

Balancing Personalization and Protocol Fidelity: A Protocol-Guided Dialogue Architecture for LLMs within Cognitive Behavioral Therapeutic Interventions^{*}

B.C.M. Devilee^{1,*}, Bin Yu¹, Simone M. de Droog¹ and Somaya Ben Allouch^{1,2}

¹Amsterdam University of Applied Sciences, Amsterdam, the Netherlands, Wibautstraat 2-4 1091 GM Amsterdam, The Netherlands

²University of Amsterdam, Amsterdam, the Netherlands, Spui 12, 1012 WX Amsterdam, The Netherlands

Abstract

This paper proposes a safer approach to using large language models (LLMs) in behavior change and mental health support systems. LLM-based systems can respond inconsistently when rules and guidelines are included only as instructional prompts within the LLM itself. Instructional prompts influence the dialogue but cannot enforce strict rules. The risk is that the model will drift from the intended intervention, thereby undermining the mechanism of change. This position paper argues that therapeutic protocols should be implemented externally in a system-level orchestration layer that governs what the LLM can do at each phase of an intervention. Established frameworks for designing and evaluating digital BCSS provide a strong conceptual foundation and set standards for describing system structure and intent. However, they offer limited guidance on how to technically enforce behavioral constraints in LLM-based systems. We introduce Protocol-Guided Dialogue, in which an explicitly defined therapeutic protocol structures the interaction through phases, goals, progression criteria, and decision rules, while the LLM generates the dialogue within these constraints. Using Competitive Memory Training (COMET) as an illustrative case, we show how phased protocol logic can be translated into system-level controls, including examples of what the LLM can and cannot do and say in each phase.

Keywords

Mental health, adolescents, GenAI, LLM, LLM-generated dialogue, COMET, CBT, protocol fidelity,

1. Introduction

Adolescent mental health faces increasing challenges worldwide [1, 2]. Many mental health problems begin before the age of 14 and continue into adolescence and early adulthood [3]. Low self-esteem and social avoidance are common in this group and may reinforce one another, limiting positive social experiences and undermining self-confidence [4, 5]. Persistent low self-esteem during adolescence has been identified as a predictor of later depressive symptoms and broader mental health difficulties [4, 5].

Despite the availability of effective interventions, most of which are based on Cognitive Behavioral Therapy (CBT), many adolescents still do not receive adequate support. Limited therapist availability, long waiting lists, and systemic barriers to care leave a substantial number of young people without timely help [2, 6, 7]. Evidence-based interventions include Compassion-Focused Therapy, Ecological Momentary Intervention [8], and Competitive Memory Training (COMET) [9].

At the same time, young people are already turning to digital tools for support. Emerging evidence shows that individuals, including those experiencing anxiety, depression, or crisis, are using Large Language Models (LLMs) as informal mental health support tools [10]. Chatbots are available at any

BCSS 2026 The 14th International Workshop on Behavior Change Support Systems, March 10, 2026, Hakodate, Japan.

^{*} Corresponding author

^{*}B.C.M. Devilee

✉ b.c.m.devilee@hva.nl (B.C.M. Devilee); b.yu@hva.nl (B. Yu); s.m.de.droog@hva.nl (S. M. d. Droog); s.ben.allouch@hva.nl (S. B. Allouch)

ORCID 0009-0004-5173-7982 (B.C.M. Devilee); 0000-0002-3128-7441 (B. Yu); 0000-0002-2899-7143 (S. M. d. Droog);

0000-0002-3520-4016 (S. B. Allouch)



© 2026 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

time, which makes them accessible and low-threshold forms of support. Recent advances in LLM technology have further expanded the possibilities for digital mental health interventions [11].

Generative models produce more natural, context-sensitive dialogue and create new possibilities for personalizing tone and framing in ways that earlier rule-based systems cannot achieve. This creates new opportunities to make digital interventions more engaging and personally relevant. However, this flexibility also introduces risks. LLM-generated dialogue is adaptive and probabilistic. The same question can lead to slightly different answers depending on context. In therapeutic settings where interventions depend on careful sequencing and clearly defined progression, variability can become consequential. A response that appears supportive can deviate from the intended therapeutic structure [11, 12].

For this reason, many existing digital mental health systems continue to rely on rule-based or pre-scripted dialogue structures [11]. These systems follow predefined flows and decision trees. They provide predictability and control, which is beneficial in psychological and psychiatric care: therapeutic protocols are structured and phased for a reason. Certain techniques are introduced only when specific conditions are met. Phasing acts as a safety mechanism and a quality safeguard. It prevents premature or inappropriate intervention steps and makes treatment delivery transparent and evaluable [11, 13]. Many evidence-based therapies allow tailoring within predefined structures [13]. Also, modular or process-based approaches show that flexibility can be aligned with clinical responsibility [14, 15]. Yet rule-based systems often feel rigid and less natural, especially for young people accustomed to digital systems that employ natural, user-tailored dialogue [16, 12].

Therefore, the main challenge is to design digital mental health interventions that employ natural dialogue to provide engaging, personalized interactions while maintaining therapeutic integrity and protocol fidelity. Protocol fidelity refers to the degree to which an intervention is delivered as intended by its developers, including adherence to prescribed phases, therapeutic components, sequencing, and progression criteria.

Established frameworks for designing digital Behavior Change Support Systems (BCSS), most notably the Persuasive Systems Design (PSD) model [17], provide a strong conceptual basis for shaping interactive behavior change interventions and have set the standard for BCSS designers to describe software structure, function, and intention. However, they offer limited guidance on how behavioral constraints should be enforced within the technical architecture of large language model-based systems. FAITA [18], an evaluation framework for LLM use in BCSS, emphasizes the importance of data control and system safety. However, because it functions as an assessment tool, it does not specify how these requirements should be implemented architecturally within an intervention.

This position paper argues that therapeutic phase logic should not be embedded solely within the LLM. Instead, it should be externalized into a system-level orchestration layer that explicitly governs phase transitions, permissible actions, and progression criteria, in this way, generative dialogue can remain flexible, while the therapeutic structure remains stable and enforceable.

2. The Case for Externalized Protocol Control

LLMs are typically governed through prompt engineering, tone guardrails, and output filtering. These techniques aim to steer responses and prevent undesirable content. However, such governance operates primarily at the level of surface output and does not provide enforceable control over therapeutic process logic. Large language models do not apply fixed rules. They generate responses probabilistically based on the evolving conversational context. Instructions embedded in prompts influence generation but do not function as binding constraints. With each new user exchange, the interpretation of prior instructions may shift. In longer or repeated interactions, models can reinterpret or gradually deprioritize earlier guidance. Over time, this conversational drift may lead to deviations from the originally intended structure [14, 16, 12]. In protocol-based therapeutic interventions, such deviations are consequential. Most evidence-based treatments rely on carefully sequenced exercises, controlled emotional activation, and clearly defined progression. Small changes in timing, framing, or intensity may appear minor

at the level of individual responses but can undermine the intended mechanism of change across phases [16, 12, 15, 13]. This structural mismatch motivates our position: protocol and process rules in LLM-enabled behavior change support systems (BCSS) should not be embedded solely within the generative model. Instead, therapeutic phase logic must be externalized into a system-level orchestration layer where state transitions, allowable actions, and progression criteria are explicitly represented and enforceable. Established BCSS frameworks, such as the Persuasive Systems Design (PSD) model [17], provide conceptual guidance for designing behavior change systems. However, they offer limited direction on how to embed therapeutic rules within LLM-based architectures [11, 14, 12]. To address this gap, this paper conceptualizes Protocol-Guided Dialogue as a system-level architecture that separates language generation from protocol control.

3. Treatment Fidelity and Architectural Governance

Treatment fidelity research emphasizes that intervention execution must be monitored to preserve internal validity and interpretability. Variability in delivery can undermine outcomes without being visible at the surface level [15]. The LLM effectively functions as a component of the “care provider”, generating the moment-to-moment delivery of the intervention [11, 16]. Without structured governance, model variability can lead to therapeutic drift: deviations in timing and intensity may undermine the integrity of the intervention, [12, 15, 19]. If the core components of an intervention are phase-dependent, fidelity cannot be reduced to producing “safe-sounding” or “empathetic” text. It requires adherence to phases and progression criteria to preserve the intended mechanism of change. Existing governance strategies, including prompt engineering, retrieval augmentation, and output filtering, operate at the content level. They do not encode therapeutic phase logic or progression criteria in an enforceable manner [11, 14, 12, 19].

They shape surface characteristics or remove disallowed outputs, but they do not guarantee sustained phase coherence or protocol sequencing over time. LLMs lack explicit state-transition mechanisms, making internal decision processes difficult to inspect or control [14]. Implicit constraint enforcement within prompts makes it impossible to guarantee phase coherence, progression control, and evaluability over time [11, 16, 19]. This risks a non-transparent decision layer within BCSS architectures.

The consequences of this architectural opacity are increasingly visible in empirical evaluations of youth-facing generative psychotherapy chatbots and companion systems. They report high accessibility and conversational fluency, yet weak therapeutic grounding, limited risk monitoring, and inconsistent crisis management [16, 12, 15].

This discrepancy highlights a central governance risk: conversational competence can create a strong impression of therapeutic adequacy without ensuring adherence to the intended mechanism of change. Individual responses may appear empathic and contextually appropriate, but fidelity failures and safety drift can accumulate across turns without being immediately visible at the surface level [16, 20, 15]. Taken together, these points establish that the clinical reliability of LLM-enabled BCSS depends on enforceable protocol logic on a system-level orchestration layer. Without such architectural separation, guarantees of phase integrity, safety, and evaluability remain fragile. This separation aligns with emerging evaluation frameworks for AI-enabled mental health care, which emphasize structured oversight, accountability mechanisms, and explicit control [14, 18].

4. Operationalization: Using COMET as an Illustrative Case

This section presents the formal design output of the proposed approach: a structured translation of phased protocol logic into system-level governance rules for generative dialogue. We present a structured mapping from therapeutic phases (illustrated using COMET) to phase-appropriate generative behavior, and a governance framework that includes protocol state tracking, phase-specific allowable actions, formal transition criteria, and structured logging for evaluation and controllability.

4.1. Phased Protocol Logic: COMET as a Design Reference

Competitive Memory Training (COMET) serves as the reference model for this translation [21]. COMET is a protocol-based cognitive behavioral intervention with seven sequential phases, each defined by specific goals and technique. It is aimed at reducing negative self-evaluations, and the intervention activates and reinforces positive self-related representations. This careful phasing is central to COMET’s effectiveness; progression depends on meeting phase-specific criteria. COMET is used here as an analytical reference for reasoning not as a direct implementation template. Table I summarizes the seven phases, their primary functions, and associated techniques.

Table 1
COMET: Phases, Functions, and Techniques

Phase	Function	Technique
1. Identification	Identifying negative self-representations	Formulation of self-beliefs; problem context
2. Positive self-representation	Building competing memory networks	Positive autobiographical memories; self-talk
3. Strengthening	Enhancing emotional salience	Imagery; sensory cues; music
4. Embodiment	Affective grounding	Body posture; facial expression
5. Counter-conditioning	Application under emotional challenge	Positive self-representation paired with insecurity-triggering situations
6. Attention training	Correcting attentional bias	Exercises directing attention toward positive stimuli
7. Transfer	Generalization to daily life	Practice in everyday contexts

4.2. From Therapeutic Phasing to Phase-Aware LLM Architecture

The COMET framework illustrates how phased protocol logic can be translated into Protocol-Guided Dialogue by binding LLM behavior to the active intervention phase. Responses are constrained to support the current therapeutic objective while allowing adaptive variation in language, tone, and narrative framing. Encoding protocol phases as system-level control logic shifts governance from content-level filtering to process-level regulation [13]. The system-level orchestration layer defines explicit phase boundaries and progression conditions within which the LLM may operate. Generative flexibility is preserved, but only within therapeutically permitted ranges. Evidence-based phased interventions can thus inform the design of constrained, phase-aware LLM architectures in BCSS. Table 2 summarizes the relationship between COMET phases and their implications for LLM behavior.

Table 2
COMET Principles and Implications for LLM Behavior

Function	Technique	LLM Behavior
Identifying negative self-representations	Formulation of self-beliefs; problem context	Provoke baseline self-beliefs and situations; reflect and summarize; avoid reframing; no exposure or challenge.
Building competing memory networks	Positive autobiographical memories; self-talk	Elicit recall of positive memories; scaffold self-talk; keep focus on evidence-based positive cues; avoid linking to difficult triggers.
Enhancing emotional salience	Imagery; sensory cues; music	Intensify imagery prompts and sensory detail; reinforce affect; maintain containment; monitor distress.
Affective grounding	Body posture; facial expression	Prompt posture and mimic cues; anchor to the current positive state; keep language concrete and brief.
Application under emotional challenge	Positive self-representation paired with insecurity-triggering situations	Introduce insecurity-triggering situations only if criteria are met; apply stop rules for escalation.
Correcting attentional bias	Exercises directing attention toward positive stimuli	Guide attention to adaptive cues; interrupt rumination loops; use short repeated prompts.
Generalization to daily life	Practice in everyday contexts	Provoke examples from daily contexts; structure reflection; consolidate learning; define completion criteria.

4.3. Translating Phased Protocol Logic into System Architecture

To implement protocol governance in LLM-enabled BCSS, it's important to specify control at different levels. Because LLMs compute the probability of the next token given all previous tokens in the context window they do not reason about goals, roles, or intentions. After each token the LLM selects the most likely continuation according to patterns learned during training. The LLM has no internal representation of "therapy," "safety," or "protocol fidelity" as stable commitments. These concepts only exist as statistical regularities in language. If the context suggests that a refusal is likely, the model produces a refusal. If the context shifts, the probability distribution shifts with it. In architectural terms role stability cannot be assumed. The governance level 0 therefore fixes the functional role externally. It defines what kind of system the LLM is allowed to instantiate before generation begins. Governance Level 1 Phase governance is to let the system know which intervention phase is active, so it is possible to set the rules for permitted actions within that phase. Level 2 I the state governance makes it possible to make decisions to go to the next phase. Level 3 is built in for monitoring and controllability to enforce logging of decisions and transitions. The mechanisms were derived directly from the structure of the COMET protocol. Each therapeutic phase defines a specific psychological objective, a set of allowed techniques, and clear boundaries regarding what must not occur. Table 2 outlines these governance layers and shows how design choices affect intervention progress, allowed behaviors, and evaluation.

Table 3
Levels of Governance, and Enforced Behavior

Governance Level	Governance focus	Mechanism	Influences	Enforces
Level 0: Global Therapeutic Role	What kind of system is this?	Allowed intervention functions, personalization scope boundaries, and global stop rules	All phases	Domain containment, and therapeutic identity
Level 1: Phase Governance	What is allowed in this phase?	Phase-specific permissible actions and containment rules	LLM output within active phase	Phase coherence
Level 2: State Governance	When may progression occur?	Protocol state tracking and transition criteria	Phase transitions	Sequencing integrity
Level 3: Audit and Logging	How is this monitored?	Structured logging of decisions and transitions	Observability and controllability	Fidelity evaluation

This layered governance structure enables structured scaffolding, controlled repetition, and phase-consistent progression. By separating role definition, phase constraints, and state transitions, the architecture preserves therapeutic sequencing while allowing contextual responsiveness. At the same time, structured logging and explicit state control enhance system observability and accountability, consistent with established principles in behavior change theory and BCSS design [13, 15, 22].

5. Counterarguments and Design Challenges

5.1. More Complexity

One may argue that finite-state or tightly scripted systems already provide clear state control and are easier to validate in clinical contexts. Indeed, many healthcare conversational agents rely on rule-based or scripted dialogue management precisely because such systems offer predictability and clearer validation pathways [12]. Introducing system-level orchestration for generative interaction increases architectural complexity and computational overhead. However, this additional complexity reflects a shift in functional scope. Generative systems without explicit orchestration already exhibit latent complexity due to probabilistic behavior and opaque internal decision processes [14, 18]. Externalizing protocol logic into an orchestration layer restructures the BCSS into explicitly auditable and controllable components. In this way, orchestration functions as a strategy for managing complexity: it enhances interpretability and governance of system behavior.

5.2. Limited Flexibility

A second concern targets the enforcement of guardrails. Excessive external guardrails can limit generative flexibility and reduce contextual adaptivity when applied rigidly. Empirical findings in the LLM literature indicate that imposing a single dominant external selection criterion does not necessarily yield the highest quality response, which suggests that the rigid application of an external constraint can exclude contextually appropriate but less conforming outputs [19]. Conversational competence and perceived responsiveness are important predictors of engagement in youth facing digital mental health systems [16, 23]. Systems that feel overly directive may undermine trust and sustained participation. However, therapeutic autonomy is not equivalent to unrestricted conversational freedom. Behavior change theory emphasizes guided progression, scaffolding, and structured sequencing as necessary conditions for safe and effective intervention [20]. Large Language Models, Chorpita Modularity Design Application 2005, Wanniarachchi Personalization Variables Digital 2025. The

orchestration layer should therefore operate as meta-adaptive governance: it regulates when and how guardrails are invoked to preserve therapeutic coherence while maintaining contextual responsiveness.

5.3. Governance versus clinical effectiveness

Beyond questions of architectural trade-offs, a more fundamental concern addresses whether structural governance ensures therapeutic benefit. Architectural safeguards can enhance transparency and fidelity, but they do not in themselves guarantee clinical effectiveness [13, 15]. Therapeutic success depends on additional factors, including tailoring, contextual relevance, and mechanism activation [15]. The proposed orchestration architecture should therefore be understood as a necessary structural condition for clinical reliability, rather than as sufficient evidence of therapeutic efficacy.

6. Conclusion and Recommendations for Future Research

The development of LLM has expanded opportunities for tailoring mental health BCSS while simultaneously introducing risks related to process fidelity, safety, and accountability. Protocol-Guided Dialogue addresses this tension by externalizing the LLM's therapeutic phase logic, decision criteria, and guardrails into explicit systemlevel governance. In doing so, it reframes architectural control as a structural condition for clinically responsible generative systems. Looking ahead, research on AI-enabled mental health interventions must move beyond optimizing conversational quality alone and address how therapeutic intent is formally encoded, monitored, and evaluated over time. As generative models increasingly mediate care-related interactions, the transparency of intervention logic and the observability of process progression become central design concerns. Advancing this agenda requires interdisciplinary collaboration among clinicians, behavior change researchers, interaction designers, and AI system developers to develop shared architectural standards for governing generative therapeutic interactions.

Declaration on Generative AI

Generative AI tools were used in a supportive capacity during the preparation of this paper. Elicit was used to provide an overview of key scientific publications. ChatGPT and Grammarly were used to improve linguistic clarity, ensure consistency in terminology throughout the manuscript, and perform grammar and spelling checks. After using these tools, the authors reviewed and edited the content as needed and take full responsibility for the publication's content.

References

- [1] A. Tarasenko, G. Josy, H. Minnis, J. Hall, A. Danese, J. Y. Lau, S. Cortese, A. Stringaris, C. Redlich, D. Ougrin, Mental health of children and young people in the WHO Europe region, *The Lancet Regional Health - Europe* 57 (2025) 101459. doi:10.1016/j.lanepe.2025.101459.
- [2] Comprehensive Mental Health Action Plan 2013-2030, Technical Report, World Health Organization, Geneva, 2021.
- [3] H. Gu, P. Zhang, J. Li, The effect of self-esteem on depressive symptoms among adolescents: The mediating roles of hope and anxiety, *Humanities and Social Sciences Communications* 11 (2024) 1–6. doi:10.1057/s41599-024-03249-1.
- [4] K. Langford, K. McMullen, L. Bridge, L. Rai, P. Smith, K. A. Rimes, A cognitive behavioural intervention for low self-esteem in young people who have experienced stigma, prejudice, or discrimination: An uncontrolled acceptability and feasibility study, *Psychology and Psychotherapy: Theory, Research and Practice* 95 (2022) 34–56. doi:10.1111/papt.12361.
- [5] M. Masselink, E. Van Roekel, A. J. Oldehinkel, Self-esteem in Early Adolescence as Predictor of Depressive Symptoms in Late Adolescence and Early Adulthood: The Mediating Role

- of Motivational and Social Factors, *Journal of Youth and Adolescence* 47 (2018) 932–946. doi:10.1007/s10964-017-0727-z.
- [6] On My Mind Promoting, Protecting and Carings for Children’s Mental Health, Technical Report, UNICEF, My Mind, My and others, 2021.
- [7] C. Barbui, J. Alonso, D. Chisholm, S. Evans-Lacko, R. C. Keynejad, L. Lazeri, N. Miah, Z. Valuckiene, C. Gastaldon, Mental health service coverage and gaps among adults in Europe: A systematic review, *The Lancet Regional Health - Europe* 57 (2025) 101458. doi:10.1016/j.lanepe.2025.101458.
- [8] U. Reininghaus, M. Daemen, M. R. Postma, A. Schick, I. Hoes-van Der Meulen, N. Volbragt, D. Nieman, P. Delespaul, L. De Haan, M. Van Der Pluijm, J. J. F. Breedvelt, M. Van Der Gaag, R. Lindauer, J. R. Boehnke, W. Viechtbauer, D. Van Den Berg, C. Bockting, T. Van Amelsvoort, Transdiagnostic Ecological Momentary Intervention for Improving Self-Esteem in Youth Exposed to Childhood Adversity: The SELFIE Randomized Clinical Trial, *JAMA Psychiatry* 81 (2024) 227. doi:10.1001/jamapsychiatry.2023.4590.
- [9] K. Korrelboom, T. IJdema, A. Karreman, M. Van Der Gaag, The Effectiveness of Transdiagnostic Applications of Competitive Memory Training (COMET) on Low Self-Esteem and Comorbid Depression: A Meta-analysis of Randomized Controlled Trials, *Cognitive Therapy and Research* 46 (2022) 532–543. doi:10.1007/s10608-021-10286-6.
- [10] T. Rousmaniere, S. B. Goldberg, J. Torous, Large language models as mental health providers, *The Lancet Psychiatry* 13 (2026) 7–9. doi:10.1016/S2215-0366(25)00269-X.
- [11] Z. Guo, A. Lai, J. H. Thygesen, J. Farrington, T. Keen, K. Li, Large Language Models for Mental Health Applications: A Systematic Review, *JMIR mental health* 11 (2024) e57400.
- [12] L. Laranjo, A. G. Dunn, H. L. Tong, A. B. Kocaballi, J. Chen, R. Bashir, D. Surian, B. Gallego, F. Magrabi, A. Y. S. Lau, E. Coiera, Conversational agents in healthcare: A systematic review, *Journal of the American Medical Informatics Association* 25 (2018) 1248–1258. doi:10.1093/jamia/ocy072.
- [13] B. F. Chorpita, E. L. Daleiden, J. R. Weisz, Modularity in the design and application of therapeutic interventions, *Applied and Preventive Psychology* 11 (2005) 141–156. doi:10.1016/j.appsy.2005.05.002.
- [14] M. Rahsepar Meadi, T. Sillekens, S. Metselaar, A. Van Balkom, J. Bernstein, N. Batelaan, Exploring the Ethical Challenges of Conversational AI in Mental Health Care: Scoping Review, *JMIR Mental Health* 12 (2025) e60432. doi:10.2196/60432.
- [15] A. J. Bellg, B. Borrelli, B. Resnick, J. Hecht, D. S. Minicucci, M. Ory, G. Ogedegbe, D. Orwig, D. Ernst, S. Czajkowski, Treatment Fidelity Workgroup of the NIH Behavior Change Consortium, Enhancing Treatment Fidelity in Health Behavior Change Studies: Best Practices and Recommendations From the NIH Behavior Change Consortium., *Health Psychology* 23 (2004) 443–451. doi:10.1037/0278-6133.23.5.443.
- [16] K. Sobowale, D. K. Humphrey, S. Y. Zhao, Evaluating Generative AI Psychotherapy Chatbots Used by Youth: Cross-Sectional Study, *JMIR Mental Health* 12 (2025) e79838. doi:10.2196/79838.
- [17] H. Oinas-Kukkonen, A foundation for the study of behavior change support systems, *Personal and Ubiquitous Computing* 17 (2013) 1223–1235. doi:10.1007/s00779-012-0591-5.
- [18] A. Golden, E. Aboujaoude, The Framework for AI Tool Assessment in Mental Health (FAITA - Mental Health): A scale for evaluating AI -powered mental health tools, *World Psychiatry* 23 (2024) 444–445. doi:10.1002/wps.21248.
- [19] X. Chen, R. Aksitov, U. Alon, J. Ren, K. Xiao, P. Yin, S. Prakash, C. Sutton, X. Wang, D. Zhou, Universal Self-Consistency for Large Language Model Generation (2023). doi:10.48550/arXiv.2311.17311. arXiv:2311.17311.
- [20] E. M. Bender, T. Gebru, A. McMillan-Major, S. Shmitchell, On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?, in: *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, ACM, Virtual Event Canada, 2021, pp. 610–623. doi:10.1145/3442188.3445922.
- [21] K. Korrelboom, COMET voor negatief zelfbeeld. Handleiding voor therapeuten, *Protocolen voor*

de GGZ, Springer Nature, 2023.

- [22] H. Oinas-Kukkonen, M. Harjumaa, A Systematic Framework for Designing and Evaluating Persuasive Systems, in: H. Oinas-Kukkonen, P. Hasle, M. Harjumaa, K. Segerståhl, P. Øhrstrøm (Eds.), *Persuasive Technology*, volume 5033, Springer, 2008, pp. 164–176. doi:10.1007/978-3-540-68504-3_15.
- [23] V. U. Wanniarachchi, C. Greenhalgh, A. Choi, J. R. Warren, Personalization variables in digital mental health interventions for depression and anxiety in adolescents and youth: A scoping review, *Frontiers in Digital Health* 7 (2025) 1500220. doi:10.3389/fdgth.2025.1500220.