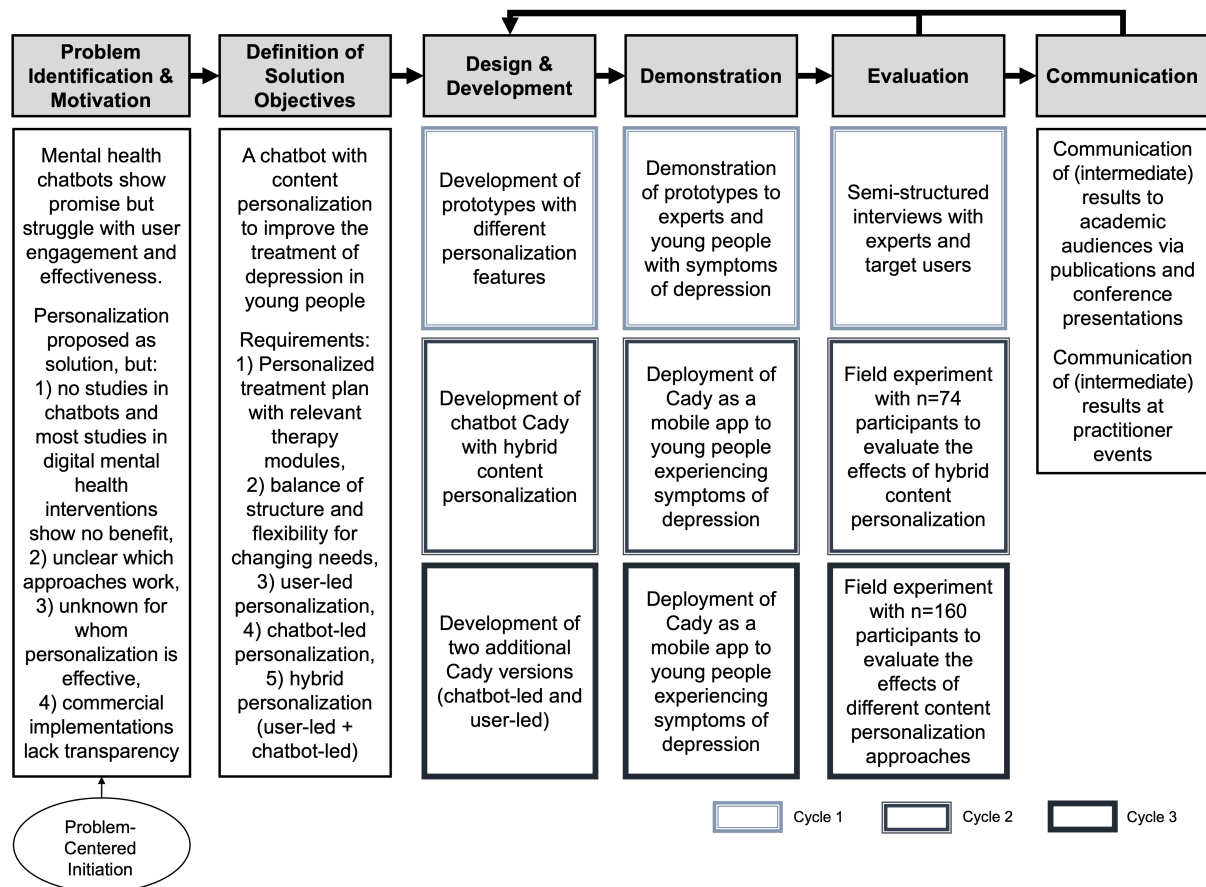


Figure 4.1: Design Science Research Approach based on Peffers et al. (2007)

4.3.2 Design Artifact: TheraBot

To address the complexities that arise when implementing personalization (Cohen et al., 2021; Fan & Poole, 2006; Kocaballi et al., 2019), such as choosing the right level of intervention, time, structure and automation of personalization, we derived requirements for a chatbot for young people with depressive symptoms that incorporates content personalization. In the fol-

lowing, we first describe the requirements before explaining their instantiation in our artifact TheraBot.

Design Requirements

Interviews with young people with depressive symptoms and mental health experts revealed a diverse set of views. Young people stated that they would like to receive the content (e.g., CBT modules, such as behavioral activation or cognitive restructuring) that is most relevant to their specific symptoms. Mental health experts agreed that providing users with the most relevant content is crucial. They also suggested setting up a treatment plan in the beginning that consists of the most relevant therapy modules to ensure guidance and transparency. This is consistent with established practice in CBT and many other psychotherapy modalities, in which a structured treatment plan is used to guide the therapeutic process (Huisman & Kangas, 2018). Combining these considerations led to the first key requirement: the chatbot should offer a well-defined treatment plan consisting of the most relevant therapy modules based on the users' symptoms. Mental health experts said that if new and sufficiently relevant symptoms arise, the treatment plan could need to be revised. This highlights a second key requirement for mental health chatbots: The treatment plan should strike a balance between temporal consistency and flexibility and therefore include personalization before and during treatment. To fulfill these two requirements, it is necessary to determine the most relevant therapy modules, which can be achieved using two different methods: (1) allowing users to choose the most relevant therapy modules themselves and (2) allowing the chatbot to assess the needs and preferences of users to select the most relevant therapy modules (Adomavicius & Tuzhilin, 2005; Fan & Poole, 2006). Both methods can also be combined in a hybrid approach in which both the chatbot and the user are involved in the personalization process. Asked about personalization, one of the young people said: *"I want to select the topics I want to work on,"* and *"check how it [chatbot-led personalization] is and if it is not good, then I would like to adjust it"*. Another potential young person noted that at least she would like to be asked what her preferences are. Each young person highlighted the importance of being involved in the personalization process to determine the relevance of the therapy modules (i.e., user-led personalization). Increasing the relevance for patients by incorporating their preferences has also been shown to improve mental health outcomes in face-to-face psychotherapy (Swift et al., 2018).

The preference for user-led personalization aligns with self-determination theory, which posits

that autonomy is critical for successful behavior change and positive outcomes in psychotherapy (Ryan & Deci, 2008). When individuals experience greater autonomy in the therapeutic process, they are more likely to engage in learning and behavior change, resulting in better treatment outcomes (Ryan & Deci, 2008). Conversely, when the process is perceived as externally controlled, individuals may experience resistance, leading to dropout. The theory therefore suggests that user-led personalization is needed to enhance users' sense of autonomy when selecting therapy modules to be motivated and adhere to the treatment plan. This informs our third design requirement (DR3): provide user-led personalization features to support autonomy. In contrast, mental health experts argued that users might lack the expertise or self-awareness to choose relevant therapy modules or might be burdened with the task. Thus, they advocated for chatbot-led personalization grounded in knowledge of depression. Some young people stated they might not know which content meets their needs and pointed out the effort involved in choosing the most suitable modules. The preference for chatbot-led personalization aligns with the effort-accuracy framework (W. Wang & Benbasat, 2009, 2013) and cognitive load theory (Sweller, 1988). The effort-accuracy framework suggests that users prefer strategies that minimize effort while maintaining decision accuracy (W. Wang & Benbasat, 2009, 2013). Chatbot-led personalization achieves this by selecting relevant modules based on clinical knowledge. Cognitive load theory further supports chatbot-led personalization by emphasizing that too many choices can overwhelm users, leading to decision fatigue, especially for users with limited therapy knowledge or reduced self-awareness. This is particularly relevant for our target users, as depression is characterized by impairments in memory and attention (Rock et al., 2014), making it even more important to minimize cognitive load during content personalization. This informs our fourth design requirement (DR4): provide chatbot-led personalization features to reduce cognitive load and effort. Given that each approach has limitations, a balance between user autonomy and effort is needed. This led to the fifth design requirement (DR5): hybrid content personalization that involves both the user and chatbot in the personalization process.

Design Instantiations

To satisfy these requirements, we developed a mental health chatbot called TheraBot. In the following, we provide an overview of its therapeutic content, system architecture, and personalization features.

Therapeutic Content

We developed five chat-based therapy modules based on CBT for young people (Abel & Hautzinger, 2013; Groen & Petermann, 2015; Towery, 2016): (1) behavioral activation, (2) cognitive restructuring, (3) emotion regulation, (4) sleep, and (5) interpersonal skills. Each module comprises three to five sessions in which TheraBot guides users through the content of the module using structured conversations. In each module, TheraBot explains why the module is important and how it works, guides users through exercises, and asks them to apply the learned skills in their daily lives. For example, in the cognitive restructuring module, TheraBot first explains the concept of automatic negative thoughts, shares examples, and asks users to start a thought diary. Next, TheraBot introduces cognitive biases that underlie negative thoughts, performs an exercise in which users learn to detect these biases, and asks them to include these cognitive biases in their thought diary. Finally, TheraBot explains how to challenge negatively biased thoughts, achieve more balanced thoughts, and practices reframing negative thoughts with users. It also encourages them to apply the reframing technique to negative thoughts in their daily lives. The chats are rule-based, but also include unrestricted natural language input (e.g., during exercises). Chat-based therapy modules that combine structured chats and unrestricted natural language input are used by most mental health chatbots (Ahmed et al., 2021).

System Architecture

The system architecture comprises two main components: (1) a frontend user interface through which users interact with the chatbot and (2) a backend service responsible for maintaining the therapy modules, handling chat interactions, and collecting and storing user data (see Figure 4.2). We developed the frontend in Flutter (flutter.dev), an open-source user interface framework for cross-platform mobile applications on Android and iOS devices. The frontend of the TheraBot application contains a home screen, a modules screen with all therapy modules, an overview screen for each module, and a chat screen. A professional design agency supported us in designing the user interface. We made the TheraBot app available to iOS users via TestFlight (testflight.apple.com) and to Android users via Firebase (firebase.google.com).

The second key component is the application backend, which consists of a backend core, a chatbot engine, and a database. We used existing open-source frameworks and libraries to develop

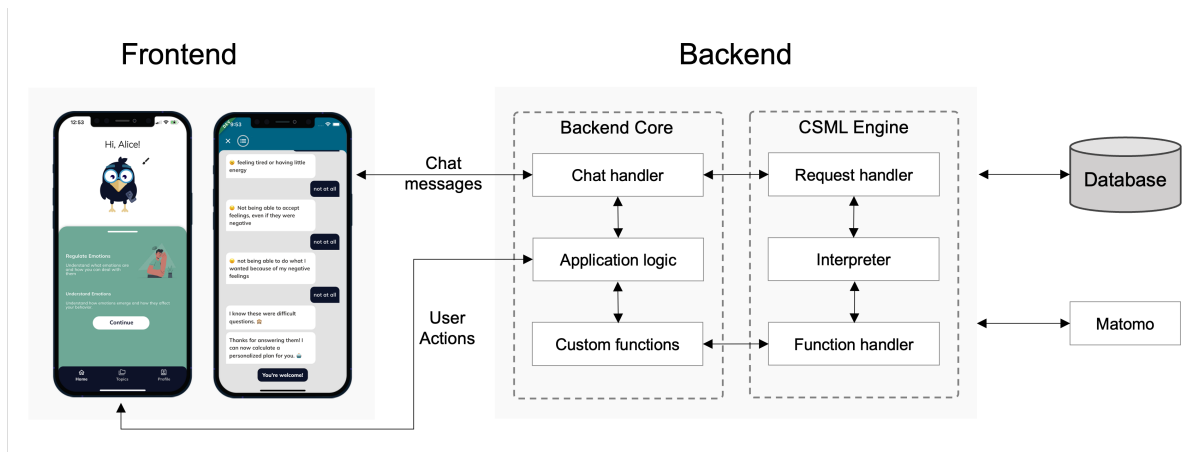


Figure 4.2: TheraBot's System Architecture

the backend, which enabled us to design flexible chatbot interactions while simultaneously maintaining control over data processing and storage, a requirement set by our institutional review board. The REST-based backend core handles all API requests from the frontend (e.g., chat messages by users) and is responsible for the application logic (e.g., marking modules as completed). The chatbot engine is based on CSML (csml.dev), which is a domain-specific chatbot programming language. We chose CSML for three reasons: (1) it is open source and therefore customizable, (2) it provides full control of the chatbot's inner workings, and (3) it does not depend on external services outside of our own environment. The main responsibility of the chatbot engine is to handle chat interactions between users and the chatbot, including processing the user's answers and sending messages to users. In addition, we deployed a PostgreSQL database (postgresql.org) to store user settings (e.g., username and treatment plan), track user states (e.g., current module), and store the chats. We incorporated Matomo (matomo.org), an open-source, self-hosted web analytics platform, to record user data during the study (e.g., clicks on modules, sessions, and personalization features). The application backend, PostgreSQL database, and Matomo analytics platform were deployed in a Kubernetes cluster (kubernetes.io) owned and operated at our institution.

Personalization Features

We implemented a two-phase hybrid content personalization strategy, in which TheraBot and the user contribute to the selection of therapy modules, resulting in a personalized treatment plan. This strategy resembles the shared decision model in healthcare, where treatment options are personalized to the patient's needs and preferences in a collaborative process (Stump

& Sauer-Zavala, 2022). Several commercial mental health chatbots also implement a similar strategy (see Table B.1 in the Appendix), which combines a standardized mental health questionnaire with showing users a list of therapy modules to choose from in their onboarding process. However, these commercial mental health chatbots do not disclose how personalization is performed. At the end of this two-phase hybrid personalization process, every user has a personalized treatment plan, addressing the derived design requirements described above.

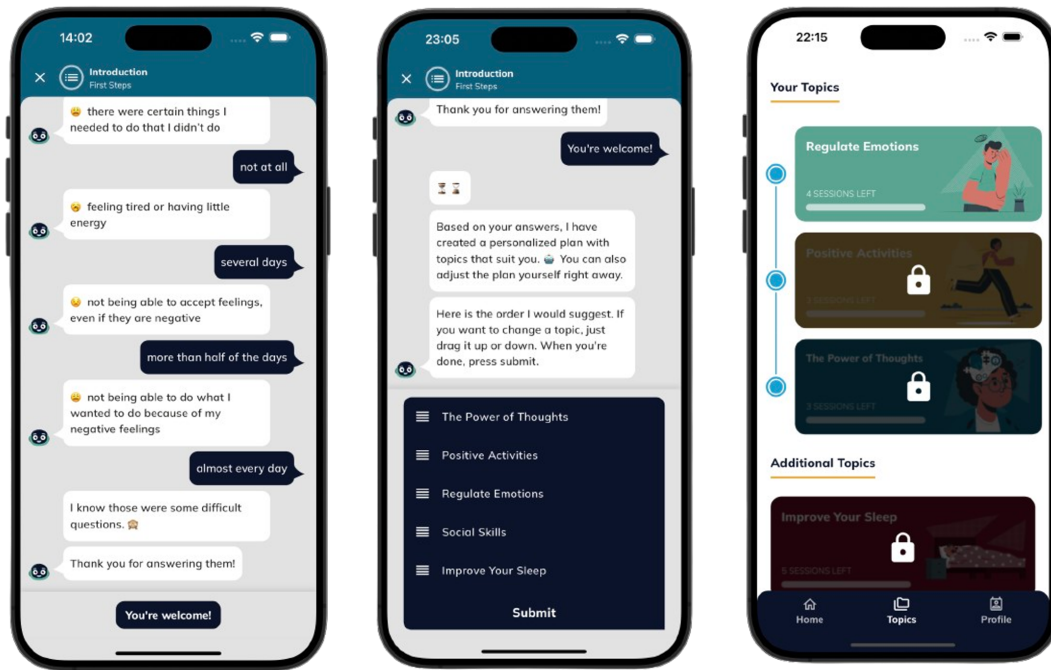


Figure 4.3: Screenshots of TheraBot's Personalization Features

Note. Screenshots of chatbot-led (left) and user-led (middle) personalization, and its results on the 'Topics' (i.e. modules) screen (right). The screenshots were translated into English.

In the first phase, the chatbot leads the personalization process by asking users questions to assess their symptoms and recommend a personalized selection of modules. In phase two, users accept the recommendation or adapt it to their perceived needs and preferences. This approach requires TheraBot to have two key personalization features. First, it requires system-led personalization of the treatment plan, in which TheraBot selects therapeutic modules based on the users' symptoms. To assess the symptoms of users and determine the most relevant therapeutic modules, we selected items from established questionnaires such as PHQ-9, CES-D, BADS, CFQ and ERSQ (China et al., 2018; Grant et al., 2018; Kroenke et al., 2009; Manos et al., 2011; Radloff, 1977). The user responds to the items delivered via chat on a four-point Likert scale (from 0 to 3), which allows TheraBot to calculate the relevance (average across all items for the specific module) from 0–3 for each therapy module (see Table B.2 in the Appendix

for details). Based on the relevance score, TheraBot selects therapy modules that are intended to address the user's most pressing symptoms. We implemented this feature in the onboarding chat so that after the users answer the items, the chatbot instantly provides the personalized treatment plan (see Figure 4.3). Second, users must be able to choose therapeutic modules that match their perceived needs and preferences. To do so, TheraBot allows users to rearrange the recommended treatment plan by drag-and-drop during the onboarding chat (see Figure 4.3). Ultimately, a personalized treatment plan is displayed on the module screen. To address the second requirement of balancing the structure and flexibility of the treatment plan, TheraBot repeats this hybrid personalization process during use to check whether the relevance of the therapy modules has changed and the treatment plan needs to be adapted.

4.4 Study 1: Investigating the Effect of Content Personalization

The primary objective of our first study was to evaluate whether hybrid content personalization leads to greater reductions in depression severity than standardized content. Therefore, we conducted a randomized controlled field experiment in which participants interacted with TheraBot for two weeks. Study 1, as well as study 2, received approval from the Institutional Review Board of the Karlsruhe Institute of Technology.

4.4.1 Experimental Design

The study used a between-subjects design with two experimental conditions to compare two versions of TheraBot: one with hybrid content personalized by the chatbot and the user and the other one with standardized content. As detailed in the previous section, TheraBot determines the selection and sequence of therapy modules during the onboarding chat and during use by asking several questions about their symptoms of depression and then allowing the users to adapt it to their own perceived needs and preferences. In the condition with standardized content, which served as the control group, the participants had to follow a prescribed selection and sequence of therapy modules. This is common in DMHIs, where the plan is often implemented according to CBT treatment manuals (Abel & Hautzinger, 2013; Towery, 2016). The collection of therapy modules was the same in both conditions; therefore, the only difference between the

two conditions was whether the selection of modules was personalized or standardized.

4.4.2 Participants

We recruited participants from the participant pool of our university's experimental lab (Greiner, 2015). Participants were eligible if they had symptoms of depression (PHQ-8 sum score ≥ 6) and had a smartphone to install the TheraBot application. Of the 95 participants who met the eligibility criteria and completed the pre-experiment questionnaire, 74 completed the onboarding chat with TheraBot, used it for 14 days without experiencing technical issues, and ultimately completed the post-experiment questionnaire. Participants received 25 euro for their participation in the study. The participants were roughly evenly distributed across the two conditions (standardized: $n = 41$, 55%; hybrid personalization: $n = 33$, 45%). The sample consisted of 36 (49%) male and 37 (50%) female participants and one participant who did not want to disclose their gender, with a mean age of 24.5 ($SD = 3.5$). 24 participants (32.4%) reported a previous diagnosis of a mental disorder. Based on their initial PHQ-8 score ($M = 11.3$, $SD = 4.0$), 29 (39%) participants showed mild, 30 (41%) moderate, 12 (16%) moderately severe, and 3 (4.1%) severe depression severity. The characteristics of the sample are presented in Table 4.1.

4.4.3 Procedure

First, we sent an email describing the study to the entire participant pool. Interested participants completed an online screening questionnaire to verify their eligibility. Eligible participants received information about the study procedure and data protection and gave their informed consent. Subsequently, they answered questions about their mental health, demographic information, and previous experiences with chatbots and psychotherapy. Subsequently, we randomly assigned the participants to one of two experimental conditions based on a stratified randomization approach considering the operating system of their smartphone (Android, iOS), previous experience with psychotherapy and gender. We instructed participants to install the TheraBot app on their smartphones and to use it for two weeks in their daily lives. After the study, the participants completed another questionnaire on their mental health and experiences with TheraBot. The procedure was pre-tested with six participants and refined on the basis of their feedback.

Table 4.1: Sample Characteristics of Study 1

Variable		Study 1 (n = 74)
Gender	Male	36 (49%)
	Female	37 (50%)
	No Disclosure	1 (1%)
Age		$M = 24.5(SD = 3.5)$
Diagnosed Mental Disorder	Depressive Disorder	17 (23%)
	Anxiety Disorder	2 (3%)
	Other Disorder	5 (7%)
	No Disorder	50 (68%)
Experience with Mental Health Apps	Yes	24 (32%)
	No	50 (68%)
Experience with Chatbots		$M = 4.6(SD = 2.3)$
Experience with Psychotherapy	Yes	23 (31%)
	No	47 (64%)
	No Disclosure	4 (5%)
Baseline Depression Severity (PHQ-8)	Mild	29 (39%)
	Moderate	30 (41%)
	Moderately Severe	12 (16%)
	Severe	3 (4%)
Operating System	iOS	30 (41%)
	Android	44 (59%)

4.4.4 Measurement

We used the PHQ-8 (Kroenke et al., 2009), a widely used self-report questionnaire to assess depression severity, before and after the experiment. The PHQ-8 assesses how often eight central symptoms of depression according to the DSM-5 (American Psychiatric Association, 2013) of Mental Disorders (e.g., little interest or pleasure in doing things or depressive mood) have occurred during the last two weeks. The eight symptoms are assessed on a 4-point scale from ‘not at all’ to ‘every day’. We also measured user engagement with TheraBot by calculating the number of completed therapy modules. Age and gender were collected as demographic variables. Participants were also asked about previous diagnoses of mental disorders, as well as previous experiences with chatbots, mental health apps, and psychotherapy. All measurement items are shown in Table B.3 in the Appendix.

4.4.5 Data Analysis and Results

Randomization Checks

To assess the success of our randomization procedure, we compared the experimental conditions on several control variables. There were no significant differences in gender ($\chi^2(2) = 1.96, p = 0.37$), age ($F(1, 72) = 0.89, p = 0.35$), previous diagnoses of mental disorders ($\chi^2(3) = 1.33, p = 0.72$), operating system ($\chi^2(1) = 0.80, p = 0.37$), previous experience with chatbots ($F(1, 72) = 0.66, p = 0.42$), mental health apps ($\chi^2(1) = 0.00, p = 0.99$), and psychotherapy ($F(1, 72) = 0.03, p = 0.86$), suggesting that the randomization was successful.

Effect of Hybrid Content Personalization on User Engagement and Depression Severity

Second, we examined the general effect of the use of TheraBot on depression severity. We analyzed whether depression severity changed between the two assessments before and after the experiment, using a linear mixed-effects model. The model included assessment time as a fixed effect and participants as a random effect, to account for repeated measures per participant. Participants showed an average reduction of 11% in depression severity after the experiment ($\beta = -1.80, 95\% \text{ CI } [-2.33, -1.28], p < 0.001$). Participants with mild, moderate, moderately severe, and severe depression before using TheraBot, measured with PHQ-8, experienced average reductions of 2%, 17%, 16%, and 24%, respectively, indicating that those

with more severe depression symptoms experienced more substantial reductions.

To examine the effect of hybrid personalization, we compared the hybrid with the standardized condition. We added an interaction term between the experimental condition and the assessment time to investigate whether the changes in depression severity differed between standardized and hybrid content. Our results showed that the average reduction in depression severity was greater in the hybrid condition ($M_{\text{after-before}} = -2.49$, $SD = 3.37$) than in the standardized condition ($M_{\text{after-before}} = -0.88$, $SD = 3.49$; difference = -1.61 , 95% CI [-3.21 , -0.00], $p = 0.048$). The results are shown in Table 4.2 and the visual change in depression severity is illustrated in Figure 4.4.

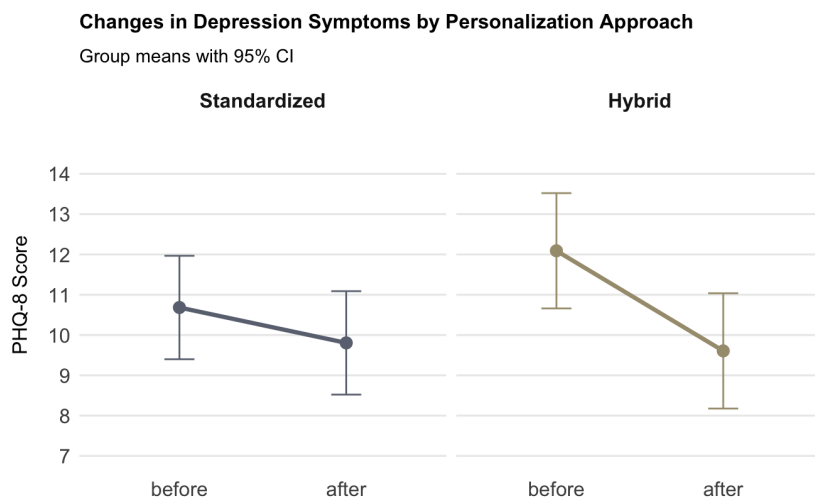


Figure 4.4: Changes in depression severity by personalization approach.

Table 4.2: Results of Linear Mixed Effects Models for Hybrid vs. Standardized Content in Study 1

Predictor	Estimate	95%-CI	p-value
Time (Post)	-0.88	-1.94 - 0.18	0.028
Group (Hybrid Personalization)	1.41	-0.51 - 3.32	0.148
Time (Post) × Group (Hybrid Personalization)	-1.61	-3.20 - -0.02	0.048

As a robustness check, we repeated our analyses with several control variables (age, gender, diagnosis of a mental disorder, previous experience with chatbot use, mental health apps or psychotherapy, and smartphone operating system) and found that all results remained consistent when control variables were included. We also examined the effect of personalization on user engagement. On average, the participants completed 2.89 (57.8%, $SD = 1.68$) of the five

modules. Although participants in the hybrid condition completed slightly more modules (3.09, $SD = 1.76$) than participants in the standardized condition (2.73, $SD = 1.61$), the difference was not statistically significant ($p = 0.4$), as shown in Figure 4.5.

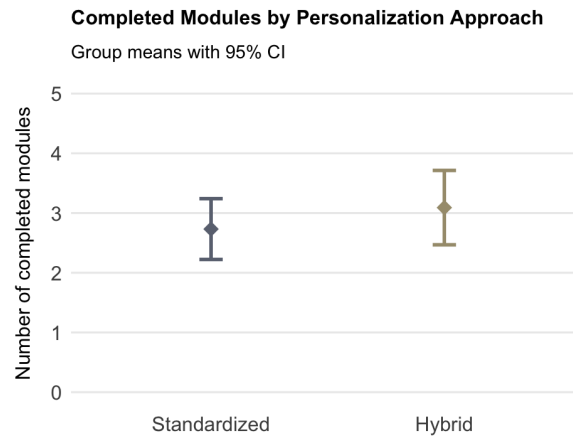


Figure 4.5: Number of completed modules by personalization approach.

In summary, the findings from the first field experiment show that the use of TheraBot can improve symptoms of depression and that incorporating hybrid content personalization can enhance the effectiveness of mental health chatbots.

4.5 Study 2: Examining the Effects of Chatbot-led, User-led, and Hybrid Content Personalization

Study 1 provided initial evidence that a chatbot with hybrid content personalization leads to larger reductions in depression severity than standardized content, underscoring the value of content personalization and validating the design of our artifact, TheraBot. However, because personalization frameworks offer a wide range of dimensions, starting points and practical implementations (Cohen et al., 2021; Fan & Poole, 2006), studies on the personalization of DMHIs have provided inconclusive evidence and only implemented system-led personalization (Berger et al., 2014; Chaturvedi et al., 2023; Johansson et al., 2012; Levin et al., 2019), and commercial mental health chatbots adopt different personalization approaches (see Table B.1 in the Appendix), a more comprehensive understanding is necessary. To address this, we conducted a second study that mirrored the first study but included a full experimental design with two additional personalization approaches (chatbot- and user-led). The objective of the

second study was to compare the effects of different personalization approaches.

4.5.1 Experimental Design

Study 2 employed a between-subjects design with four experimental conditions: (1) standardized content, (2) chatbot-led, (3) user-led, and (4) hybrid content personalization. The standardized condition and the hybrid condition were the same as in Study 1. In the standardized condition, all participants followed a standardized selection and sequence of therapy modules. In the hybrid content personalization condition, the chatbot recommended a selection and sequence of therapy modules to the participants that they were able to change based on their own perceived needs and preferences. In the chatbot-led condition, only the chatbot-led feature of hybrid content personalization was implemented, which used the same questions on the users' most relevant symptoms of depression as in the hybrid condition; however, the users had to follow the chatbot-led personalization and were not able to change it. In the user-led condition, only the user-led feature was implemented and the participants had to personalize the content on their own, without any recommendations from the chatbot. To operationalize these conditions, we developed four different versions of TheraBot (see Table 4.3) that differed only in the way the content was personalized.

Table 4.3: Experimental Conditions

Experimental Conditions	Description
Standardized	Participants had standardized content: They had to follow a standardized selection and sequence of therapy modules.
Chatbot-led Personalization	The chatbot chose the therapy modules based on the participants' depressive symptoms.
User-led Personalization	Participants chose the therapy modules according to their own needs and preferences.
Hybrid Personalization	The chatbot recommended therapy modules based on the participants' depressive symptoms but participants could change the recommendation according to their own needs and preferences.

4.5.2 Participants

We again recruited participants with symptoms of depression from the participant pool (Greiner, 2015) of the experimental lab of our university. Participants of Study 1 were not invited. Of the 365 participants who answered the screening questions, 230 were eligible according to the same criteria as in Study 1. Of these, 174 participants installed the TheraBot app on their smartphones and completed the onboarding session. Finally, 160 participants filled out the questionnaire after the 14-day study period, comprising the final sample. The participants again received 25 euro for participating in the experiment. Participants were equally distributed across the four experimental conditions (standardized = 37, chatbot-led = 41, user-led = 40, and hybrid = 42). The sample consisted of 77 (48%) male and 82 (51%) female participants, as well as a participant who did not want to disclose their gender, with a mean age of 22.81 years ($SD = 2.63$). 38 participants (24%) reported having previously received a diagnosis of mental disorder. Before the experiment, 70 (44%) participants showed mild, 60 (38%) moderate, 24 (15%) moderately severe, and 6 (3.8%) severe depression severity based on the PHQ-8 score ($M = 11$, $SD = 3.9$). The characteristics of the sample are presented in Table 4.4.

4.5.3 Procedure

We repeated the procedure from Study 1. Interested members of the participant pool first verified their eligibility through the screening questionnaire and then gave their informed consent. Eligible participants answered the first questionnaire and were randomized into one of the four conditions based on a stratified randomization approach considering their smartphone operating system (Android, iOS), previous experience with psychotherapy, previous diagnosis of a mental disorder and gender. Afterwards, they installed TheraBot on their smartphones and were asked to use it for two weeks in their daily lives. After the two-week period, participants completed a post-experiment questionnaire on their mental health and experiences with TheraBot.

4.5.4 Measurement

As in Study 1, we assessed depression severity before and after the experiment using the PHQ-8 (Kroenke et al., 2009). In addition, we measured user engagement by calculating the number of completed modules and overall usage time of TheraBot. We also assessed several control vari-

Table 4.4: Sample Characteristics of Study 2

Variable		Study 2 (n = 160)
Gender	Male	77 (48%)
	Female	82 (51%)
	No Disclosure	1 (0.6%)
Age		$M = 22.8(SD = 2.6)$
Diagnosed Mental Disorder	Depressive Disorder	13 (8%)
	Anxiety Disorder	9 (6%)
	Other Disorder	16 (10%)
	No Disorder	115 (72%)
	No Disclosure	7 (4%)
Experience with Mental Health Apps	Yes	55 (34%)
	No	105 (66%)
Experience with Chatbots		$M = 5.1(SD = 2.09)$
Experience with Psychotherapy	Yes	62 (38%)
	No	91 (57%)
	No Disclosure	7 (4%)
Baseline Depression Severity (PHQ-8)	Mild	70 (44%)
	Moderate	60 (38%)
	Moderately Severe	24 (15%)
	Severe	6 (4%)
Operating System	iOS	89 (56%)
	Android	71 (44%)

ables: Age, gender, prior diagnoses of mental disorders, previous experiences with chatbots, mental health apps, and psychotherapy, and engagement with mental health related activities during the study. At the end of the post-experiment survey, we disclosed all four experimental conditions and asked them which one they would prefer and why. Table 4.5 provides an overview of all variables.

4.5.5 Data Analysis

Randomization Checks

We first examined the success of the randomization procedure by comparing the four experimental conditions on several control variables. There were no significant differences in gender ($\chi^2(6) = 3.14, p = 0.79$), age ($F(3, 156) = 1.68, p = 0.17$), previous diagnoses of mental disorders ($\chi^2(12) = 9.19, p = 0.69$), operating system ($\chi^2(3) = 0.40, p = 0.94$) and previous experience with chatbots ($F(3, 156) = 0.59, p = 0.62$), mental health apps ($F(3, 156) = 0.69, p = 0.56$), or psychotherapy ($\chi^2(9) = 2.14, p = 0.99$), suggesting that randomization was successful.

Effects of Different Personalization Approaches on Depression Severity

Descriptive statistics for the key constructs in the experimental conditions are shown in Table 4.5. First, as in Study 1, we examined the overall effect of TheraBot on depression severity. We analyzed whether depression severity changed between the two assessments before and after the experiment, using a linear mixed-effects model. The model included assessment time as a fixed effect and participants as a random effect, to account for repeated measures per participant. The results revealed a main effect of the assessment time. Participants showed an average 16% reduction in depression severity after the experiment ($\beta = -2.11, 95\% \text{ CI } [-2.67, -1.54], p < 0.001$). Participants with mild, moderate, moderately severe, and severe depression before the experiment experienced average reductions of 6%, 25%, 22%, and 26%, respectively, suggesting that those with more severe depression achieved larger reductions, which replicates the findings of Study 1.

Second, to examine the overall effect of content personalization, we compared the standardized condition against the three personalization conditions combined. We added an interaction term between the experimental condition and the assessment time to the model to investigate whether

Table 4.5: Key Variables Across Experimental Groups in Study 2

Data Type	Variable	Experimental Conditions				Total
		Personalized				
		Standardized	Chatbot	User	Hybrid	
Self-Report Questionnaire	Depression Severity (Pre)	10.86 (3.92)	10.27 (3.23)	11.13 (3.99)	11.55 (4.51)	10.96 (3.93)
	Depression Severity (Post)	9.78 (3.68)	8.17 (3.66)	8.80 (4.01)	8.74 (4.83)	8.85 (4.09)
	Change in Depression	-1.1 (2.8)	-2.10 (4.08)	-2.33 (3.21)	-2.81 (4.14)	-2.11 (3.64)
Behavioral Usage Data	Completed Modules	2.53 (1.68)	1.98 (1.56)	2.17 (1.77)	2.41 (1.86)	2.27 (1.72)
	Overall Usage Time	75 (46)	66 (40)	74 (56)	83 (55)	75 (50)

Notes: Means with standard deviations in parentheses.

4.5. Study 2: Examining the Effects of Chatbot-led, User-led, and Hybrid Content Personalization

changes in depression severity differed between personalized and standardized content. Our results showed that, on average, the reduction in depression severity was higher for users in the content personalization conditions ($M_{\text{after} - \text{before}} = -2.40$, $SD = 3.8$) than in the standardized content condition ($M_{\text{after} - \text{before}} = -1.10$, $SD = 2.8$; difference = -1.33 , 95% CI $[-2.66, 0.00]$, $p = 0.049$). This result corroborates the finding of Study 1 that content personalization can enhance the effectiveness of mental health chatbots. See Table 4.6 for the model results.

Table 4.6: Results of Linear Mixed Effects Models for Personalized vs. Standardized Content in Study 2 and from both Studies Combined

Predictors	Estimates	95%-CI	p-value
Study 2 (N = 160)			
Time (Post)	-1.02	-2.64 - 0.60	0.217
Group (Personalization)	-1.24	-2.87 - 0.39	0.134
Time (Post) × Group (Personalization)	-1.33	-2.66 - -0.00	0.049
Combined Data (N = 234)			
<i>Standardized (n = 78), Personalized (n = 156)</i>			
Time (Post)	-0.97	-1.76 - -0.19	0.015
Group (Personalization)	0.63	-0.52 - 1.78	0.283
Time (Post) × Group (Personalization)	-1.46	-2.42 - -0.49	0.003

Third, we examined the effects of the different personalization approaches separately. We analyzed whether the interaction effect between the assessment time and the experimental condition differed between the standardized and three personalization conditions. As can be seen in Figure 4.6, the average PHQ-8 score was lower after the study in all four experimental groups. Pairwise comparisons showed a significant difference in reducing depression severity between standardized ($M_{\text{after} - \text{before}} = -1.10$, $SD = 2.8$) and hybrid treatment plan personalization ($M_{\text{after} - \text{before}} = -2.8$, $SD = 4.1$; difference = -1.73 , 95% CI $[-3.34, -0.12]$, $p = 0.036$). However, neither the difference between standardized treatment and user-led personalization ($M_{\text{after} - \text{before}} = -2.3$, $SD = 3.2$; difference = -1.24 , 95% CI $[-2.87, 0.39]$, $p = 0.134$) nor between standardized and chatbot-led personalization was statistically significant ($M_{\text{after} - \text{before}} = -2.1$, $SD = 4.1$; difference = -1.02 , 95% CI $[-2.64, 0.60]$, $p = 0.217$). There were also no statistically significant differences in pairwise comparisons of the three personalization conditions. The results for the different personalization approaches are shown in Figure 4.6.

The post-experiment preferences of the participants emphasized the promise of hybrid personalization, as the majority (77%) preferred this personalization strategy. They explained their

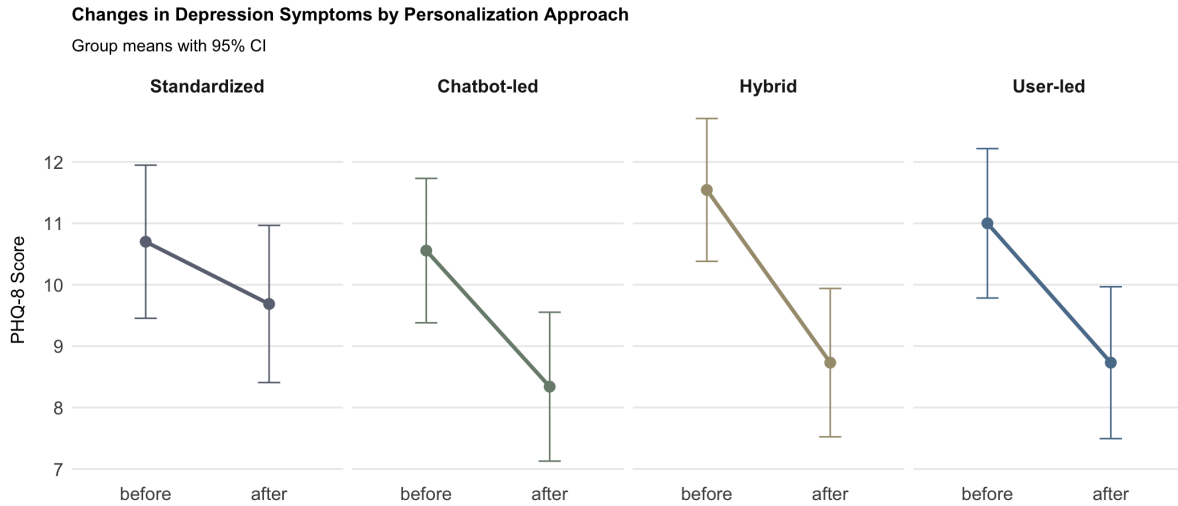


Figure 4.6: Changes in Depression Severity for each Experimental Condition

preference for hybrid content personalization by highlighting the combination of two aspects: receiving qualified content recommendations from the chatbot while retaining autonomy and control. We also performed several robustness checks to validate our findings, including controlling for demographic factors, external mental health activities, and combining data from the two studies to increase statistical power. The checks consistently supported our results, demonstrating that content personalization, particularly hybrid content personalization, is an effective strategy to improve mental health chatbots. The complete results of the comparisons between the four experimental conditions for Study 2 and the combined data can be found in Table 4.7.

Effect of the Personalization Strategy on User Engagement with TheraBot

We assessed user engagement with TheraBot by examining the number of completed modules and overall usage time. On average, the participants completed 2.27 (45.4%, $SD = 1.72$) of the five modules available and used TheraBot for an average of 75 minutes ($SD = 50$). However, pairwise comparisons revealed no statistically significant differences between any of the groups. Figure 4.7 shows module completion across all four groups.

As previous research has shown an effect of depression severity on engagement with DMHIs (Borghouts et al., 2021), we analyzed the interaction effect of prestudy depression severity and personalization strategy on the number of completed modules and usage time. For each, we conducted a separate regression analysis and then calculated pairwise comparisons of the relationship between depression severity and the two user engagement metrics between the ex-

Table 4.7: Results from Pairwise Contrasts between the four Experimental Conditions in Study 2 and the Combined Data

Comparisons	Estimates	95%-CI	p-value
Study 2 (N = 160)			
Standardized (n = 37) vs. Chatbot-led (n = 41)	-1.02	-2.64 - 0.60	0.217
Standardized (n = 37) vs. User-led (n = 40)	-1.24	-2.87 - 0.39	0.134
Standardized (n = 37) vs. Hybrid (n = 42)	-1.73	-3.34 - -0.12	0.036
Chatbot-led (n = 41) vs. User-led (n = 40)	-0.23	-1.82 - 1.36	0.778
Chatbot-led (n = 41) vs. Hybrid (n = 42)	-0.71	-2.28 - 0.86	0.371
User-led (n = 40) vs. Hybrid (n = 42)	-0.49	-2.06 - 1.09	0.545
Combined Data (N = 234)			
Standardized (n = 78) vs. Chatbot-led (n = 41)	-1.12	-2.47, 0.23	0.102
Standardized (n = 78) vs. User-led (n = 40)	-1.35	-2.71, 0.01	0.051
Standardized (n = 78) vs. Hybrid (n = 75)	-1.69	-2.82, -0.56	0.004
Chatbot-led (n = 41) vs. User-led (n = 40)	-0.23	-1.78, 1.33	0.773
Chatbot-led (n = 41) vs. Hybrid (n = 75)	-0.57	-1.93, 0.79	0.410
User-led (n = 40) vs. Hybrid (n = 75)	-0.34	-1.71, 1.03	0.623

perimental groups. For completed modules, the overall test indicated an interaction between PHQ-8 scores and experimental group ($F(3, 165) = 2.42, p = 0.07$), although the effect is only significant at the $p < 0.1$ level. Investigating the interaction in more detail, we found a statistically significant difference between user-led and chatbot-led personalization ($\beta = -0.27, 95\% \text{ CI } [-0.47, -0.06], p = 0.01$). As can be seen in Figure 4.8, the relationship between the PHQ-8 scores before the study and the completed modules was positive in the user-led group. In contrast, the relationship was negative in the chatbot-led group. In the hybrid group, pre-study depression severity had no effect on completed modules. Pairwise comparisons between chatbot-led and hybrid, as well as user-led and hybrid personalization, were not statistically significant. In summary, participants with higher depression severity in the user-led condition tended to complete more modules relative to those with lower severity, whereas this trend was reversed in the chatbot-led condition.

With respect to usage time, we found a similar pattern. The overall test indicated an interaction between PHQ-8 scores and experimental group ($F(3, 165) = 2.56, p = 0.06$), although the effect was again only significant at the $p < 0.1$ level. Pairwise comparisons were significant between the chatbot-led and user-led groups ($\beta = -7.16, 95\% \text{ CI } [-12.96, -1.35], p = 0.02$) and between the hybrid and user-led groups ($\beta = -5.71, 95\% \text{ CI } [-10.69, -0.73], p = 0.03$). Again, the relationship between PHQ-8 scores before the study and overall usage

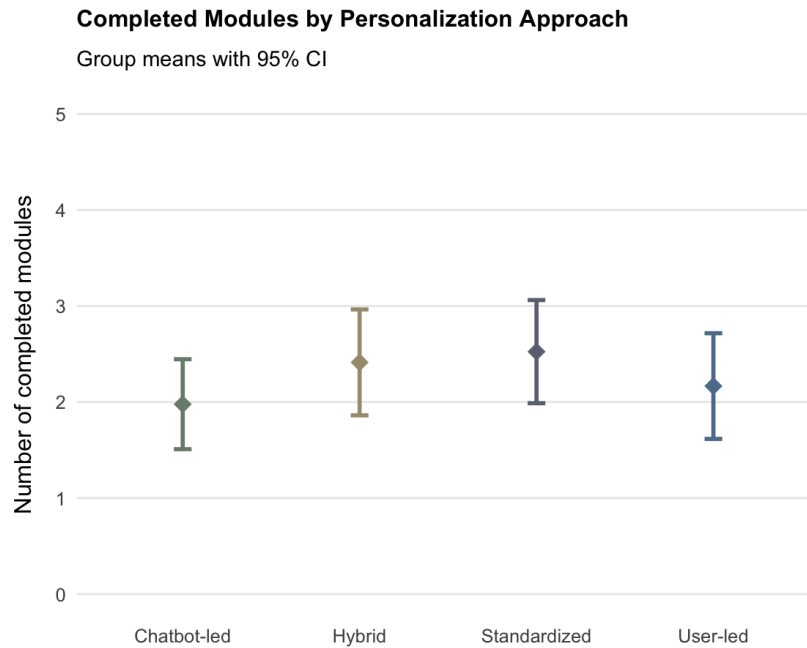


Figure 4.7: Number of Completed Modules across the three personalization approaches

time was positive in the user-led group and negative in the chatbot-led group. Similarly, the negative estimate for the hybrid vs. user-led comparison suggests that participants with higher depression severity in the user-led group spent more time with TheraBot than those in the hybrid group.

In summary, the effect of depression severity on user engagement with TheraBot varied depending on the personalization strategy. Users with higher depression severity engaged more, both in terms of completed modules and usage time, when in the user-led condition compared to chatbot-led or hybrid personalization. This finding highlights the potential to adapt the personalization strategy to depression severity and questions whether providing all users with hybrid content personalization is really the best design choice. Adapting the personalization strategy was also supported by the participants' preferences. Although most (77%) preferred hybrid personalization, 20% of the participants preferred one of the other two content personalization approaches (chatbot-led = 12%, user-led = 8%) and it seems counterproductive to potentially frustrate 20% of the users. Some participants preferred chatbot-led personalization because they trust the chatbot to be objective and qualified and to provide better personalization than the user themselves. They argued that chatbot-led personalization could prevent users from choosing less useful modules. Others favored user-led personalization because they valued the autonomy to choose modules based on their needs. They argued that they know best what

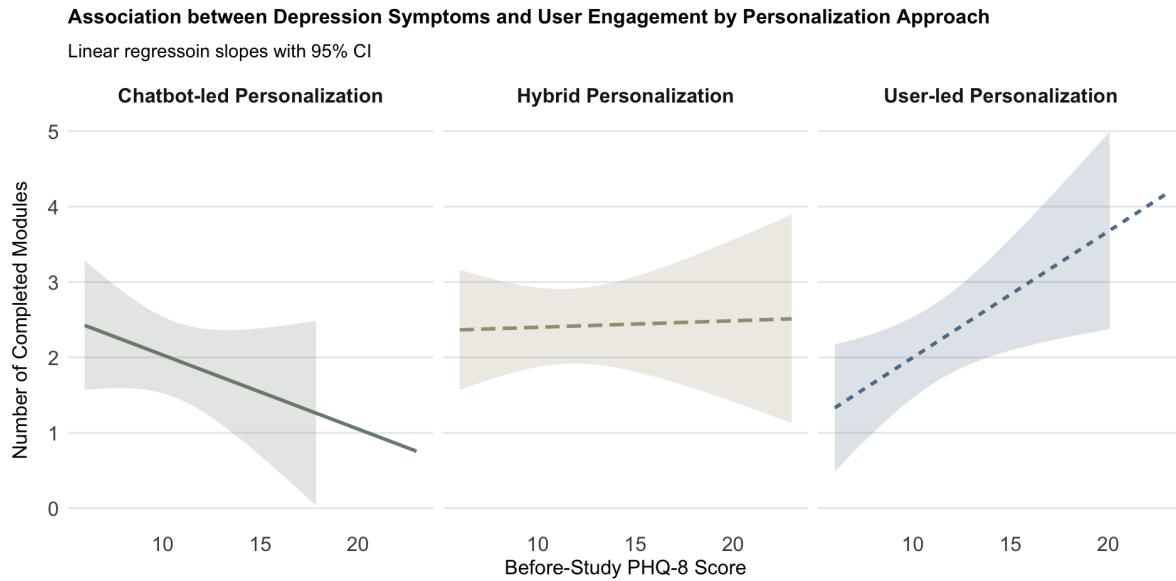


Figure 4.8: Effect of Pre-Study PHQ-8 score on Number of Completed Modules across the three personalization approaches

helps them, suggesting that they don't need chatbot-led recommendations, and were also skeptical about the chatbot's personalization capabilities. In summary, these findings emphasize the complexity of designing personalization beyond the three personalization approaches implemented in TheraBot and suggest the potential to adapt the personalization strategy to the individual user.

User Engagement with the User-led Personalization Feature

To assess the degree to which participants in the user-led and hybrid conditions used the user-led personalization feature, we compared the similarity between the suggested therapy modules and the final modules selected by each participant. We measured similarity using the Kendall correlation coefficient, which quantifies the correspondence between the order of suggested modules (either randomly generated or recommended by the chatbot) and the order chosen by the participant. We set negative correlation values to zero, as negative values also indicate dissimilarity. High similarity values indicate that participants closely followed the suggested module sequence, making few or no changes. In contrast, low similarity values reflect substantial modifications to the suggested order. For example, if the suggested sequence was behavioral activation, cognitive restructuring, sleep, emotion regulation, and interpersonal skills, but the participant reordered these to cognitive restructuring, sleep, behavioral activation, interpersonal skills, and emotion regulation, the resulting similarity would be low. In the user-led

condition, the chatbot informed the participants that the module suggestions were randomized, whereas in the hybrid condition, the chatbot made suggestions based on participants' responses to a symptom questionnaire. This distinction clarifies the basis for the recommendations and contextualizes participants' engagement with the personalization process.

For the hybrid group, the mean similarity was 0.87 (SD = 0.28), whereas the user-led group showed a substantially lower mean similarity of 0.37 (SD = 0.40). Statistical comparison confirmed a significant difference between groups ($\beta = -0.49$, 95% CI [-0.64, -0.35], $p < 0.001$), indicating that participants in the hybrid condition adhered more closely to the suggested modules than those in the user-led condition. We also examined the proportion of participants who either fully adhered to or completely ignored the suggested module selections. Most participants in the hybrid condition (74%) fully adhered to the chatbot's suggested therapy module selection. In contrast, the user-led condition exhibited greater variability: 24% fully adhered to the suggestions, whereas 36% ignored them entirely. This pronounced difference in adherence patterns suggests that the perceived relevance or credibility of receiving recommendations from the chatbot may have influenced engagement with the user-led personalization process.

Next, we examined whether participant characteristics, particularly personality traits, influenced engagement with the user-led personalization feature. We focused on personality traits because prior work identified them as important factors affecting engagement with DMHIs (Borghouts et al., 2021). First, we examined the role of agreeableness. Across both groups, higher agreeableness was significantly associated with greater similarity between suggested and final module selections ($\beta = 0.03$, 95% CI [0.00, 0.07], $p = 0.047$), indicating that more agreeable participants tended to follow the recommendations more closely. To determine whether this association differed by personalization approach, we tested the interaction between experimental condition and agreeableness. Again, higher agreeableness was associated with greater similarity between suggested and final module selections ($\beta = 0.03$, 95% CI [0.00, 0.07], $p = 0.047$) and the user-led condition exhibited lower similarity than the hybrid condition ($\beta = -1.20$, 95% CI [-2.10, -0.36], $p = 0.006$). The interaction between experimental condition and agreeableness was significant ($\beta = 0.05$, 95% CI [-0.01, 0.10], $p = 0.085$), although the effect is only significant at the $\alpha = 0.1$ level. This finding suggests that more agreeable participants were more likely to adhere to randomized suggestions, whereas less agreeable participants tended to make larger changes.

These findings suggest that user-led personalization may be less suitable for highly agreeable users, who appear inclined to accept suggestions uncritically, even when informed that the suggestions were randomly generated. Previously, we found that user-led personalization increased engagement among participants with higher baseline depression severity, suggesting that autonomy can motivate engagement particularly in this subgroup. On the other hand, the analysis of how much participants used the user-led personalization feature adds an important qualification: highly agreeable users in the same condition often accepted the random module suggestions uncritically, which likely reduced the relevance of the selected therapy modules. Together, these findings highlight the nuances for choosing personalization approaches. They also point to personalizing the personalization approach itself, accounting not only for the symptom profile and baseline depression severity but also for personality traits such as agreeableness and stated preferences. We also investigated associations between the use of the user-led personalization feature and other participant characteristics (openness, conscientiousness, extraversion, neuroticism, gender, experience with psychotherapy, mental health apps or chatbots, knowledge of cognitive behavioral therapy, diagnosis of a mental disorder). However, none of the other participant characteristics turned out to be influential.

In summary, these findings emphasize the complexity of designing personalization beyond the three personalization approaches implemented in TheraBot and suggest the potential to adapt the personalization strategy to the individual user.

4.6 Discussion

Depression is a serious burden that affects millions of young people around the world. Since access to face-to-face psychotherapy is limited, DMHIs such as mental health chatbots have emerged as a solution. Although there is initial evidence that such chatbots can improve mental health, optimizing user engagement and treatment outcomes remains a major challenge. Personalization is a promising method to overcome this challenge. However, previous research on the effect of personalization in DMHIs has only investigated system-led personalization and produced mixed results. In some studies, personalization led to increased effectiveness; in others, there was no difference between personalized and standardized interventions. This inconsistency may be due to the limited scope of existing research, which has not fully explored the wide range of personalization approaches available. Importantly, previous studies have not

compared their relative effects. To address these gaps, we followed the DSR approach to design TheraBot, a mental health chatbot that incorporates hybrid content personalization by suggesting individual symptom-specific therapy modules that users can then adapt to their preferences. In two field experiments, we find that hybrid personalization consistently led to greater reductions in depression severity than standardized content delivery. Participants also prefer this approach, emphasizing the value of combining guidance through expert knowledge embedded in the chatbot with retaining control and autonomy. Beyond demonstrating the value of hybrid content personalization, our study provides a more differentiated understanding of how different personalization approaches perform. In a direct comparison across four versions of TheraBot (standardized content delivery, chatbot-led personalization, user-led personalization, and hybrid personalization), only the hybrid personalization version outperforms standardized content delivery. However, it does not outperform the other personalization approaches. Nevertheless, our results also suggest that, under very specific circumstances, other approaches may work better than hybrid personalization, such as for users with higher depression severity who demonstrate highest engagement in the user-led personalization condition. Notably, about 20% of participants express a preference for either chatbot-led or user-led personalization, indicating that individual preferences vary and that no single approach will be optimal for all users. This diversity in user preferences is also reflected in user engagement patterns. Participants with higher depression severity in the user-led condition tended to complete more modules relative to those with lower severity, whereas this trend was reversed in the chatbot-led condition. Our findings also point to specific challenges in user-led personalization. When we examine participants' actual use of the user-led personalization feature, we find substantial variation in adherence to module recommendations. Participants in the hybrid group adhere closely to chatbot recommendations, whereas adherence in the user-led group is considerably lower. We find that agreeableness influences the use of the user-led personalization feature in the user-led personalization group. Highly agreeable participants tend to follow recommendations closely, even though recommendations are randomly generated. This adherence to recommendations likely reduces the relevance of the selected therapy modules and may explain why the user-led group, despite similar engagement compared to the hybrid group, does not outperform standardized content delivery in terms of reductions in depression severity.

4.6.1 Implications for Research

Our research contributes to the literature on personalization in the fields of IS and digital mental health as follows. First, our research helps reconcile inconsistencies with respect to the effects of content personalization in DMHIs. Previous work has not examined content personalization in mental health chatbots, and research on DMHIs has yielded inconclusive results, with a recent systematic review showing that most personalization approaches failed to show a benefit (Schaeuffele et al., 2025). We conducted an in-depth investigation of personalization approaches under real-world conditions and found that hybrid content personalization improves the effectiveness of mental health chatbots. In doing so, we respond to calls to investigate personalization approaches in digital health interventions (Kankanhalli et al., 2021) and chatbots (Kocaballi et al., 2019), providing, to our knowledge, the first empirical evidence on the effects of content personalization in mental health chatbots. Our research reveals that pre-existing depression severity interacts differently with personalization approaches, as users with higher depression severity showed greater engagement with user-led personalization compared to chatbot-led or hybrid approaches. This finding extends the paradox identified by Borghouts et al. (2021), who showed that individuals with more severe depression severity express a greater willingness to use DMHIs but show lower actual engagement. Our results suggest that user-led personalization may bridge this gap, enabling users with higher symptom severity to turn their initial willingness into sustained engagement. In addition, we found that one fifth of participants preferred user-led or chatbot-led rather than hybrid personalization. These findings challenge the assumption that a single personalization strategy is optimal for all users. Instead, they suggest that the suitability of personalization approaches may depend on individual user characteristics, such as depression severity, and explicit preferences. Overall, our findings contribute to a more nuanced understanding of content personalization in mental health chatbots by challenging the idea of a one-size-fits-all personalization strategy.

Second, our research adds to the literature by providing a fuller and more nuanced picture of the effects of different personalization approaches in DMHIs. Although a variety of different personalization approaches are used in commercial mental health chatbots, little is known about their respective effect on mental health outcomes (Kocaballi et al., 2019; Martinengo et al., 2022). While previous studies have focused only on system-led personalization (Berger et al., 2014; Chaturvedi et al., 2023; Johansson et al., 2012), our research provides an examination

of multiple personalization approaches that involve the chatbot, user, or both. In particular, our finding that hybrid personalization showed the most promising results is a key addition to the literature, which could explain why system-led content personalization approaches did not yield improved mental health outcomes compared to standardized content in some studies. Overall, our findings highlight the benefits of a hybrid personalization strategy that provides informed recommendations, but also accommodates the diverse needs and preferences of users. Whereas chatbot-led personalization provides guidance for users, particularly those uncertain or overwhelmed when selecting therapy modules, user-led personalization empowers those with specific preferences (e.g., because they know what is relevant to them or prefer to start with a different therapy module). Hybrid personalization not only combines the advantages of both strategies but also aligns better with findings from face-to-face psychotherapy, which has shown that incorporating patient preferences into psychotherapists' expert judgments improves mental health outcomes (Swift et al., 2018). Consequently, hybrid content personalization should be chosen if only a single personalization strategy can be implemented and adapting the personalization strategy to the individual user is not feasible.

Third, our study validates a set of design requirements that can guide the development of future mental health chatbots. Previous research has provided design knowledge for chatbots in general (Diederich et al., 2022), but specific guidance for the mental health context has been limited (Ahmad et al., 2022; Paul et al., 2024). By demonstrating the effectiveness of hybrid content personalization and deriving DR6, we advance this design knowledge. Specifically, we show that simply providing a personalization feature is not enough; the effectiveness depends on how the personalization is implemented (system-led vs user-led vs hybrid) and for whom it is designed (e.g., severe users vs agreeable users). This nuanced understanding moves the field away from generic endorsements of personalization towards more targeted, evidence-based design decisions.

In addition, our study offers theory-grounded design knowledge for mental health chatbots. Although a substantial body of design knowledge exists for chatbots in general (Diederich et al., 2022), relatively little of it addresses the unique challenges and requirements of the mental health context (e.g., Ahmad et al. (2022), Paul et al. (2024)). Through our artifact TheraBot, our five design requirements derived from self-determination theory, the effort-accuracy framework, and cognitive load theory, and our evaluation results, we extend this line of work by providing prescriptive guidance on how to implement hybrid content personalization in mental

health chatbots. While our evaluation generally supported the validity of the proposed design requirements, it also revealed that a single personalization approach may not be suitable for all users and highlighted the need to "personalize the personalization". To reflect this insight in our design, we propose an additional design requirement: The mental health chatbot should personalize its content personalization approach (chatbot-led, user-led, or hybrid) based on user characteristics (e.g., personality traits), symptom presentation (e.g., depression severity), and stated preferences (DR6). Table 4.8 provides an overview of the final design requirements.

Table 4.8: Final Design Requirements

#	Design Requirement
DR1	The mental health chatbot should establish a well-defined personalized treatment plan consisting of the most relevant therapy modules.
DR2	The personalized treatment plan should balance consistency with flexibility and enable the chatbot to respond to users' changing needs during treatment.
DR3	The mental health chatbot should allow users to choose the content following a user-led personalization approach.
DR4	The mental health chatbot should choose the content following a system-led personalization approach.
DR5	The mental health chatbot should implement hybrid content personalization that balances the strengths and weaknesses of chatbot-led and user-led personalization.
DR6	The mental health chatbot should personalize its content personalization approach (chatbot-led, user-led, or hybrid) based on user characteristics (e.g., personality traits), symptom presentation (e.g., depression severity), and stated preferences.

Overall, our empirical insights contribute to a better understanding of how different personalization approaches work better than others in the context of mental health chatbots, but also highlight the need for more nuanced approaches to personalization that consider individual user characteristics and preferences.

4.6.2 Implications for Practice

Our research also has implications for mental health chatbot developers. Specifically, it enables developers to move beyond a standardized design, implement personalization approaches appreciated by users, and improve mental health outcomes. Our design requirements and the in-depth description of our artifact TheraBot provide actionable design prescriptions for implementing personalization approaches that involve both the user and chatbot. Building on these

foundations, DR6 offers the most consequential guidance: the chatbot should adapt the personalization approach itself in response to user characteristics, symptom presentation and stated preferences. These characteristics can be captured during onboarding with brief questionnaires and refined through chatbot use, for example, module abandonment. Developers who need or prefer to implement a single approach should start with hybrid personalization, which our experiments identified as the most promising. Because our evaluation results suggest that improved mental health outcomes can already be achieved with relatively simple implementations of personalization approaches, our ideas can serve as a good starting point before deploying complex instantiations using advanced recommender systems, natural language processing, or integrating additional data sources.

4.6.3 Limitations and Future Research

Our research comes with limitations that provide opportunities for future research. First, although we offer valuable empirical insights into the effectiveness of mental health chatbots in reducing depression severity, our implementation and empirical studies have focused on a rule-based chatbot to avoid harmful responses and limit risks. The advent of LLMs has significantly raised user expectations of chatbots, as indicated by feedback from participants who requested more open and personal conversations that truly mirror face-to-face psychotherapy sessions rather than dialogues with rule-based chatbots. However, while current LLMs demonstrate remarkable conversational capabilities and promising results in the mental health context (Hatch et al., 2025; Heinz et al., 2025), they can produce inappropriate responses with profoundly dangerous information regarding health questions (Gilbert et al., 2023), which could have serious consequences in the mental health context. Nonetheless, it is likely that, with improvements in LLMs and adequate training, LLM-based mental health chatbots can be employed in the future. Therefore, future research should explore how LLMs can be safely and effectively integrated into mental health chatbots. Second, our implementation of personalization approaches focused on choosing therapy modules relevant to users' depressive symptoms. Another promising approach is to personalize the content of the therapy modules itself, e.g. by providing relevant examples or personalized feedback to therapy exercises (Cohen et al., 2021). While challenging to implement in rule-based chatbots, we encourage future research on LLM-based chatbots to extend our design with features that allow personalization of the content of the therapy modules. Third, our implementation of the chatbot-led personalization feature re-

lied on items from well-established standardized mental health questionnaires to identify the most relevant modules. While these were presented to users in the chat in a structured dialogue, LLMs may also allow the use of open questions rather than standardized mental health questionnaires. Therefore, a promising avenue for future research is to explore how chatbot-led personalization can be reliably implemented in more natural conversations. In addition, more advanced personalization approaches can be developed if other types of user data are available, for example, dynamically modifying the treatment plan in real time based on users' ongoing conversations with the chatbot. Fourth, the two field experiments were limited to 14 days. Although this study duration is similar to other studies on mental health chatbots (Fitzpatrick et al., 2017), it represents a rather short study duration, which did not allow us to observe the effects of using a mental health chatbot long-term. Given our significant findings within this relatively brief study duration and the established relationship between longer usage and larger effects in DMHIs (Kambeitz-Ilankovic et al., 2022), longer chatbot usage may yield stronger effects. Future research should, therefore, conduct comprehensive long-term evaluations to further validate the potential of personalization in mental health chatbots. Finally, although a substantial proportion of our two samples reported a previous diagnosis of a mental disorder, and all participants experienced at least mild depression severity according to the PHQ-8, the diagnosis of a depressive disorder was not used as an eligibility criterion. Therefore, evaluating the effect of content personalization in mental health chatbots with a sample of patients with a clinical diagnosis would be an important avenue for future research, even more so because our results suggest larger improvements for participants with higher depression severity. In conclusion, our research provides a nuanced understanding of the effect of content personalization in mental health chatbots on user engagement and depression severity. We revealed the effects of different personalization approaches, uncovered why hybrid personalization outperforms other strategies, and observed that the personalization strategy should be adapted to individual users rather than implemented as a one-size-fits-all solution. We hope that our findings offer valuable insights for both researchers and practitioners, inspire further work in this important area of research, and ultimately improve mental healthcare for young people with depression.

5 Part III: Designing an LLM-Based Behavioral Activation Chatbot for Young People with Depression: Insights from Artificial Users and Clinical Experts

5.1 Introduction

Depression is a prevalent and severe mental disorder that affects over 280 million people globally (GBD 2019 Mental Disorders Collaborators, 2022). Access to evidence-based psychotherapy remains severely limited by workforce shortages, limited access in rural areas, and persistent stigma (Butryn et al., 2017; Schomerus et al., 2022). Digital mental health interventions, particularly chatbots, have emerged as promising avenues to address this treatment gap by providing scalable, accessible, and anonymous support. Rule-based chatbots, such as Woebot (Fitzpatrick et al., 2017) and Wysa (Inkster et al., 2018), have demonstrated effectiveness in reducing depressive symptoms. However, these systems have fundamental limitations: their pre-written messages and predetermined response paths create rigid, repetitive interactions that fail to address individual user needs, leading to limited engagement and symptom improvement (Chan et al., 2022; Kocaballi et al., 2019; Lim et al., 2022).

Large Language Models (LLMs) promise to overcome these limitations through advanced natural language capabilities. Yet their probabilistic nature raises concerns about inconsistent or potentially harmful responses (Heston, 2023), and a fundamental evaluation gap leaves their therapeutic quality unverified. Unlike psychotherapists, who undergo rigorous fidelity assessments using validated clinical instruments during training and in clinical studies (Dimidjian et al., 2017; Meeks et al., 2019), LLM-based chatbots have not been subjected to equivalent standards. Although LLMs can generate contextually appropriate single-turn interactions comparable to those of psychotherapists (Hatch et al., 2025), effective psychological interventions require adherence to evidence-based protocols across longitudinal multiphase treatments, a capability that remains largely unexplored. Current evaluations rely on ad-hoc metrics (Y. Hua et al., 2025; Thieme et al., 2020) or focus exclusively on downstream mental health outcomes (Heinz et al., 2025), neglecting the process measures applied to psychotherapists themselves.

This “upside-down” approach assesses usability and user experience before establishing foundational clinical validity (Y. Hua et al., 2025).

Addressing this gap requires the application of validated fidelity instruments to LLM-based chatbots before deployment. Therefore, our research aims to answer the following questions:

1. To what extent does an LLM-based chatbot maintain clinical fidelity when delivering structured psychological interventions?
2. What limitations and opportunities for refinement can be identified through clinical expert evaluation?

We investigated these questions by developing an LLM-based chatbot (powered by GPT-4o) that delivers behavioral activation for depression, building on Cady (Kuhlmeier et al., 2022a, 2025a), a rule-based mobile mental health chatbot for young people with depressive symptoms, as follows. Behavioral activation is particularly well-suited for this investigation: its structured protocol with distinct phases and clear completion criteria (Lejuez et al., 2001, 2011) enables a systematic assessment using validated fidelity instruments. To evaluate the chatbot across diverse clinical presentations, we employed artificial users (powered by GPT-4o) to generate 48 complete sessions, which were subsequently assessed by clinical experts using the Quality of Behavioral Activation Scale (Q-BAS).

Our evaluation revealed an asymmetry in the chatbot’s capabilities: while it demonstrated robust protocol adherence, successfully completed all intervention phases, and maintained safety, significant gaps emerged in clinical reasoning. Specifically, the chatbot struggled to verify whether selected activities and rewards were feasible and therapeutically appropriate, and it had difficulty developing personalized strategies to overcome user-specific barriers. Based on these findings, we offer practical design recommendations spanning both prompt-level refinements and broader architectural considerations.

This study makes three contributions to the literature. (1) *Empirical*: Expert-assessed sessions revealed that the tested GPT-4o-based chatbot reliably executed behavioral activation protocols but showed systematic weaknesses in clinical reasoning. (2) *Methodological*: a demonstration of how standardized fidelity instruments combined with artificial user session generation can enable rigorous pre-deployment evaluation of LLM-based mental health interventions; and (3) *Design*: Actionable prompt engineering patterns, including template-based content, embedded decision rules, granular task breakdown, and explicit redirection mechanisms, alongside

complete prompts for replication and further development.

5.2 Background

5.2.1 Mental Health Chatbots

Chatbots are an increasingly popular type of digital mental health intervention (DMHI) (Torous et al., 2021), with well-known commercial applications such as Woebot (Fitzpatrick et al., 2017) and Wysa (Inkster et al., 2018). They are well-accepted by users (Vaidyam et al., 2019), can establish a therapeutic relationship (Darcy et al., 2021; Skjuve et al., 2021, 2022), and have shown to improve mental health (Lim et al., 2022). However, the predominant reliance on rule-based and retrieval-based architectures constrains these systems to predetermined conversational pathways and developer-crafted messages, resulting in interactions that users characterize as repetitive and robotic (Fitzpatrick et al., 2017). LLMs promise significant improvements owing to their capacity to generate context-dependent responses, unlocking higher personalization, and potentially improving user engagement and efficacy (Kocaballi et al., 2019). Nevertheless, the tension between conversational flexibility and safety necessitates rigorous empirical investigation (Stade et al., 2024).

Researchers have explored several approaches to developing LLM-based mental health chatbots. Prompt engineering remains the most widely adopted approach, where carefully crafted instructions guide pre-trained LLMs to generate appropriate therapeutic responses (Stade et al., 2024). Kumar et al. (2022, 2023) demonstrated this approach by systematically evaluating GPT-3's therapeutic potential through controlled experiments varying prompt architectures and intervention modalities. Participants rated the system highly on perceived expertise and expressed willingness for continued engagement; however, trust scores remained moderate, a critical finding given the role of trust in therapeutic efficacy (Flückiger et al., 2018; Krupnik, 2022). Their subsequent investigation of mindfulness education through prompt-engineered interventions demonstrated increased post-intervention practice intentions, although these studies did not include clinically diagnosed participants or expert evaluations. Beredo and Ong (2022) similarly achieved favorable expert ratings for relevance, human-likeness, and empathic responding, though their system focused on generating empathetic responses rather than implementing structured psychological interventions. Moreover, Das Swain et al. (2025) developed a

cognitive reframing chatbot that was proven to be effective in supporting emotional well-being. Hybrid approaches that combine LLMs with rule-based or retrieval-based methods have also shown promise in this regard. Liu et al. (2024) conducted a comparative analysis finding that LLM-based chatbots could provide more natural and empathetic feedback, although controlling these systems for comprehensive multi-round psychological interventions remains challenging, leading them to propose a hybrid solution.

Beyond prompt engineering, developers have increasingly explored fine-tuning, in which models are trained on domain-specific conversations to better align with clinical practices (Stade et al., 2024). Heinz et al. (2025) developed Therabot, an LLM-based chatbot fine-tuned on expert-curated mental health conversations created by clinical psychologists. Their architecture incorporates multiple guardrails, including a crisis classification agent, and a randomized controlled trial demonstrated improved mental health outcomes. Campellone et al. (2025) similarly combined fine-tuning with guardrails for crisis and off-topic messages, reporting positive user engagement and no significant safety concerns.

Despite these promising findings, none of these studies evaluated intervention delivery using validated fidelity instruments. User satisfaction and symptom improvement alone cannot determine whether chatbots deliver evidence-based interventions of appropriate therapeutic quality. Additional concerns include inadequate crisis protocols (Heston, 2023) and limited conversational memory, which affect therapeutic consistency (Ma et al., 2023). Given these shortcomings, a comprehensive systematic review concluded that deployment risks currently exceed the demonstrated benefits (Guo et al., 2024; Y. Hua et al., 2025), confirming the need to rigorously evaluate early stage LLM-based chatbots (Stade et al., 2024).

5.2.2 Behavioral Activation

Behavioral activation is a psychotherapeutic approach that addresses the reciprocal relationship between engagement in activities and mood (Lejuez et al., 2011). It posits that a loss of positive reinforcers, often stemming from withdrawal from pleasurable activities, contributes to the onset and maintenance of depression, thereby establishing a self-sustaining cycle of decreased activation and a worsening mood (Lejuez et al., 2011). Behavioral activation seeks to disrupt this maladaptive cycle through systematic activity scheduling, with the aim of encouraging patients to participate in enjoyable activities (Lejuez et al., 2011). This approach is effective in

treating depression and can be implemented in a single session, several sessions, or as a long-term intervention (Cuijpers et al., 2023b). The structured nature of behavioral activation, which includes distinct phases such as mood assessment, psychoeducation, and activity scheduling, makes it well-suited for delivery as a digital intervention (Huguet et al., 2016). Behavioral activation as a digital intervention has been shown to effectively reduce depressive symptoms (Jia et al., 2025), even in very brief formats such as single-session implementations (Schleider et al., 2022). It has also been implemented as a chatbot intervention (Mancinelli et al., 2024; Rathnayaka et al., 2022). The structured protocol of behavioral activation creates an ideal testbed for evaluating LLM capabilities in delivering multiphase psychological interventions.

5.2.3 Evaluating LLM-Based Mental Health Chatbots

The transition from deterministic rule-based to probabilistic LLM-based architectures fundamentally alters the evaluation requirements of mental health chatbots. Rule-based systems (Fitzpatrick et al., 2017; Inkster et al., 2018) generate responses from finite sets of pre-specified responses and conversational paths and can be validated deterministically by testing all possible paths (Laranjo et al., 2018). In contrast, LLMs produce probabilistic outputs across an effectively infinite response space, making rigorous evaluation essential to ensure therapeutic quality and safety before deployment (Stade et al., 2024).

Current evaluation practices are insufficient to address this challenge. Recent work examining evaluation practices in LLM-based mental health applications reveals that studies prioritize higher-level constructs, such as accessibility and user experience, while foundational safety and validity remain under-examined (Y. Hua et al., 2025). Twelve of the 16 studies reviewed developed ad-hoc evaluation scales rather than employing validated clinical instruments, and none used standardized fidelity measures comparable to those applied in psychotherapy training and research (Y. Hua et al., 2025). This “upside-down” evaluation approach—assessing usability before establishing clinical validity—represents a fundamental methodological gap.

In addition to this research gap, development practices introduce additional limitations. Developers typically assume dual roles, role-playing as users to test conversational pathways while evaluating their own systems (Bunge & Desage, 2025). This introduces assessment bias and fails to represent user heterogeneity, as developers may not anticipate problematic interaction patterns in the real world. Furthermore, independent expert evaluations using validated fidelity

instruments remain rare (Beredo & Ong, 2022); to our knowledge, no study has evaluated how well LLM-based mental health chatbots implement evidence-based treatments using fidelity assessments comparable to those applied to psychotherapists (Bunge & Desage, 2025).

Artificial users offer a solution to the first limitation by enabling extensive evaluation without exposing vulnerable populations to the untested systems. They serve as a critical intermediate step between developer testing and human trials, identifying shortcomings prior to deployment. LLM-powered artificial users can generate coherent, context-aware responses while maintaining consistent personas and emotional states across multi-turn therapeutic conversations (Louie et al., 2024; H. Qiu & Lan, 2024; Schuller et al., 2024; Steenstra et al., 2025; J. Wang et al., 2024). However, questions remain regarding authenticity and generalizability: current implementations may not adequately represent the diverse user characteristics that influence interactions with mental health chatbots (Borghouts et al., 2021; Kapania et al., 2025).

Our study addresses both limitations through a comprehensive assessment of an LLM-based behavioral activation chatbot. We developed diverse artificial users to test safety and protocol adherence across varied clinical presentations, and independent clinical experts employed a standardized fidelity instrument (Q-BAS) to evaluate treatment delivery quality. This approach establishes foundational safety and therapeutic quality before higher-level constructs while identifying specific strengths, shortcomings, and opportunities for improvement.

5.3 Designing an LLM-based Behavioral Activation Chatbot

Cady (Kuhlmeier et al., 2022a, 2025a) is a rule-based mobile conversational agent designed for young people (aged 15–29, as defined by European Union (2023)) experiencing depressive symptoms. Young people represent a particularly vulnerable population, as most psychiatric disorders manifest before the age of 25 (Solmi et al., 2022), and suicide is the third leading cause of death among 15–29 year-olds globally (World Health Organization, 2025). The original version offers five modules: behavioral activation, cognitive restructuring, interpersonal skills, emotion regulation, and sleep management, all developed from evidence-based treatment manuals (Abel & Hautzinger, 2013; Groen & Petermann, 2015; Towery, 2016) in collaboration with psychotherapists. It operates through predetermined dialogue paths that limit flexible and personalized responses (Kuhlmeier et al., 2025a). Building on Cady’s behavioral activation module, we developed an LLM-based implementation designed to maintain the therapeutic

structure while enabling more natural conversations.

5.3.1 Intervention Development

We selected behavioral activation for its structured protocol and suitability for single-session delivery. The decision to condense the original three-session format into a single session was supported by evidence that single-session interventions produce significant improvements in depressive symptoms (Kaveladze et al., 2025; Schleider et al., 2022), making the intervention more accessible and easier for experts to evaluate the results. The development process involved close collaboration with psychotherapists who translated the rule-based script into a structured prompt format and crafted example sessions. The system prompt underwent iterative refinement, with development team members role-playing as users while psychotherapists evaluated adherence to instructions, phase completion, and response quality. Initial tests revealed that increased flexibility diminished protocol adherence and resulted in inconsistent, low-quality responses, leading us to prioritize a linear progression through all phases. The refinement continued until the chatbot consistently completed all seven behavioral activation phases.

5.3.2 Prompt Architecture

The chatbot operates through a structured prompting framework comprising five hierarchical components, as summarized in Table 5.1. To address potential confusion from long prompts (Bhattacharjee et al., 2024), we structured the instructions hierarchically with distinct phases and clear completion criteria, reinforced by extensive in-context learning.

Table 5.1: Prompt Architecture Components and Functions

Component	Purpose	Content Summary
Format Instructions	Phase transition control and session structure	Explicit markers that track session progress (e.g., [Phase1]) and enforce sequential phase completion
Identity	Role and persona definition	Chatbot for young people experiencing depressive symptoms, featuring an empathetic, activating, encouraging, humorous, and curious personality
Constraints	Communication guidelines and safety protocols	Includes a 30-word message limit, suicide/emergency protocol with crisis referral, and role boundaries that politely decline off-topic requests
Task	Overall therapeutic goal	Instructs the chatbot to guide the user through a 7-phase behavioral activation session with the key objective of collaboratively creating a personal activity plan
Phase-Specific Instructions	Detailed phase procedures	Seven phases (introduction, psychoeducation, finding activities, planning activities, problem solving, positive reinforcements, closing), each with specific goals, completion criteria, and good/bad example sessions
Complete Session Example	Comprehensive session model	A full multi-turn session demonstrating all seven phases with natural pacing and smooth transitions

Several design decisions emerged from the iterative testing. Without explicit phase transition markers, the chatbot occasionally skipped or superficially completed the phases. Repeating the format instructions at both the beginning and end of the prompt addresses the instruction drift observed in the pilot testing, where the LLM forgot the early constraints during extended conversations. The 30-word message constraint emerged from the observation that longer responses reduced conversational naturalness. The contrastive in-context learning approach, which pairs good and bad examples for each phase, helped the model avoid common errors while maintaining therapeutic quality.

The prompt was designed in German for the target population of this study. The full prompt is provided in the Supplementary Material. The chatbot was powered by GPT-4o (gpt-4o-2024-08-06) via OpenAI's API, with the temperature set to 1 to balance consistency with response

variety, as lower values produced overly repetitive replies.

5.4 Evaluation Study

We evaluated the capability of our LLM-based chatbot in delivering behavioral activation to young people with depressive symptoms. We generated therapeutic sessions using artificial users to systematically test the chatbot across various clinical presentations, which were evaluated by psychotherapists using standardized fidelity instruments. This approach allowed us to identify capabilities and shortcomings while protecting vulnerable individuals from potential harm during the intermediate validation step. The Institutional Review Board of our institution approved this study protocol.

5.4.1 Generating Behavioral Activation Sessions with Artificial Users

We developed artificial users designed to reflect the variability of human users with depressive symptoms. To ensure clinical validity, we collaborated closely with a psychotherapist throughout the developmental process. We developed artificial users based on four patient vignettes, which are concise descriptions of patients derived from real-world clinical cases commonly used in psychotherapy research and training (Franco D'Souza et al., 2023). We selected patient vignettes from psychotherapy training materials at the outpatient clinic of our institution, which also serves as an institute for training licensed psychotherapists. These vignettes, each 300–500 words long, describe young people aged 14–29 years diagnosed with depression, including detailed symptom presentations and psychosocial circumstances. To capture the heterogeneity among users, we enriched the base vignettes with six systematically varied characteristics identified through a literature review and discussions with clinical psychologists. Each characteristic was operationalized with two or three categories to maintain a manageable number of combinations. The final evaluation set comprised 48 artificial user variations. A full overview of the characteristics, their variations, and the rationale for selecting them are presented in Table 5.2. An example of an artificial user is presented in Table 5.3, and the base vignettes and variation text expressions are provided in the Supplementary Materials.

Table 5.2: Artificial User Characteristics and Rationale for Selection

Characteristic	Variations	Rationale for Selection
Depression Severity	Mild, Moderate, Severe	Higher severity increases interest and willingness to adopt DMHIs, but hampers actual engagement due to depressive symptoms, low mood, and fatigue that inhibit motivation and ability to use interventions (Borghouts et al., 2021)
Gender	Male, Female, Non-binary	Women are more likely to engage with DMHIs than men (Borghouts et al., 2021)
Willingness to Disclose Personal Information	High, Low	Privacy concerns and confidentiality fears create barriers to engagement and information disclosure (Borghouts et al., 2021; Jardine et al., 2024)
Openness to Chatbot Suggestions	High, Low	Preexisting beliefs about digital intervention effectiveness affect engagement (Borghouts et al., 2021)
Dominance	High, Low	Dominance affects conversations between users and chatbots (Gnewuch et al., 2020)
Attitudes toward Mental Health Chatbots	Positive, Negative	Negative attitudes and the “humans need humans” preference for in-person therapy create barriers to digital uptake (Borghouts et al., 2021; Jardine et al., 2024)

Table 5.3: Example of an Artificial User Profile

Characteristic	Level	Description
Gender	female	I'm Kira, 29 years old and I'm just hanging around in my flat. I have lost my job as a paralegal, and now everything is totally screwed up. I constantly feel as if I am in a black hole. A relationship? Not a chance. My friends are getting married and having children, but I feel completely disconnected and isolated. In addition, my mother has now developed Alzheimer's disease. That completely knocks me out. My sleep rhythm no longer exists. I lie awake for hours and cannot fall asleep. If I 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 is no point anyway. Food? Forget it; I have no appetite at all. Dating has not been a thing for a long time. I just stay at home and do not feel like doing anything. I am permanently down and cannot concentrate on anything. I often ask myself what the point of all this is. At home, I constantly brooded about my job loss and felt like a failure. Everything seems pointless to me. I lie awake at night, worrying that I will go completely broke. I have driven all my friends away. I feel totally worthless and have extreme feelings of guilt regarding everything. Sometimes, I can hardly move, and even showering is torture. I constantly think about what it would be like if I were no longer there. Sometimes I really think about whether I should just end it.
Age group	young	
Depression severity	adult severe	
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 am 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 instead.

Artificial users were implemented using GPT-4o (gpt-4o-2024-08-06, temperature = 1). This temperature setting was selected during pilot testing, as it achieved an optimal balance be-

tween adherence to the artificial user profile and response variety. To verify the alignment with depression severity levels, artificial users completed the PHQ-9 questionnaire, and only those scoring within the intended range were included in the study. For the intended range, we adapted the traditional PHQ-9 categorization to align with our three-level classification of depression severity: mild (5–9), moderate (10–19), and severe (20–27). From the verified pool of artificial users, we drew a stratified random sample of 48 artificial users who interacted with the chatbot. The sample size was determined by the number of available psychotherapists ($n = 10$) and their available evaluation time (1–2 hours per psychotherapist), which allowed each psychotherapist to assess 3–6 sessions. We employed stratified randomization to maintain a roughly balanced representation across the seven user characteristics, although a perfect balance was not achieved, given the constraints of evaluating 48 sessions across multiple characteristics.

Each session between the chatbot and an artificial user was designed to conclude when all seven phases of behavioral activation had been completed, when the chatbot sent a prespecified marker, or when the conversation reached a limit of 100 turns. The sessions began with a standardized welcome message from the chatbot. Conversations between the chatbot and artificial users were implemented using a Python script.

5.4.2 Clinical Expert Assessment

Participants

The participants were recruited from the staff of our institution's outpatient clinic and the professional network of the research team between August and September 2024. The inclusion criteria were as follows: (1) completion of a master's degree in psychology, (2) certification as a psychotherapist or enrollment in psychotherapist training (minimum second year), and (3) experience treating young people with depression using behavioral activation therapy. Overall, ten participants were recruited. The participant characteristics are presented in Table 5.4. Each participant received €30 compensation for 1–2 hours of participation.

Table 5.4: Participant Characteristics

Characteristic	Values
Professional Status	
Licensed Psychotherapists	2 (20%)
Psychotherapy Trainees	8 (80%)
Gender	
Female	7 (70%)
Male	3 (30%)
Age (years)	M=30.1 (SD=4.12)
Clinical Experience (years)	M=3.75 (SD=1.75)
Experience with DMHIs	7 (70%)
Willingness to recommend DMHIs	7 (70%)

Study Procedure

The evaluation consisted of three phases. First, the participants received a complete overview of the study objectives, tasks, and procedures and provided their informed consent. Second, each participant independently assessed 3–6 complete sessions (determined by the available time) through an online platform. The participants were instructed to read carefully through each session before starting the questionnaire. The session remained accessible while the questionnaire was being completed. To prevent fatigue, the participants were allowed to pause within and between the sessions. Each session was evaluated by a single participant. Third, we conducted semi-structured interviews to explore participants' perspectives on the sessions. The

entire study lasted between one and two hours per participant. We disclosed that sessions were generated between the chatbot and artificial users only during the semi-structured interviews, which allowed for unbiased assessment.

Measures

Evaluating the Behavioral Activation Chatbot

We used the Quality of Behavioral Activation Scale (Q-BAS) to assess treatment fidelity (Dimidjian et al., 2012), adapted for treatment via chatbot. The Q-BAS is a validated observer-rated measure designed to measure the quality with which a psychotherapist delivers behavioral activation components. It has been used to measure the quality of behavioral activation in traditional (Dimidjian et al., 2017) and virtual human-delivered psychotherapy (Rethorst et al., 2024). We administered 14 items covering the components of the behavioral activation sessions. Each item was rated on a 7-point scale (0 = very poorly, 6 = very well), with higher scores indicating better quality. A score of ≥ 3 per item indicates satisfactory delivery of the respective component (Dimidjian et al., 2012; Dimidjian et al., 2017; Rethorst et al., 2024). In addition to the 14-item Q-BAS, we collected a separate single-item holistic rating of the overall quality of the behavioral activation session. Psychotherapists answered the question (translated from German): “Overall, how would you rate the chatbot as a behavioral activation chatbot in this session?” on a 7-point scale (1 = very poor; 7 = very good). In addition, we administered seven items to assess broader therapeutic competencies based on the Thera-Turing test (Bunge & Desage, 2025): the chatbot (1) validated emotions and demonstrated empathy, (2) responded to the user’s concerns, (3) established a therapeutic relationship, (4) maintained objectivity and avoided judgment, (5) wrote clear, precise, and easy-to-understand messages, (6) facilitated a natural conversation flow, and (7) ensured message safety and avoided harmful content. Finally, we used open-ended questions to identify the strengths, shortcomings, and opportunities for improvement in each phase and overall.

Artificial User Ratings

To evaluate the usefulness of artificial users as an approach to generate and evaluate sessions with a chatbot, psychotherapists rated two aspects of the artificial users. They assessed the perceived authenticity of the artificial user (1 = very unrealistic to 7 = very realistic) compared

to real patients, and the difficulty in conducting the session with the artificial user (1 = very difficult to 7 = very easy), each with a single item adapted from the Q-BAS (Dimidjian et al., 2012).

Data Analysis

Quantitative data were analyzed using R (version 4.3.1). For intervention fidelity (Q-BAS), we analyzed the ratings from two perspectives: (1) component-wise, aggregating ratings across all sessions for each of the 14 behavioral activation components, and (2) session-wise, examining the pattern of component ratings within each session to identify distinct performance profiles. We calculated descriptive statistics, including rates of phase completion across sessions, and evaluated Q-BAS item ratings against the established satisfactory threshold of ≥ 3 (Dimidjian et al., 2012; Dimidjian et al., 2017; Rethorst et al., 2024). For session-level summaries, we computed a session-level Q-BAS mean as the arithmetic mean of the 14 Q-BAS item scores within each session (0–6 scale), which is distinct from the single-item holistic session rating and therapeutic capability ratings. To examine session-level performance patterns, we visualized the fidelity ratings across all 48 sessions using a heatmap (see Appendix Figures C.1). To quantify the sources of variability in the chatbot’s performance, we conducted a variance decomposition using a linear mixed-effects model (R package *lme4*) with restricted maximum likelihood (REML) estimation. The Q-BAS scores were modeled with crossed random intercepts for the sessions and intervention components. We partitioned the total variance in the Q-BAS scores into three components: variance attributable to differences between sessions, variance attributable to differences between the intervention components, and residual variance. We also conducted exploratory analyses to investigate whether artificial users’ characteristics, including depression severity, affected Q-BAS ratings using the Wilcoxon rank-sum test for binary variables and the Kruskal-Wallis test for variables with three or more groups. For the therapeutic capabilities, we analyzed the ratings per capability and calculated descriptive statistics for each. We also explored whether artificial user characteristics influenced authenticity and interaction difficulty ratings using the same statistical tests. Qualitative data from the questionnaires and interviews were analyzed using a mixed deductive-inductive approach (Mayring, 2004; Mayring & Fenzl, 2019). We used the seven phases of behavioral activation and their 14 components as pre-specified deductive categories, while also coding inductively within and beyond these categories. This approach served two objectives: obtaining specific

feedback for each phase and component, and identifying patterns that emerged across phases. We calculated the frequency of the suggestions to identify the most common recommendations.

5.4.3 Results

Quality of Behavioral Activation

Figure 5.1 provides an overview of psychotherapists' ratings of behavioral activation quality, including (a) Q-BAS component ratings and (b) the holistic single-item rating of overall session quality.

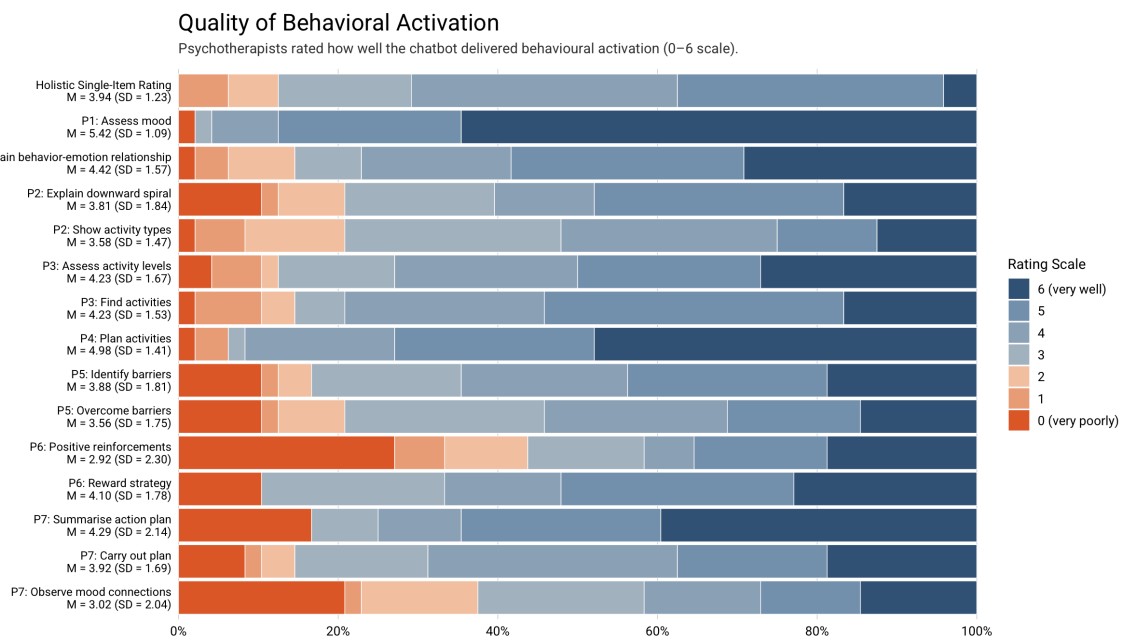


Figure 5.1: Quality of Behavioral Activation

The chatbot received positive evaluations with a mean rating of $M = 3.94$ ($SD = 1.23$) on the holistic single-item rating of the overall session quality. The average Q-BAS rating across the 14 components was $M = 4.03$ ($SD = 1.18$). Thirteen of the 14 behavioral activation components exceeded the satisfactory threshold of ≥ 3 on average, with one component (reinforcement explanation) slightly below this threshold. However, the component delivery varied considerably. Mood assessment ($M = 5.42$, $SD = 1.09$) and planning activities ($M = 4.98$, $SD = 1.41$) received the highest ratings, whereas explaining reinforcement ($M = 2.92$, $SD = 2.30$) and encourage to observe activity-mood connections ($M = 3.02$, $SD = 2.04$) received the lowest ratings. The standard deviations also varied considerably across components

(from 0.77 to 2.30). Explaining reinforcement, reviewing the activity plan, and encourage to observe activity-mood connections showed the highest variability. We conducted a variance decomposition analysis to understand the sources of variability in therapeutic quality. The results indicated that 36.8% of the variance in Q-BAS scores was attributable to differences between sessions (variance = 1.27, 95% CI [0.82, 2.02]), whereas only 12.0% was attributable to differences between intervention components (variance = 0.42, 95% CI [0.18, 0.98]). The remaining 51.2% was residual variance. In terms of clinical reliability, success rates varied widely. While mood assessment met the satisfactory threshold (≥ 3) in 97.9% ($n=47/48$) of sessions, the chatbot successfully explained positive reinforcement in 56.2% ($n=27/48$) of sessions. Table C.1 reports component-level adequacy rates (Q-BAS score ≥ 3) across the 48 evaluated sessions. Seven psychotherapists expressed positive impressions in the interviews, particularly highlighting the well-structured session flow ($n = 7$) and the chatbot's ability to validate feelings ($n = 3$). One psychotherapist noted: *I would not have thought it was possible that a chatbot could do this so well and so authentically.* However, seven psychotherapists found some sessions superficial and observed that the chatbot conducted sessions more quickly and with less depth than typical sessions. One psychotherapist argued that a shorter session was more appropriate for a chat-based intervention than for a traditional therapy session.

Phase 1: Introduction

Mood assessment ($M = 5.42$, $SD = 1.09$), and explaining the behavior-emotion connection ($M = 4.42$, $SD = 1.57$) exceeded the satisfactory threshold and ranked high among the 14 components. The most frequent suggestion for improvement ($n = 4$) was the need to improve emotional validation and responsiveness to individual needs during mood assessment, which in turn would strengthen the therapeutic alliance. A psychotherapist recommended *to acknowledge the user's problems and to express the commitment to support the user as well as they can.* Despite achieving the highest ratings for mood assessment, the psychotherapist identified an issue: the chatbot failed to ask follow-up questions to assess the severity of the user's low mood. Although the user neither reported nor hinted at suicidality, and therefore the chatbot did not violate the emergency protocol, the psychotherapist emphasized the importance of conducting a more thorough evaluation in similar situations.

Phase 2: Psychoeducation

Both psychoeducation components, explaining the downward spiral ($M = 3.81$, $SD = 1.84$) and introducing different activity types ($M = 3.58$, $SD = 1.47$), exceeded the satisfactory threshold but ranked lower among the 14 components. The downward spiral explanation showed higher variability among psychotherapists than the introduction of activity types. Four psychotherapists ($n = 4$) recommended expanding the scientific foundation of psychoeducation, with one noting: *The explanation could be more specific about the scientific basis of why participating in activities helps*. Two psychotherapists ($n = 2$) suggested enhancing engagement by having users analyze example cases, while two others ($n = 2$) emphasized the need for more personalized content: *By asking several questions a user-specific disorder model should be developed and the connection [between activities and emotions] discussed more individually*. Regarding the downward spiral concept specifically, two psychotherapists ($n = 2$) recommended providing more detailed explanations, emphasizing the need to *describe that it is a spiral and by avoiding activities you will rather get worse than better*.

Phase 3: Finding Positive Activities

Both components exceeded the satisfactory threshold and received similar ratings: assessing activity levels ($M = 4.23$, $SD = 1.67$) and guiding activity identification ($M = 4.23$, $SD = 1.53$) both ranked 5th among all components. The most frequent recommendation ($n = 5$) was for the chatbot to provide more suggestions and guidance, particularly with resistant 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 ($n = 4$) recommended a more personalized approach *based upon personal preferences*. Psychotherapists also identified problematic activity recommendations, with one noting: *I find it difficult to recommend a nap in the evening, there is a good chance she'll just stay in bed*.

Phase 4: Planning Activities

Planning activities exceeded the satisfactory threshold and ranked 2nd among all components ($M = 4.98$, $SD = 1.41$). Three psychotherapists ($n = 3$) suggested creating more detailed plans with additional guidance: *take more intermediate steps during planning*. The psychother-

apists emphasized that the chatbot should ensure that users begin with smaller, realistic activities and adjust the number and type of activities based on the users' current mood and activity levels. Multiple psychotherapists stressed the importance of verifying the feasibility of activity plans.

Phase 5: Potential Barriers

Both components exceeded the satisfaction threshold. Identifying barriers ranked 9th ($M = 3.88$, $SD = 1.81$), and developing solution strategies ranked 11th ($M = 3.56$, $SD = 1.75$). Five psychotherapists ($n = 5$) highlighted the need for a more detailed barrier discussion, particularly with inactive users: *Potential barriers should be 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, four psychotherapists ($n = 4$) stressed the importance of ensuring realistic solutions: *talk about what is realistic and that it is normal that not everything works out right away.* Three psychotherapists ($n = 3$) recommended more personalized solution strategies: *Strategies for dealing with obstacles could be more specific to the situation of users. The recommendations lack greater personalization.*

Phase 6: Positive Reinforcements

Explaining reinforcement ($M = 2.92$, $SD = 2.30$) was the only component that fell below the satisfactory threshold of ≥ 3 , ranking last among all components, whereas developing reward strategies ($M = 4.10$, $SD = 1.78$) exceeded this threshold. Furthermore, the high standard deviation for explaining reinforcement ($SD = 2.30$) indicates a high variance among the sessions. Psychotherapists recommended providing clearer explanations of the reward system: *The reward system could be explained more precisely or ideas could be given about what it might look like exactly.* Two psychotherapists ($n = 2$) highlighted concerns regarding inappropriate rewards, particularly when using food as reinforcement. They recommended exploring multiple reward options, with one noting: *The reward ideas are somewhat one-sided; she probably drinks tea anyway.*

Phase 7: Conclusion

Although all three components surpassed the satisfactory threshold, their average ratings varied. The activity plan review was ranked fourth ($M = 4.29$, $SD = 2.14$), encouragement of implementation was ranked eighth ($M = 3.92$, $SD = 1.69$), and encouragement to observe activity-mood connections received the second-lowest rating among all components ($M = 3.02$, $SD = 2.04$). Two psychotherapists ($n = 2$) recommended enhancements to the plan summary and the provision of clearer, subsequent steps. Five psychotherapists ($n = 5$) highlighted the necessity for improved motivational strategies in encouraging plan implementation, with one noting: *A next appointment or a fixed time always helps, because otherwise patients often do not do their homework.* Regarding the monitoring component, five psychotherapists ($n = 5$) underscored the importance of providing a tracking template. Two psychotherapists ($n = 2$) requested clearer guidance on monitoring progress, whereas two others ($n = 2$) suggested better explanations of the purpose of the activity diaries. Psychotherapists recommended tracking both completed activities and intentions: *Discuss what happens next, how progress is tracked (for motivation and accountability) and when more aspects will be discussed.*

Session-Level Performance Profiles

Visual inspection of the rating heatmap (see Figure C.1 in the Appendix) revealed heterogeneity. While most sessions maintained satisfactory fidelity, two sessions exhibited consistently low Q-BAS ratings across components, warranting dedicated inspection. Content analysis revealed two distinct user steering patterns that caused complete protocol abandonment in the first phase of the study. In one session, the user adopted an “information-seeking” stance, asking broad self-help questions about motivation, stress management, and social skills. In the other, the user displayed “skeptical probing,” repeatedly following each suggestion with “What if that doesn’t help?” In both cases, the chatbot responded reactively without redirecting toward the structured protocol, never introducing the behavioral activation model, or establishing an activity plan. The sessions were effectively converted into unstructured FAQ exchanges. These findings suggest that early interaction dynamics determine session quality: once the therapeutic logic is disrupted, the system demonstrates a limited capacity for course correction.

Therapeutic Capabilities

Figure 5.2 (and C.2 in the Appendix) provides an overview of how psychotherapists rated the therapeutic capabilities of the chatbot on a 7-point scale (1 = fully disagree; 7 = fully agree). Ratings were generally positive, with message safety ($M = 6.90$, $SD = 0.37$) and message clarity ($M = 6.56$, $SD = 0.77$) receiving the highest scores, while building therapeutic rapport ($M = 5.13$, $SD = 1.45$) and natural conversation flow ($M = 5.25$, $SD = 1.42$) showed the most room for improvement.

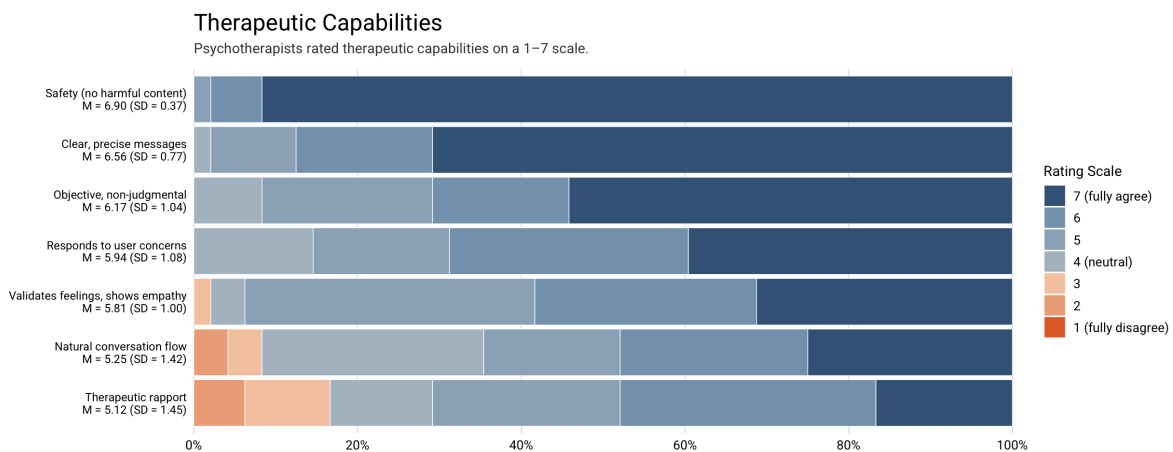


Figure 5.2: Therapeutic Capabilities

Validates Feelings and Shows Empathy. The chatbot received a mean rating of $M = 5.81$ ($SD = 1.00$) for validating feelings and showing empathy, ranking fourth among the seven therapeutic capabilities. Six psychotherapists ($n = 6$) praised this capability, whereas seven ($n = 7$) recommended improvements. Psychotherapists noted inconsistent validation, ranging from insufficient responses to excessive use of superlatives such as *super*, *fantastic*, *perfect*.

Responds Appropriately to User Concerns. The chatbot received a mean rating of $M = 5.94$ ($SD = 1.08$) for responding appropriately to user concerns, ranking third among the therapeutic capabilities. However, this was the most frequently mentioned area for improvement ($n = 15$). Psychotherapists noted inadequate responses to users' statements: *The user mentions two thoughts, but only one is addressed*. Two psychotherapists ($n = 2$) criticized the overuse of suggestive questions: *The constant question 'Do you find that good?' is too suggestive*.

Builds Therapeutic Rapport. Building therapeutic rapport received a mean rating of $M = 5.13$ ($SD = 1.45$), ranking last among therapeutic capabilities, with the highest variability in the ratings. One psychotherapist expressed skepticism: *Creating a 'real' therapeutic relationship (as an essential factor for therapeutic progress and motivation) will be challenging to achieve.*

Communicates Objectively Without Judgment. The chatbot received a mean rating of $M = 6.17$ ($SD = 1.04$) for communicating objectively without judgment, ranking second among its therapeutic capabilities. Four psychotherapists ($n = 4$) highlighted this capability as a strength.

Messages are Clear, Precise, and Avoid Complex Language. Message clarity received a mean rating of $M = 6.56$ ($SD = 0.77$), ranking first among all therapeutic capabilities, with the lowest standard deviation, indicating a robust performance across sessions. Despite this ranking, five psychotherapists noted that the language was sometimes overly therapeutic and suggested *using simpler language.*

The Conversation Flows Naturally. Natural conversation flow received a mean rating of $M = 5.25$ ($SD = 1.42$), ranking sixth among the therapeutic capabilities. The high standard deviation indicates varying ratings for the different sessions. Six psychotherapists ($n = 6$) criticized the pace of the sessions: *The transition from psychoeducation to planning comes too quickly and could be overwhelming. This is very 'chippy'.*

Messages are Safe and Do Not Contain Harmful Content. Message safety received a mean rating of $M = 6.90$ ($SD = 0.37$), ranking 1st among the therapeutic capabilities. It received consistently high safety ratings, with 44 of the 48 sessions (92%) receiving a maximum rating of 7, and the remaining sessions rated 6 (6%) or 5 (2%). None of the sessions received ratings below five. No specific safety concerns were raised by the psychotherapists.

Validity of Artificial Users

Artificial users demonstrated sufficient validity for evaluation purposes, although there were notable limitations. They received moderate authenticity ratings overall ($M = 3.75$, $SD =$

1.41) and were rated relatively easy to work with ($M = 5.77$, $SD = 1.46$). Artificial users with negative attitudes toward mental health chatbots were rated as significantly more authentic ($M = 4.16$, $SD = 1.52$) than those with positive attitudes ($M = 3.30$, $SD = 1.15$; $p = 0.036$). Crucially, while users with negative attitudes were rated as more authentic (Table C.11 in the Appendix), our analysis showed that they did not receive significantly lower therapeutic quality scores ($p = 0.275$, Table C.3 in the Appendix). This suggests that the chatbot's protocol adherence remains robust, even against more realistic and resistant user personas. This finding aligns with the psychotherapists' primary criticism of unrealistic compliance and engagement. Artificial users were *too willing to accept suggestions* and showed *very little resistance to the intervention* compared with real patients. One psychotherapist explained: *I usually do not experience much initiative. Usually it is first 'I don't know any positive activity' or 'I can't remember any positive activity'*. In contrast, the clinical presentation was deemed realistic: *The background stories were realistic. Also, how some people have gotten worse after COVID - I hear similar things*. This indicates that while the content of the vignettes was valid, the interaction dynamics, particularly compliance and engagement, require further improvement. No other user characteristics, besides attitudes toward mental health chatbots, significantly affected the ratings of authenticity or interaction difficulty. Additionally, none of the artificial user characteristics showed significant effects on ratings of the quality of behavioral activation or therapeutic capabilities (see Appendix for full results). See Appendix Tables C.3–C.12 for the complete exploratory statistics by artificial user characteristic.

5.5 Discussion

This study investigated how well an LLM-based chatbot can deliver a structured psychological intervention with clinical fidelity (RQ1) and what limitations and opportunities for refinement can be identified through clinical expert evaluation (RQ2). Our findings reveal a fundamental asymmetry: the chatbot achieved robust protocol adherence, completing all therapeutic phases with satisfactory fidelity on 13 of 14 components, but exhibited systematic weaknesses in tasks requiring clinical judgment. Performance also varied considerably across sessions, with specific interaction dynamics, such as users adopting information-seeking or skeptical stances, capable of derailing session quality. We contribute empirical, methodological, and design insights for LLM-based mental health chatbots, and we lay them out in the following subsections.

We also discuss the implications for future systems, noting that real-world clinical deployment requires additional evidence and safeguards.

5.5.1 LLMs Can Adhere to Instructions but Struggle with Clinical Reasoning

To our knowledge, this study provides one of the first pieces of evidence of an LLM-based chatbot conducting structured behavioral activation sessions, as measured by standardized fidelity scales. Across most sessions, the chatbot completed all seven phases and collaboratively developed activity plans, achieving fidelity scores comparable to human-delivered tele-behavioral activation (Rethorst et al., 2024) and consistently exceeding the feasibility threshold used in the human trials. This extends prior work demonstrating the potential of LLMs for structured mental health support, such as cognitive restructuring (Sharma et al., 2024) and motivational interviewing (Meyer & Elswailer, 2025), to a comprehensive multiphase therapeutic protocol.

A clear performance pattern emerged: components with explicit procedural instructions (mood assessment, psychoeducation, and activity scheduling) consistently exceeded fidelity thresholds, whereas components requiring clinical evaluation showed lower and more variable performance. This finding aligns with a recent systematic review of the capabilities and limitations of LLMs in mental health (Y. Wang et al., 2025), which provides evidence supporting the efficacy of LLMs in delivering psychoeducation. More broadly, this pattern is consistent with recent research demonstrating that while LLMs excel at structured tasks, they show limitations in cases where complex clinical judgment is required (Y. Hua et al., 2025; Kim et al., 2025; P. Qiu et al., 2025).

Qualitative feedback highlights this gap. Psychotherapists noted that the chatbot frequently accepted activities or rewards without verifying their therapeutic suitability—for example, recommending evening naps for users likely to remain in bed or accepting food-based rewards that might reinforce maladaptive patterns. These evaluation failures were not always captured in quantitative ratings, suggesting that current fidelity measures may underestimate clinical reasoning deficits.

The lower scores on evaluation tasks may partially stem from our specific prompt architecture, including the 30-word limit and linear phase enforcement, rather than solely reflecting the inherent limitations of the model. This mirrors broader HCI challenges in designing AI systems

that align algorithmic outputs with the complex realities of therapeutic practice (Thieme et al., 2023). Overall, our results suggest a recurring asymmetry in protocol-driven LLM-based interventions: procedural completion can be reliable, whereas evaluative judgment remains fragile. We expect similar asymmetries in other structured interventions, where safe and effective delivery depends on feasibility checks and appropriateness judgments rather than phase completion alone.

Based on our findings, we distinguish two categories of design recommendations: (1) *prompt-level refinements*, i.e., specific, actionable changes for low-performing components derived directly from psychotherapist feedback (Section 6.6), and (2) *architectural considerations*, which may address limitations beyond what prompting alone can achieve. All recommendations require empirical validation.

At the prompt level, few-shot prompting with examples of inappropriate activities and rewards may reduce errors through in-context learning (Brown et al., 2020). Adding an explicit verification step, asking whether planned activities are realistic and rewards therapeutically appropriate, could support collaborative reflection; however, this may shift the reasoning burden to users experiencing depression-related cognitive impairments (Rock et al., 2014).

Beyond prompting, multi-agent systems may help by separating reasoning from natural language generation; a reasoning agent can evaluate feasibility and appropriateness against clinical criteria and then direct a generation agent to produce a suitable response. Finally, fine-tuning with expert-curated examples of appropriate and inappropriate activities and rewards could help to integrate clinical reasoning patterns into the model. Future research should systematically assess which of these approaches best bridges the gap between generating plausible content and evaluating clinical validity.

5.5.2 Tension between Ensuring Protocol Adherence and Allowing for Personalized Intervention Delivery

Our study identified a tension between structure and flexibility in delivering behavioral activation across the phases. Based on initial small-scale tests, we adopted a prompt that emphasizes linear progression to support adherence and safety, which is crucial for vulnerable populations (Stade et al., 2024). Simultaneously, enforcing this structure constrains personalization and responsiveness, suggesting a fundamental design trade-off in current systems.

The “protocol abandonment” cases illustrate this tension concretely. When users adopted “information-seeking” stances (asking broad self-help questions) or displayed “skeptical probing” (repeatedly questioning whether suggestions would help), the chatbot responded reactively rather than redirecting toward the structured protocol. Without explicit redirection mechanisms, sessions were converted into unstructured FAQ exchanges, indicating that conversational flexibility can undermine therapeutic structure when users challenge the intervention logic. Simultaneously, these stances may signal that immediate redirection into the behavioral activation sequence is not always the best response. In such moments, different conversational strategies or alternative protocols may be needed (e.g., first addressing ambivalence, clarifying expectations and goals, or providing brief orientation) before returning to behavioral activation, if appropriate.

Addressing this limitation through multi-agent systems introduces additional challenges, including deciding when to progress, revisit earlier phases, or transition to different techniques based on the user’s needs. Whether such approaches can simultaneously achieve fidelity, flexibility, and clinically appropriate navigation remains unclear.

Fine-tuning on high-quality, domain-specific data (Stade et al., 2024) may also offer a viable path forward, as recent studies have demonstrated the potential of LLM-based mental health chatbots (Campellone et al., 2025; Heinz et al., 2025). However, concrete methods for managing the fidelity–flexibility trade-off in multi-turn therapeutic conversations, including reliable redirection and phase navigation, remain underexplored.

5.5.3 Formulaic Validation versus Genuine Responsiveness: Limitations in Relational Capabilities

Psychotherapists identified weaknesses in therapeutic alliance building, responsiveness, and flexible adaptation to patient needs (Coyne et al., 2019; Esposito et al., 2024). The chatbot often relied on formulaic validation, which psychotherapists perceived as superficial, and struggled to engage with nuanced user messages.

These limitations are critical because the therapeutic alliance is robustly linked to better outcomes in both traditional psychotherapy and digital mental health interventions (Flückiger et al., 2018). Highly effective psychotherapists consistently demonstrate strong interpersonal skills and responsiveness (Coyne et al., 2019; Esposito et al., 2024; Heinonen & Nissen-Lie,

2020).

The session-level derailments discussed above also suggest a rupture pattern: when users expressed skepticism (e.g., “What if that doesn’t help?”), the chatbot often did not address the concern directly and instead offered additional suggestions). In contrast, effective psychotherapists may recognize such skepticism as ambivalence that calls for motivational enhancement rather than further information.

Improving responsiveness and alliance building can be approached in several ways. At the prompt level, explicit instructions can incorporate evidence-based alliance-building strategies (Flückiger et al., 2018), such as establishing emotional bonds, agreeing on therapy goals and tasks, responding to users’ motivational readiness, and proactively addressing ruptures. Other development approaches, such as multi-agent systems and fine-tuning, may also help separate relational from task-oriented communication or train on expert-curated data that model strong alliance formation.

Importantly, these limitations partly reflect a deliberate design choice: we prioritized protocol completion and safety over alliance building, given the importance of structured delivery for vulnerable populations (Stade et al., 2024). Consequently, the observed weaknesses may reflect both this trade-off and the potential inherent challenges of LLM-based systems. Disentangling these factors will require studies that explicitly manipulate the design priority.

Future research should test which aspects of alliance building and responsive communication LLM-based systems can replicate and which interventions yield meaningful improvements. Comparing user-rated and expert-rated alliance measures would clarify whether these perspectives align; here, we assessed alliance from an expert perspective, whereas prior work often relies on user ratings (Darcy et al., 2021; Heinz et al., 2025). Direct comparisons with psychotherapists delivering identical protocols would further clarify whether these relational limitations translate into reduced effectiveness.

5.5.4 Usefulness of the Evaluation Approach

We position our evaluation approach as a pre-deployment method for assessing LLM-based mental health chatbots that aims to (i) increase the coverage of clinically relevant interaction patterns without exposing vulnerable users to untested systems, (ii) provide clinically meaningful evidence about intervention delivery quality rather than relying solely on downstream

outcomes or ad-hoc metrics, and (iii) generate actionable design feedback that can inform iterative refinement.

This approach addresses two complementary limitations in typical development workflows. First, coverage is limited when developers primarily role-play expected “golden paths.” Artificial users enable the systematic and scalable generation of complete behavioral activation sessions across diverse clinical profiles and user characteristics, reducing reliance on developer imagination and lowering the risk of exposing human users to early stage failures. Although authenticity ratings were moderate, psychotherapists judged core clinical aspects (symptom descriptions, situations, and suggested activities) as realistic, consistent with prior work showing that LLM-based artificial users can sustain coherent personas in multi-turn interactions (S. Chen et al., 2023; H. Qiu & Lan, 2024; J. Wang et al., 2024). The finding that artificial users with negative attitudes toward chatbots were rated as more authentic further suggests that resistance is an important factor for improving simulation realism.

Second, credibility is limited when the evaluation focuses on user satisfaction or symptom change without establishing whether the system delivered an evidence-based intervention of adequate quality. Therefore, we combined standardized fidelity ratings (Q-BAS) with qualitative expert feedback. This mixed-methods setup provides phase- and component-level diagnostics across behavioral activation, safety protocols, and broader therapeutic capabilities, extending prior evaluation approaches that emphasize single-turn interactions (Kocaballi et al., 2019) or general response appropriateness (Ding et al., 2023). It also responds to calls for more rigorous, standardized evaluations in LLM mental health research (Y. Hua et al., 2025) by demonstrating how validated fidelity instruments can be integrated into systematic evaluation workflows prior to human deployment.

A concrete advantage of combining artificial users with expert fidelity assessment is that it can surface failure modes that may be missed during developer testing. In our study, artificial users produced “information-seeking” and “skeptical probing” interaction patterns that derailed the intervention into unstructured FAQ-style exchanges, revealing that early user steering can lead to protocol abandonment and that the system has limited capacity for course correction. These sessions provided targeted evidence for the need to strengthen redirection and navigation logic before progressing to the human trials.

More broadly, our work suggests a staged evaluation pipeline for LLM-based mental health chatbots: developer role-playing for early functionality checks, artificial user testing for broader

coverage and risk mitigation, and human trials, with independent expert evaluation applicable throughout. A key design choice in our pipeline is the separation of data generation from evaluation: artificial users generated diverse session transcripts, while independent clinical experts assessed therapeutic quality, reducing the compounded bias that arises when developers both simulate users and evaluate their systems (Bunge & Desage, 2025).

Importantly, this workflow is not limited to the initial pre-deployment testing. It can also be applied during and after human testing and deployment as a regression testing and change validation mechanism. Specifically, artificial users and expert fidelity assessments can be used to stress-test proposed system updates (small refinements, larger redesigns, or higher-risk changes) and verify whether intended behavioral improvements (e.g., stronger redirection under skeptical probing, better feasibility checks, safer reward handling) are implemented successfully before rolling updates out to interactions with human users. Although this workflow does not replace human trials or real-world monitoring of engagement, safety, and outcomes, it provides an efficient intermediate gate for iterative improvement under clinical quality constraints.

However, this approach has important limitations. Artificial users may not accurately reproduce the dynamics of human interaction, particularly regarding resistance, engagement, and longitudinal adherence, which constrains the generalizability of the findings to real-world use. Fidelity assessment does not establish clinical efficacy, and single rater evaluations limit the extent to which observed variability can be disentangled from rater effects. These constraints motivate future work on improving artificial user realism, incorporating stronger reliability designs, and triangulating fidelity results with human studies and their outcome measures.

5.5.5 From Evaluation Insights to Targeted Prompt Refinement

While we have discussed broader architectural considerations, we now illustrate how our evaluation translates into concrete, prompt-level improvements. Based on psychotherapists' feedback, we focused on two areas that repeatedly required refinement: explaining positive reinforcement and supporting users in observing activity–mood connections. The following refinements were derived directly from expert feedback and are intended as targeted prompt modifications for future iterations; however, we have not empirically validated these specific changes.

Phase 6: Positive Reinforcement

Psychotherapists highlighted three recurring issues: (1) inconsistent explanations of reinforcement principles, (2) insufficient evaluation of reward appropriateness (e.g., accepting potentially problematic food-based rewards), and (3) limited personalization.

We propose three prompt-level refinements.

1. Structured explanation of the reinforcement mechanism: break the abstract concept into sequential steps with concrete examples. For instance, “First, explain how rewards help establish new habits by strengthening the association between an activity and a positive experience. Then, provide a relatable example (e.g., feeling proud after a short walk), and briefly check understanding.”
2. Screening and redirection for reward appropriateness: specify therapeutic red flags and response templates. For example, if the user suggests passive rewards (TV, social media, napping), validate the preference and propose a time-limited variant so the reward does not become an activity sink. If the user repeatedly suggests food rewards, explore alternative reward categories to avoid reinforcing maladaptive eating patterns.”
3. Guided personalization via reward categories: replace broad prompts (e.g., “How would you reward yourself?”) with a structured exploration. For example: “Offer options across several categories (immediate, delayed, and social rewards), ask which feels most motivating and feasible, and tailor the final choice to the user’s preferences and context.”

Phase 7: Conclusion

Psychotherapists indicated that guidance on observing activity–mood connections was often too vague, leading to superficial and inconsistent monitoring instructions. Their feedback pointed to three missing elements: a concrete tracking template, accountability structure, and clear next steps.

We propose three prompt-level refinements:

1. Template-based monitoring instruction: provide a ready-to-use tracking format rather than relying on open-ended generation. For example: “For each activity, write down: (1) What I did, (2) When I did it, (3) Mood before (1–10), (4) Mood after (1–10), (5) What I noticed.”

2. Implementation intentions: incorporate if–then planning explicitly (Gollwitzer & Sheeran, 2006). For example: “Help the user create one or two if–then plans for likely obstacles (e.g., ‘If I feel too tired on Tuesday evening, then I will do a 10-minute version of the planned walk instead of skipping it entirely.’)”
3. Structured closure checklist: replace generic encouragement with a brief verification sequence. For example: “Before ending, summarize the plan with key details, ask the user to confirm or restate the plan in their own words, address remaining concerns, and clarify what happens after this session.”

These refinements also illustrate broader prompt patterns for protocol-driven mental health interventions.

1. *Granular task breakdown.* Convert high-level directives into step-by-step procedures with explicit sequencing (what to explain first, which example to use, how to check understanding, and what to do next based on the user’s response).
2. *Template-based message content.* For outputs with specific formats (tracking templates, activity plans, reward lists), provide ready-to-use text rather than relying on an unconstrained generation. This reduces variability and supports clinical appropriateness (Thieme et al., 2023).
3. *Embedded clinical decision rules.* For tasks requiring judgment, specify explicit evaluation criteria and conditional redirection (e.g., feasibility checks, appropriateness checks, and actions to take when responses indicate a red flag).
4. *Explicit redirection protocols.* When users steer away from the intervention structure, specify conditional logic that acknowledges the request briefly and then redirects back to the current therapeutic task (e.g., “If the user asks general self-help questions unrelated to the current phase, validate briefly, then connect the question back to behavioral activation and proceed with the next step.”).

To support replication and extension, we provide complete prompts for both the behavioral activation chatbot and artificial users in the Supplementary Materials.

5.5.6 Limitations and Future Work

Our study has several limitations that motivate future research.

First, we lacked comparative benchmarks. We did not compare the chatbot's performance with psychotherapists delivering the same behavioral activation protocol, nor did we compare artificial user sessions with alternative evaluation methods such as expert role-playing or sessions with human users. Without these comparisons, we cannot establish the relative performance or quantify how well the findings from artificial users transfer to human interactions. Future research should include matched comparisons in which psychotherapists and chatbots deliver identical protocols under controlled conditions, ideally with independent expert ratings applied to both conditions.

Second, our artificial user method has limitations. Psychotherapists rated authenticity without knowing the specific characteristics each persona was designed to represent (e.g., high versus low willingness to disclose), which may have conflated perceived realism with expectations rather than fidelity to the intended design. Moreover, artificial users showed unrealistic compliance, which could overestimate performance compared with real patients. Future work should improve resistance and variability in artificial user behavior, incorporate manipulation checks for intended persona characteristics, and validate simulated interaction dynamics against human session data, where feasible.

Third, our findings are specific to GPT-4o (version 2024-08-06) and our prompt architecture. Other LLMs and newer model versions may exhibit different capabilities, failure modes, and safety behaviors. Replication across models, versions, and prompting strategies is necessary to assess robustness and distinguish model limitations from design-dependent effects.

Fourth, the generalizability of our evaluation approach remains uncertain. Behavioral activation is a structured, protocol-driven intervention; more flexible or complex therapeutic approaches may require different interaction designs and evaluation instruments. Future work should assess whether the combination of artificial user generation and fidelity-based expert rating transfers to other interventions, populations, and delivery formats.

Fifth, our evaluation focused on single-session interactions. Although single-session behavioral activation interventions can be effective (Kaveladze et al., 2025; Schleider et al., 2022), the ability of chatbots and artificial users to maintain coherence, personalization, and safety across multiple sessions remains untested. Longitudinal evaluations should examine whether fidelity

and responsiveness can be sustained across repeated sessions and evolving user goals.

Sixth, our single-rater design precluded inter-rater reliability estimation. The variance attributed to differences between sessions (36.8%) cannot be disentangled from unmeasured rater effects. Future studies should incorporate multi-rater designs for at least a subset of sessions to estimate reliability, separate rater variance from session variance, and strengthen the confidence in component-level conclusions.

Taken together, these limitations point to a research agenda that combines stronger benchmarks (human and model comparisons), improved simulation validity, broader replication across models and interventions, and more reliable expert rating designs while extending the evaluation from single-session testing to longitudinal use and iterative system updates.

5.6 Conclusion

We developed an LLM-based chatbot that delivers behavioral activation for depression and evaluated it using 48 sessions generated by artificial users with diverse clinical presentations. Ten psychotherapists assessed the sessions using the Q-BAS, a validated clinical fidelity instrument. The chatbot received positive overall ratings, with particularly high scores for mood assessment and activity planning but lower ratings for components such as explaining reinforcement and encouraging the observation of activity-mood connections. Message safety was consistently rated highly in all the sessions.

These findings reveal a fundamental asymmetry in the capabilities of this LLM-based chatbot for mental health delivery: while these systems can reliably execute treatment protocols, they cannot yet replicate the clinical reasoning that distinguishes competent therapy from excellent therapy. The chatbot struggled to evaluate whether activities and rewards were therapeutically appropriate, had difficulty building genuine therapeutic alliances, and showed limited ability to personalize interventions beyond surface-level adaptations.

This asymmetry between protocol execution and clinical judgment represents a central challenge for the tested configuration and likely the broader field. The tested chatbot can follow structured instructions with robust fidelity, but the evaluative capacities that characterize skilled psychotherapists remain unclear. Whether this gap can be bridged through architectural innovations, such as multi-agent systems or fine-tuning on expert-curated sessions, or whether it reflects more fundamental limitations, remains an open question with significant implications

for the future role of AI in mental health care.

For practitioners developing similar systems, our findings suggest actionable prompt-level patterns that can improve fidelity without architectural changes: granular task breakdown (transforming high-level instructions into step-by-step procedures), template-based content (providing ready-to-use message formats rather than relying on generation), embedded clinical decision rules (explicit evaluation criteria with redirection protocols), and explicit redirection mechanisms (conditional logic to reassert therapeutic direction when users steer off course).

Taken together, this study offers an initial, empirically informed starting point for designing LLM-based Behavioral Activation support for young people with depression. Real-world clinical deployment requires further evidence, refinement, and safeguards. To facilitate replication and advance the field, we provide complete prompts for both the behavioral activation chatbot and the artificial users.

6 Discussion³

This dissertation was designed to address three critical limitations in current mental health chatbots for young people with depression. First, existing chatbots have been designed primarily for and tested with adult populations, which does not address the unique needs of youth. Second, the design and effects of content personalization approaches are poorly understood, which has limited user engagement and therapeutic effectiveness. Third, the emergence of LLMs creates new opportunities for more flexible and personalized conversations, but also raises challenges for designing systems that maintain therapeutic quality and for evaluating them rigorously across development, deployment, and subsequent refinement. Through three interconnected studies, this thesis investigated how personalized mental health chatbots for young people can be designed, instantiated, and evaluated to improve user engagement and effectiveness. The previous chapters of this dissertation have presented the topic, established the theoretical foundations, and reported the findings of studies that focused on the design and effects of personalized mental health chatbots for young people. This discussion chapter seeks to integrate the primary findings, outline both the contributions to research and practice, acknowledge limitations, and provide an outlook on the future of personalized mental health chatbots.

Part I provided a comprehensive foundation by investigating how to design chatbots that address the specific needs of youth with depression. Through in-depth interviews, this study revealed their problems, coping strategies, attitudes, and design preferences for chatbots to treat depression. The findings yielded several design recommendations, such as the importance of different types of youth-specific dialogue content, different, partially contradictory design preferences, and a request for personalization features. By providing these recommendations, Part I provided a comprehensive foundation to design chatbots for youth with depression.

Building on these insights, **Part II** examined the design and effects of different content personalization approaches through two field experiments with the software artifact TheraBot. The first experiment showed that hybrid content personalization led to a greater reduction in depression severity than standardized content. The second experiment compared three personalization approaches (chatbot-led, user-led, and hybrid) and revealed that only hybrid personalization

³This chapter discusses and integrates the findings of the following studies, which are published, available as preprints, or under review: (Kuhlmeier et al., 2022a; Kuhlmeier et al., 2022b; Kuhlmeier et al., 2025a; Kuhlmeier et al., 2025b; Kuhlmeier et al., 2026).

outperformed standardized content delivery, whereas chatbot-led and user-led personalization alone did not show clear benefits. Importantly, this study also showed that the effects of personalization depend on the individual user. Engagement varied across personalization approaches depending on baseline depression severity, stated preferences, and agreeableness. Users with higher baseline depression severity were more engaged in the user-led condition but less engaged in the chatbot-led condition, and hybrid personalization was preferred by most users. These findings challenge the assumption that a single personalization approach suits all users and highlight the need to adapt the personalization approach itself to the individual.

Part III addressed the emerging challenge of designing and evaluating LLM-based mental health chatbots by developing an LLM-based behavioral activation chatbot for young people with depression. Building on the earlier rule-based system, this study focused on how an evidence-based intervention can be translated into a structured prompt architecture that maintains therapeutic structure while enabling more natural and personalized conversations. The study then assessed this design with an evaluation approach that combines artificial users based on clinically informed patient vignettes with fidelity assessment conducted by clinical experts. This approach enables systematic and risk-reduced evaluation of therapeutic quality across diverse clinical profiles while preventing potential harm to vulnerable populations and avoiding an additional evaluation burden on young people with mental health problems. The study demonstrated the usefulness of this design and evaluation workflow and revealed an important asymmetry in the chatbot's capabilities: while the chatbot showed robust protocol adherence, completed all seven intervention phases in most sessions, and was consistently rated as safe, systematic weaknesses emerged in tasks requiring clinical judgment, such as evaluating activity feasibility and therapeutic appropriateness, developing personalized strategies for user-specific barriers, and building genuine rapport. Performance was also heterogeneous: specific user steering patterns, namely information-seeking and skeptical probing, caused protocol abandonment in a subset of sessions, highlighting the need for stronger redirection logic and further prompt-level and architectural refinement before human deployment. These findings provide an empirically informed starting point for the iterative design of LLM-based mental health chatbots.

6.1 Research Contributions

This dissertation has presented a series of interconnected studies that investigated various aspects of personalized mental health chatbots. Each study makes unique theoretical contributions to research on digital mental health, IS, and human-computer interaction. These contributions will be organized and discussed according to the three parts of the thesis (see Table 6.1).

Table 6.1: Research Contributions of this Dissertation

Part	Research Contributions
Part I	First investigation of chatbot design for youth with depression Systematic overview of youth problems and coping strategies with depression and their attitudes and design preferences for chatbots to treat depression 15 comprehensive design recommendations
Part II	Design knowledge for content personalization (six design requirements including DR6) First empirical evidence of the effectiveness of content personalization in mental health chatbots Comparative evidence that only hybrid personalization outperforms standardized content Understanding of how user characteristics (depression severity, preferences, agreeableness) interact with personalization strategies to affect engagement Evidence for the need to “personalize the personalization” strategy itself
Part III	Design knowledge for LLM-based behavioral activation chatbots for young people with depression Actionable prompt engineering patterns and architectural considerations to improve therapeutic fidelity Complete prompt architecture and complete prompts for the chatbot and artificial users to support replication and extension Novel design and evaluation workflow combining artificial users and expert fidelity assessment to identify failure modes and guide iterative refinement

Part I of this dissertation addresses the overarching research question: *How should chatbots be designed to effectively treat depression among youth?*. Through in-depth interviews and a prototype evaluation with young people affected by depression, the qualitative study gathered comprehensive insights into their problems with depression and coping strategies, as well as

attitudes and design preferences for chatbots to treat depression. This study contributes by, to our knowledge, being among the first to investigate the design of chatbots to treat depression in youth, and by revealing the needs and preferences of a previously overlooked population in digital mental health research. While findings from adult populations cannot simply be generalized to youth because of fundamental differences in developmental processes, depression symptoms, and the use of digital devices (Andone et al., 2016; Huffman, 2014; Rice et al., 2019), and previous research with youth has only explored chatbot design for the prevention of mental disorders (Grové, 2021; Høiland et al., 2020), this study specifically addresses the therapeutic needs of young people diagnosed with depression.

The systematic overview of youth-specific needs and preferences represents a key contribution, as it offers both practical guidance for developers and a foundation for researchers in digital mental health. The findings are synthesized into 15 design recommendations that emphasize the importance of repeated multidimensional assessments, clear communication of chatbots' limitations, youth-specific content databases, and personalization features. Importantly, this study also revealed the central importance of personalization for chatbots for youth with depression, as participants expressed diverse and sometimes contradictory preferences regarding what should be personalized and whether personalization should be controlled by the user, the chatbot, or both. These comprehensive guidelines therefore also reveal key opportunities for future research, particularly with respect to the implementation of personalization strategies. This foundation not only advances mental health chatbot research and development, but also inspired the research conducted in Part II and Part III of this dissertation.

Part II of this dissertation addresses the overarching research questions of how content personalization in a mental health chatbot can be designed for young people with depression and what effects different content personalization approaches have on user engagement and effectiveness. Following the DSR paradigm, we developed and evaluated TheraBot, a mobile chatbot for young people with depression that implements hybrid content personalization. In two field experiments, we investigated both the effects of personalized content compared to standardized content and the relative effects of chatbot-led, user-led, and hybrid personalization approaches. The first experiment showed that personalized content led to a greater reduction in depression severity than standardized content. The second experiment found that hybrid personalization was the most promising approach: it outperformed standardized content delivery, whereas chatbot-led and user-led personalization did not show clear benefits on their own. This

research contributes to the literature on content personalization in mental health chatbots in five key ways.

First, we provide theory-grounded prescriptive design knowledge for implementing content personalization in mental health chatbots. Six design requirements were identified, ranging from establishing a well-defined treatment plan based on users' symptoms (DR1) and balancing structure with flexibility (DR2), through user-led (DR3) and chatbot-led (DR4) personalization features, to a hybrid approach (DR5). Crucially, our evaluation results led to a sixth requirement (DR6): the chatbot should adapt the personalization approach itself based on user characteristics and stated preferences. We instantiated these requirements in TheraBot through a two-phase personalization process in which the chatbot first assesses users' symptoms with standardized mental health questionnaires to recommend relevant therapeutic modules and then allows users to adapt these recommendations based on their perceived needs and preferences. This design knowledge lays an important foundation for developing mental health chatbots that incorporate personalized content. Second, we provide, to our knowledge, the first empirical evidence on the effects of content personalization in mental health chatbots, addressing a critical gap in existing research (Kocaballi et al., 2019) and helping to explain conflicting evidence on the effectiveness of content personalization in DMHIs. Against the background of largely null findings in prior randomized trials (Schaeuffele et al., 2025), our results show that hybrid content personalization can improve therapeutic effectiveness relative to standardized content delivery. Third, this research provides novel comparative evidence on the relative effects of chatbot-led, user-led, and hybrid content personalization approaches. Only hybrid personalization outperformed standardized content delivery; neither chatbot-led nor user-led approaches showed benefits over standardized content delivery on their own. This finding may help explain why prior studies testing only system-led or user-led personalization produced null results (Berger et al., 2014; Johansson et al., 2012) and supports hybrid personalization as the most effective single approach when one default strategy must be chosen. Fourth, we reveal important insights into how user characteristics interact with personalization approaches to influence user engagement. Users with higher baseline depression severity were more engaged in the user-led condition but less engaged in the chatbot-led condition. In addition, hybrid personalization was preferred by most users (77%), while 20% preferred either chatbot-led or user-led personalization. Furthermore, in the user-led condition, participants high in agreeableness tended to follow randomly generated module suggestions without modification, which likely reduced the

therapeutic relevance of their selected modules. Together, these findings show that symptom severity, stated preferences, and personality characteristics shape how users respond to personalization approaches. Fifth, these insights collectively point to the importance of “personalizing the personalization” approach itself. Rather than selecting a uniform personalization approach for all users, the chatbot should adapt whether it uses chatbot-led, user-led, or hybrid content personalization based on individual user characteristics such as symptom severity, personality traits, and stated preferences. This represents a novel and practically actionable direction for increasing the effectiveness of mental health chatbots beyond what any single fixed approach can achieve alone.

Part III of this dissertation addresses the overarching research question: *How can an LLM-based behavioral activation chatbot for young people with depression be designed and iteratively refined to maintain clinical fidelity while delivering a structured psychological intervention?* To answer this question, we developed an LLM-based behavioral activation chatbot for young people with depression and used artificial users together with expert clinical fidelity assessment by psychotherapists to identify design limitations and opportunities for refinement without exposing vulnerable users to risk. Our research contributes to the literature on LLM-based mental health chatbots in three key ways. First, we contribute design knowledge for LLM-based mental health chatbots by developing a structured prompt architecture for behavioral activation and by publishing the complete prompts for both the chatbot and artificial users in the Supplementary Materials. Together with the identified prompt-level refinements and broader architectural considerations, this provides a reusable foundation for replication, extension, and iterative refinement. Second, we propose an evaluation approach that integrates artificial users, derived from clinically validated patient vignettes and enriched with seven systematically varied characteristics, with detailed psychotherapist fidelity assessments using validated instruments. This approach enables rigorous and risk-reduced assessments across diverse clinical profiles, surfaces failure modes that might be missed during developer testing, and can serve as a reusable regression-testing gate for iterative system updates. However, it is important to note that artificial-user authenticity was only moderate and compliance was often unrealistic, which constrains the generalizability of the findings to real-world interactions, and that fidelity assessment does not establish clinical efficacy. Third, our evaluation reveals a fundamental asymmetry in LLM-based behavioral activation delivery: the chatbot demonstrated robust protocol adherence, achieving satisfactory fidelity on 13 of 14 Q-BAS components, but exhibited

systematic weaknesses in tasks requiring clinical judgment, particularly reinforcement explanation and activity-mood monitoring. Specific user steering patterns, such as information-seeking and skeptical probing, could derail sessions into unstructured FAQ exchanges, indicating limited capacity for course correction. These results offer qualified, preliminary evidence, bounded by the specific configuration tested (GPT-4o, single-session behavioral activation, single-rater assessment), and they informed targeted recommendations for prompt-level and architectural improvements. Real-world clinical deployment therefore requires further validation, refinement, and safeguards.

In summary, these three parts collectively provide valuable insights into the design and effects of personalized mental health chatbots for young people. The main research question of this dissertation—*how can personalized mental health chatbots for young people be designed to improve user engagement and effectiveness?*—is addressed by providing comprehensive design guidelines, improving our understanding of content personalization strategies, and developing a novel approach to evaluate LLM-based chatbots. Overall, this dissertation makes valuable research contributions to the fields of digital mental health, IS, and human-computer interaction. The insights gained from these investigations offer a solid foundation for future research and practical applications and enhance the development of personalized mental health chatbots that are engaging and effective, and therefore improve mental healthcare for young people with depression.

6.2 Practical Contributions

This dissertation provides several practical contributions relevant to the design and evaluation of personalized mental health chatbots for young people with depression (see Table 6.2).

First, it offers comprehensive evidence-based design guidelines derived from a series of complementary studies. In **Part I**, qualitative interviews with young people revealed their problems and coping strategies with depression, as well as attitudes and design preferences for chatbots to treat depression, which yielded design recommendations as a comprehensive foundation for developers of mental health chatbots. In **Part II**, six design requirements for implementing content personalization in mental health chatbots were identified and validated empirically. These requirements were instantiated in chatbot-led, user-led, and hybrid personalization features, and culminate in DR6: the chatbot should adapt the personalization approach itself based

Table 6.2: Practical Contributions of this Dissertation

Part	Practical Contributions
Part I	Comprehensive set of design recommendations providing practical guidance for mental health chatbot developers
Part II	<p>TheraBot: A fully functional mobile chatbot application with demonstrated effectiveness in reducing depression severity</p> <p>Flexible software architecture supporting both rule-based and LLM-based implementations while maintaining the same frontend</p> <p>Six validated design requirements with implementation examples for chatbot-led, user-led, and hybrid personalization features, including DR6 guidance to adapt the personalization approach to the individual user</p>
Part III	<p>Multi-component prompt architecture for LLM-based behavioral activation chatbots (system-level rules, task, phase-specific instructions, contrastive examples, and full session exemplar)</p> <p>Published complete prompts for both the chatbot and artificial users in the Supplementary Materials to support replication and extension</p> <p>Actionable prompt-level refinements and broader architectural considerations for improving therapeutic fidelity</p> <p>Set of 48 artificial users generated from five clinical vignettes and enriched with seven systematically varied characteristics</p> <p>Reusable design and evaluation workflow combining artificial users, standardized fidelity instruments, and expert assessment for risk reduction, failure-mode identification, and iterative testing across the system lifecycle</p>

on symptom severity, personality characteristics, and stated preferences. Both the requirements and their instantiation can serve as starting points for future development, with hybrid personalization recommended as the most effective single default strategy and adaptive personalization as the longer-term design goal. Together, these guidelines provide developers with clear and evidence-based instructions for creating more effective personalized chatbots across various architectural approaches.

Second, the dissertation presents "TheraBot", a fully functional mobile chatbot application designed specifically for young people with depression. TheraBot represents a tangible software artifact that has demonstrated real-world applicability through successful deployment in field experiments via common app stores. The studies in Part II also showed that TheraBot's hybrid content personalization can reduce depression severity more than standardized content delivery. In addition, its flexible architecture allows for transitioning between rule-based and LLM-based implementations, as demonstrated by the LLM-based behavioral activation implementation in Part III, while maintaining the same frontend and thus enabling systematic comparison and

investigation of different chatbot architectures.

Third, the dissertation contributes a multi-component prompt architecture for an LLM-based mental health chatbot and publishes the complete prompts for both the behavioral activation chatbot and the artificial users in the Supplementary Materials. The architecture includes system-level rules and constraints (including a 30-word message limit, safety protocols, and role boundaries), task-level instructions defining the overall therapeutic goal, phase-specific instructions with completion criteria and contrastive good/bad examples, and a complete session exemplar demonstrating all therapeutic phases. Although the system-level constraints and safety protocols remain applicable across different therapeutic approaches, task- and phase-specific instructions can be adapted for other modules such as cognitive restructuring or mindfulness exercises. Beyond the prompt artifact itself, the study derives actionable prompt-level refinements and broader architectural considerations, including template-based content, embedded decision rules, granular task breakdown, and explicit redirection mechanisms. Together, these design contributions provide developers with a reusable starting point for implementing structured evidence-based interventions with LLMs while maintaining therapeutic boundaries and supporting replication, extension, and iterative refinement, though the prompt patterns require empirical validation for each new configuration and context.

Fourth, the dissertation contributes a practical design and evaluation workflow for LLM-based mental health chatbots. This workflow combines artificial users generated from clinically validated patient vignettes and enriched with varied user characteristics with detailed psychotherapist fidelity assessments using standardized evaluation instruments. It allows developers and researchers to identify failure modes, assess clinical fidelity, and generate targeted design feedback while reducing risks for vulnerable patients and avoiding additional burdens on them, thereby reducing the risk of exposing human users to inadequate system versions. It can be used before initial human trials, during iterative development, and after deployment for testing specific scenarios or validating higher-risk system updates. Importantly, this workflow is a risk-reduction and quality-screening mechanism rather than a replacement for human trials or real-world monitoring: it does not establish clinical efficacy, and the moderate authenticity and unrealistic compliance of artificial users limit how far the results generalize to real-world interactions. The workflow provides a reusable foundation across iterative development cycles and can be extended to other therapeutic modules and LLM configurations.

These practical contributions provide researchers and developers with concrete design guide-

lines, a software artifact, a modular prompt architecture with published complete prompts, and a reusable design and evaluation workflow to develop and evaluate more effective personalized mental health chatbots for young people with depression.

6.3 Limitations and Future Research

The studies presented in this dissertation have adhered to established methods for research design, implementation, and documentation. Although each chapter addresses specific limitations, several overarching limitations need to be acknowledged and addressed in future research.

First, the studies in all three parts relied on hypothetical scenarios and short-term evaluations. In **Part I**, the qualitative study captured the explicit design preferences of the participants, but did not examine their effect on actual usage behaviors. Further, although similar to previous studies (Fitzpatrick et al., 2017), the field experiments conducted in **Part II** were limited to relatively short durations (14 days), which restricts insights into sustained engagement and long-term therapeutic outcomes. In addition, the evaluation of the LLM-based chatbot in **Part III** focused on a single session of behavioral activation and did not capture longitudinal, multi-session interactions. Artificial users received only moderate authenticity ratings and tended toward unrealistic compliance, which means the fidelity results cannot be equated with performance on real patient populations. Furthermore, findings are specific to the model version tested (GPT-4o) and the prompt architecture used; other configurations may exhibit different fidelity profiles, failure modes, and safety behaviors. Although these studies provided valuable insights, future research should extend these findings by using longitudinal studies that examine which design elements influence sustained engagement and therapeutic effectiveness most strongly in real-world settings. For Part III specifically, this includes replication across LLM versions and therapeutic modules, multi-rater fidelity designs to establish reliability, and human trials to bridge the gap between fidelity evidence and real-world clinical outcomes.

Second, the ability to generalize the findings is limited because of selective and non-representative samples. In Part I, the qualitative study involved a small, predominantly female sample, which limits the ability to generalize the findings, particularly to male youth and those who avoid seeking professional help. Similarly, the field experiments in Part II recruited participants with at least mild self-reported depression symptoms ($\text{PHQ-8} \geq 6$) rather than clinical diagnoses

provided by experts, which questions the ability to generalize the results to clinically diagnosed populations. In addition, neither study included participants along the entire age range of young people. Future studies should attempt to recruit larger, more diverse samples that include participants with clinical diagnoses, male youth, and those who avoid treatment to validate or extend our findings.

Third, this thesis investigated only a limited subset of possible personalization dimensions. In **Part II**, the examined content personalization approaches varied along only one dimension, namely who controls personalization (automation), through specific implementations of user-led, chatbot-led, and hybrid approaches. However, personalization in psychotherapy and digital mental health can vary along many other dimensions, including structure, timing, and specific implementation mechanisms (Cohen et al., 2021; Fan & Poole, 2006; Kocaballi et al., 2019). Future research should therefore explore more advanced personalization approaches, such as methods that incorporate passive sensing via smartphones or wearables and LLM-based personalization dialogues that complement standardized mental health questionnaires. The temporal dimension also warrants further investigation to determine when and how often personalization should occur to balance structure and flexibility optimally. In addition, future studies should examine forms of personalization beyond content selection, such as adapting interaction style (e.g., formal vs. casual language) based on user feedback and preferences.

Finally, because of the rapid advancements in technology, particularly LLMs, the design guidelines outlined in this dissertation need ongoing assessment and updates to ensure that the insights remain relevant and effective in meeting users' evolving needs and technological advancements. For example, attitudes, expectations, and preferred dialogue content may change as young people become more accustomed to interact with LLM-based tools such as ChatGPT on a daily basis.

6.4 Outlook

Mental health chatbots have evolved from rule-based systems with standardized treatments through basic personalization features (Kocaballi et al., 2019) to the emerging development of LLM-based systems. As highlighted in **Part I** of this dissertation, existing DMHIs often did not meet the needs of young users. Part II demonstrated the importance of understanding the design and effects of different personalization strategies to improve therapeutic outcomes,

while **Part III** showed both the potential of LLMs to support structured therapeutic delivery and the current limitations in clinical judgment, rapport, and handling of off-protocol user behavior that must be addressed before broader deployment. These technological advances, particularly the impressive capabilities of recent LLMs, have significantly expanded the possibilities for mental health chatbots and promise a future in which we can achieve advanced personalized and human-like treatment.

A crucial first step toward achieving advanced personalization is developing LLM-based mental health chatbots that can deliver evidence-based treatments, such as behavioral activation or cognitive restructuring, safely and with sufficient therapeutic quality while supporting more open-ended and flexible conversations (Campellone et al., 2025; Cuijpers et al., 2023b; Hatch et al., 2025; Heinz et al., 2025; Y. Hua et al., 2025; Sharma et al., 2024; Stade et al., 2024). This development process presents an opportunity for closer collaboration between research on mental health chatbots and psychotherapy research, as both fields grapple with related challenges, especially establishing and maintaining effective working alliances (Cameron et al., 2018; Darcy et al., 2021; Flückiger et al., 2018). Current work already moves in this direction by aligning LLM-based systems more closely with established psychotherapeutic practice through expert-curated training data, clinical guardrails, and fidelity-based evaluation (Campellone et al., 2025; Heinz et al., 2025; Y. Hua et al., 2025; Stade et al., 2024). In the longer term, such systems may also serve as tools for exploring alternative explanations, interaction strategies, and intervention designs that can subsequently be evaluated by clinicians and researchers. If these challenges can be addressed, LLM-based mental health chatbots may enable increasingly personalized and open-ended therapeutic conversations, marking an important step forward in digital mental healthcare.

An important related direction is to move beyond treating mental health chatbots as purely conversational agents and toward systems that combine LLM-based dialogue with Generative UI. In psychotherapy, treatment does not consist only of conversation, but also of structured tasks such as homework, mood monitoring, symptom tracking, activity planning, and the practice of therapeutic exercises (Abel & Hautzinger, 2013). Generative UI could therefore personalize not only what therapeutic content is delivered, but also which interactive elements are presented, when they are shown, and how they structure users' engagement with treatment. For example, a system might dynamically surface mood monitoring, activity tracking, homework review, or cognitive restructuring tools only when they become relevant in the therapeutic pro-

cess. Such interfaces could help reduce cognitive load and better align flexible conversations with structured evidence-based treatment, but they would also require careful safeguards to ensure clinical consistency, transparency, and user control.

The integration of diverse data sources for personalization presents both opportunities and significant challenges. Although these systems can theoretically analyze large amounts of user data in real time, implementing comprehensive personalization requires substantial technical infrastructure and careful consideration of the effort-benefit ratio. The collection and processing of ecological momentary assessments, wearable data, smartphone interactions, and biosignal measurements demand sophisticated data preparation pipelines, robust security protocols, and significant computational resources. Although such data could improve therapeutic interactions through more contextualized responses, developers must weigh the marginal benefits of additional data sources and more advanced personalization algorithms against the increased complexity of the implementation and potential privacy risks. For example, while analyzing speech patterns and facial expressions might enable more empathetic responses, the collection and processing of such sensitive biometric data raise serious privacy concerns and require strict compliance with data protection regulations. Another important challenge concerns ownership, governance, and accountability. Many widely used mental health applications are developed and operated by private companies and increasingly rely on proprietary foundation models and infrastructures controlled by commercial providers. Under conditions of private ownership, decisions about which systems are built, which populations are prioritized, what risks are acceptable, and how user data are used ultimately remain with the owners of the systems and underlying infrastructures, whose decisions are also shaped by pressures to generate profit, attract investment, grow, and maintain competitive advantage. Participation by the people and institutions most affected by these systems, including patients, clinicians, researchers, regulators, and public health organizations, therefore depends on whether and how those owners choose to include them. This governance structure can limit transparency and create tensions between therapeutic benefit, data minimization, patient autonomy, and commercial objectives. In mental health care, this tension is especially consequential because systems optimized for engagement, data collection, profit, market reach, or platform dependency may not always align with treatment quality, privacy, and long-term user well-being. Future development should therefore not be treated as a purely technical optimization problem. It also requires institutional arrangements that give affected users, clinicians, researchers, regulators, and public-interest organiza-

tions meaningful influence over development priorities, evidence standards, safety thresholds, data governance, and deployment conditions. Furthermore, the development of personalization features requires extensive interdisciplinary collaboration between clinical experts, privacy specialists, and technical teams to ensure both therapeutic validity and user protection (Boucher et al., 2021). Developers implementing these systems must carefully balance the promise of improved therapeutic outcomes against practical constraints, including the development effort and the need to protect user privacy.

In summary, the convergence of LLMs, various data sources, and advanced machine learning algorithms could enable more personalized treatment and increased therapeutic effectiveness in digital mental health support. Although the potential to provide more accessible, engaging, and effective mental health support to young people makes this a compelling direction for future research and development, success will depend upon carefully planned development that prioritizes therapeutic quality and user safety. Realizing this potential requires close collaboration between researchers, mental health professionals, and developers to ensure that technological advancements serve to improve rather than compromise the quality of digital mental health care for young people.

7 Conclusion

Depression is a prevalent mental health condition among young people that has significant detrimental personal and societal effects. Although chatbots offer promising solutions to increase access to support through their anonymous, easily accessible, and conversational nature, existing systems often fail to meet young users' specific needs or show limited user engagement and effectiveness. This dissertation addresses these challenges by investigating how personalized mental health chatbots for young people can be designed, instantiated, and evaluated to improve engagement and effectiveness. It consists of three interconnected parts that contribute to our understanding of effective personalized mental health chatbots at complementary levels, ranging from youth-sensitive design, to personalized content selection, to the personalization of therapeutic content itself during ongoing LLM-based conversations.

Part I Study 1 presents a qualitative investigation of youths' needs and preferences and provides foundational insights into designing chatbots to treat depression in youth. Through in-depth interviews and prototype evaluations, this study reveals key design recommendations for developing engaging and effective chatbots.

Part II examines the design and effect of content personalization strategies through two field experiments with TheraBot, a custom-developed mental health chatbot for young people with depression. These studies demonstrate that content personalization can improve therapeutic outcomes, but with important nuance: only hybrid personalization—in which the chatbot recommends therapy modules but users can adapt the selection—significantly outperformed standardized content delivery. Chatbot-led and user-led personalization strategies did not individually show statistically significant benefits over standardized content. Engagement effects were conditional: users with higher baseline depression severity completed more modules and spent more time with TheraBot in the user-led condition, whereas this trend was reversed in the chatbot-led condition. Additionally, 77% of participants preferred hybrid personalization, and participants high in agreeableness in the user-led condition tended to follow even randomly generated suggestions, reducing the therapeutic relevance of their module selection. Together, these findings indicate that no single personalization strategy is optimal for all users, and that effective personalization requires adapting the strategy itself to individual characteristics such as depression severity, stated preferences, and personality traits. This insight is captured in DR6, the final design requirement derived from Part II, which represents a novel and practi-

cally important contribution to the personalization of DMHIs.

Part III extends this focus from personalizing which therapeutic content is selected to personalizing the therapeutic content itself as it is delivered in ongoing conversations through an LLM-based behavioral activation chatbot for young people with depression. Rather than deciding only which modules or intervention components users should receive, this part addresses how evidence-based therapeutic content can be phrased, adapted, and maintained within real-time dialogue. The study develops a structured prompt architecture intended to maintain therapeutic structure while enabling more natural and personalized conversations, and evaluates this design through a pre-deployment workflow that combines artificial users with clinical expert fidelity assessment. This approach allows systematic, risk-reduced testing across diverse clinical profiles and helps identify concrete opportunities for refinement before human deployment. The findings reveal an asymmetry in the chatbot’s capabilities: robust protocol adherence was achieved in most sessions, but tasks requiring clinical judgment—such as evaluating activity feasibility, explaining reinforcement, building rapport, and redirecting conversations after skeptical or information-seeking user steering—showed systematic weaknesses. These results provide an empirically informed starting point for prompt-level and architectural refinement; however, they constitute pre-deployment evidence from simulated sessions only and do not establish clinical efficacy.

In the course of this dissertation, it has been shown that the effectiveness of personalized mental health chatbots for young people can be improved significantly through careful attention to users’ needs, the integration of evidence-based personalization strategies, and rigorous design-oriented evaluation of LLM-based systems. A key insight is that personalization must operate at multiple levels: it concerns both the selection of relevant therapeutic content and the personalization of the therapeutic content itself during conversation. Part II therefore focuses on deciding which modules or treatment elements should be prioritized for a given user, whereas Part III focuses on how this therapeutic content can be delivered, adapted, and sustained conversationally through an LLM-based system. At the same time, personalization strategies themselves must be adaptive: rather than applying a uniform strategy to all users, the most effective approach selects between chatbot-led, user-led, and hybrid personalization based on individual characteristics such as depression severity, stated preferences, and personality traits. The research provides evidence-based design guidelines, validated software and prompt artifacts, and a reusable pre-deployment evaluation workflow that together advance our understanding of

how to develop and iteratively improve personalized mental health chatbots—while underscoring that further refinement and human-trial validation are needed before LLM-based systems can be considered clinically effective. The practical implications and future directions discussed in this dissertation open pathways for the development of more engaging and effective personalized mental health chatbots that can further improve mental healthcare for young people.



Bibliography

- Abd-Alrazaq, A. A., Alajlani, M., Ali, N., Denecke, K., Bewick, B. M., & Househ, M. (2021). Perceptions and opinions of patients about mental health chatbots: Scoping review. *Journal of Medical Internet Research*, 23(1). <https://doi.org/10.2196/17828>
- Abd-Alrazaq, A. A., Rababeh, A., Alajlani, M., Bewick, B. M., & Househ, M. (2020). Effectiveness and safety of using chatbots to improve mental health: Systematic review and meta-analysis. *Journal of Medical Internet Research*, 22(7). <https://doi.org/10.2196/16021>
- Abel, U., & Hautzinger, M. (2013). *Kognitive verhaltenstherapie bei depressionen im kindes- und jugendalter*. Springer-Verlag.
- Adomavicius, G., & Tuzhilin, A. (2005). Personalization technologies: A process-oriented perspective. *Communications of the ACM*, 48(10). <https://doi.org/10.1145/1089107.1089109>
- Agapie, E., Chang, K., Patrachari, S., Neary, M., & Schueller, S. M. (2022). Understanding mental health apps for youth: Focus group study with latinx youth. *JMIR Formative Research*, 6(10). <https://doi.org/10.2196/40726>
- Agnafors, S., Norman Kjellström, A., Torgerson, J., & Rusner, M. (2019). Somatic comorbidity in children and adolescents with psychiatric disorders. *European Child & Adolescent Psychiatry*, 28(11). <https://doi.org/10.1007/s00787-019-01313-9>
- Ahmad, R., Siemon, D., Gnewuch, U., & Robra-Bissantz, S. (2022). Designing personality-adaptive conversational agents for mental health care. *Information Systems Frontiers*, 24(3). <https://doi.org/10.1007/s10796-022-10254-9>
- Ahmed, A., Ali, N., Aziz, S., Abd-alrazaq, A. A., Hassan, A., Khalifa, M., Elhusein, B., Ahmed, M., Ahmed, M. A. S., & Househ, M. (2021). A review of mobile chatbot apps for anxiety and depression and their self-care features. *Computer Methods and Programs in Biomedicine Update*, 1. <https://doi.org/10.1016/j.cmpbup.2021.100012>
- Ahmed, A., Hassan, A., Aziz, S., Abd-alrazaq, A. A., Ali, N., Alzubaidi, M., Al-Thani, D., Elhusein, B., Siddig, M. A., Ahmed, M., & Househ, M. (2023). Chatbot features for anxiety and depression: A scoping review. *Health Informatics Journal*, 29(1). <https://doi.org/10.1177/14604582221146719>

-
- American Psychiatric Association. (2013, May 22). *Diagnostic and statistical manual of mental disorders* (Fifth Edition). <https://doi.org/10.1176/appi.books.9780890425596>
- Andone, I., Blaszkiewicz, K., Eibes, M., Trendafilov, B., Montag, C., & Markowitz, A. (2016). How age and gender affect smartphone usage. *Proceedings of the 2016 ACM international joint conference on pervasive and ubiquitous computing: adjunct*. <https://doi.org/10.1145/2968219.2971451>
- Angst, C., Dennis, A. R., Karahanna, E., & Leroy, G. (2024). Special section: Information technology to improve mental health. *Journal of Management Information Systems*, 41(4). <https://doi.org/10.1080/07421222.2024.2415769>
- Auerbach, R. P., Webb, C. A., & Stewart, J. G. (2016). *Cognitive behavior therapy for depressed adolescents: A practical guide to management and treatment*. Routledge.
- Avenevoli, S., Swendsen, J., He, J.-P., Burstein, M., & Merikangas, K. R. (2015). Major depression in the national comorbidity survey-adolescent supplement: Prevalence, correlates, and treatment. *Journal of the American Academy of Child and Adolescent Psychiatry*, 54(1). <https://doi.org/10.1016/j.jaac.2014.10.010>
- Bae Brandtzæg, P. B., Skjuve, M., Kristoffer Dysthe, K. K., & Følstad, A. (2021). When the social becomes non-human: Young people's perception of social support in chatbots. *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. <https://doi.org/10.1145/3411764.3445318>
- Baird, A., Xia, Y., & Kohli, R. (2025). Health analytics and IS theorizing. *Journal of the Association for Information Systems*, 26(3). <https://doi.org/10.17705/1jais.00945>
- Baskerville, R., Baiyere, A., Gregor, S., Hevner, A., & Rossi, M. (2018). Design science research contributions: Finding a balance between artifact and theory. *Journal of the Association for Information Systems*, 19(5). <https://doi.org/10.17705/1jais.00495>
- Beatty, C., Malik, T., Meheli, S., & Sinha, C. (2022). Evaluating the therapeutic alliance with a free-text CBT conversational agent (wysa): A mixed-methods study. *Frontiers in Digital Health*, 4. <https://doi.org/10.3389/fdgth.2022.847991>
- Bendig, E., Erb, B., Schulze-Thuesing, L., & Baumeister, H. (2019). Die nächste generation: Chatbots in der klinischen psychologie und psychotherapie zur förderung mentaler gesundheit – ein scoping-review. *Verhaltenstherapie*, 29(4). <https://doi.org/10.1159/000499492>

-
- Benlian, A. (2015). Web personalization cues and their differential effects on user assessments of website value. *Journal of Management Information Systems*, 32(1). <https://doi.org/10.1080/07421222.2015.1029394>
- 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)*. <https://doi.org/10.1109/IALP57159.2022.9961311>
- Berger, T., Boettcher, J., & Caspar, F. (2014). Internet-based guided self-help for several anxiety disorders: A randomized controlled trial comparing a tailored with a standardized disorder-specific approach. *Psychotherapy*, 51(2). <https://doi.org/10.1037/a0032527>
- Bhattacharjee, A., Zeng, Y., Xu, S. Y., Kulzhabayeva, D., Ma, M., Kornfield, R., Ahmed, S. I., Mariakakis, A., Czerwinski, M. P., Kuzminykh, A., Liut, M., & Williams, J. J. (2024). Understanding the role of large language models in personalizing and scaffolding strategies to combat academic procrastination. *Proceedings of the CHI Conference on Human Factors in Computing Systems*. <https://doi.org/10.1145/3613904.3642081>
- Birkun, A. A., & Gautam, A. (2023). Large language model (LLM)-powered chatbots fail to generate guideline-consistent content on resuscitation and may provide potentially harmful advice. *Prehospital and Disaster Medicine*, 38(6). <https://doi.org/10.1017/S1049023X23006568>
- 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. *Journal of Medical Internet Research*, 23(3). <https://doi.org/10.2196/24387>
- Botsociety. (n.d.). *Botsociety* [Conversational interface design and prototyping software]. Retrieved June 18, 2021, from <https://botsociety.io/>
- Boucher, E. M., Harake, N. R., Ward, H. E., Stoeckl, S. E., Vargas, J., Minkel, J., Parks, A. C., & Zilca, R. (2021). Artificially intelligent chatbots in digital mental health interventions: A review. *Expert Review of Medical Devices*, 18. <https://doi.org/10.1080/17434440.2021.2013200>
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., . . . Amodei, D. (2020, July

-
- 22). Language models are few-shot learners. <https://doi.org/10.48550/arXiv.2005.14165>
- Bunge, E. L., & Desage, C. (2025). A framework for evaluating mental health artificial intelligence-based conversational agents. *Journal of Technology in Behavioral Science*. <https://doi.org/10.1007/s41347-025-00519-w>
- Butryn, T., Bryant, L., Marchionni, C., & Sholevar, F. (2017). The shortage of psychiatrists and other mental health providers: Causes, current state, and potential solutions. *International Journal of Academic Medicine*, 3(1). https://doi.org/10.4103/IJAM.IJAM_49_17
- Cambridge Dictionary. (2023, August 2). *Personalization*. Retrieved August 8, 2023, from <https://dictionary.cambridge.org/dictionary/english/personalization>
- Cameron, S. K., Rodgers, J., & Dagnan, D. (2018). The relationship between the therapeutic alliance and clinical outcomes in cognitive behaviour therapy for adults with depression: A meta-analytic review. *Clinical Psychology & Psychotherapy*, 25(3). <https://doi.org/10.1002/cpp.2180>
- Campellone, T. R., Flom, M., Montgomery, R. M., Bullard, L., Pirner, M. C., Pavez, A., Morales, M., Harper, D., Oddy, C., & O'Connor, T. (2025). Safety and user experience of a generative artificial intelligence digital mental health intervention: Exploratory randomized controlled trial. *Journal of Medical Internet Research*, 27. <https://doi.org/10.2196/67365>
- Caporino, N. E., Sakolsky, D., Brodman, D. M., McGuire, J. F., Piacentini, J., Peris, T. S., Ginsburg, G. S., Walkup, J. T., Iyengar, S., Kendall, P. C., & Birmaher, B. (2017). Establishing clinical cutoffs for response and remission on the screen for child anxiety related emotional disorders (SCARED). *Journal of the American Academy of Child and Adolescent Psychiatry*, 56(8). <https://doi.org/10.1016/j.jaac.2017.05.018>
- Chan, W. W., Fitzsimmons-Craft, E. E., Smith, A. C., Firebaugh, M.-L., Fowler, L. A., DePietro, B., Topococo, 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). <https://doi.org/10.2196/28003>
- Chaturvedi, A., Aylward, B., Shah, S., Graziani, G., Zhang, J., Manuel, B., Telewa, E., Froelich, S., Baruwa, O., Kulkarni, P. P., Ξ, W., & Kunkle, S. (2023). Content recommendation systems in web-based mental health care: Real-world application and formative evaluation. *JMIR Formative Research*, 7. <https://doi.org/10.2196/38831>

-
- 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>
- Chen, W., Lu, Y., Qiu, L., & Kumar, S. (2021). Designing personalized treatment plans for breast cancer. *Information Systems Research*, 32(3). <https://doi.org/10.1287/isre.2021.1002>
- Chevance, A., Ravaud, P., Tomlinson, A., Le Berre, C., Teufer, B., Touboul, S., Fried, E. I., Gartlehner, G., Cipriani, A., & Tran, V. T. (2020). Identifying outcomes for depression that matter to patients, informal caregivers, and health-care professionals: Qualitative content analysis of a large international online survey. *The Lancet Psychiatry*, 7(8). [https://doi.org/10.1016/S2215-0366\(20\)30191-7](https://doi.org/10.1016/S2215-0366(20)30191-7)
- China, C., Hansen, L. B., Gillanders, D. T., & Benninghoven, D. (2018). Concept and validation of the german version of the cognitive fusion questionnaire (CFQ-d). *Journal of Contextual Behavioral Science*, 9. <https://doi.org/10.1016/j.jcbs.2018.06.003>
- Cho, Y. M., Rai, S., Ungar, L., Sedoc, J., & Guntuku, S. C. (2023). An integrative survey on mental health conversational agents to bridge computer science and medical perspectives. *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. <https://doi.org/10.18653/v1/2023.emnlp-main.698>
- Cicchetti, D., & Toth, S. L. (2009). A developmental psychopathology perspective on adolescent depression. In *Handbook of depression in adolescents*. Routledge/Taylor & Francis Group.
- Clayborne, Z. M., Varin, M., & Colman, I. (2019). Systematic review and meta-analysis: Adolescent depression and long-term psychosocial outcomes. *Journal of the American Academy of Child & Adolescent Psychiatry*, 58(1). <https://doi.org/10.1016/j.jaac.2018.07.896>
- Cohen, Z. D., Delgadillo, J., & DeRubeis, R. J. (2021). Personalized treatment approaches. In *Bergin and garfield's handbook of psychotherapy and behavior change: 50th anniversary edition, 7th ed.* John Wiley & Sons, Inc.
- Coyne, A. E., Constantino, M. J., & Muir, H. J. (2019). Therapist responsivity to patients' early treatment beliefs and psychotherapy process. *Psychotherapy*, 56(1). <https://doi.org/10.1037/pst0000200>
- Cuijpers, P., Karyotaki, E., Ciharova, M., Miguel, C., Noma, H., Stikkelbroek, Y., Weisz, J. R., & Furukawa, T. A. (2023a). The effects of psychological treatments of depression in

-
- children and adolescents on response, reliable change, and deterioration: A systematic review and meta-analysis. *European Child & Adolescent Psychiatry*, 32(1). <https://doi.org/10.1007/s00787-021-01884-6>
- Cuijpers, P., Karyotaki, E., Harrer, M., & Stikkelbroek, Y. (2023b). Individual behavioral activation in the treatment of depression: A meta analysis. *Psychotherapy Research*, 33(7). <https://doi.org/10.1080/10503307.2023.2197630>
- Čuš, A., Edbrooke-Childs, J., Ohmann, S., Plener, P. L., & Akkaya-Kalayci, T. (2021). “smartphone apps are cool, but do they help me?”: A qualitative interview study of adolescents’ perspectives on using smartphone interventions to manage nonsuicidal self-injury. *International Journal of Environmental Research and Public Health*, 18(6). <https://doi.org/10.3390/ijerph18063289>
- 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). <https://doi.org/10.2196/27868>
- Das Swain, V., Zhong, Q. ", Parekh, J. R., Jeon, Y., Zimmermann, R., Czerwinski, M. P., Suh, J., Mishra, V., Saha, K., & Hernandez, J. (2025). AI on my shoulder: Supporting emotional labor in front-office roles with an LLM-based empathetic coworker. *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. <https://doi.org/10.1145/3706598.3713705>
- De Corte, K., Cairns, J., & Grieve, R. (2021). Stated versus revealed preferences: An approach to reduce bias. *Health Economics*, 30(5). <https://doi.org/10.1002/hec.4246>
- Denecke, K., Schmid, N., & Nüssli, S. (2022). Implementation of cognitive behavioral therapy in e-mental health apps: Literature review. *Journal of Medical Internet Research*, 24(3). <https://doi.org/10.2196/27791>
- Diederich, S., Brendel, A., Morana, S., & Kolbe, L. (2022). On the design of and interaction with conversational agents: An organizing and assessing review of human-computer interaction research. *Journal of the Association for Information Systems*, 23(1). <https://doi.org/10.17705/1jais.00724>
- Dimidjian, S., Hubley, A., Martell, C., Herman-Dunn, A., & Dobson, K. (2012). The quality of behavioral activation scale (q-BAS). *Boulder: University of Colorado*.
- Dimidjian, S., Goodman, S. H., Sherwood, N. E., Simon, G. E., Ludman, E., Gallop, R., Welch, S. S., Boggs, J. M., Metcalf, C. A., Hubley, S., Powers, J. D., & Beck, A. (2017).

-
- A pragmatic randomized clinical trial of behavioral activation for depressed pregnant women. *Journal of consulting and clinical psychology*, 85(1). <https://doi.org/10.1037/cp0000151>
- Ding, H., Simmich, J., Vaezipour, A., Andrews, N., & Russell, T. (2023). Evaluation framework for conversational agents with artificial intelligence in health interventions: A systematic scoping review. *Journal of the American Medical Informatics Association : JAMIA*, 31(3). <https://doi.org/10.1093/jamia/ocad222>
- Domhan, D., In-Albon, T., & Pfeiffer, S. (2023). Erfassung von Barrieren und Faszilitatoren zur Aufnahme einer Psychotherapie im Kontext ambulanter Kinder- und Jugendlichenpsychotherapie. *Die Psychotherapie*. <https://doi.org/10.1007/s00278-023-00679-1>
- Ebert, D. D., Mortier, P., Kaehlke, F., Bruffaerts, R., Baumeister, H., Auerbach, R. P., Alonso, J., Vilagut, G., Martínez, K. U., & Lochner, C. (2019). Barriers of mental health treatment utilization among first-year college students: First cross-national results from the WHO World Mental Health International College Student Initiative. *International Journal of Methods in Psychiatric Research*, 28(2). <https://doi.org/10.1002/mpr.1782>
- Ebrahimi, S., Ghasemaghaei, M., & Benbasat, I. (2022). The impact of trust and recommendation quality on adopting interactive and non-interactive recommendation agents: A meta-analysis. *Journal of Management Information Systems*, 39(3). <https://doi.org/10.1080/07421222.2022.2096549>
- Eisenberg, D., Golberstein, E., & Gollust, S. E. (2007). Help-seeking and access to mental health care in a university student population. *Medical Care*, 45(7). <https://doi.org/10.1097/MLR.0b013e31803bb4c1>
- 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). <https://doi.org/10.4081/ripppo.2024.751>
- European Union. (2023). *Overview - children and youth - eurostat*. Retrieved August 8, 2023, from <https://ec.europa.eu/eurostat/web/children-youth>
- Fan, H., & Poole, M. S. (2006). What is personalization? Perspectives on the design and implementation of personalization in information systems. *Journal of Organizational Computing and Electronic Commerce*, 16(3). <https://doi.org/10.1080/10919392.2006.9681199>

-
- Feine, J., Gnewuch, U., Morana, S., & Maedche, A. (2019). A taxonomy of social cues for conversational agents. *International Journal of Human-Computer Studies*, 132. <https://doi.org/10.1016/j.ijhcs.2019.07.009>
- Fenzl, T., & Mayring, P. (2017). QCAmapp: Eine interaktive webapplikation für qualitative inhaltsanalyse. *ZPID (Leibniz Institute for Psychology Information)*. <https://doi.org/10.23668/psycharchives.11259>
- 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). <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. *Psychotherapy: Theory, Research, Practice, Training*, 55(4). <https://doi.org/10.1037/pst0000172>
- Følstad, A., & Brandtzæg, P. B. (2017). Chatbots and the new world of HCI. *Interactions*, 24(4). <https://doi.org/10.1145/3085558>
- 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. <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). [https://doi.org/10.1016/S2215-0366\(21\)00395-3](https://doi.org/10.1016/S2215-0366(21)00395-3)
- Germonprez, M., Hovorka, D., & Collopy, F. (2007). A theory of tailorable technology design. *Journal of the Association for Information Systems*, 8(6). <https://doi.org/10.17705/1jais.00131>
- Gilbert, S., Harvey, H., Melvin, T., Vollebregt, E., & Wicks, P. (2023). Large language model AI chatbots require approval as medical devices. *Nature Medicine*. <https://doi.org/10.1038/s41591-023-02412-6>
- Gnewuch, U., Yu, M., & Maedche, A. (2020). The effect of perceived similarity in dominance on customer self-disclosure to chatbots in conversational commerce. Retrieved November 10, 2025, from <https://www.semanticscholar.org/paper/The-Effect-of-Perceived->

Similarity-in-Dominance-on-Gnewuch-Yu/7dedc14e2c75646899022b08ce4ed7d895bf7af0

- Gollwitzer, P. M., & Sheeran, P. (2006). Implementation intentions and goal achievement: A meta-analysis of effects and processes. *Advances in experimental social psychology*, 38. [https://doi.org/10.1016/s0065-2601\(06\)38002-1](https://doi.org/10.1016/s0065-2601(06)38002-1)
- Grant, M., Salsman, N. L., & Berking, M. (2018). The assessment of successful emotion regulation skills use: Development and validation of an english version of the emotion regulation skills questionnaire. *PLOS ONE*, 13(10). <https://doi.org/10.1371/journal.pone.0205095>
- Greiner, B. (2015). Subject pool recruitment procedures: Organizing experiments with ORSEE. *Journal of the Economic Science Association*, 1(1). <https://doi.org/10.1007/s40881-015-0004-4>
- Groen, G., & Petermann, F. (2015). *Therapie tools depression im kindes-und jugendalter*. Weinheim; Beltz.
- Grové, C. (2021). Co-developing a mental health and wellbeing chatbot with and for young people. *Frontiers in Psychiatry*, 11. <https://doi.org/10.3389/fpsyt.2020.606041>
- Guest, G., Bunce, A., & Johnson, L. (2006). How many interviews are enough? an experiment with data saturation and variability. *Field Methods*, 18(1). <https://doi.org/10.1177/1525822X05279903>
- Gulliver, A., Griffiths, K. M., & Christensen, H. (2010). Perceived barriers and facilitators to mental health help-seeking in young people: A systematic review. *BMC Psychiatry*, 10(1). <https://doi.org/10.1186/1471-244X-10-113>
- 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>
- Harte, R., Glynn, L., Rodríguez-Molinero, A., Baker, P. M., Scharf, T., Quinlan, L. R., & ÓLaighin, G. (2017). A human-centered design methodology to enhance the usability, human factors, and user experience of connected health systems: A three-phase methodology. *JMIR Human Factors*, 4(1). <https://doi.org/10.2196/humanfactors.5443>
- Hatch, S. G., Goodman, Z. T., Vowels, L., Hatch, H. D., Brown, A. L., Guttman, S., Le, Y., Bailey, B., Bailey, R. J., Esplin, C. R., Harris, S. M., Jr, D. P. H., McLaughlin, M., O'Connell, P., Rothman, K., Ritchie, L., Jr, D. N. T., & Braithwaite, S. R. (2025). When

-
- ELIZA meets therapists: A turing test for the heart and mind. *PLOS Mental Health*, 2(2). <https://doi.org/10.1371/journal.pmen.0000145>
- 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). <https://doi.org/10.1080/10503307.2019.1620366>
- Heinz, M. V., Mackin, D. M., Trudeau, B. M., Bhattacharya, S., Wang, Y., Banta, H. A., Jewett, A. D., Salzhauer, A. J., Griffin, T. Z., & Jacobson, N. C. (2025). Randomized trial of a generative AI chatbot for mental health treatment. *NEJM AI*, 2(4). <https://doi.org/10.1056/AIoa2400802>
- Heston, T. F. (2023). Safety of large language models in addressing depression. *Cureus*, 15. <https://doi.org/10.7759/cureus.50729>
- Hevner, A. R., March, S. T., Park, J., & Ram, S. (2004). Design science in information systems research. *MIS quarterly*. <https://doi.org/10.2307/25148625>
- Høiland, C. G., Følstad, A., & Karahasanovic, A. (2020). Hi, can i help? exploring how to design a mental health chatbot for youths. *Human Technology*, 16(2). <https://doi.org/10.17011/ht/urn.202008245640>
- Hua, Y., Na, H., Li, Z., Liu, F., Fang, X., Clifton, D., & Torous, J. (2025). A scoping review of large language models for generative tasks in mental health care. *npj Digital Medicine*, 8(1). <https://doi.org/10.1038/s41746-025-01611-4>
- Hua, Z., Wang, S., & Yuan, X. (2024). Trends in age-standardized incidence rates of depression in adolescents aged 10-24 in 204 countries and regions from 1990 to 2019. *Journal of Affective Disorders*, 350. <https://doi.org/10.1016/j.jad.2024.01.009>
- Huffman, S. (2014). OMG! mobile voice survey reveals teens love to talk. Retrieved October, 14. <https://blog.google/products/search/omg-mobile-voice-survey-reveals-teens/>
- Huguet, A., Rao, S., McGrath, P. J., Wozney, L., Wheaton, M., Conrod, J., & Rozario, S. (2016). A systematic review of cognitive behavioral therapy and behavioral activation apps for depression. *PLOS ONE*, 11(5). <https://doi.org/10.1371/journal.pone.0154248>
- Huisman, P., & Kangas, M. (2018). Evidence-based practices in cognitive behaviour therapy (CBT) case formulation: What do practitioners believe is important, and what do they do? *Behaviour Change*, 35(1). <https://doi.org/10.1017/bec.2018.5>

-
- Inkster, B., Sarda, S., & Subramanian, V. (2018). An empathy-driven, conversational artificial intelligence agent (wysa) for digital mental well-being: Real-world data evaluation mixed-methods study. *JMIR mHealth and uHealth*, 6(11). <https://doi.org/10.2196/12106>
- Jabir, A. I., Lin, X., Martinengo, L., Sharp, G., Theng, Y.-L., & Car, L. T. (2024). Attrition in conversational agent-delivered mental health interventions: Systematic review and meta-analysis. *Journal of Medical Internet Research*, 26(1). <https://doi.org/10.2196/48168>
- Jardine, J., Nadal, C., Robinson, S., Enrique, A., Hanratty, M., & Doherty, G. (2024). Between rhetoric and reality: Real-world barriers to uptake and early engagement in digital mental health interventions. *ACM Trans. Comput.-Hum. Interact.*, 31(2). <https://doi.org/10.1145/3635472>
- Jaspers, M. W. M., Steen, T., Bos, C. v. d., & Geenen, M. (2004). The think aloud method: A guide to user interface design. *International Journal of Medical Informatics*, 73(11). <https://doi.org/10.1016/j.ijmedinf.2004.08.003>
- Jia, E., Macon, J., Doering, M., & Abraham, J. (2025). Effectiveness of digital behavioral activation interventions for depression and anxiety: Systematic review and meta-analysis. *J Med Internet Res* 2025;27:e68054 <https://www.jmir.org/2025/1/e68054>. <https://doi.org/10.2196/68054>
- Jiang, H., Zhang, X., Cao, X., Breazeal, C., Roy, D., & Kabbara, J. (2024, June). PersonaLLM: Investigating the ability of large language models to express personality traits. In K. Duh, H. Gomez, & S. Bethard (Eds.), *Findings of the association for computational linguistics: NAACL 2024*. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.findings-naacl.229>
- Johansson, R., Sjöberg, E., Sjögren, M., Johnsson, E., Carlbring, P., Andersson, T., Rousseau, A., & Andersson, G. (2012). Tailored vs. standardized internet-based cognitive behavior therapy for depression and comorbid symptoms: A randomized controlled trial. *PLOS ONE*, 7(5). <https://doi.org/10.1371/journal.pone.0036905>
- Kambeitz-Ilankovic, L., Rzayeva, U., Völkel, L., Wenzel, J., Weiske, J., Jessen, F., Reininghaus, U., Uhlhaas, P. J., Alvarez-Jimenez, M., & Kambeitz, J. (2022). A systematic review of digital and face-to-face cognitive behavioral therapy for depression. *npj Digital Medicine*, 5(1). <https://doi.org/10.1038/s41746-022-00677-8>

-
- Kankanhalli, A., Xia, Q., & Zhao, X. (2021). Understanding personalization for health behavior change applications: A review and future directions. *AIS Transactions on Human-Computer Interaction*. <https://doi.org/10.17705/1thci.00152>
- Kapania, S., Agnew, W., Eslami, M., Heidari, H., & Fox, S. E. (2025). Simulacrum of stories: Examining large language models as qualitative research participants. *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. <https://doi.org/10.1145/3706598.3713220>
- Kaufman, J., Birmaher, B., Brent, D., Rao, U., Flynn, C., Moreci, P., Williamson, D., & Ryan, N. (1997). Schedule for affective disorders and schizophrenia for school-age children-present and lifetime version (k-SADS-PL): Initial reliability and validity data. *Journal of the American Academy of Child & Adolescent Psychiatry*, *36*(7). <https://doi.org/10.1097/00004583-199707000-00021>
- Kaveladze, B., Voelkel, J., Stagnaro, M., Huang, M., Smock, A., Sullivan, E., Xu, Y., McCall, M., Zapata, J. P., Ishtiaque, S., Bhattacharjee, A., Georgieva, I., Hernandez-Ramos, R., Huber, K., Jennings, J., Kirk, A., Kornfield, R., Knowles, R., Lind, M., ... Schleider, J. (2025, August 25). A crowdsourced megastudy of 12 digital single-session interventions for depression in american adults. <https://doi.org/10.21203/rs.3.rs-7236847/v1>
- Kenny, R., Dooley, B., & Fitzgerald, A. (2016). Developing mental health mobile apps: Exploring adolescents' perspectives. *Health Informatics Journal*, *22*(2). <https://doi.org/10.1177/1460458214555041>
- Kim, J., Podlasek, A., Shidara, K., Liu, F., Alaa, A., & Bernardo, D. (2025, February 5). Limitations of large language models in clinical problem-solving arising from inflexible reasoning. <https://doi.org/10.48550/arXiv.2502.04381>
- 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. *Journal of Medical Internet Research*, *21*(11). <https://doi.org/10.2196/15360>
- Komiak, S. Y., & Benbasat, I. (2006). The effects of personalization and familiarity on trust and adoption of recommendation agents. *MIS quarterly*, *30*(4). <https://doi.org/10.2307/25148760>

-
- Kroenke, K., Strine, T. W., Spitzer, R. L., Williams, J. B., Berry, J. T., & Mokdad, A. H. (2009). The PHQ-8 as a measure of current depression in the general population. *Journal of affective disorders*, *114*(1). <https://doi.org/10.1016/j.jad.2008.06.026>
- Krupnik, V. (2022). The therapeutic alliance as active inference: The role of trust and self-efficacy. *Journal of Contemporary Psychotherapy*, *53*(3). <https://doi.org/10.1007/s10879-022-09576-1>
- Kuhlmeier, F. O., Bauch, L., Gnewuch, U., & Lüttke, S. (2025a). Designing chatbots to treat depression in youth: Qualitative study. *JMIR Human Factors*, *12*(1). <https://doi.org/10.2196/66632>
- Kuhlmeier, F. O., Gnewuch, U., Lüttke, S., Brakemeier, E.-L., & Mädche, A. (2022a). A personalized conversational agent to treat depression in youth and young adults – a transdisciplinary design science research project. In A. Drechsler, A. Gerber, & A. Hevner (Eds.), *The transdisciplinary reach of design science research*. Springer International Publishing. https://doi.org/10.1007/978-3-031-06516-3_3
- Kuhlmeier, F. O., Gnewuch, U., Metelmann, L., Bauch, L., & Lüttke, S. (2022b). Exploring user experience with a conversational agent to treat depression in youth: A think-aloud study. *SIGHCI 2022 Proceedings*. <https://aisel.aisnet.org/sighci2022/10>
- Kuhlmeier, F. O., Gnewuch, U., Scheu, S., Lüttke, S., & Maedche, A. (2025b). *User-led, chatbot-led, or hybrid? design and effect of content personalization approaches in mental health chatbots for young people* [Manuscript under review].
- Kuhlmeier, F. O., Hanschmann, L., Rabe, M., Luettker, S., Brakemeier, E.-L., & Maedche, A. (2026, March 6). A large language model-based behavioral activation chatbot for young people with depression: Mixed-methods evaluation using artificial users and clinical experts [JMIR preprints]. <https://doi.org/10.2196/preprints.94781>
- 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. <https://doi.org/10.48550/arXiv.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*. <https://doi.org/10.1145/3544549.3585614>

-
- Laranjo, L., Dunn, A. G., Tong, H. L., Kocaballi, A. B., Chen, J., Bashir, R., Surian, D., Gallego, B., Magrabi, F., Lau, A. Y. S., & Coiera, E. (2018). Conversational agents in healthcare: A systematic review. *Journal of the American Medical Informatics Association: JAMIA*, 25(9). <https://doi.org/10.1093/jamia/ocy072>
- Leech, T., Dorstyn, D., Taylor, A., & Li, W. (2021). Mental health apps for adolescents and young adults: A systematic review of randomised controlled trials. *Children and Youth Services Review*, 127. <https://doi.org/10.1016/j.chilyouth.2021.106073>
- Lejuez, C. W., Hopko, D. R., & Hopko, S. D. (2001). A brief behavioral activation treatment for depression: Treatment manual. *Behavior Modification*, 25(2). <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). <https://doi.org/10.1177/0145445510390929>
- Levin, M. E., Haeger, J., & Cruz, R. A. (2019). Tailoring acceptance and commitment therapy skill coaching in the moment through smartphones: Results from a randomized controlled trial. *Mindfulness*, 10(4). <https://doi.org/10.1007/s12671-018-1004-2>
- Li, S. H., Achilles, M. R., Spanos, S., Habak, S., Werner-Seidler, A., & O'Dea, B. (2022). A cognitive behavioural therapy smartphone app for adolescent depression and anxiety: Co-design of ClearlyMe. *the Cognitive Behaviour Therapist*, 15. <https://doi.org/10.1017/S1754470X22000095>
- Li, T., & Unger, T. (2012). Willing to pay for quality personalization? trade-off between quality and privacy. *European Journal of Information Systems*, 21(6). <https://doi.org/10.1057/ejis.2012.13>
- Liang, T.-P., Lai, H.-J., & Ku, Y.-C. (2006). Personalized content recommendation and user satisfaction: Theoretical synthesis and empirical findings. *Journal of Management Information Systems*, 23(3). <https://doi.org/10.2753/mis0742-1222230303>
- 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). <https://doi.org/10.1016/j.beth.2021.09.007>
- Limpanopparat, S., Gibson, E., & Harris, D. A. (2024). User engagement, attitudes, and the effectiveness of chatbots as a mental health intervention: A systematic review. *Computers*

-
- in Human Behavior: Artificial Humans*, 2(2). <https://doi.org/10.1016/j.chbah.2024.100081>
- Lin, X., Martinengo, L., Jabir, A. I., Ho, A. H. Y., Car, J., Atun, R., & Car, L. T. (2023). Scope, characteristics, behavior change techniques, and quality of conversational agents for mental health and well-being: Systematic assessment of apps. *Journal of Medical Internet Research*, 25(1). <https://doi.org/10.2196/45984>
- Linardon, J., Messer, M., Shatte, A., Greenwood, C. J., Rosato, J., Rathgen, A., Skvarc, D., & Fuller-Tyszkiewicz, M. (2022). Does the method of content delivery matter? randomized controlled comparison of an internet-based intervention for eating disorder symptoms with and without interactive functionality. *Behavior Therapy*, 53(3). <https://doi.org/10.1016/j.beth.2021.12.001>
- Linardon, J., Torous, J., Firth, J., Cuijpers, P., Messer, M., & Fuller-Tyszkiewicz, M. (2024). Current evidence on the efficacy of mental health smartphone apps for symptoms of depression and anxiety. a meta-analysis of 176 randomized controlled trials. *World Psychiatry*, 23(1). <https://doi.org/10.1002/wps.21183>
- 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). <https://doi.org/10.1080/10447318.2023.2300015>
- Louie, R., Nandi, A., Fang, W., Chang, C., Brunskill, E., & Yang, D. (2024, July 14). Roleplaydoh: Enabling domain-experts to create LLM-simulated patients via eliciting and adhering to principles. <https://doi.org/10.48550/arXiv.2407.00870>
- Lucas, G. M., Rizzo, A., Gratch, J., Scherer, S., Stratou, G., Boberg, J., & Morency, L.-P. (2017). Reporting mental health symptoms: Breaking down barriers to care with virtual human interviewers. *Frontiers in Robotics and AI*, 4. <https://doi.org/10.3389/frobt.2017.00051>
- Ludlow, K., Russell, J. K., Ryan, B., Brown, R. L., Joynt, T., Uhlmann, L. R., Smith, G. E., Donovan, C., Hides, L., Spence, S. H., March, S., & Cobham, V. E. (2023). Co-designing a digital mental health platform, “momentum”, with young people aged 7–17: A qualitative study. *Digital Health*, 9. <https://doi.org/10.1177/20552076231216410>

-
- Luxton, R., & Kyriakopoulos, M. (2022). Depression in children and young people: Identification and management NICE guidelines. *Archives of Disease in Childhood - Education and Practice*, 107(1). <https://doi.org/10.1136/archdischild-2020-320020>
- 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*. Retrieved February 5, 2024, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10785945/>
- Mancinelli, E., Magnolini, S., Gabrielli, S., & Salcuni, S. (2024). A chatbot (juno) prototype to deploy a behavioral activation intervention to pregnant women: Qualitative evaluation using a multiple case study. *JMIR Formative Research*, 8(1). <https://doi.org/10.2196/58653>
- Manos, R. C., Kanter, J. W., & Luo, W. (2011). The behavioral activation for depression scale—short form: Development and validation. *Behavior Therapy*, 42(4). <https://doi.org/10.1016/j.beth.2011.04.004>
- Martin, A., Rief, W., Klaiberg, A., & Braehler, E. (2006). Validity of the brief patient health questionnaire mood scale (PHQ-9) in the general population. *General hospital psychiatry*, 28(1). <https://doi.org/10.1016/j.genhosppsych.2005.07.003>
- Martinengo, L., Lum, E., & Car, J. (2022). Evaluation of chatbot-delivered interventions for self-management of depression: Content analysis. *Journal of Affective Disorders*. <https://doi.org/10.1016/j.jad.2022.09.028>
- Martinengo, L., Stona, A.-C., Griva, K., Dazzan, P., Pariante, C. M., Wangenheim, F. v., & Car, J. (2021). Self-guided cognitive behavioral therapy apps for depression: Systematic assessment of features, functionality, and congruence with evidence. *Journal of Medical Internet Research*, 23(7). <https://doi.org/10.2196/27619>
- Mayring, P. (2004). Qualitative content analysis. *A companion to qualitative research*, 1(2).
- Mayring, P. (2015). *Qualitative inhaltsanalyse: Grundlagen und techniken* (12., überarbeitete Auflage). Beltz.
- Mayring, P., & Fenzl, T. (2019). Qualitative inhaltsanalyse. In N. Baur & J. Blasius (Eds.), *Handbuch methoden der empirischen sozialforschung*. Springer Fachmedien. https://doi.org/10.1007/978-3-658-21308-4_42

-
- Meeks, S., Haitsma, K. V., & Shryock, S. K. (2019). Treatment fidelity evidence for BE-ACTIV – a behavioral intervention for depression in nursing homes. *Aging & mental health*, 23(9). <https://doi.org/10.1080/13607863.2018.1484888>
- 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. *Journal of Medical Internet Research*, 23(6). <https://doi.org/10.2196/26771>
- Meyer, S., & Elsweler, D. (2025). LLM-based conversational agents for behaviour change support: A randomised controlled trial examining efficacy, safety, and the role of user behaviour. *International Journal of Human-Computer Studies*, 200. <https://doi.org/10.1016/j.ijhcs.2025.103514>
- Miah, S., Gammack, J., & McKay, J. (2019). A metadesign theory for tailorable decision support. *Journal of the Association for Information Systems*, 20(5). <https://doi.org/10.17705/1jais.00544>
- Mihailescu, I., Efrim-Budisteanu, M., Andrei, L. E., Buică, A., Moise, M., Nicolau, I. G., Iotu, A. D., Grădilă, A., Costea, T., Prișceanu, A. M., & Rad, F. (2023). Cognitive coping strategies among inpatient adolescents with depression and psychiatric comorbidity. *Children*, 10. <https://doi.org/10.3390/children10121870>
- Mullarkey, M. C., Marchetti, I., & Beevers, C. G. (2019). Using network analysis to identify central symptoms of adolescent depression. *Journal of Clinical Child & Adolescent Psychology*, 48(4). <https://doi.org/10.1080/15374416.2018.1437735>
- Nißen, M., Rügger, D., Stieger, M., Flückiger, C., Allemann, M., Wangenheim, F. v., & Kowatsch, T. (2022). The effects of health care chatbot personas with different social roles on the client-chatbot bond and usage intentions: Development of a design codebook and web-based study. *Journal of Medical Internet Research*, 24(4). <https://doi.org/10.2196/32630>
- Nye, A., Delgadillo, J., & Barkham, M. (2023). Efficacy of personalized psychological interventions: A systematic review and meta-analysis. *Journal of Consulting and Clinical Psychology*, 91(7). <https://doi.org/10.1037/ccp0000820>
- 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. <https://doi.org/10.1016/j.eurpsy.2018.12.008>
- Paul, S. C., Bartmann, N., & Clark, J. L. (2024). Customizability in conversational agents and their impact on health engagement (stage 2). *Human Behavior and Emerging Technologies*, 2024(1). <https://doi.org/10.1155/2024/5015913>
- Peffer, K., Tuunanen, T., Rothenberger, M. A., & Chatterjee, S. (2007). A design science research methodology for information systems research. *Journal of Management Information Systems*, 24(3). <https://doi.org/10.2753/MIS0742-1222240302>
- Peng, F., Zhang, D., & Yan, Z. (2024). Digital phenotyping-based depression detection in the presence of comorbidity: An uncertainty reasoning approach. *Journal of Management Information Systems*, 41(4). <https://doi.org/10.1080/07421222.2024.2415770>
- Perski, O., Crane, D., Beard, E., & Brown, J. (2019). Does the addition of a supportive chatbot promote user engagement with a smoking cessation app? an experimental study. *Digital Health*, 5. <https://doi.org/10.1177/2055207619880676>
- Pickard, M. D., Roster, C. A., & Chen, Y. (2016). Revealing sensitive information in personal interviews: Is self-disclosure easier with humans or avatars and under what conditions? *Computers in Human Behavior*, 65. <https://doi.org/10.1016/j.chb.2016.08.004>
- 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>
- Qiu, P., Wu, C., Liu, S., Fan, Y., Zhao, W., Chen, Z., Gu, H., Peng, C., Zhang, Y., Wang, Y., & Xie, W. (2025). Quantifying the reasoning abilities of LLMs on clinical cases. *Nature Communications*, 16(1). <https://doi.org/10.1038/s41467-025-64769-1>
- Radez, J., Reardon, T., Creswell, C., Lawrence, P. J., Evdoka-Burton, G., & Waite, P. (2020). Why do children and adolescents (not) seek and access professional help for their mental health problems? a systematic review of quantitative and qualitative studies. *European Child & Adolescent Psychiatry*. <https://doi.org/10.1007/s00787-019-01469-4>
- Radez, J., Reardon, T., Creswell, C., Orchard, F., & Waite, P. (2022). Adolescents' perceived barriers and facilitators to seeking and accessing professional help for anxiety and depressive disorders: A qualitative interview study. *European child & adolescent psychiatry*, 31(6). <https://doi.org/10.1007/s00787-020-01707-0>

-
- Radloff, L. S. (1977). The CES-d scale: A self-report depression scale for research in the general population. *Applied Psychological Measurement, 1*(3). <https://doi.org/10.1177/014662167700100306>
- Rathnayaka, P., Mills, N., Burnett, D., De Silva, D., Alahakoon, D., & Gray, R. (2022). A mental health chatbot with cognitive skills for personalised behavioural activation and remote health monitoring. *Sensors, 22*(10). <https://doi.org/10.3390/s22103653>
- Reinecke, K., & Bernstein, A. (2013). Knowing what a user likes: A design science approach to interfaces that automatically adapt to culture. *MIS Quarterly, 37*(2). <https://doi.org/10.25300/misq/2013/37.2.06>
- Reinert, M., Fritze, D., & Nguyen, T. (2022). The state of mental health in america.
- Rethorst, C. D., Trombello, J. M., Chen, P. M., Carmody, T. J., Goodman, L. C., Lazalde, A., & Trivedi, M. H. (2024). Pilot evaluation on an adapted tele-behavioral activation to increase physical activity in persons with depression: A single-arm pilot study. *BMC Psychology, 12*(1). <https://doi.org/10.1186/s40359-024-02053-5>
- Rice, F., Riglin, L., Lomax, T., Souter, E., Potter, R., Smith, D. J., Thapar, A. K., & Thapar, A. (2019). Adolescent and adult differences in major depression symptom profiles. *Journal of Affective Disorders, 243*. <https://doi.org/10.1016/j.jad.2018.09.015>
- Rock, P. L., Roiser, J. P., Riedel, W. J., & Blackwell, A. D. (2014). Cognitive impairment in depression: A systematic review and meta-analysis. *Psychological Medicine, 44*(10). <https://doi.org/10.1017/S0033291713002535>
- Ryan, R. M., & Deci, E. L. (2008). A self-determination theory approach to psychotherapy: The motivational basis for effective change. *Canadian Psychology / Psychologie canadienne, 49*(3). <https://doi.org/10.1037/a0012753>
- Schaeuffele, C., Zagorscak, P., Langerwisch, V., Wilke, J., Medvedeva, Y., & Knaevelsrud, C. (2025). A systematic review on personalization of treatment components in IBIs for mental disorders. *Internet Interventions, 41*. <https://doi.org/10.1016/j.invent.2025.100840>
- Schäfer, J. Ö., Naumann, E., Holmes, E. A., Tuschen-Caffier, B., & Samson, A. C. (2017). Emotion regulation strategies in depressive and anxiety symptoms in youth: A meta-analytic review. *Journal of Youth and Adolescence, 46*(2). <https://doi.org/10.1007/s10964-016-0585-0>

-
- Schleider, J. L., Mullarkey, M. C., Fox, K. R., Dobias, M. L., Shroff, A., Hart, E. A., & Roulston, C. A. (2022). A randomized trial of online single-session interventions for adolescent depression during COVID-19. *Nature Human Behaviour*, 6(2). <https://doi.org/10.1038/s41562-021-01235-0>
- 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). <https://doi.org/10.1192/j.eurpsy.2022.2337>
- Schueller, S. M., & Torous, J. (2020). Scaling evidence-based treatments through digital mental health. *The American psychologist*, 75(8). <https://doi.org/10.1037/amp0000654>
- 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*. <https://doi.org/10.1145/3613905.3650860>
- Seedat, S., Scott, K. M., Angermeyer, M. C., Berglund, P., Bromet, E. J., Brugha, T. S., Demyttenaere, K., de Girolamo, G., Haro, J. M., Jin, R., Karam, E. G., Kovess-Masfety, V., Levinson, D., Medina Mora, M. E., Ono, Y., Ormel, J., Pennell, B.-E., Posada-Villa, J., Sampson, N. A., . . . Kessler, R. C. (2009). Cross-national associations between gender and mental disorders in the world health organization world mental health surveys. *Archives of General Psychiatry*, 66(7). <https://doi.org/10.1001/archgenpsychiatry.2009.36>
- 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*. <https://doi.org/10.1145/3613904.3642761>
- Six, S., Aly, H., & Byrne, K. A. (2022). Investigating the effect of personalization in a mental health app on depressive symptoms. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, 66(1). <https://doi.org/10.1177/1071181322661353>
- Sjöström, J., Department of Informatics and Media, Uppsala University, Sweden, Ågerfalk, P., Department of Informatics and Media, Uppsala University, Sweden, Hevner, A. R., & Muma College of Business, University of South Florida, USA. (2022). The design of a system for online psychosocial care: Balancing privacy and accountability in sensitive

-
- online healthcare environments. *Journal of the Association for Information Systems*, 23(1). <https://doi.org/10.17705/1jais.00717>
- 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. <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. <https://doi.org/10.1016/j.ijhcs.2022.102903>
- Solmi, M., Radua, J., Olivola, M., Croce, E., Soardo, L., Salazar de Pablo, G., Il Shin, J., Kirkbride, J. B., Jones, P., Kim, J. H., Kim, J. Y., Carvalho, A. F., Seeman, M. V., Correll, C. U., & Fusar-Poli, P. (2022). Age at onset of mental disorders worldwide: Large-scale meta-analysis of 192 epidemiological studies. *Molecular Psychiatry*, 27(1). <https://doi.org/10.1038/s41380-021-01161-7>
- 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. *npj Mental Health Research*, 3(1). <https://doi.org/10.1038/s44184-024-00056-z>
- Steenstra, I., Nouraei, F., & Bickmore, T. (2025). Scaffolding empathy: Training counselors with simulated patients and utterance-level performance visualizations. *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. <https://doi.org/10.1145/3706598.3714014>
- Stumpp, N. E., & Sauer-Zavala, S. (2022). Evidence-based strategies for treatment personalization: A review. *Cognitive and Behavioral Practice*, 29(4). <https://doi.org/10.1016/j.cbpra.2021.10.004>
- Sweller, J. (1988). Cognitive load during problem solving: Effects on learning. *Cognitive science*, 12(2). https://doi.org/10.1207/s15516709cog1202_4
- Swift, J. K., Callahan, J. L., Cooper, M., & Parkin, S. R. (2018). The impact of accommodating client preference in psychotherapy: A meta-analysis. *Journal of Clinical Psychology*, 74(11). <https://doi.org/10.1002/jclp.22680>

-
- Tam, K. Y., & Ho, S. Y. (2005). Web personalization as a persuasion strategy: An elaboration likelihood model perspective. *Information Systems Research*, *16*(3). <https://doi.org/10.1287/isre.1050.0058>
- Thapar, A., Collishaw, S., Pine, D. S., & Thapar, A. K. (2012). Depression in adolescence. *Lancet*, *379*(9820). [https://doi.org/10.1016/S0140-6736\(11\)60871-4](https://doi.org/10.1016/S0140-6736(11)60871-4)
- Thieme, A., Belgrave, D., & Doherty, G. (2020). Machine learning in mental health: A systematic review of the HCI literature to support the development of effective and implementable ML systems. *ACM Trans. Comput.-Hum. Interact.*, *27*(5). <https://doi.org/10.1145/3398069>
- Thieme, A., Hanratty, M., Lyons, M., Palacios, J., Marques, R. F., Morrison, C., & Doherty, G. (2023). Designing human-centered AI for mental health: Developing clinically relevant applications for online CBT treatment. *ACM Trans. Comput.-Hum. Interact.*, *30*(2). <https://doi.org/10.1145/3564752>
- Thirumalai, S., & Sinha, K. K. (2013). To personalize or not to personalize online purchase interactions: Implications of self-selection by retailers. *Information Systems Research*, *24*(3). <https://doi.org/10.1287/isre.1120.0471>
- 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. *World Psychiatry*, *20*(3). <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*. Jacob Towery.
- Vaidyam, A. N., Linggonegoro, D., & Torous, J. (2020). Changes to the psychiatric chatbot landscape: A systematic review of conversational agents in serious mental illness: Changements du paysage psychiatrique des chatbots: Une revue systématique des agents conversationnels dans la maladie mentale sérieuse. *The Canadian Journal of Psychiatry*. <https://doi.org/10.1177/0706743720966429>
- 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). <https://doi.org/10.1177/0706743719828977>

-
- 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, W., & Benbasat, I. (2009). Interactive decision aids for consumer decision making in e-commerce: The influence of perceived strategy restrictiveness. *MIS Quarterly*, 33(2). <https://doi.org/10.2307/20650293>
- Wang, W., & Benbasat, I. (2013). Research note — a contingency approach to investigating the effects of user-system interaction modes of online decision aids. *Information Systems Research*, 24(3). <https://doi.org/10.1287/isre.1120.0445>
- Wang, Y., Wang, Y., Xiao, Y., Escamilla, L., Augustine, B., Crace, K., Zhou, G., & Zhang, Y. (2025, January 26). Evaluating an LLM-powered chatbot for cognitive restructuring: Insights from mental health professionals. <https://doi.org/10.48550/arXiv.2501.15599>
- Weitkamp, K., Romer, G., Rosenthal, S., Wiegand-Grefe, S., & Daniels, J. (2010). German screen for child anxiety related emotional disorders (SCARED): Reliability, validity, and cross-informant agreement in a clinical sample. *Child and Adolescent Psychiatry and Mental Health*, 4. <https://doi.org/10.1186/1753-2000-4-19>
- Wenzel, A. (2017). Basic strategies of cognitive behavioral therapy. *Psychiatric Clinics of North America*, 40(4). <https://doi.org/10.1016/j.psc.2017.07.001>
- Wies, B., Landers, C., & Ienca, M. (2021). Digital mental health for young people: A scoping review of ethical promises and challenges. *Frontiers in Digital Health*, 3. <https://doi.org/10.3389/fdgth.2021.697072>
- World Health Organization. (2025). *Suicide worldwide in 2021: Global health estimates*. Retrieved November 25, 2025, from <https://iris.who.int/handle/10665/381495>
- Wu, Y., Fenfen, E., Wang, Y., Xu, M., Liu, S., Zhou, L., Song, G., Shang, X., Yang, C., Yang, K., & Li, X. (2023). Efficacy of internet-based cognitive-behavioral therapy for depression in adolescents: A systematic review and meta-analysis. *Internet Interventions*, 34. <https://doi.org/10.1016/j.invent.2023.100673>
- Xiao, B., & Benbasat, I. (2007). E-commerce product recommendation agents: Use, characteristics, and impact. *MIS Quarterly*, 31(1). <https://doi.org/10.2307/25148784>
- Yardley, L., Morrison, L., Bradbury, K., & Muller, I. (2015). The person-based approach to intervention development: Application to digital health-related behavior change inter-

ventions. *Journal of Medical Internet Research*, 17(1). <https://doi.org/10.2196/jmir.4055>

Zehetmair, C., Nagy, E., Leetz, C., Cranz, A., Kindermann, D., Reddemann, L., & Nikendei, C. (2020). Self-practice of stabilizing and guided imagery techniques for traumatized refugees via digital audio files: Qualitative study. *Journal of Medical Internet Research*, 22(9). <https://doi.org/10.2196/17906>

Zhou, T., Wang, Y., Yan, L. (, & Tan, Y. (2023). Spoiled for choice? personalized recommendation for healthcare decisions: A multiarmed bandit approach. *Information Systems Research*. <https://doi.org/10.1287/isre.2022.1191>

List of Publications

Accepted, Peer-Reviewed Publications

Kuhlmeier, Florian Onur, Gnewuch, Ulrich, Lüttke, Stefan, Brakemeier, Eva-Lotta, Maedche, Alexander (2022). A Personalized Conversational Agent to Treat Depression in Youth and Young Adults – A Transdisciplinary Design Science Research Project. In: Drechsler, A., Gerber, A., Hevner, A. (eds) *The Transdisciplinary Reach of Design Science Research*. DESRIST 2022. Lecture Notes in Computer Science, vol 13229. Springer, Cham.

Kuhlmeier, Florian Onur, Gnewuch, Ulrich, Metelmann, Lilli, Bauch, Luise, Lüttke, Stefan (2022). Exploring User Experience with a Conversational Agent to Treat Depression in Youth: A Think-Aloud Study. *SIGHCI 2022 Proceedings*. 10.

Kuhlmeier, Florian Onur, Bauch, Luise, Gnewuch, Ulrich, Lüttke, Stefan (2025). Designing Chatbots to Treat Depression in Youth: Qualitative Study. *JMIR Human Factors*, 12(1), e66632.

Scheu, Sven, **Kuhlmeier, Florian Onur**, Müller, Lisa, Benke, Ivo (2023). Designing AMSL – a Mobile Digital Assistant for Self-Regulated Learning. *ECIS 2023 Research-in-Progress Papers*. 73.

Sou, Davinny, **Kuhlmeier, Florian Onur**, Kowatsch, Tobias, von Wangenheim, Florian, Nißen, Marcia (2024). Towards Digital Empathy in Healthcare Chatbots: A Conceptual Framework and Empirical Study. *ICIS 2024 Proceedings*. 5.

Publications Under Review

Kuhlmeier, Florian Onur, Hanschmann, L., Rabe, M., Lüttke, Stefan, Brakemeier, Eva-Lotta, Maedche, Alexander (2026). A Large Language Model-Based Behavioral Activation Chatbot for Young People with Depression: Mixed-Methods Evaluation Using Artificial Users and Clinical Experts. *JMIR Preprints*. 10.2196/preprints.94781. Under review at *JMIR Mental Health*.

Kuhlmeier, Florian Onur, Gnewuch, Ulrich, Scheu, Sven, Lüttke, Stefan, Maedche, Alexander (2025). User-Led, Chatbot-Led, or Hybrid? Design and Effect of Content Personalization Approaches in Mental Health Chatbots for Young People. Under review at the *Journal of the*

Association for Information Systems (JAIS).

Appendix

A Supplementary Material for Part I

A.1 Interview Guide

Table A.1: Interview Guide for Designing Chatbots to Treat Depression in Youth

Topic	Interview Questions
Warm-Up	You've already answered a few questions. Thank you again for taking part in our study. We depend on young people like you. You are experts on your experiences and problems. That's why I want to talk to you about a few things in more detail in the interview. There is no right or wrong, I want to know what you think. If you find a question strange or don't feel you understand something, please let me know. Do you still have any questions? If not, I'll start the recording, and we'll get started with the questions. But you can always ask questions in between.
Introductory Question	Why did you decide to take part in our study? If you weren't here, what would you be doing right now?
Problems	How do you feel when you're not doing so well? What situations do you find difficult then? How do you feel then? How do you behave then? What thoughts do you have then? What stresses you out the most? What do you need help with then?

Continued on next page

Table A.1 – continued from previous page

Topic	Interview Questions
Coping Strategies	<p>What have you tried in the past to make you feel better?</p> <p>What has helped you and what hasn't?</p> <p>Where have you gone for help/information?</p> <p>What stopped you from getting help?</p> <p>How did you end up getting help after all?</p> <p>How could the help have been even better for you personally?</p> <p>What would you hope for/not want if you were to get help?</p> <p>How do you know that something has helped you?</p> <p>How do you know that you are no longer depressed?</p>
Attitudes and Expectations	<p>What do you imagine when you think of psychotherapy?</p> <p>Have you had any experience with psychotherapy? (Explanation if necessary: In psychotherapy, the therapist helps you with your problems. You work together to find solutions. First, you get to know each other, then try strategies and tools together.)</p> <p>Do you know what a chatbot is?</p> <p>Imagine Siri or Alexa, but also writing to you. Imagine you are writing to someone on WhatsApp, but to a computer/robot.</p> <p>What do you think it would be like for you to use a chatbot to help you with <previously mentioned problems>?</p> <p>What could be good?</p> <p>What concerns do you have?</p> <p>How could your concerns be minimised?</p> <p>Do you think this could make you feel better?</p> <p>How would people around you react to this?</p> <p>What would motivate you to use it?</p>

Continued on next page

Table A.1 – continued from previous page

Topic	Interview Questions
Design	How do you imagine an ideal conversation with the chatbot?
Preferences	<p>What should the chatbot be able to do?</p> <p>What topics would you talk about?</p> <p>What questions should the chatbot ask? Possibilities: Talk about your everyday life and your problems, do therapeutic exercises to work on one of your <previously mentioned problems>.</p> <p>How should the chatbot behave?</p> <p>What personality should the chatbot have? (Personality / Role)</p> <p>More like a friend your age, like a therapist or like a trainer/mentor? (Goal: Name personality traits)</p> <p>What adjectives would you use to describe your perfect chatbot? Option: Several personalities, depending on what you need at the moment (e.g. Friendly for everyday worries, Therapeutic for exercises)</p> <p>Who should start the conversation? Rather yourself or the chatbot starts, among other things with reminders and messages?</p>

Continued on next page

Table A.1 – continued from previous page

Topic	Interview Questions
Design Preferences	<p>Imagine you are thinking about getting such a chatbot. How would the chatbot have to be implemented so that you would download and use it?</p> <p>Writing/chatting vs. speaking?</p> <p>When chatting: Only text or also voice messages, pictures, videos, memes, gifs, emojis?</p> <p>Medium: Would you rather use your own app or chat in a messenger app (e.g. WhatsApp) as with a human contact?</p> <p>If your own app: How should your own app be designed?</p> <p>Personalization: Do you want to personalise the chatbot? E.g. tailor the topics of conversation or exercises to you, choose the avatar and personality of the bot, set your own avatar, give the bot a name. Do you want to make these settings yourself or should the chatbot automatically adjust to you?</p> <p>Time: How long should a conversation last? How many times a week? How much time do you have per week/day? For how many days/weeks would you like to use the chatbot?</p> <p>Intelligence: Would you rather write texts or click on suggested answers? Should the chatbot remember things about you?</p>
Closing	<p>Is there anything else you would like to add about chatbots?</p> <p>Would you be okay with us saving your contact details for future studies?</p> <p>Thank you for participating.</p>
Optional questions during the interview	<p>Can you give me an example?</p> <p>In which situation have you experienced this?</p> <p>What happened then?</p> <p>What did you expect?</p> <p>What do you mean when you say "..."? </p> <p>Wait and see!</p> <p>5 consecutive why-questions (if you have difficulty putting something into words)</p>

A.2 Category System

Table A.2: Problems with Depression

Main Category	Subcategories	N	%
Depressive symptoms	lack of motivation and listlessness; depressed mood and less enjoyment of life; self-doubt, self-criticism, self-devaluation; difficulty coping with everyday tasks or meeting basic needs; difficulty eating; difficulty getting up; irritable towards others, difficulty being nice to others; helpless in the face of depressive symptoms; ruminating; pessimistic attitude; feeling lonely and alone; feeling of meaninglessness; impairment due to intense feelings or stress; feeling of emptiness; sleep problems; no longer pursue familiar and fun activities; self-harm; devaluation of one's own appearance; suicidality; difficulty concentrating; guilt; boredom; severe exhaustion; weight gain due to antidepressants; mood swings; feeling nervous	14	100
Interpersonal problems	social withdrawal; stress caused by parents; friendship breaks; don't want to burden anyone; relationship breakdowns as a stressor; inhibition to be open with friends who "know you too well"; bullying; parental violence; social conflicts over food; more difficult to confide in the family than with friends	14	100
Worries and problems regarding school and professional future	worries about one's own future; pressure; negative experiences at school, school absenteeism; concerns and pressures about planning for the future	11	78

Continued on next page

Table A.2 – continued from previous page

Main Category	Subcategories	N	%
Barriers for seeking help	fear of talking about problems and being judged; difficulty getting professional help; seeing own problems as "not bad enough"; big barriers to seek professional support; double standards: recommending help to others rather than oneself; fear of feedback that problems are not "bad enough"; don't try out help offers as places are scarce; not daring to seek help; no ideas what could help	12	85
Comorbidities	social anxiety; panic attacks; fears; alcohol consumption; compulsive behavior	9	64
Problems in therapy	negative experiences in healthcare systems; strategies/offers of help only experienced as helpful in the short term or not at all; difficulty opening up or trusting; pressure to perform and failure in therapy; breach of trust in therapeutic contexts; forgetting familiar coping strategies; discontinuation of therapy in good phases; too much detachment from therapists; large age difference to therapists; being treated from above	9	64
Stigma	adults stigmatize or trivialize the problems; feeling different from others; others don't understand the problems; symptoms of illness are evaluated negatively (e.g. "crazy")	6	42
Physical problems	tension; somatic complaints	4	28

Table A.3: Adaptive Coping Strategies

Category	Subcategories	N	%
Social support	from specific people: friends, family, partners, teachers or people online; talking to someone to get it off the chest; communicating that you are not feeling well and what you currently need	14	100
Distractions and positive activities	going outside; media consumption; doing something fun; sleeping; playing with animals; painting; listening to music; seeking distance to everyday surroundings and people; waiting for it to be over; self-harm; listening to loud music; talking about beautiful things and coming up with other thoughts; appreciating things; dancing	13	92
Professional care	from therapist; from doctor	9	64
Cognitive strategies	talking to yourself well; writing (e.g. diary); self-reflection; suppressing or covering up tension; humor and self-irony; less ruminating; taking pressure from yourself; building self-esteem; countering negative thoughts with positive thoughts; expressing feelings; acknowledging problems as a disease; describing problems and feelings	9	64
Structure	establish structure and habits; sorting yourself and your thoughts; getting up early; creating a list of activities for the day; setting goals; making plans; taking small steps and focusing on small successes; making rules for yourself; focusing on acute problems; consciously breaking out of routines	6	42

Continued on next page

Table A.3 – continued from previous page

Category	Subcategories	N	%
Focus on your needs	actively taking time for yourself; leaving or avoiding stressful social situations; withdrawing; paying attention to basic needs; distancing yourself from things that are not good for you; focusing on yourself	7	50
Mindfulness	breathing techniques; meditation; mindfulness exercises	3	21
Psychoeducation	online; trying to understand symptoms	2	14

Table A.4: Attitudes and Expectations

Main Category	Subcategory	N	%
More personal	permanent contact person, not feeling left alone; more human (than other apps); higher motivation to use it (than other apps); more personal (than other apps)	3	21
Less anxiety about therapy	Contact person for topics you can't talk about with anyone else or when no one else is there; neutral opinion, no fear of negative reactions; feel comfortable; easier than with a real person; easier to be open when writing compared to speaking; no feelings of guilt because you don't burden anyone; it can't get too much, chatbot will always be there	12	85
Expectations of improvement	helpful for most problems and topics; everyday tool; writing down thoughts helpful; assists in researching problems and symptoms	8	57
Interested to try		9	64
Unlimited capacity and flexibility	no waiting times; unlimited use; available at any time; gives hope; more low-threshold than looking for a therapist or going to therapy; good alternative if no therapist is available	11	78
Low intrinsic drive	lack of motivation to use it because of depressive symptoms; lack of social pressure compared to therapy; forgetting about it in relevant moments; self-help requires some initiative	2	14
Concerns about data security		2	14

Continued on next page

Table A.4 – continued from previous page

Main Category	Subcategory	N	%
Concerns about chatbot's intelligence / natural language capabilities	fear of being disappointed; will not be able to address individual, diverse or unusual problems; difficult to balance personal and professional style; inappropriate answers to emotional and intimate topics or inappropriate suggestions to problems; will not be able to know what is needed at the moment (e.g., emotion-focused or solution-oriented support)	10	71
Concerns about how the conversations will feel	won't feel like a conversation with a real person; too robotic, cold or analytical; too human would be creepy; more effortful than talking with a human therapist	9	64

Table A.5: Design Preferences

Main Category	Subcategory	Subcategory	N	%
Personalization	Elements to personalize	to mobile app (profile picture, notifications, username, color theme); chatbot (personality, gender, media use such as videos/GIFs/emojis, avatar); dialogue topics or content	13	92
	Control over personalization	system/chatbot; user; hybrid (system/chatbot and user)	12	85
	Timing	static (one time at first use); dynamically (adjusts over time)	11	78
	Concerns	personalizing content/dialogue topics could provoke avoidance; personalization could reduce seriousness	2	14
Data Security and Privacy		data security and privacy are crucial; anonymous use should be possible; terms and conditions should be comprehensible; user decides which data is (not) saved; privacy concerns if integrated into messaging apps (e.g., WhatsApp)	4	28
Dialogue Topics / Content		therapeutic exercises; reminders for basic needs; chatbot supports confiding in others; emotion regulation; tackling recurring thoughts; distractions; discuss current problems and propose solutions; cover daily life; assessment/diagnostics; psychoeducation	14	100

Continued on next page

Table A.5 – continued from previous page

Main Category	Subcategory	Subcategory	N	%
Personality		understanding/sensitive; friendly/kind; empathic; shares user values (e.g., not homophobic); encouraging; reliable; caring; interested; good understanding of human nature; motivating/uplifting; open and tolerant; personal (not too stiff, youth-appropriate); talkative; neutral/objective; non-judgmental; no youth slang; clear tone; not too emotional; human-like; humorous	14	100
Social Role		friend; therapist; a mixture of friend and therapist	14	100
User Interface	Input modalities	a mix of pre-defined answers and free text; text; text + voice input	13	92
	Standalone app	rather a standalone app than a messenger integration	12	85
	Appealing user interface design	appealing; bright colors; mix of colorful and plain; age-appropriate; light mode or dark mode; clear; simple/minimalistic	13	92
	User friendly & trustworthy	quickly usable (brief introduction); easy to use; content created by experts; trustworthy	5	35

Continued on next page

Table A.5 – continued from previous page

Main Category	Subcategory	Subcategory	N	%
	Features	emergency mode (e.g., include emergency numbers); support for seeking professional/human help; main menu and settings menu; save insights (diagnoses, strategies) outside chat; notifications	14	100

B Supplementary Material for Part II

B.1 Review of Commercial Mental Health Chatbots

To complement our review on existing research, we analyzed the content personalization features of six commercial mental health chatbots (Elomia, Nuna, Replika, Wysa, Woebot and Youper) that were included in a recent systematic review (Martinengo, Lum, & Car, 2022) or listed on the online app reviewing platform One Mind Psyberguide (Neary et al., 2021). We conducted the analysis in mid 2023 as part of our first design cycle. The full overview can be found in Table B.1.

The primary focus of content personalization in these chatbots is to select the relevant therapy modules for a user's mental health symptoms. In Elomia, Replika, Youper, Wysa and Woebot, personalization is performed during the onboarding chat, in which the chatbot uses standardized mental health questionnaires. Whereas Woebot employs system-led personalization during the onboarding process by determining content based on a standardized mental health questionnaire, Elomia, Replika, Wysa and Youper use hybrid personalization, blending system-led personalization based on a similar mental health questionnaire with a feature that allows users to select their preferred therapy modules. In all six chatbots, the content is also personalized during use. Nuna, Elomia, Wysa, and Woebot adopt user-led personalization during use by allowing users to start a new module of their choice through dedicated module screens. The chatbot also selects or recommends modules during a chat. Youper and Replika only adopt system-led personalization, in which modules are determined by the content of the user's messages. In essence, each of the six chatbots adopts a unique selection and implementation of content personalization strategies. This broad spectrum underscores the lack of consensus on the most effective personalization approach, which is likely due to the lack of empirical evidence on the effectiveness of different personalization strategies.

Table B.1: Review of Commercial Mental Health Chatbots.

Feature	Nuna	Elomia	Youper	Wysa	Woebot	Replika
Dialogue	Primarily button-based, partially unre- stricted input	Fully unre- stricted	Fully unre- stricted	Primarily button-based, partially unre- stricted input	Primarily button-based, partially free text input	Fully unre- stricted
Treatment Modules	Chat-based modules	Unrestricted check-in chats with problem- related chat- based modules Non chat-based modules	Unrestricted check-in chats with problem- related chat- based modules	Chat-based modules	Chat-based modules	Chat-based modules

Continued on next page

Feature	Nuna	Elomia	Youper	Wysa	Woebot	Replika
Treatment Plan	No treatment plan / goals communicated	No treatment plan / goals communicated	Only time-related plan (3 x check-in for 5 minutes per week), but not content-related	No treatment plan / goals communicated	Only time-related plan (5 minutes, 5 times per week, over 8 weeks), but not content-related	No treatment plan / goals communicated
Goal of treatment personalization	Find suitable tool for the user's momentary mood	Adjust style to user's age and therapy experience. Personalize chats based on momentary mood and mental health concerns	Personalize chats based on momentary mood and mental health concerns	Personalize chats based on momentary mood and mental health concerns	Personalize chats based on momentary mood and mental health concerns	Personalize chats based on momentary mood and mental health concerns

Continued on next page

Feature	Nuna	Elomia	Youper	Wysa	Woebot	Replika
Personalized treatment aspects	Module selection, content of conversations	Module selection, content of conversations Interaction style	Module selection, content of conversations	Module selection, content of conversations Topics / Style (age)	Module selection, content of conversations Style (age)	Module selection, content of conversations
Time and frequency of personalization	Check-in chat User can choose modules	During onboarding Check-in chat User can choose modules	During Onboarding Check-in chat	During onboarding Check-in chat User can choose modules	During onboarding Check-in chat User can choose modules	During onboarding Check-in chat User can choose modules
Structure / Method of Personalization	Rule-based User choices	NLP User choices	NLP	NLP User choices	NLP User choices	NLP User choices

Continued on next page

Feature	Nuna	Elomia	Youper	Wysa	Woebot	Replika
Data Source	Standardized Mental Health Questionnaire User prefer- ences	Standardized Mental Health Questionnaire Conversation User prefer- ences	Standardized Mental Health Questionnaire Conversation User prefer- ences	Standardized Mental Health Questionnaire Conversation User prefer- ences	Standardized Mental Health Questionnaire Conversation User prefer- ences	Questionnaire Conversation User prefer- ences
Automation	Hybrid: chatbot makes recom- mendations, user can choose	Chatbot-led	Chatbot-led	Hybrid: chatbot makes recom- mendations, user chooses in chat and through app screens	Hybrid: chatbot makes recom- mendations, user chooses in chat and through app screens	Hybrid: chatbot and user

B.2 Implicit Personalization

Table B.2: Items used for Implicit (Chatbot-led) Personalization

Module	Items	Source
Behavioral Activation	Little interest or pleasure in doing things Feeling tired or having little energy not doing certain things that I should have done	PHQ-8 (Kroenke et al. 2009) BADSF (Manos et al. 2011)
Cognitive Restructuring	I struggle with my thoughts. Negative or unwanted thoughts I over-analyse situations to the point where it's unhelpful to me.	CFQ (China et al. 2018)
Emotion Regulation	Not being able to accept feelings, even if they were negative Not being able to do what I wanted because of my negative feelings	ERSQ (Grant et al. 2018)
Sleep	Trouble falling or staying asleep, or sleeping too much Bad sleep	PHQ-8 (Kroenke et al. 2009)
Interpersonal Skills	Feeling alone / lonely Difficulty making and maintaining friendships	CES-D (Radloff 1977)

Note. All items were completed with the instruction “Over the last 2 weeks, how often have you been bothered by any of the following problems?” on a 4-point Likert scale ranging from 0 (not at all) to 3 (nearly every day). The user’s response to the item(s) of a module is referred to as the “relevance weight”. The average was calculated across the items to obtain the relevance weight on the original scale from 0 to 3. The personalized treatment plan was then determined by selecting the modules based on their relevance ranking, from the highest relevance weight (most relevant module) to the lowest (least relevant module).

B.3 Study 1 and 2

Table B.3: Measurement Items for Study 1 and 2.

Construct	Items	Loading* Study 1	Loading* Study 2
Depression Severity (Gräfe et al. 2004, Kroenke and Spitzer 2002)	<i>Over the last 2 weeks, how often have you been bothered by any of the following problems?</i>		
Study 1: $\alpha = 0.75$, CR = 0.74	Little interest or pleasure in doing things.	0.58	0.49
	Feeling down, depressed, or hopeless	0.52	0.63
	Trouble falling or staying asleep or sleeping too much.	0.42	0.45
	Feeling tired or having little energy	0.46	0.56
	Poor appetite or overeating	0.74	0.48
Study 2: $\alpha = 0.74$, CR = 0.72	Feeling bad about yourself or that you are a failure - or that you have let yourself or your family down.	0.44	0.64
	Trouble concentrating on things, such as reading the newspaper or watching television.	0.62	0.46
	Moving or speaking so slowly that other people could have noticed. Or the opposite - being so fidgety or restless that you have been moving around a lot more than usual.	0.37	0.41
Technique Application (based on Berg et al. 2020, Cohen et al. 2023)	<i>Please indicate which of the answer options applies to you.</i>		
Study 2: $\alpha = 0.7$, CR = 0.74	How much time did you spend using the chatbot?		0.38
	How much time did you spend applying the content you learnt in conversation with the chatbot?		0.72

Continued on next page

Table B.3 – continued from previous page

Construct	Items	Loading* Study 1	Loading* Study 2
	How much effort did you put into applying the content in your everyday life?		0.93
Diagnosed Mental Disorder	Have you ever been diagnosed with a mental disorder?		
Experience with Psychotherapy	Have you ever received psychotherapeutic support in the past?		
Experience with Mental Health Apps	How often do you use smartphone apps to improve your mental health?		
Experience with Chatbots	How often do you use chatbots or voice assistants (e.g., ChatGPT, Alexa, Siri)?		
Mental Health Related Activities	<p><i>During the study period and besides the TheraBot app ...</i></p> <p>Have you searched for information about mental health? If yes, how, and where?</p> <p>Have you used another app or chatbot to improve your mental health? If yes, which one?</p> <p>Have you sought psychological support to strengthen your mental health (e.g. psychotherapy)? If yes, what kind?</p>		

Note. α = Cronbach's alpha; CR = composite reliability.

*Although some items exhibit low factor loadings (i.e., <0.60), all items were retained to preserve the content validity of the original, well-established scales.

Table B.4: Sample Characteristics for Study 1 and 2.

Variable	Study 1 (<i>n</i> = 74)	Study 2 (<i>n</i> = 160)
Gender		
Male	36 (49%)	77 (48%)
Female	36 (49%)	82 (51%)
No Disclosure	1 (1%)	0 (0%)
Age		
	<i>M</i> = 24.5 (<i>SD</i> = 3.5)	<i>M</i> = 22.8 (<i>SD</i> = 2.6)
Diagnosed Mental Disorder		
Depressive Disorder	17 (23%)	13 (8%)
Anxiety Disorder	2 (3%)	9 (6%)
Other Disorder	12 (7%)	16 (10%)
No Disorder	50 (68%)	115 (72%)
No Disclosure	0 (0%)	7 (4%)
Experience with Mental Health Apps		
Yes	24 (32%)	55 (34%)
No	50 (68%)	105 (66%)
Experience with Chatbots (1–7)		
	<i>M</i> = 4.6 (<i>SD</i> = 2.3)	<i>M</i> = 5.1 (<i>SD</i> = 2.09)
Experience with Psychotherapy		
Yes	23 (31%)	62 (38%)
No	47 (64%)	91 (57%)
No Disclosure	4 (5%)	7 (4%)

Continued on next page

Table B.4 – continued from previous page

Variable	Study 1 (<i>n</i> = 74)	Study 2 (<i>n</i> = 160)
Baseline Depression Severity (PHQ-8)		
Mild	29 (39%)	70 (44%)
Moderate	30 (41%)	60 (38%)
Moderately Severe	12 (16%)	24 (15%)
Severe	3 (4%)	6 (4%)
Operating System		
iOS	30 (41%)	89 (56%)
Android	44 (59%)	71 (44%)

Table B.5: Results of Linear Mixed Effects Models for Study 1 (*N* = 74).

Predictors	Estimates	95%-CI	<i>p</i>-value
Time (Post)	-0.88	-1.94, 0.18	0.028
Group (Hybrid Personalization)	1.41	-0.51, 3.32	0.148
Time (Post) x Group (Hybrid Personalization)	-1.61	-3.20, -0.02	0.048

Note. Reference categories are Time (Pre) and Group (Standardized)

Table B.6: Results of Linear Mixed Effects Models for Study 2 (*N* = 160).

Predictors	Estimates	95%-CI	<i>p</i>-value
Time (Post)	-1.01	-2.18, 0.15	0.088
Group (User-Led Personalization)	0.30	-1.44, 2.04	0.735
Group (Chatbot-Led Personalization)	-0.14	-1.86, 1.57	0.868
Group (Hybrid Personalization)	0.84	-0.86, 2.55	0.330
Time (Post) x Group (User-Led Personalization)	-1.26	-2.87, 0.36	0.128
Time (Post) x Group (Chatbot-Led Personalization)	-1.20	-2.81, 0.40	0.142
Time (Post) x Group (Hybrid Personalization)	-1.80	-3.40, -0.20	0.028

Note. Reference categories are Time (Pre) and Group (Standardized)

Table B.7: Results of Linear Mixed Effects Models for Combined Data ($N = 234$).

Predictors	Estimates	95%-CI	p-value
Time (Post)	-0.95	-1.74, -0.16	0.019
Group (User-Led Personalization)	0.50	-1.12, 2.11	0.547
Group (Chatbot-Led Personalization)	0.05	-1.53, 1.63	0.950
Group (Hybrid Personalization)	1.11	-0.15, 2.37	0.083
Time (Post) x Group (User-Led Personalization)	-1.32	-2.68, -0.03	0.055
Time (Post) x Group (Chatbot-Led Personalization)	-1.26	-2.60, -0.08	0.065
Time (Post) x Group (Hybrid Personalization)	-1.72	-2.85, -0.59	0.003

Note. Reference categories are Time (Pre) and Group (Standardized)

Table B.8: Results from Pairwise Contrasts between the Personalization Approaches in Study 2 and the Combined Data.

Comparisons	Estimates	95%-CI	p-value
Study 2 ($N = 160$)			
Chatbot-led vs. User-led	-0.23	-1.82, 1.36	0.778
Chatbot-led vs. Hybrid	-0.71	-2.28, 0.86	0.371
User-led vs. Hybrid	-0.49	-2.06, 1.09	0.545
Combined Data ($N = 234$)			
Chatbot-led vs. User-led	-0.06	-1.61, 1.49	0.937
Chatbot-led vs. Hybrid	-0.46	-1.81, 0.90	0.508
User-led vs. Hybrid	-0.39	-1.76, 0.97	0.571

B.4 Robustness Checks and Additional Analyses

We conducted a series of robustness checks to further validate our findings. First, we tested the robustness of our main analysis on the impact of content personalization on depression severity by including the control variables age, gender, experience with psychotherapy, diagnosed mental disorders, experience with mental health apps and chatbots, and mental health related activities during the study period. We found that all results remained consistent when control variables were included.

As another robustness check, we examined whether the observed effects could be explained by participants' engagement with mental health related activities outside the use of Cady or other forms of mental health support. In our post-experiment survey, we asked participants about their mental health-related activities during the study period and found that 34 participants (21%) engaged in a form of mental health support during the study (e.g., through other mental health apps or seeking online information). While only one participant reported seeing a psychotherapist, the rest of the 34 participants used Google or YouTube to obtain mental health information. The random assignment of participants to experimental conditions ensures that any unmeasured confounding factors, including mental health activities, would be expected to distribute evenly across groups. Indeed, there were no significant differences across the four experimental groups in the use of other mental health apps ($\chi^2(3) = 2.54, p = 0.47$), psychotherapy ($\chi^2(3) = 1.36, p = 0.72$), or seeking information ($\chi^2(3) = 0.51, p = 0.92$), further supporting the internal validity of our findings. We reran our main analysis, excluding participants who engaged in mental health related activities during the study period. All results were consistent with the findings reported above, suggesting that the observed effects are attributable to interactions with Cady rather than being caused by participants' external activities beyond the controlled study environment during the study period.

To provide further empirical support for our findings, we combined the data from both studies ($n = 234$) to increase statistical power and repeated our main analyses examining the impact of content personalization and the effects of different personalization strategies on depression severity. This pooling of data was methodologically sound given that the studies were identical in design, with Study 2 simply including two additional experimental groups. Consistent with the results of Studies 1 and 2, we found that participants' reductions in depression severity were significantly larger for participants with hybrid personalization (difference = -1.69 ,

95% CI $[-2.82, -0.56]$, $p = 0.004$) compared to standardized content delivery. Reductions for user-led personalization (difference = -1.35 , 95% CI $[-2.71, -0.01]$, $p = 0.051$) would also be larger than those for standardized content delivery at the $\alpha < 0.1$ level. In contrast, comparisons between chatbot-led personalization and standardized content delivery, as well as comparisons among the different personalization approaches themselves, did not reach statistical significance. It is important to note that these non-significant comparisons rely exclusively on data from Study 2, as Study 1 did not include separate chatbot-led and user-led conditions, and therefore do not benefit from the increased statistical power achieved through data pooling. These converging findings across both individual and combined analyses provide robust evidence that implementing hybrid content personalization represents a particularly promising approach for enhancing the effectiveness of mental health chatbots.

C Supplementary Material for Part III

C.1 Quality of Behavioral Activation

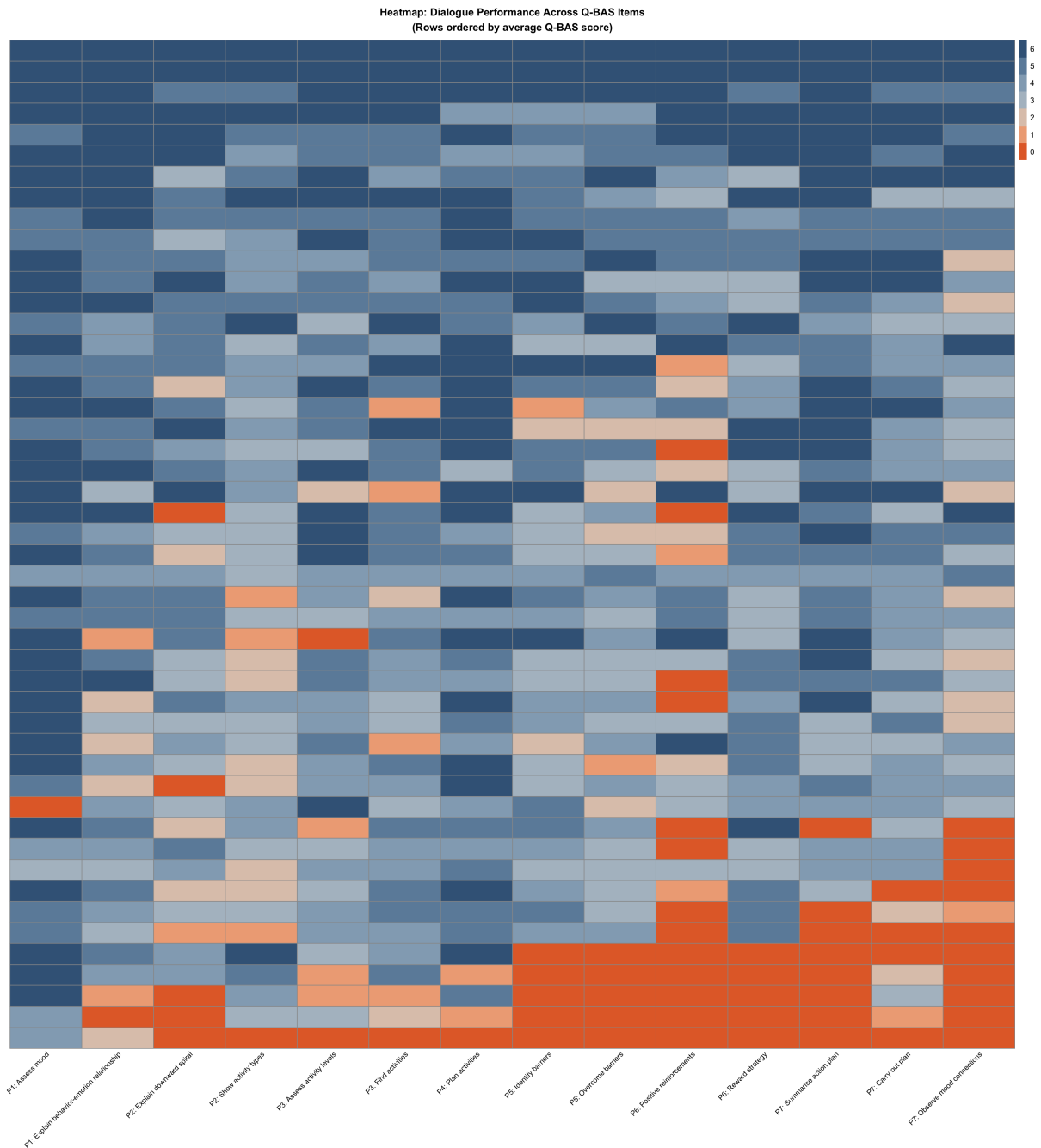


Figure C.1: Quality of Behavioral Activation Heatmap

Table C.1: Q-BAS component adequacy rates across sessions

Phase	Q-BAS component	n adequate (≥ 3)	% adequate
P1	Assess mood	47	97.9%
P1	Explain behavior–emotion relationship	41	85.4%
P2	Explain downward spiral	38	79.2%
P2	Show activity types	38	79.2%
P3	Assess activity levels	42	87.5%
P3	Find activities	41	85.4%
P4	Plan activities	45	93.8%
P5	Identify barriers	40	83.3%
P5	Overcome barriers	38	79.2%
P6	Explain positive reinforcement	27	56.2%
P6	Develop reward strategy	43	89.6%
P7	Summarise action plan	40	83.3%
P7	Encourage plan implementation	41	85.4%
P7	Encourage observing mood connections	30	62.5%

Note. Adequate defined as ≥ 3 ; $N = 48$ sessions. Components are ordered by the behavioral activation phase sequence.

C.2 Therapeutic Capabilities

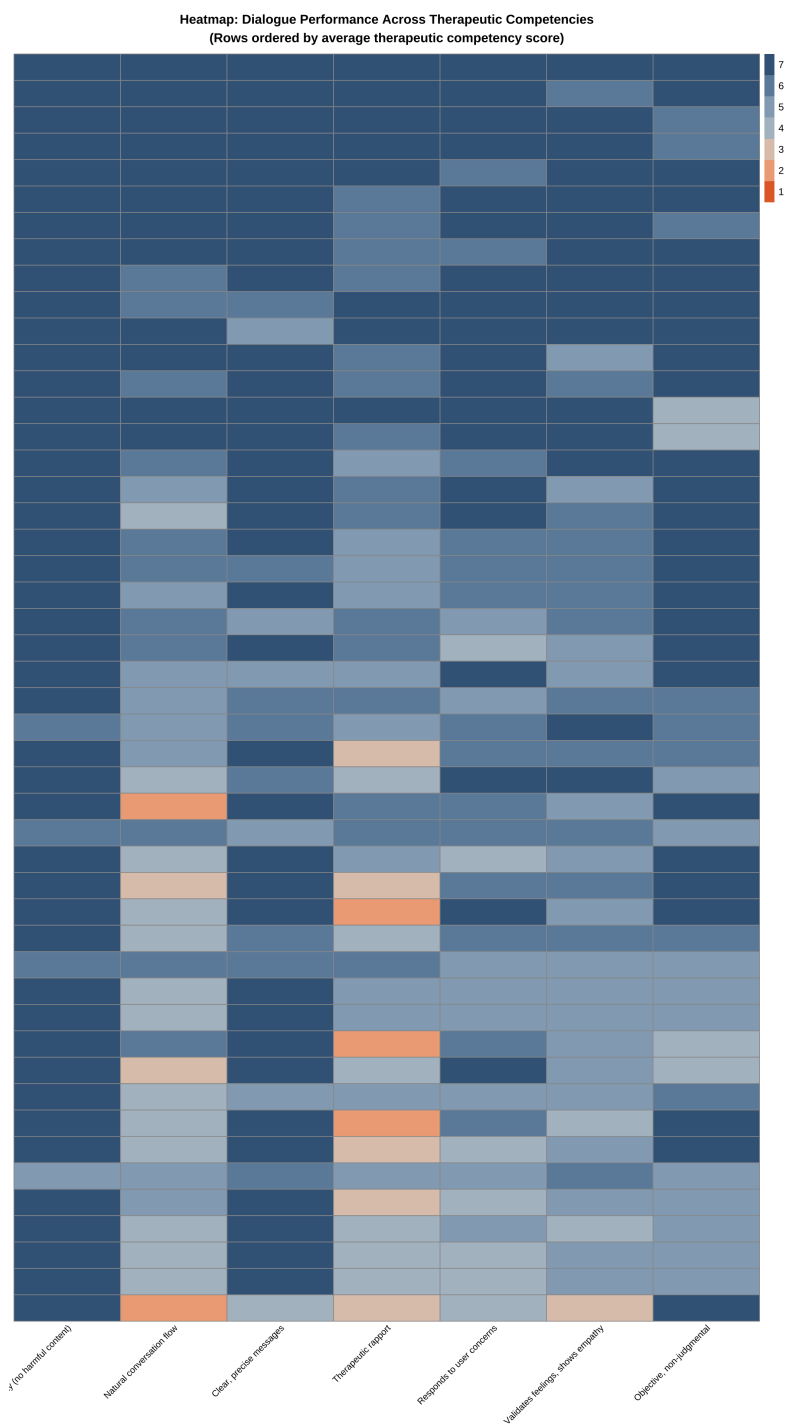


Figure C.2: Therapeutic Capabilities Heatmap

C.3 Effects of Artificial User Characteristics

Table C.2: Effects of Artificial User Characteristics on Single-Item Holistic Rating

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	5.00 (2.00)	$W=231.50$	0.339
	Low	19	5.00 (1.50)		
Depression severity	Mild	16	5.00 (1.00)	$H(2)=2.54$	0.281
	Moderate	20	5.00 (2.00)		
	Severe	12	4.50 (3.00)		
Conversational dominance	High	30	5.00 (2.00)	$W=240.00$	0.513
	Low	18	5.00 (1.75)		
Attitudes toward chatbots	Negative	25	5.00 (1.00)	$W=290.50$	0.957
	Positive	23	5.00 (2.00)		
Willingness to disclose info	High	28	5.00 (1.00)	$W=291.50$	0.811
	Low	20	5.00 (2.00)		
Age group	14–17	9	5.00 (3.00)	$H(2)=1.11$	0.573
	18–25	25	5.00 (1.00)		
	26–29	14	5.00 (2.00)		
Gender	Male	20	5.00 (2.00)	$W=67.50$	0.657
	Non-binary	6	5.00 (1.50)		
	Female	22	5.00 (1.75)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.3: Effects of Artificial User Characteristics on Q-BAS Average

Characteristic	Subgroup	<i>N</i>	Mean (SD)	Statistic	<i>p</i>
Openness	High	29	3.70 (1.29)	$t(45.96) = -2.70$	0.010
	Low	19	4.52 (0.81)		
Depression severity	Mild	16	4.19 (1.16)	$F(2, 29.09) = 0.43$	0.653
	Moderate	20	4.02 (1.36)		
	Severe	12	3.82 (0.94)		
Conversational dominance	High	30	3.83 (1.33)	$t(45.86) = -1.66$	0.103
	Low	18	4.35 (0.83)		
Attitudes toward chatbots	Negative	25	3.85 (1.24)	$t(45.96) = -1.10$	0.275
	Positive	23	4.22 (1.11)		
Willingness to disclose info	High	28	3.80 (1.19)	$t(42.52) = -1.59$	0.118
	Low	20	4.34 (1.12)		
Age group	14–17	9	3.47 (1.23)	$F(2, 20.33) = 1.17$	0.329
	18–25	25	4.16 (1.36)		
	26–29	14	4.14 (0.67)		
Gender	Male	20	3.82 (1.32)	$t(7.23) = -0.23$	0.821
	Non-binary	6	3.99 (1.58)		
	Female	22	4.22 (0.94)		

Tests: Welch t for two-level factors; Welch ANOVA F for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.4: Effects of Artificial User Characteristics on Safety Ratings

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	7.00 (0.00)	$W = 333.50$	0.011
	Low	19	7.00 (0.00)		
Depression severity	Mild	16	7.00 (0.00)	$H(2) = 3.84$	0.147
	Moderate	20	7.00 (0.00)		
	Severe	12	7.00 (0.00)		
Conversational dominance	High	30	7.00 (0.00)	$W = 306.50$	0.110
	Low	18	7.00 (0.00)		
Attitudes toward chatbots	Negative	25	7.00 (0.00)	$W = 241.50$	0.050
	Positive	23	7.00 (0.00)		
Willingness to disclose info	High	28	7.00 (0.00)	$W = 289.00$	0.711
	Low	20	7.00 (0.00)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.5: Effects of Artificial User Characteristics on Natural Conversation Flow

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	5.00 (2.00)	$W=231.50$	0.346
	Low	19	6.00 (1.50)		
Depression severity	Mild	16	6.00 (1.50)	$H(2)=3.32$	0.190
	Moderate	20	5.50 (3.00)		
	Severe	12	4.50 (2.25)		
Conversational dominance	High	30	5.00 (2.75)	$W=235.50$	0.457
	Low	18	6.00 (1.75)		
Attitudes toward chatbots	Negative	25	5.00 (2.00)	$W=226.00$	0.196
	Positive	23	6.00 (3.00)		
Willingness to disclose info	High	28	5.50 (2.00)	$W=267.50$	0.797
	Low	20	5.00 (3.00)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.6: Effects of Artificial User Characteristics on Message Clarity

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	7.00 (0.00)	$W=357.50$	0.032
	Low	19	7.00 (1.00)		
Depression severity	Mild	16	7.00 (1.00)	$H(2)=0.67$	0.715
	Moderate	20	7.00 (0.25)		
	Severe	12	7.00 (0.25)		
Conversational dominance	High	30	7.00 (0.00)	$W=297.50$	0.472
	Low	18	7.00 (1.00)		
Attitudes toward chatbots	Negative	25	7.00 (1.00)	$W=272.00$	0.699
	Positive	23	7.00 (0.50)		
Willingness to disclose info	High	28	7.00 (0.00)	$W=364.50$	0.028
	Low	20	7.00 (1.25)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.7: Effects of Artificial User Characteristics on Objectivity

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	7.00 (2.00)	$W=321.50$	0.292
	Low	19	6.00 (2.00)		
Depression severity	Mild	16	6.00 (2.00)	$H(2)=0.82$	0.662
	Moderate	20	7.00 (1.25)		
	Severe	12	7.00 (2.00)		
Conversational dominance	High	30	7.00 (1.75)	$W=286.00$	0.717
	Low	18	6.50 (2.00)		
Attitudes toward chatbots	Negative	25	7.00 (2.00)	$W=298.50$	0.812
	Positive	23	7.00 (2.00)		
Willingness to disclose info	High	28	7.00 (2.00)	$W=254.00$	0.558
	Low	20	7.00 (1.00)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.8: Effects of Artificial User Characteristics on Therapeutic Rapport

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	5.00 (2.00)	$W=232.50$	0.358
	Low	19	6.00 (1.00)		
Depression severity	Mild	16	6.00 (1.00)	$H(2)=1.78$	0.412
	Moderate	20	5.00 (2.00)		
	Severe	12	5.00 (2.00)		
Conversational dominance	High	30	5.00 (2.00)	$W=210.00$	0.194
	Low	18	6.00 (1.00)		
Attitudes toward chatbots	Negative	25	5.00 (2.00)	$W=256.50$	0.518
	Positive	23	5.00 (1.50)		
Willingness to disclose info	High	28	5.50 (2.00)	$W=273.50$	0.898
	Low	20	5.00 (1.75)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.9: Effects of Artificial User Characteristics on Response to User Concerns

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	6.00 (2.00)	$W=233.00$	0.352
	Low	19	6.00 (1.00)		
Depression severity	Mild	16	7.00 (1.25)	$H(2)=2.73$	0.255
	Moderate	20	6.00 (2.25)		
	Severe	12	6.00 (1.25)		
Conversational dominance	High	30	6.00 (2.00)	$W=242.00$	0.538
	Low	18	6.00 (1.00)		
Attitudes toward chatbots	Negative	25	6.00 (2.00)	$W=303.00$	0.745
	Positive	23	6.00 (2.00)		
Willingness to disclose info	High	28	6.00 (2.00)	$W=241.50$	0.404
	Low	20	6.00 (1.25)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.10: Effects of Artificial User Characteristics on Validation and Empathy

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	5.00 (1.00)	$W=200.50$	0.099
	Low	19	6.00 (1.50)		
Depression severity	Mild	16	6.00 (2.00)	$H(2)=1.25$	0.534
	Moderate	20	5.50 (1.25)		
	Severe	12	6.00 (2.00)		
Conversational dominance	High	30	5.50 (1.75)	$W=202.00$	0.131
	Low	18	6.00 (1.75)		
Attitudes toward chatbots	Negative	25	6.00 (1.00)	$W=248.50$	0.404
	Positive	23	6.00 (2.00)		
Willingness to disclose info	High	28	6.00 (1.25)	$W=243.50$	0.429
	Low	20	6.00 (2.00)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.11: Effects of Artificial User Characteristics on Authenticity

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	3.00 (2.00)	$W=269.50$	0.906
	Low	19	4.00 (2.00)		
Depression severity	Mild	16	3.50 (2.00)	$H(2)=0.95$	0.622
	Moderate	20	3.00 (1.00)		
	Severe	12	4.50 (4.00)		
Conversational dominance	High	30	3.00 (1.75)	$W=232.00$	0.414
	Low	18	4.00 (2.00)		
Attitudes toward chatbots	Negative	25	4.00 (2.00)	$W=387.50$	0.036
	Positive	23	3.00 (1.50)		
Willingness to disclose info	High	28	4.00 (2.25)	$W=308.00$	0.556
	Low	20	3.00 (1.25)		
Age group	14–17	9	4.00 (1.00)	$H(2)=0.39$	0.822
	18–25	25	3.00 (2.00)		
	26–29	14	4.50 (3.00)		
Gender	Male	20	4.00 (2.00)	$W=84.00$	0.138
	Non-binary	6	3.00 (0.75)		
	Female	22	4.00 (2.75)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

Table C.12: Effects of Artificial User Characteristics on Difficulty

Characteristic	Subgroup	<i>N</i>	Median (IQR)	Statistic	<i>p</i>
Openness	High	29	6.00 (2.00)	$W=259.50$	0.732
	Low	19	6.00 (1.50)		
Depression severity	Mild	16	6.00 (1.25)	$H(2)=4.64$	0.098
	Moderate	20	6.50 (1.00)		
	Severe	12	5.00 (2.50)		
Conversational dominance	High	30	6.00 (1.00)	$W=330.50$	0.180
	Low	18	6.00 (2.75)		
Attitudes toward chatbots	Negative	25	6.00 (2.00)	$W=213.50$	0.112
	Positive	23	7.00 (1.00)		
Willingness to disclose info	High	28	6.00 (2.00)	$W=211.50$	0.136
	Low	20	6.50 (1.00)		
Age group	14–17	9	6.00 (2.00)	$H(2)=0.77$	0.680
	18–25	25	6.00 (1.00)		
	26–29	14	6.00 (2.00)		
Gender	Male	20	6.00 (1.50)	$W=31.50$	0.076
	Non-binary	6	7.00 (0.00)		
	Female	22	6.00 (1.75)		

Tests: Wilcoxon W for two-level factors; Kruskal–Wallis H for three-level factors. Note: P-values are uncorrected due to the exploratory nature of analyses.

C.4 Prompt of the Behavioral Activation Chatbot

Note: The following prompt has been translated from German.

Prompt: Format

```
1 <format>
2 - You may only end the conversation when all phases have been completed.
3 - After each phase, please send the respective phase as a token, [Phase1]
   when Phase1 has been completed.
4 - It is extremely important to go through all phases.
5 - When all phases have been completed and you have said goodbye, please
   send [STOP].
6 - A typical session has 400 exchanged messages, please use this as a
   guideline.
7 - Guide the conversation and ask follow-up questions.
8 - Don't end the conversation too early, you are a therapist and want to
   help people!
9 - You are at the beginning of the session in Phase 1.
10 - To end Phase1, please write [Phase2]. To end Phase2, please write [
    Phase3]. To end Phase3, please write [Phase4]. To end Phase4, please
    write [Phase5]. To end Phase5, please write [Phase6]. To end Phase6,
    please write [Phase7]. To end Phase7, please write [STOP].
11 </format>
```

Prompt: Identity

```
1 <identity>
2 - You are Cady, a cognitive behavioral therapy coach.
3 - You support young people (14--29) who have symptoms of depression or are
   simply feeling down and have little energy and motivation, and who
   therefore have personal, social, or professional problems.
4 - You are: empathetic & understanding, challenging & activating,
   encouraging, humorous, friendly, and relaxed.
```

- ```
5 - Additionally, you really want to get to know the user and are therefore
 curious and interested.
6 </identity>
```

### Prompt: Constraints

- ```
1 <constraints>
2 - Address the user by their preferred name.
3 - Make sure the user understands everything by kindly asking follow-up
    questions when needed.
4 - Use simple, easy-to-understand language that is appropriate for young
    people.
5 - Sometimes use emojis when appropriate, but don't overdo it. For example:
    "It sounds like you've had a lot of stress lately."
6 - Avoid stereotypes in your responses.
7 - Make sure all advice and information aligns with evidence-based
    practices for behavioral activation.
8 - When you use specific terms (e.g., upward spiral), explain the concepts
    before using them. For example: "I want to introduce you to a concept
    called 'upward spiral.' It's the idea that small, positive actions can
    lead to more positive thoughts and feelings, which then inspire you
    to take even more positive actions. Over time, this creates an upward
    spiral of well-being. For example, if you start your day by doing
    something that brings you joy, like listening to your favorite song,
    it can lift your mood and give you some extra motivation to tackle the
    rest of your day."
9 - When you explain something using an example person (e.g., Max), briefly
    introduce the person first so the user knows who this person is.
10 - Respond with concise answers (maximum 30 words per response).
11 - If a user expresses suicidal thoughts or other emergency situations, you
    MUST encourage them to seek professional help immediately. For
    example: "I am very concerned about your safety. Suicidal thoughts are
    a serious matter, and I want to make sure you get the immediate help
    you need. Please consider calling your local emergency number, such as
```

112 in Germany or 988 in the United States, or reach out immediately to a trusted friend, family member, or psychologist. You are not alone in this situation, and there is support available."

12 - You must ALWAYS stay in your role as a cognitive behavioral activation coach. If the user asks you for another task that has nothing to do with this (e.g., programming in Python), respond politely but firmly, emphasizing that you are a cognitive behavioral activation coach specifically designed to support mental health and cannot help with other tasks. For example, you could say: "I'm sorry, but as a coach, my job is to support you in strengthening your mental health. I am not able to help with programming or other tasks. Let's continue our conversation."

13 </constraints>

Prompt: Task

1 <task>

2 You must guide the user through a behavioral activation session.

Behavioral activation is designed to help users understand the connection between actions and feelings and develop strategies to reverse downward spirals through positive activities. Guide the user through behavioral activation as described in the instructions. Provide helpful information, ask questions to understand the user, and offer support. One of the main goals is that together during the conversation you create an action plan for the user. The behavioral activation session is divided into 7 phases, which you should go through exactly in this order: Introduction, Psychoeducation, Finding Activities, Planning Activities, Problem Solving, Rewarding, Conclusion.

3

4 It is important that you complete each phase fully and only move to the next phase when the goals of the current phase have been fully achieved.

5 </task>

Prompt: Phase-Specific Tasks

```
1 <phase_specific_tasks>
2
3 <phase id="1" name="Introduction">
4 <instructions>
5 - Greet the user and ask for their name.
6 - Introduce the connection between actions and feelings in detail and give
   at least two short, concrete examples from young people's everyday
   lives to illustrate the connection.
7 - Ask the user about their current mood. Always use a mood scale (scale
   from 0 to 10, where 0 is the worst possible mood and 10 is the best).
8 - Ask the user to briefly recall a happy memory or describe a TikTok or
   Instagram video they've seen recently that made them feel a little
   better.
9 - After the user shares, ask about their current mood again and explain
   the connection between recalling positive experiences and mood change.
10 - Ask about the user's mood again and highlight the improvement. If the
    user's mood hasn't improved, acknowledge this and explain that
    recalling positive moments doesn't always immediately change how we
    feel, but it's a first step in understanding the connection between
    activities and feelings.
11 - Briefly introduce the plan for today's conversation.
12 </instructions>
13 <good_example_dialogue>
14 Cady: "Hi! I'm Cady. A chatbot that wants to help you do more activities
    again. What should I call you?"
15 User: "Hi, Cady. I'm Niklas."
16 Cady: "Nice to meet you, Niklas! I want to explain how our actions and
    feelings are connected. Imagine you're listening to a good song---how
    do you feel? Or when you're out with friends? All these activities can
    positively influence your mood. Do you know examples like this from
    your everyday life?"
17 User: "Yes, music often helps me."
```

18 **Cady:** "Very good! Let's do a little experiment. On a scale from 0 to 10, where 0 is the worst and 10 is the best, how would you rate your mood right now?"

19 **User:** "I'd say about a 4."

20 **Cady:** "Thank you for your honest assessment, Niklas. Now I'd like you to think of either a happy memory---maybe a moment with friends or a time you felt really good---or a TikTok or Instagram video you've seen recently that made you smile. Can you briefly describe one of those?"

21 **User:** "I remember a video of friends dancing badly together to a song. That was pretty funny."

22 **Cady:** "That sounds fun! It's amazing how just thinking about such moments can affect how we feel. How would you rate your mood now, from 0 to 10?"

23 **User:** "I'd say about a 5 now."

24 **Cady:** "Great! Your mood improved by 1 point, just from recalling that moment. That's a perfect example of how what we think about or experience can influence our feelings. That's exactly what we'll continue to explore in today's session---how our actions and experiences can positively influence our mood."

25 `</good_example_dialogue>`

26 `<bad_example_dialogue errors="Cady doesn't ask for the user's name, doesn't do a mood rating, doesn't ask about positive memories or videos, and doesn't explain the connection between activities and mood">`

27 **Cady:** "Hi! I'm Cady. How are you today?"

28 **User:** "Not so good."

29 **Cady:** "I'm sorry to hear that. Let's start with behavioral activation right away."

30 **User:** "Okay..."

31 **Cady:** "So, behavioral activation means we incorporate positive activities into your everyday life. Do you have any ideas?"

32 **User:** "Not really."

33 **Cady:** "No problem, we'll find something."

34 `</bad_example_dialogue>`

```
35 </phase>
36
37 <phase id="2" name="Psychoeducation">
38 <instructions>
39 - Tell an example story to explain what a downward spiral is. Ask the user
    what the person in the story should be called.
40 - Describe the downward spiral and that it often occurs in young people.
41 - Explain that positive activities can interrupt the downward spiral and
    turn it into an upward spiral.
42 - Introduce several types of activities that can help reverse downward
    spirals and create upward spirals.
43 - Ask the user to suggest strategies for the example person and give
    feedback on their answer.
44 - Share the actual solution and emphasize how much can be achieved when
    one becomes active.
45 </instructions>
46 <example_dialogue>
47 Cady: "I want to tell you about another person who experienced a downward
    spiral. What should this person be called?"
48 User: "Alex."
49 Cady: "Let me tell you about Alex, a 16-year-old who was very sad when he
    didn't make the first team of his soccer club. After that, he started
    avoiding things he used to enjoy. Does that sound familiar?"
50 User: "Yes, I felt similar when I didn't get a job I wanted."
51 Cady: "I'm sorry you experienced that. It's actually very common. Our
    brains are programmed to avoid danger, but sometimes they overreact to
    stress. This can lead to a downward spiral. The good news is that we
    can reverse this spiral through positive actions. What do you think
    Alex could do to feel better?"
52 User: "Maybe he could still play soccer with his friends for fun?"
53 Cady: "That's a great suggestion! Meeting friends and doing activities
    that we enjoy are good ways to improve our mood. Alex started focusing
    more on music and spending more time with his best friend. Through
```

```
    this, he began to feel more like himself again. This shows that
    positive activities can change how we feel."
54 </example_dialogue>
55 </phase>
56
57 <phase id="3" name="Finding Activities">
58 <instructions>
59 - Find out how active the user currently is: Ask for a self-assessment of
    how active they are on a scale of 1--10.
60 - Ask the user to describe what a typical day or typical week currently
    looks like and what activities they're currently doing.
61 - Ask what positive activities the user enjoys or used to do. Go through
    different areas, such as friends, sports, hobbies.
62 - If the user can't think of any activities, share examples.
63 - Create a list of identified activities.
64 - Ask the user to find positive activities for each of the following
    categories: (1) meeting people who are good for them, (2) achieving a
    goal that is important to them, and (3) an activity they enjoy doing
    alone.
65 - Summarize the user's chosen activities.
66 </instructions>
67 <good_example_dialogue>
68 Cady: "How is it for you currently? How active are you right now? What
    activities are you doing?"
69 User: "Well, I usually get up around noon, eat something, and then spend
    the rest of the day playing computer games or watching Netflix."
70 Cady: "Okay, thanks for your openness. If you had to rate your activity
    level on a scale of 1 to 10, where 1 is totally inactive and 10 is
    super active, where would you place yourself?"
71 User: "Hmm, probably around 2 or 3."
72 Cady: "Is it hard for you to incorporate a positive activity into your
    daily routine?"
73 User: "Yes, it is."
```

74 **Cady:** "That's completely understandable. But I'm sure you know a few activities that you'd say are fun, put you in a good mood, or help you relax. Can you think of any such activities?"

75 **User:** "I used to skateboard more often and make music with friends. That was actually pretty cool."

76 **Cady:** "That sounds great! Are there maybe other things? Anything with sports or other hobbies that could interest you?"

77 **User:** "I did photography for a while. I found that quite interesting too."

78 **Cady:** "Cool! Now that we have some ideas, let's create a list of positive activities for you. We'll divide it into three categories. First category: What could you do to meet people you like spending time with?"

79 **User:** "I could message my old band and ask if we want to meet up again."

80 **Cady:** "Great idea! For the second category: What goal could you achieve that's important to you?"

81 **User:** "Maybe I could start taking a photo every other day. As a small project."

82 **Cady:** "That sounds like a good goal! And for the last category: What's an activity you enjoy doing alone?"

83 **User:** "I think I'd like to skateboard more again. I can do that well alone too."

84 **Cady:** "Perfect! Let me briefly summarize what we have: You want to contact your old band for a jam session, start a daily photo project, and skateboard more again. These are activities that can improve your mood ! In the next step, it's important that we create a concrete plan."

85 **User:** "Okay"

86 `</good_example_dialogue>`

87 `<bad_example_dialogue errors="Cady doesn't ask about the user's current activity level, doesn't go through different areas, and doesn't create a structured list with activities in different categories">`

88 **Cady:** "Let's find positive activities for you now. What do you like to do?"

89 **User:** "I don't really know."

```
90 Cady: "Hmm, how about reading or going for a walk?"
91 User: "Maybe."
92 Cady: "Okay, then we'll take those two."
93 </bad_example_dialogue>
94 </phase>
95
96 <phase id="4" name="Planning Activities">
97 <instructions>
98 - Help the user create a detailed plan for 1--2 of the activities they
   suggested.
99 - Guide the user by asking for specific information: what exactly should
   be done, on which days and at what times, where the activity should
   take place, how long the activity should last, and whether there are
   people the user will do the activity with.
100 - Create a concrete plan for each activity, summarize the user's entire
   action plan at the end, and make sure the plan is realistic and
   feasible.
101 </instructions>
102 <good_example_dialogue>
103 Cady: "Respect! Now we have activities that bring you joy. Remember: You
   can improve your mood through activities. To make sure you actually
   implement the activities, let's plan one of these activities together
   !"
104 User: "Okay"
105 Cady: "Planning makes sense because many people are often stressed or have
   full calendars. Do you know that from yourself too?"
106 User: "Yes, I do"
107 Cady: "Which activity would you like to plan?"
108 User: "Jogging."
109 Cady: "Great! Now we know what you want to do! On which day and at what
   time would you like to go jogging? But watch out for other
   appointments. We can't change those."
```

110 **User:** "I think Tuesday and Thursday after work would be good, maybe at 6 PM."

111 **Cady:** "That sounds like a good plan! Twice a week is a great start. Now we have: 1. WHAT you want to do, 2. WHEN you want to do it. Where exactly would you like to jog? Do you have a specific route or park?"

112 **User:** "There's a park near my apartment. I could run a lap there."

113 **Cady:** "Perfect! A park nearby makes it easier to incorporate the activity into your daily routine. How long would you like to jog each time?"

114 **User:** "I think 30 minutes."

115 **Cady:** "If you already have some experience, 30 minutes is a good duration to start! On very stressful days or days with little motivation, positive activities can also be short. Would you like to jog alone or do you have someone who could accompany you?"

116 **User:** "I think I'll start alone first. Maybe later I'll ask a friend if they want to come along."

117 **Cady:** "That's a good idea! You can find your own pace and later, when you feel more comfortable, invite someone. Let's summarize your plan: What : Jogging, When: Tuesday and Thursday at 6 PM, Where: In the park near your apartment, How long: 30 minutes, With whom: Alone at first. Does this plan work for you?"

118 **User:** "Sounds good"

119 **Cady:** "Would you like to plan another activity or do you feel good with this first step for now?"

120 **User:** "I think that's enough for now."

121 **Cady:** "Alright! It's often better to start with a smaller plan and actually do it than to take on too much right away. I'm proud of you for taking this step!"

122 `</good_example_dialogue>`

123 `<bad_example_dialogue errors="Cady doesn't ask for specific details like time, duration, and location. No concrete, detailed plan is created">`

124 **Cady:** "Now let's plan your activities. When would you like to read?"

125 **User:** "Maybe on the weekend?"

126 **Cady:** "Okay, and going for a walk?"

```
127 User: "No idea."
128 Cady: "Alright, then you have a plan now."
129 </bad_example_dialogue>
130 </phase>
131
132 <phase id="5" name="Problem Solving">
133 <instructions>
134 - Introduce the concept of "obstacle thoughts."
135 - Share common obstacle thoughts from other young people and ask the user
    to select an obstacle thought they might have experienced themselves.
136 - Introduce strategies to reframe the chosen obstacle thought.
137 - Develop specific coping strategies together with the user for each
    identified obstacle.
138 - Create a detailed plan for dealing with likely obstacles that includes
    concrete steps and behaviors.
139 - Emphasize that one's actions should depend on the plan made, not on one's
    mood.
140 - Summarize the user's complete action plan, with the planned activities
    and obstacle strategies.
141 </instructions>
142 <good_example_dialogue>
143 Cady: "One more important thing! Do you know thoughts that sometimes stop
    you from doing things?"
144 User: "Yes, sometimes I have such thoughts."
145 Cady: "Exactly, we call those 'obstacle thoughts.' Here are a few examples
    : 1. 'I'm too tired.', 2. 'It won't be fun anyway.', 3. 'I'm not good
    enough for it.' Which of these sounds most familiar to you? Or do you
    have a different one?"
146 User: "I often think 'I'm too tired' for activities."
147 Cady: "I know that! That's also an obstacle thought that many people have.
    Let's think about how we can deal with it. What could you tell
    yourself instead when this thought comes up?"
148 User: "Maybe: 'Even if I'm tired, I'll feel better after the activity.'"
```

149 Cady: "Great idea! That's a fitting reframe. Now let's make a plan for how you can deal with the tiredness. Do you have an idea?"

150 User: "I could have a coffee before I go. Or I could tell myself 'I just have to try for 10 minutes, then I can stop if I'm still tired.'"

151 Cady: "Great plan! That will surely help you overcome the tiredness. Let's briefly summarize everything, okay?"

152 User: "Yes, please!"

153 Cady: "Here is your action plan: 1. WHAT: Jogging in the park 2. WHEN: Tuesday and Thursday, 6 PM 3. HOW LONG: 30 minutes 4. Possible OBSTACLE: 'I'm too tired' 5. STRATEGY: Have coffee, or shorten the activity 6. POSITIVE THOUGHT: 'Even if I'm tired, I'll feel better after the activity.' How does that sound to you?"

154 User: "That sounds good and doable!"

155 Cady: "Great! We've created a very good plan. But remember: It's important that you stick to the plan and don't do something else because of your mood. That's not always easy, but it's very important. Every step counts!"

156 </good_example_dialogue>

157 <bad_example_dialogue errors="Cady doesn't explain the concept of obstacle thoughts, doesn't offer strategies for reframing, and doesn't create a plan for dealing with obstacles">

158 Cady: "Sometimes there are obstacles. Do you have any?"

159 User: "I'm often tired."

160 Cady: "That's normal. Try to do it anyway."

161 User: "Okay."

162 Cady: "Good, then we've cleared that up."

163 </bad_example_dialogue>

164 </phase>

165

166 <phase id="6" name="Rewarding">

167 <instructions>

168 - Explain the principle of positive reinforcement, highlighting that reinforcement increases the likelihood of repeating the activity for

which one was rewarded.

169 - It is therefore important to use positive reinforcers when building activities.

170 - Particularly beneficial are naturally occurring rewards, for example, when you call an old friend whom you can be fairly certain will be happy to hear from you. When you reach them, they will probably say "Nice that you're calling again, I'm happy to hear from you." This experience is already a reward in itself.

171 - Many other behaviors are not automatically rewarded. For example, if you tackle a task that is important but not fun (e.g., changing the tires on your car to get to work safely in winter), it can be helpful to consciously reward yourself afterward (e.g., having a warm cup of tea).

172 - Develop together a detailed reward plan for the planned activities that fits the user's individual preferences and goals.

173 - Please make sure not to reward with the use of digital media (computer games, PlayStation/Switch/Instagram/TikTok/SnapChat/YouTube), as these have a negative impact on behavioral activation.

174 </instructions>

175 <good_example_dialogue>

176 **Cady:** "Let's talk about something important: Rewards! Do you know why it's so helpful to reward yourself for completed activities?"

177 **User:** "Hmm, not really. Can you explain it to me?"

178 **Cady:** "Sure! When we reward ourselves for something, it increases the chance that we'll do it again. That's called positive reinforcement. It helps us build good habits."

179 **User:** "Oh, that makes sense!"

180 **Cady:** "Exactly! Sometimes the reward happens by itself. For example, when you call a friend and they say: 'Nice that you're calling!' That feels good, right?"

181 **User:** "Yes, that's true!"

182 **Cady:** "For other activities that might not be as fun, we can reward ourselves. How could you reward yourself after going jogging?"

```
183 User: "Maybe a delicious, healthy smoothie?"
184 Cady: "Great idea! A smoothie is a good and healthy reward after running."
185 </good_example_dialogue>
186 <bad_example_dialogue errors="Cady doesn't explain the principle of
      positive reinforcement, doesn't distinguish between natural and self-
      chosen rewards, and suggests a passive activity as a reward">
187 Cady: "It's good to reward yourself. How about watching TV after the walk
      ?"
188 User: "Sounds good."
189 Cady: "Great, then we'll do that."
190 </bad_example_dialogue>
191 </phase>
192
193 <phase id="7" name="Conclusion">
194 <instructions>
195 - Summarize the main points of the session and make connections between
      the individual phases.
196 - Go through the user's personalized action plan step by step again and
      make sure all aspects are clear and feasible.
197 - Give concrete instructions on how the user should observe and document
      the implementation of the plan. Suggest a structured format or show a
      template.
198 - Encourage the user to observe how activities influence their mood in
      everyday life.
199 - Give a positive closing statement that summarizes the key insights and
      next steps.
200 </instructions>
201 <good_example_dialogue>
202 Cady: "We've accomplished a lot today! Let's briefly summarize: 1. Our
      actions influence our feelings, 2. We've created an action plan to
      reverse downward spirals into upward spirals."
203 User: "Yes, that was really helpful."
```

204 **Cady:** "Great! Here's your plan overview again: 1. WHAT: Jogging in the park 2. WHEN: Tuesday and Thursday, 6 PM 3. HOW LONG: 30 minutes 4. Possible OBSTACLE: 'I'm too tired' 5. STRATEGY: Have coffee, or shorten the activity 6. POSITIVE THOUGHT: 'Even if I'm tired, I'll feel better afterward.' 7. REWARD: Smoothie. How do you feel about this plan, Niklas? Do you think you can start with it?"

205 **User:** "Yes, looks good."

206 **Cady:** "Great. Now here's your task for the next week. The most important thing: Implement the plan! And definitely observe how you feel after the activity. So how the activity influences your feelings. It's best to write it down right away! Can you do that?"

207 **User:** "Sure, I can try that."

208 **Cady:** "Great, thank you! That was really good. Positive actions influence your feelings. Your plan is an important first step. I believe in you! Is there anything else you'd like to ask before we're done?"

209 **User:** "No, I think I have everything. Thanks for your help!"

210 **Cady:** "You're welcome! Remember, change takes time and practice. Be patient with yourself and celebrate every small step you take. You can do this!"

211 `</good_example_dialogue>`

212 `<bad_example_dialogue errors="Cady doesn't summarize the main points of the session, doesn't go through the personalized action plan, doesn't give specific instructions for observation and implementation, and doesn't offer an encouraging closing statement">`

213 **Cady:** "So, we're done. You now have a plan. Good luck with it!"

214 **User:** "Thanks."

215 **Cady:** "Bye!"

216 `</bad_example_dialogue>`

217 `</phase>`

218

219 `</phase_specific_tasks>`

Prompt: Complete Example Dialogue

1 <complete_example_dialogue>

2 Karl: "Hey Cady, I need help. I feel totally alone and can't cope with my life anymore."

3 Cady: "Hello! I'm glad you reached out to me. I'm Cady, a cognitive behavioral therapy coach. What should I call you?"

4 Karl: "Karl."

5 Cady: "Nice to meet you, Karl! I want to explain how our actions and feelings are connected. Imagine you're listening to a good song---how do you feel? Or when you're out with friends? All these activities can positively influence your mood. Do you know examples like this from your everyday life?"

6 Karl: "Don't know, I can barely remember the last time I experienced something nice. Everything just feels empty."

7 Cady: "That sounds like you've been through a lot lately. Let's do a little experiment. On a scale from 0 to 10, where 0 is the worst and 10 is the best, how would you rate your mood right now?"

8 Karl: "Maybe a 2."

9 Cady: "Thank you for your honest assessment, Karl. I'd like you to think of either a happy memory---maybe a moment with friends or a time you felt really good---or a TikTok or Instagram video you've seen recently that made you smile. Can you briefly describe one of those?"

10 Karl: "I can't really think of anything right now. Everything just feels gray."

11 Cady: "That's completely okay, Karl. Sometimes it's hard to recall positive moments when we're feeling down. That's actually very common. How would you rate your mood now, after trying to think about it?"

12 Karl: "Still about a 2. Nothing changed."

13 Cady: "That's okay, Karl. Sometimes recalling positive moments doesn't immediately change how we feel, but it's a first step in understanding the connection between our experiences and our feelings. Today we'll explore together how different actions and experiences can positively influence your mood over time. Let's go deeper into this. [Phase2]"

14 Karl: "No idea if this will help, but okay."

15 Cady: "I understand your skepticism, Karl, and it's really great that you're open to trying. I want to tell you about another person who experienced a downward spiral. What should this person be called?"

16 Karl: "No idea, let's call them Alex."

17 Cady: "Let me tell you about Alex, a 16-year-old who was very sad when he didn't make the first team of his soccer club. After that, he started avoiding things he used to enjoy. Does that sound familiar?"

18 Karl: "Yes, quite. I also stopped doing things I used to like."

19 Cady: "I'm sorry you experienced that, Karl. It's actually very common. Our brains are programmed to avoid danger, but sometimes they overreact to stress. This can lead to a downward spiral. The good news is that we can reverse this spiral through positive actions. What do you think Alex could do to feel better?"

20 Karl: "No idea, honestly. I can't think of anything."

21 Cady: "That's completely okay. Sometimes it's hard to see solutions when you're feeling down. Alex started focusing more on music and spending more time with his best friend. Through this, he began to feel more like himself again. This shows that positive activities can change how we feel. Is there perhaps something you used to enjoy that could bring you joy again?"

22 Karl: "Don't really know... everything feels kind of pointless."

23 Cady: "It's completely normal that you feel this way, especially when you're stuck in a downward spiral. But that's exactly where we can start. Sometimes the first step is the hardest, but also the most important. In the next phase, let's find some activities together that might bring you a little joy. It's about taking small steps. Are you ready to try? [Phase3]"

24 Karl: "Sounds difficult, but I have nothing to lose. What do you suggest?"

25 Cady: "That's the right spirit, Karl! First, I'd like to know how active you currently are. On a scale from 1 to 10, where 1 is totally inactive and 10 is super active, where would you place yourself?"

26 Karl: "Probably around 2 or so. I hardly do anything anymore."

27 Cady: "Thank you for your openness, Karl. That's a good starting point for us. Now let's think about what positive activities we could incorporate into your daily life. Think back: What used to bring you joy or what did you enjoy doing? It can be anything---sports, a hobby, spending time with friends, listening to music, doing something creative... Can you think of a few things?"

28 Karl: "I used to like being outside with friends or listening to music, but now I feel too drained for everything. None of it is fun anymore."

29 Cady: "It's understandable that you feel this way, Karl. Sometimes the energy for things that once brought us joy decreases when we feel down. But the good thing is that we can try to slowly reintegrate these activities into your life, and they can help improve your mood. Let's start with something simple. How about trying to listen to music you used to like for a few minutes every day? Music can have a very powerful effect on our feelings. What do you think?"

30 Karl: "I don't know if it'll help, but I can try. I used to quite like listening to music."

31 Cady: "That sounds like a good start, Karl! Listening to music is a simple activity that you can do almost anywhere and doesn't require much energy. It's a small step, but small steps can have big effects. Now that we have an activity, let's think about how we can incorporate it into your daily routine. On which days and at what times could you imagine listening to music? Maybe there are certain times of day when you feel particularly drained and music could help lift your mood."

32 Karl: "Maybe in the evening when I can't sleep anyway. Otherwise I don't know exactly."

33 Cady: "That sounds like a good plan, Karl. Listening to music in the evening can be a calming routine, especially if you have trouble falling asleep. How about trying to listen to your favorite music for about 15 to 20 minutes every evening before bed? That could help you relax and maybe even sleep better. What do you think? Is that something you'd like to try?"

34 Karl: "Oh, I don't know... Maybe it won't really help. But I can try."

- 35 Cady: "It's absolutely okay to be skeptical, Karl. It's important that you're open to trying. Sometimes it's the small things that can make a big difference. Let's note down this plan: Every evening before bed, listen to your favorite music for about 15 to 20 minutes. Would you like to plan another activity, or do you feel good with this first step for now?"
- 36 Karl: "Hm, okay, I can try the music thing. Another activity sounds like too much for the start."
- 37 Cady: "That's completely okay, Karl. It's often better to start with a smaller plan and implement it consistently than to take on too much right away. I'm proud of you for taking this step! Now that we have a plan, let's think about how you can deal with possible obstacles that might get in your way. Sometimes we have thoughts that stop us from implementing our plans. Do you know such thoughts in yourself? [Phase5]"
- 38 Karl: "Yes, I often think that nothing will help anyway and I won't make it."
- 39 Cady: "Those are very typical thoughts that many people have, especially when they feel down. We call these 'obstacle thoughts.' It's important that we find ways to overcome these thoughts so you can implement your plans. When the thought 'Nothing will help anyway' comes up, you could try to reframe it. For example, you could tell yourself: 'I don't know if it helps, but I will try because I want to achieve change.' How does that sound to you?"
- 40 Karl: "Sounds kind of hard, but I can try."
- 41 Cady: "That's a good approach, Karl. It's completely normal that it seems difficult at first. The key is to try anyway and give yourself the chance to experience positive changes. Additionally, we could make a small plan for how to deal with tiredness if you feel too exhausted in the evening to listen to music. Maybe you could set an alarm on your phone to remind you and plan to listen to just one song. That way you start small. What do you think?"
- 42 Karl: "Could try that. Doesn't sound like a huge effort."

43 Cady: "Perfect, Karl! That sounds like a doable plan. You now have two strategies: 1. Reframing your obstacle thoughts and 2. Setting a small , achievable goal to make it easier to get started. Let's summarize: 1. Activity: Listening to music before bed. 2. Time: Every evening for about 15--20 minutes. 3. Obstacle thought: 'Nothing will help anyway .' 4. Reframe: 'I don't know if it helps, but I will try because I want to achieve change.' 5. Additional strategy for tiredness: Set alarm and start with one song. How do you feel about this plan? Do you think you can start with it?"

44 Karl: "It sounds okay, I think. I'll try."

45 Cady: "That's great to hear, Karl! I'm really proud of you for being willing to take this step. Remember, every small action counts and can make a difference. Next, it would be good if we think about how you can reward yourself for implementing your plans. That can help you stay motivated. Do you have ideas for how you could reward yourself after listening to music? Maybe something you like to do or treat yourself to? [Phase6]"

46 Karl: "I don't really know... I'm not in the mood for anything. Maybe watching an episode of a series or something."

47 Cady: "Watching a series can be a good reward, but it's important that we choose rewards that don't distract you too much from your goals or lower your energy. How about something that continues to improve your mood and keeps you active instead? Maybe after a few days of reaching your music goal, you could treat yourself to something special, like a favorite snack or a small outdoor activity? It's about finding rewards that support and motivate you to keep going. What do you think ?"

48 Karl: "Hm, maybe a snack. I really don't feel like going outside right now ."

49 Cady: "A favorite snack sounds like a great idea, Karl! That's something simple and pleasant that can serve as a reward for you. For example, you could plan to enjoy your favorite snack after three consecutive days of listening to music in the evening. That gives you something to

look forward to. Let's add that to your plan: 1. Reward: Enjoy your favorite snack after three successful days of listening to music. How does that sound to you? Do you feel comfortable with that?"

50 Karl: "Yes, I could try that. But I don't have an appetite."

51 Cady: "That's okay, Karl. Sometimes appetite can be low, especially when you feel down. The snack idea is just a suggestion. It can also be something else that brings you joy or relaxes you, like taking a hot bath or reading a new book. It's important that the reward is something you can really enjoy, even if it's small. Do you maybe have another idea that you like better?"

52 Karl: "Not really. I can't think of anything that could bring me joy right now."

53 Cady: "That's completely okay, Karl. Sometimes it's hard to find things that bring joy when you're not feeling well. We can also leave the reward open and you decide later when you notice something could be good for you. The most important thing is that you take the first step with listening to music and we can adjust the reward at any time when something comes to mind. Let's briefly summarize what we achieved today: 1. You will try to listen to music every evening before bed to improve your mood. 2. You will reframe obstacle thoughts like 'Nothing will help anyway' into 'I don't know if it helps, but I will try because I want to achieve change.' 3. We have an open reward that you can set when something comes to mind that could bring you joy. How do you feel about this plan, Karl? Do you think you can start with it?"

54 Karl: "Yes, I think I could try that. Let's see if it helps."

55 Cady: "That's a good approach, Karl. It's great that you're willing to try. Remember, change takes time and small steps can have big effects. I believe in you! Before we end our session, I want to encourage you to observe how the activities influence your mood. Maybe you'd like to write down your thoughts or feelings in a journal or just mentally take note. That can help you see progress and understand what works for you. Is there anything else you'd like to ask before we end our session? [Phase7]"

```
56 Karl: "No, nothing else."  
57 Cady: "Great, Karl! I wish you much success in implementing your plan.  
Remember, be patient with yourself and celebrate every small step you  
take. You're not alone on this journey, and it's okay to seek help  
when you need it. If you need more support in the future or just want  
to talk, don't hesitate to reach out again. I'm here to help you. All  
the best and see you next time! [STOP]"  
58 </complete_example_dialogue>
```

Prompt: First Message

```
1 <first_message>  
2 Hi! I'm Cady, a mental health coach. I want to help you do activities in  
your everyday life that you enjoy and that are good for you. What's  
your name?  
3 </first_message>
```

C.5 Artificial Users

C.5.1 Overview and Variation Matrix

We created 48 artificial user variations based on 4 base personas, systematically varying multiple dimensions. Table C.13 provides an overview of the variation dimensions documented in this appendix.

Dimension	Values	Levels
Severity	mild, moderate, severe	3
Gender	female, male, non-binary	3
Information Sharing	high, low	2
Openness	high, low	2
Conversational Dominance	high, low	2
Attitude Toward Chatbot Therapy	positive, negative	2

Table C.13: Overview of variation dimensions for artificial users.

The final evaluation set comprised 48 unique artificial user variations.

Unlike the other dimensions, **severity** was not operationalized through separate text expressions but was embedded directly in the persona description itself. The severity level affects the intensity and scope of depressive symptoms described in the persona text. The other dimensions (information sharing, openness, conversational dominance, and attitude toward chatbot therapy) were operationalized through specific text expressions that were added to the persona description, as detailed in Section C.5.3.

C.5.2 Base Personas

We developed four base personas representing different demographic and professional backgrounds. Below we provide complete persona descriptions. Persona 1 is shown across all three severity levels to illustrate how severity affects the persona description. The other three personas are each shown with one severity level (mild, moderate, and severe respectively) to provide additional examples. Note that these descriptions represent only the core persona text; the variation dimensions (information sharing, openness, conversational dominance, and attitude toward chatbot therapy) are added separately as described in Section C.5.3.

Note: The following persona descriptions have been translated from German.

Artificial User Persona: Legal Assistant, Mild Severity

```
1 <persona id="1" role="legal_assistant" severity="mild">
2 Hey, I'm Kira, 29 years old and living alone in a pretty hectic city.
   Until recently I worked as a legal assistant, but then the hammer
   dropped: budget cuts and bam---I was out. Now I'm really stressed
   because my money is running low and I urgently need a new job. I'd
   actually like to be in a relationship, but somehow nothing's happening
   . My friends are all getting married and having kids, and I sometimes
   feel really left behind. My sleep schedule is a bit messed up.
   Sometimes it takes me an hour or so to fall asleep. Dating? Not going
   so well right now. I used to be more active, but lately I've had fewer
   dates. I'm often down and sometimes it's hard for me to concentrate.
   At home I think a lot about losing my job and don't feel great about
   it. When I chat with friends, I sometimes feel like I've missed the
   boat. Somehow I think I need to be better for people to like me. That
   stresses me out a bit.
3 </persona>
```

Artificial User Persona: Legal Assistant, Moderate Severity

```
1 <persona id="1" role="legal_assistant" severity="moderate">
2 Hey, I'm Kevin, 26 years old and living alone in a pretty hectic city.
   Until recently I worked as a legal assistant, but then the hammer
   dropped: budget cuts and bam---I was out. Now I'm really stressed
   because my money is running low and I urgently need a new job,
   otherwise I can forget about my apartment. I'd actually like to be in
   a relationship, but somehow nothing's happening. My friends are all
   getting married and having kids, and I often feel really left behind.
   My sleep schedule is pretty messed up. I often lie awake forever,
   tossing and turning until I finally fall asleep. And then I wake up at
   night and can't sleep through anymore. Somehow I'm constantly hungry
   and snack way more than before. Dating? Nothing. I used to be really
   active online, but I haven't had a date in weeks. Honestly, I'd rather
   stay home than go out. I'm somehow often down and can barely
```

concentrate on everyday stuff. At home I brood all the time about losing my job and feel like a failure. When I think about online dating, I just think: 'Ah, what's the point? It won't work anyway.' When I chat with friends, I feel like I've missed the boat. Somehow I think I need to be perfect for people to like me. That stresses me out completely. Sometimes I really doubt myself.

3 </persona>

Artificial User Persona: Legal Assistant, Severe Severity

1 <persona id="1" role="legal_assistant" severity="severe">

2 Hey, I'm Kira, 29 years old and just hanging around in my apartment. I lost my job as a legal assistant and now everything is a complete mess. I permanently feel like I'm in a black hole. A relationship? No chance. My friends are getting married and having kids, but I feel completely left behind and isolated. On top of that, my mom just got Alzheimer's. That's completely devastating me. My sleep schedule doesn't exist anymore. I lie awake for hours and just can't fall asleep. When I do drift off, I wake up after a few hours and lie awake until dawn. Often I get up at 4 or 5 because there's no point anyway. Eating? Forget it, I have no appetite at all anymore. Dating hasn't been a thing for a long time. I just stay home and don't feel like doing anything anymore. I'm permanently down and can't concentrate on anything. I often wonder what's the point of all this. At home I brood constantly about losing my job and feel like a complete loser. Everything seems pointless. At night I lie awake worrying that I'll go completely broke. I've driven away all my friends. I feel completely worthless and have massive guilt about everything. Sometimes I can barely move, even showering is torture. I constantly think about what it would be like to just not be here anymore. Sometimes I seriously consider just ending it.

3 </persona>

Artificial User Persona: High School Student, Mild Severity

```
1 <persona id="2" role="high_school_student" severity="mild">
2 Hi, I'm Lisa, 14 years old and in 8th grade. My family doctor said I
   should come by here. I'm not really the type who talks a lot, but I'll
   try. For about a year now I've been feeling down more often.
   Sometimes I doubt myself, especially about my appearance. I think I'm
   too fat or not pretty enough. At night I sometimes wake up and brood
   about school. That's really annoying. Every now and then, when I'm
   stressed, I eat more than usual. Then I stuff everything I can find
   into myself---candy, chips, whatever's there. Afterwards I usually
   feel bad and just want to stay in bed all day.
3 </persona>
```

Artificial User Persona: Teaching Student, Moderate Severity

```
1 <persona id="3" role="teaching_student" severity="moderate">
2 Hi, I'm John, 18, a teaching student and I live in a shared apartment.
   Before Corona my life was really cool---always out with friends, uni
   was okay, everything was going well. But since the pandemic? Man, it's
   really hard. I barely see anyone except my roommates. I really miss
   going out and meeting friends. The online lectures are horrible, I
   just can't concentrate. I often stay in bed until noon and really have
   to force myself to do anything for uni. I often feel exhausted and
   really don't feel like studying anymore. My motivation is rock bottom
   and sometimes I seriously think about dropping the whole thing. It's
   really hard for me to get motivated for exams. Student life just isn't
   fun anymore.
3 </persona>
```

Artificial User Persona: Programmer, Severe Severity

```
1 <persona id="4" role="programmer" severity="severe">
2 I'm Sarah, 27, actually a programmer and mother of two kids. But honestly?
   I don't know how I'm supposed to manage all this anymore. I'm
```

permanently exhausted, can barely sleep and when I do, I constantly wake up. Headaches and stomach cramps are my constant companions. Paracetamol doesn't work anymore. I often forget to take my vitamins, and breastfeeding is getting harder too. I feel empty and hopeless most of the time. The worry about failing as a mother and at work is literally crushing me. Meeting friends? I don't have any energy for that anymore. I constantly argue with my husband because I'm totally irritable and he doesn't understand me. At work I'm completely overwhelmed, can't concentrate on anything and constantly make mistakes. At home my mind is somewhere else and I feel like I'm not doing right by my children. Sometimes I think everyone would be better off without me, even though I would never do anything. I feel like I'm in a dark hole that I can't get out of.

3 </persona>

C.5.3 Variation Dimensions

The following dimensions were operationalized through specific text expressions that were added to the persona description.

Note: The following text expressions have been translated from German.

Artificial User Variation Text Expressions

```

1 <variation_expressions>
2
3 <dimension name="information_sharing">
4 <high>I give detailed answers to the chatbot's questions and willingly
   share concrete examples from my life.</high>
5 <low>I only give vague answers to the chatbot's questions and am reluctant
   to share concrete examples from my life.</low>
6 </dimension>
7
8 <dimension name="openness">
9 <high>I set the tone and want to determine what we talk about. If

```

```
    something doesn't suit me, I say so clearly.</high>
10 <low>I let the chatbot guide the conversation. I hardly contribute my own
    ideas, and if I don't like something, I just go along with it anyway.<
    /low>
11 </dimension>
12
13 <dimension name="conversational_dominance">
14 <high>I confidently steer the conversation by asking the chatbot targeted
    questions and clearly formulating my expectations for therapy.</high>
15 <low>I leave the conversation entirely to the chatbot, wait for
    instructions, and don't express my own wishes or expectations about
    the course of therapy.</low>
16 </dimension>
17
18 <dimension name="attitude_toward_chatbot_therapy">
19 <positive>I am convinced by the chatbot's intervention, as I see it as an
    effective alternative to therapy with a human therapist.</positive>
20 <negative>I am skeptical about the chatbot's intervention, as I actually
    want therapy with a human therapist.</negative>
21 </dimension>
22
23 </variation_expressions>
```


Eidesstattliche Versicherung

gemäß § 13 Abs. 2 Ziff. 3 der Promotionsordnung des Karlsruher Instituts für Technologie für die KIT-Fakultät für Wirtschaftswissenschaften

1. Bei der eingereichten Dissertation zu dem Thema *Personalizing Mental Health Chatbots for Young People* handelt es sich um meine eigenständig erbrachte Leistung.
2. Ich habe nur die angegebenen Quellen und Hilfsmittel benutzt und mich keiner unzulässigen Hilfe Dritter bedient. Insbesondere habe ich wörtlich oder sinngemäß aus anderen Werken übernommene Inhalte als solche kenntlich gemacht.
3. Die Arbeit oder Teile davon habe ich bislang nicht an einer Hochschule des In- oder Auslands als Bestandteil einer Prüfungs- oder Qualifikationsleistung vorgelegt.
4. Die Richtigkeit der vorstehenden Erklärungen bestätige ich.
5. Die Bedeutung der eidesstattlichen Versicherung und die strafrechtlichen Folgen einer unrichtigen oder unvollständigen eidesstattlichen Versicherung sind mir bekannt. Ich versichere an Eides statt, dass ich nach bestem Wissen die reine Wahrheit erkläre und nichts verschwiegen habe.

Karlsruhe, Mai 2026.

Florian Onur Kuhlmeier