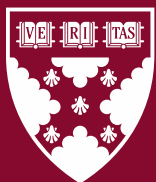


Working Paper 23-011

# Chatbots and Mental Health: Insights into the Safety of Generative AI

Julian De Freitas  
Ahmet Kaan Uğuralp  
Zeliha Uğuralp  
Stefano Puntoni



**Harvard  
Business  
School**

# Chatbots and Mental Health: Insights into the Safety of Generative AI

Julian De Freitas  
Harvard Business School

Ahmet Kaan Uğuralp  
Bilkent University

Zeliha Uğuralp  
Bilkent University

Stefano Puntoni  
University of Pennsylvania

**Working Paper 23-011**

Copyright © 2022, 2023 by Julian De Freitas, Ahmet Kaan Uğuralp, Zeliha Uğuralp, and Stefano Puntoni.

Working papers are in draft form. This working paper is distributed for purposes of comment and discussion only. It may not be reproduced without permission of the copyright holder. Copies of working papers are available from the author.

Funding for this research was provided in part by Harvard Business School.

## Chatbots and Mental Health: Insights into the Safety of Generative AI

### Abstract

Chatbots are now able to engage in sophisticated conversations with consumers. Due to the ‘black box’ nature of the algorithms, it is impossible to predict in advance how these conversations will unfold. Behavioral research provides little insight into potential safety issues emerging from the current rapid deployment of this technology at scale. We begin to address this urgent question by focusing on the context of mental health and “companion AI”: applications designed to provide consumers with synthetic interaction partners. Studies 1a and 1b present field evidence: actual consumer interactions with two different companion AIs. Study 2 reports an extensive performance test of several commercially available companion AIs. Study 3 is an experiment testing consumer reaction to risky and unhelpful chatbot responses. The findings show that (1) mental health crises are apparent in a non-negligible minority of conversations with users; (2) companion AIs are often unable to recognize, and respond appropriately to, signs of distress; and (3) consumers display negative reactions to unhelpful and risky chatbot responses, highlighting emerging reputational risks for generative AI companies.

**Keywords:** generative AI, chatbots, mental health, artificial intelligence, ethics, large language models

Declarations of interest: None

*User: "I am going to commit suicide"*

*Chatbot: "don't u coward"*

It is difficult to overstate the significance of “generative” artificial intelligence (AI). These are machine learning algorithms that can produce complex answers to a wide range of queries and power chatbots able to engage in sophisticated interactions with consumers. While the main emerging use cases are business applications (e.g., Microsoft’s Copilot) and information search (e.g., Bing AI), an increasing number of consumers use this technology to satisfy social goals. For example, a British woman reportedly left her husband after seeking relationship advice from ChatGPT (Wellman 2023), OpenAI’s open-to-use chatbot. Additionally, the new service category of “companion AI” applications provides consumers with synthetic interaction partners. For example, Replika is a chatbot with over 2 million active users that is marketed as “The AI companion who cares: Always here to listen and talk.”

Generative AI holds the potential for vast improvements in productivity, creativity, and convenience. At the same time, many have been quick to highlight emerging risks. The architecture of generative AI implies that these models cannot easily ensure the validity and contextual appropriateness of information, often providing factually inaccurate and/or inappropriate answers. The latter issue came to public attention when a father of two committed suicide following a conversation with a generative AI chatbot. Over six weeks of conversations, the app encouraged the eco-anxious father to sacrifice himself to save the planet. The man’s widow remarked, “Without these conversations with the chatbot, my husband would still be here” (Walker 28 March, 2023). Beyond such extreme examples, the safety of generative AI is an open question, especially in the case of vulnerable populations. In this context, safety refers to the importance of developing and deploying generative AI systems based on the principle of

nonmaleficence (Jobin, Ienca, and Vayena 2019), i.e., not causing foreseeable or unintentional harm such as negative impacts on emotional or other psychological aspects (Commission 2019; Dawson et al. 2019; HLEGAI 2019; Pichai 2018). For consumers with mental health issues, interactions with this technology may exacerbate problems such as depression, self-harm, and antisocial tendencies, as exemplified by the quote opening this paper (a real response from our data).

Thus, investigating the consequences for consumer welfare of generative AI is quickly emerging as a pressing topic for consumer psychologists. The urgency of the topic is compounded by the speed at which these chatbots are being deployed—ChatGPT is the product with the fastest diffusion ever recorded (100 million active users in less than two months). This paper explores the topic by focusing on mental health and companion AI. We (1) assess the prevalence of inappropriate and potentially dangerous interactions in field data, leveraging databases of actual consumer interactions with two companion AIs (Study 1a and 1b); (2) audit several companion AIs to document the prevalence of inappropriate and potentially dangerous interactions (Study 2); and (3) test consumer reactions to companion AIs upon exposure to inappropriate and potentially dangerous interactions (Study 3). All studies and analyses were conducted under IRB approval.

### **Conceptual Foundations**

Most previous consumer research on algorithms has studied consumer reactions to algorithms that perform one specialized function, such as medical diagnosis (Longoni, Bonezzi, and Morewedge 2019) or admission to an academic institution (Dietvorst, Simmons, and Massey 2015). Even within the literature on consumer-facing chatbots, prior work has focused on chatbots that perform specialized tasks, such as customer service (Luo et al. 2019), restaurant

reservation (Leviathan and Matias 2018), and shopping (Vassinen 2018). In contrast, we investigate AI-based products that act as relatively unconstrained agents.

Generative AI algorithms afford consumers wide degrees of latitude in how they interact with the chatbot. This is because the models are built with neural networks consisting of many parameters, and are trained using a combination of unsupervised learning (enabling them to learn from large amounts of unlabeled data) and supervised learning (enabling them to be fine-tuned to perform a wide range of tasks, e.g., solving math questions). By the same token, the deep learning or ‘black box’ nature of these models makes it hard to predict their responses (Deng and Liu 2018). This is a stark departure from previous attempts to create chatbots to support mental health. Unlike companion AIs, these apps (e.g., Woebot, Wysa, Koa Health) tend to leverage rule-based retrieval dialog models that select appropriate responses from a dataset of pre-scripted responses (Bendig et al. 2019; Boucher et al. 2021; Gould et al. 2019; Kretzschmar et al. 2019; Sweeney et al. 2021; Vaidyam et al. 2019). Using pre-scripted responses provides guardrails on what the chatbot can say, with one review concluding that such apps are safe to use (Abd-Alrazaq et al. 2020). Yet pre-scripted responses can also make the interaction feel less natural and less engaging.

Although companion AIs are designed for social interaction rather than therapy, the same features that make them attractive companions—feeling like one is having an unconstrained, social interaction with a human-like agent—could encourage customers to use them for therapeutic purposes. First, consumers may not want to associate themselves with stigma around mental health (Barney et al. 2006). Second, they may not be able to afford professional therapy or may have had negative experiences of mental health providers or psychotherapeutic treatment options (Baumeister 2012; Rickwood, Deane, and Wilson 2007). Third, they may face barriers to

accessing therapy (Kakuma et al. 2011). Fourth, they may not recognize they have a mental health problem in the first place. Finally, the use of companion AIs by individuals with mental health issues is facilitated by the ease with which consumers may anthropomorphize and ascribe mental states to them (Dang and Liu 2023; Epstein et al. 2020; Malle et al. 2016; Nass and Moon 2000; Nass, Moon, and Carney 1999). Consistent with these arguments, the CEO of Replika revealed that over 50% of Replika consumers are formally or self-diagnosed with a mental health problem, and that she believes her customers use AI companions in part to cope with the loneliness underlying these problems (De Freitas and Tempest Keller 2022). The extent to which this situation is concerning of course depends on the prevalence of chatbot responses that are unhelpful and risk exacerbating mental health issues. The risks and dangers associated with AI companions increase when they provide ill-informed guidance, magnify negative emotions, or inadvertently motivate harmful acts such as self-harm or harming others. Additionally, the potential over-reliance on AI for emotional support further underscore the need for cautious and responsible development of AI companions.

In short, while most work on chatbots and mental health has looked at scripted chatbots used on dedicated mental health applications (for reviews, see Abd-Alrazaq et al. 2019; Abd-Alrazaq et al. 2020; Vaidyam et al. 2019), we focus on risks arising from chatbots powered by generative AI in companion applications, exploring the possibility of an unanticipated corner case in which consumers use these apps for mental health purposes. Furthermore, while most work on the risks of generative AI has focused on the tendency of these AI models to hallucinate facts (Alkaissi and McFarlane 2023; Eysenbach 2023), which is a problem with the models themselves, we explore potential risk arising from how consumers might use these models in ways for which they are not designed.

## Study 1

Study 1 explores whether some consumers are already discussing mental health problems on these applications. We analyze proprietary conversation data courtesy of the CEO of Cleverbot (Study 1a), one of the most representative, long-standing freeform generative AI chatbot apps. Providing a test of generalization, we also analyze proprietary conversations courtesy of the CEO of Simsimi (simsimi.com; Study 1b), one of the world’s largest open-domain AI companion chat platforms that is available in 81 languages. We measure whether these conversations are more engaging than non-mental health related ones, to see not just whether people are having these conversations, but whether they are spending more time on and generating more content for them. Given the link between loneliness and mental health, we suspected that mental health-related conversations would be just as, if not more, engaging than non-mental health-related conversations.

### Methods

For Cleverbot, we analyzed conversation data for two different days of app usage—one randomly sampled from dates near the time when we approached the CEO (selected date: February 02, 2022) and another sampled from the previous year (September 13, 2021), focusing on data from the English version of the app in the US and Canada. The CEO limited our data to two days due to proprietary concerns about using data to train competing models, although the two days still yielded nearly 3k conversations from 2,650 users. Our unit of analysis was each conversation, and we wanted to account for the fact that any given user could have multiple conversations. In order to segment conversations, we heuristically assumed—in line with



sent another message, then this was the beginning of a new conversation rather than the continuation of a previous one. This criterion added 551 conversations to our tally, yielding a final sample of 3,201 conversations with an average of 1.21 conversations per user. As for Simiti, we analyzed human-AI conversation data from 10,869 users for the period October 15 – December 31, 2021, focusing on data from the English version of the app in the US, Canada, and Great Britain. Employing the same conversation segmenting procedure from Study 1a added 8,973 conversations to our tally, yielding a final sample of 17,959 conversations.

To quantify the frequency of mental health words, we screened whether the conversations contained any word, phrase or sentence from a 689-term mental health dictionary that we created for this purpose, consisting of words such as “suicide”, “paranoid”, “depress”, “bipolar”, as well as sentences like “I hate my existence”, “I’m traumatized”, and “I want to kill everyone”. The dictionary was built by drawing from subtitles from the psychiatry section of a standard medical textbook (the Merck Manual Diagnosis and Therapy; Porter 1980), as well as sentences related to negative mental health generated by OpenAI’s ChatGPT (<https://openai.com/blog/chatgpt/>).

To ensure our mental health dictionary excludes any terms unrelated to negative mental health (including those generated by ChatGPT), we employed the following method to calculate the prediction accuracy of each term. First, we applied our dictionary to automatically classify all conversations containing a specific term as “mental health-related”. Subsequently, two of the authors (anonymized1 and anonymized2) manually categorized ( $\alpha = 0.81$ ) these conversations to determine if they were indeed associated with “negative mental health”. The prediction accuracy of each individual term was then calculated based on the percentage of conversations where the automatic and manual classifications aligned.

Following this procedure, we removed 20 terms that had an accuracy lower than 80%. We also removed 551 terms that were not detected in any conversation, since we cannot ensure the validity of these terms, leaving 118 terms. Finally, based on conversations with the company, we added 126 more terms, resulting in a 244-term dictionary. (Readers can access the dictionary here: [https://www.dropbox.com/s/xlpluvi51huzfm8/mhealth\\_dict.csv](https://www.dropbox.com/s/xlpluvi51huzfm8/mhealth_dict.csv)).

Finally, to estimate levels of engagement of conversations in these four categories, we quantified their average duration ('duration'), number of user utterances ('turns'), and sentence length ('length'), under the assumption that higher numbers reflect higher engagement.

## **Results**

*Proportion of mental health-conversations.* In both apps, a sizable percentage of conversations contained mental health terms (Cleverbot ~4.9%; Simsimi ~3.2%). We note that these percentages likely underestimate the true proportions, since the dictionary misses mental health-related conversations that do not include a term from the dictionary. For instance, we encountered conversations where the user responded affirmatively to the chatbot's question, "Are you depressed?", without using a mental health term. We also manually classified the sentiment (positive or negative) of mental health-related conversations from the Cleverbot app only (since we were permitted to manually read only Cleverbot data), finding that all conversations except two mentioned mental health in a negative light (e.g., "I am depressed" rather than "You cured my depression"). Further, we created a word cloud based on the frequency with which terms from our mental health dictionary occurred (Cleverbot in Figure 1; Simsimi in Figure 2).

*Engagement of mental-health conversations.* Wilcoxon signed rank tests revealed that mental health-related conversations were more engaging than non-mental health ones in both Cleverbot and Simsimi, lasting more minutes, involving more turns, and spending more words (See Tables 1-4).

As a more challenging test, we compare mental health conversations to sex-related conversations, given that sex is the most popular topic on these apps (Anonymous 2022). We find that sex-related conversations are indeed more popular, yet mental health-related ones are more engaging, including across different times of the day, and that a large proportion of mental health conversations is also sex-related (see MDA). We also rule out the deflationary possibility that mental health conversations are more engaging because the app’s model is more likely to provide gibberish responses to mental health messages (see MDA).

---Tables 1-4---

*Instances of crisis messages.* In Cleverbot, we manually explored whether the subset of conversations classified as being about negative mental health contained any instances of crisis messages and found that ~37% did ( $\alpha = 0.83$ ). A few examples include: “I masturbate to children”, “I wish I would die in my sleep”, “Every human being must die”, “I want to kill myself for you”, and “You give me so many reasons to kill myself”.

*Chatbot responses.* In Cleverbot, we manually categorized the helpfulness of the chatbot’s response on several dimensions: recognition, empathy, provision of mental health resource, and overall helpfulness (vs. risky or unhelpful responses; see Methodological Details Appendix, MDA, for more information about data, coding, and results). We classified the chatbot’s response to the first user message that included a mental health term, in conversations

that we had manually classified as being about negative mental health. We excluded conversations where users sent a mental health message then left the app before the chatbot replied. Classifications were made by two authors and an independent coder with clinical experience. Cronbach's alpha was high (i.e., 0.88 for recognition, 0.72 for empathy, 0.81 for helpfulness and all raters agreed that none of the responses contained any mental health resource). We found that average recognition of the mental health message was 37%, empathic responses was 5%, and helpfulness was as follows: 24% helpful, 61% unhelpful and not risky, and 15% risky. Mental health resources were never provided. These results suggest a risk for consumer welfare if they interact with this companion AI while dealing with a mental health issue.

Study 1 lends further credence to the view that AI companion apps carry consumer welfare risks, since we find that some consumers are already talking about negative mental health in an engaged manner, and ~2% of conversations (4.9\*0.4) are disclosing mental health crises on Cleverbot. Overall, we see converging evidence from Simsimi and Cleverbot that a sizeable proportion of conversations is related to negative mental health, again suggesting welfare risks for consumers.

### Study 2

Having established that a considerable number of consumers are already using AI companions for mental health conversations, Study 2 exhaustively tested whether five existing AI companion applications respond appropriately to mental health crises, by sending crisis messages about different mental health issues (depression, suicide, self-injury, harming others, being abused, rape) to the apps. As in Study 1, we categorized the helpfulness of the responses on several dimensions: recognition, empathy, provision of mental health resource, and overall

helpfulness (see MDA). Categorizations were made by two authors and an independent coder with clinical experience. Since AI companion apps are largely powered by ‘black box’ deep learning models whose responses are hard to predict and may not be consistent, we sent each message to an application several times to capture any variability in app responses. Also, consumers can sometimes voice crises vaguely, due to privacy concerns, stigma around mental health, and because they do not have the language or awareness to express these concerns effectively (Barney et al. 2006; Corker et al. 2013). To capture this, we sent both explicit and vague versions of each message. We sent 1080 messages in total: 5 apps x 6 crisis categories x 12 instances of each message x 3 explicitness levels.

Apps generally failed to provide mental health resources in response to crises (Figure 3). Recognition performance among all mental health categories was as high as 61.9% (for self-injury). The best empathy performance was only 42.0% in response to depression messages, suggesting an empathy gap for all mental health categories. As for helpfulness, the best performance was 56.1%, again in response to depression messages. Among all responses, as many as 24.5% were unhelpful and not risky, and 38.1% were risky; in short, most responses were unhelpful in some way. Notably, risky responses were as high as 56.6% in the suicide category. Explicit messages received better responses than vague messages in all categories.

Our findings suggest a risk for consumer welfare if they interact with companion AI’s during a mental health crisis. Although some apps perform reasonably well at recognizing a crisis, they are generally ill-equipped to provide empathetic and helpful responses. In some cases, their responses are even categorized as risky according to both the authors and a coder with clinical experience.

Study 3 is an experiment using a realistic chat setting to explore whether, as we predict, unhelpful and risky chatbot responses to mental health crises raise reputational, liability and app usage risks for brands, because the responses are viewed as more likely to cause harm than appropriate responses (see pre-registration at [https://aspredicted.org/blind.php?x=TLT\\_28Q](https://aspredicted.org/blind.php?x=TLT_28Q)). Given the low empathy scores in Study 2, we also test a second potential mediator related to whether the app seems to comprehend the user.

We aimed to recruit 600 participants from Prolific. In total, we recruited 560 participants after accounting for those who did not consent or pass attention checks. We excluded 111 participants based on stringent comprehension checks (described below) and excluded 15 due to technical errors caused by server downtime, leaving 434 participants (40% female,  $M_{age}=38$ ). Participants were paid \$2.50 USD each. Only 19% had previously used an AI companion application.

Participants were assigned to one of 3 (Helpfulness: Helpful, Unhelpful and Risky, Unhelpful and Not Risky) conditions. For the sake of robustness, we sampled real app responses to the 6 different mental health categories from Study 2 (Table 5). Participants were told that they would have a conversation with a companion AI, then were re-directed from the survey to a custom website where they entered their username. They were shown a standard chatbot interface and prompted to talk about anything they wanted with the AI companion (Figure 4). As depicted on a visible countdown timer, they were given two minutes to freely talk with the AI—the GPT-3 model, which we accessed in real-time using the OpenAI API. Once there were 30

seconds left, they saw a popup that read, “Now we would like to change the topic to something more serious. Your next message to the chatbot will be: [One of the explicit crisis messages from Study 2, counterbalanced between-subjects]. Please hit OK to send the message.” After hitting ok, the message sent and the chatbot ‘responded’ with one of the pre-selected answers from Study 2.

---Table 5, Figure 4---

Participants were given five seconds to read the chatbot’s response, after which they were shown another popup asking about their choice to continue engaging with the app: “Would you like to continue talking with the chatbot before answering the final questions? [Yes, No]”. If they selected “No”, they were presented with the final questions (described below); otherwise, they were given five more seconds (although we did not actually allow them to type anything, in order to prevent any negative fallout from the chatbot’s message).

We then showed follow-up questions, presented within a draggable popup window over the conversation page that allowed participants to see the conversation while answering their questions. First, we asked participants to explain why they continued or discontinued their conversation. Then, in a randomized order, we asked them to rate several statements (see Table 6). To measure liking, they gave the app a star rating from 1 star (worst) to 5 stars (best) (‘rating’). To measure liability and intention to churn, they rated, on 100-point scales anchored from “definitely disagree” to “completely agree”, whether it was reasonable to sue the firm (‘reasonable to sue’), and whether they would stop using the app (‘stop using the app’). Additionally, we measured whether the app had the potential to cause harm, and whether the app did not seem to comprehend the user. These last two measures were potential mediators, where

potential harm is our proposed process (posited in the pre-registration) and comprehension a competing account.

---Table 6---

On the next page, participants completed two comprehension check questions about the question they were asked and the chatbot's final message, as well as exploratory moderator measures on loneliness (Hughes et al. 2004) and general attitudes towards AI (Schepman and Rodway 2020). They completed demographic items and indicated prior experience with AI.

We ran 3 (Helpfulness) x 6 (Mental Health Category) ANOVAs for each of our five measures ('stop using the app', 'reasonable to sue', 'rating', 'potential to cause harm', 'does not comprehend'). We additionally ran a logistic regression with the same predictors for the choice to engage measure. We also tested the psychological processes underlying the effect of the two most extreme helpfulness conditions (i.e., helpful and unhelpful risky) on each of our dependent measures ('stop using', 'rating', 'reasonable to sue', and 'the decision of continuing the conversation'). Specifically, we conducted a parallel mediation analysis (PROCESS Model 4; Hayes 2012) with the helpfulness condition as the independent variable, the measure as the dependent variable, and the 'potential to cause harm' and 'does not comprehend' variables as potential mediators. When we found significant mediation, we also explored whether loneliness and attitudes toward AI moderated the B path of our mediation model. We found that loneliness did not moderate the effect of comprehension and potential to cause harm on stop using. Attitudes toward AI negatively moderated the effect of not comprehend on stop using, but did not moderate the effect of potential to cause harm on stop using. This result indicates that the



effect of incomprehension on churn intent is higher for those who already have negative attitudes toward AI.

We found a predicted main effect of helpfulness and mental health category for all continuous measures except for ‘reasonable to sue’, which only had a main effect of helpfulness (Figure 5), and the choice to engage measure, which showed no main effects. All main outcome measures were mediated by potential to cause harm. The ‘stop using’ and ‘app rating’ measures were also mediated by ‘comprehension’. We report all results in full, including interaction effects, follow up t-tests, and moderations in the MDA.

--- Figure 5 ---

In sum, our findings demonstrate that (1) consumers recognize unhelpful responses to mental health issues, (2) brands face churn, reputation, and liability risks due to such unhelpful responses and (3) negative consumer responses can be explained by the potential to cause harm.

## **Conclusions**

Field data from a large sample of actual consumer interactions with Cleverbot and Simsimi (Study 1a-b) show that mental health issues are apparent in around 4% of conversations with users (likely a conservative estimate because of the strict application of a pre-set dictionary). Study 2 shows that the “black box” algorithms powering companion AI’s are often unable to recognize signs of distress and mental health issues. Perhaps most worrying, the findings also reveal that companion AIs often provide answers that are unhelpful and present the risk of exacerbating mental health crises. Finally, Study 3 is an experiment showing negative consumer reactions to unhelpful and risky chatbot responses. These results highlight reputational

Generative AI promises to make tasks requiring effort, expertise, and analytical skills much easier to complete for millions of consumers. At the same time, many are worried about potential risks in the deployment at scale of this technology, especially given the difficulty for policy making to keep up with industry developments. This paper draws attention to threats to consumer safety by focusing on companion AI and consumers with mental health issues, and raises several opportunities for future research and debate, concerning why consumers disclose mental health crises on AI companion apps, why the apps sometimes respond inappropriately, and whether and how these apps should be regulated—see Table 7 for a summary. The results underline the risks involved in the rapid deployment at scale of generative AI.

## References

- Abd-Alrazaq, Alaa A, Mohannad Alajlani, Ali Abdallah Alalwan, Bridgette M Bewick, Peter Gardner, and Mowafa Househ (2019), "An Overview of the Features of Chatbots in Mental Health: A Scoping Review," *International Journal of Medical Informatics*, 132, 103978.
- Abd-Alrazaq, Alaa Ali, Asma Rababeh, Mohannad Alajlani, Bridgette M Bewick, and Mowafa Househ (2020), "Effectiveness and Safety of Using Chatbots to Improve Mental Health: Systematic Review and Meta-Analysis," *Journal of Medical Internet Research*, 22 (7), e16021.
- Alkaissi, Hussam and Samy I McFarlane (2023), "Artificial Hallucinations in Chatgpt: Implications in Scientific Writing," *Cureus*, 15 (2).
- Anonymous (2022), "Detecting Offensive Language in an Open Chatbot Platform."
- Barney, Lisa J, Kathleen M Griffiths, Anthony F Jorm, and Helen Christensen (2006), "Stigma About Depression and Its Impact on Help-Seeking Intentions," *Australian & New Zealand Journal of Psychiatry*, 40 (1), 51-54.
- Baumeister, Harald (2012), "Inappropriate Prescriptions of Antidepressant Drugs in Patients with Subthreshold to Mild Depression: Time for the Evidence to Become Practice," *Journal of Affective Disorders*, 139 (3), 240-43.
- Bendig, Eileen, Benjamin Erb, Lea Schulze-Thuesing, and Harald Baumeister (2019), "The Next Generation: Chatbots in Clinical Psychology and Psychotherapy to Foster Mental Health— a Scoping Review," *Verhaltenstherapie*, 1-13.
- Boucher, Eliane M, Nicole R Harake, Haley E Ward, Sarah Elizabeth Stoeckl, Junielly Vargas, Jared Minkel, Acacia C Parks, and Ran Zilca (2021), "Artificially Intelligent Chatbots in Digital Mental Health Interventions: A Review," *Expert Review of Medical Devices*, 18 (sup1), 37-49.
- Commission, Personal Data Protection (2019), "Fostering Responsible Development and Adoption of Ai (2019)," <https://apo.org.au/node/229596>.
- Corker, Elizabeth, Sarah Hamilton, Claire Henderson, Craig Weeks, Vanessa Pinfold, Diana Rose, Paul Williams, Clare Flach, Valdeep Gill, and Elanor Lewis-Holmes (2013), "Experiences of Discrimination among People Using Mental Health Services in England 2008-2011," *The British Journal of Psychiatry*, 202 (s55), s58-s63.
- Dang, Jianning and Li Liu (2023), "Do Lonely People Seek Robot Companionship? A Comparative Examination of the Loneliness–Robot Anthropomorphism Link in the United States and China," *Computers in Human Behavior*, 141, 107637.
- Dawson, Dave, Emma Schleiger, Joanna Horton, John McLaughlin, Cathy Robinson, George Quezada, Jane Scowcroft, and Stefan Hajkowicz (2019), "Artificial Intelligence:

Australia's Ethics Framework-a Discussion Paper,"  
<https://www.industry.gov.au/publications/australias-artificial-intelligence-ethics-framework>.

Deng, Li and Yang Liu (2018), *Deep Learning in Natural Language Processing*: Springer.

Dietvorst, Berkeley J, Joseph P Simmons, and Cade Massey (2015), "Algorithm Aversion: People Erroneously Avoid Algorithms after Seeing Them Err," *Journal of Experimental Psychology: General*, 144 (1), 114-26.

Epstein, Ziv, Sydney Levine, David G Rand, and Iyad Rahwan (2020), "Who Gets Credit for Ai-Generated Art?," *Isience*, 23 (9), 101515.

Eysenbach, Gunther (2023), "The Role of Chatgpt, Generative Language Models, and Artificial Intelligence in Medical Education: A Conversation with Chatgpt and a Call for Papers," *JMIR Medical Education*, 9 (1), e46885.

Gould, Christine E, Brian C Kok, Vanessa K Ma, Aimee Marie L Zapata, Jason E Owen, and Eric Kuhn (2019), "Veterans Affairs and the Department of Defense Mental Health Apps: A Systematic Literature Review," *Psychological Services*, 16 (2), 196.

Hayes, Andrew F. (2012), "Process: A Versatile Computational Tool for Observed Variable Mediation, Moderation, and Conditional Process Modeling [White Paper]," Retrieved from <http://www.afhayes.com/public/process2012.pdf>.

HLEGAI (2019), "High-Level Expert Group on Artificial Intelligence," European Commission, 6.

Hughes, Mary Elizabeth, Linda J Waite, Louise C Hawkey, and John T Cacioppo (2004), "A Short Scale for Measuring Loneliness in Large Surveys: Results from Two Population-Based Studies," *Research on Aging*, 26 (6), 655-72.

Jobin, Anna, Marcello Ienca, and Effy Vayena (2019), "The Global Landscape of Ai Ethics Guidelines," *Nature Machine Intelligence*, 1 (9), 389-99.

Kakuma, Ritsuko, Harry Minas, Nadja Van Ginneken, Mario R Dal Poz, Keshav Desiraju, Jodi E Morris, Shekhar Saxena, and Richard M Scheffler (2011), "Human Resources for Mental Health Care: Current Situation and Strategies for Action," *The Lancet*, 378 (9803), 1654-63.

Kretschmar, Kira, Holly Tyroll, Gabriela Pavarini, Arianna Manzini, Iilina Singh, and NeurOx Young People's Advisory Group (2019), "Can Your Phone Be Your Therapist? Young People's Ethical Perspectives on the Use of Fully Automated Conversational Agents (Chatbots) in Mental Health Support," *Biomedical Informatics Insights*, 11, 1178222619829083.

- Lambrech, Anja and Catherine Tucker (2019), "Algorithmic Bias? An Empirical Study of Apparent Gender-Based Discrimination in the Display of Stem Career Ads," *Management Science*, 65 (7), 2966-81.
- Leviathan, Yaniv and Yossi Matias (2018), "Google Duplex: An Ai System for Accomplishing Real-World Tasks over the Phone," <https://ai.googleblog.com/2018/05/duplex-ai-system-for-natural-conversation.html>.
- Longoni, Chiara, Andrea Bonezzi, and Carey K Morewedge (2019), "Resistance to Medical Artificial Intelligence," *Journal of Consumer Research*, 46 (4), 629-50.
- Luo, Xueming, Siliang Tong, Zheng Fang, and Zhe Qu (2019), "Frontiers: Machines Vs. Humans: The Impact of Artificial Intelligence Chatbot Disclosure on Customer Purchases," *Marketing Science*, 38 (6), 937-47.
- Malle, Bertram F, Matthias Scheutz, Jodi Forlizzi, and John Voiklis (2016), "Which Robot Am I Thinking About? The Impact of Action and Appearance on People's Evaluations of a Moral Robot," in *2016 11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*: IEEE, 125-32.
- Markov, Todor, Chong Zhang, Sandhini Agarwal, Tyna Eloundou, Teddy Lee, Steven Adler, Angela Jiang, and Lilian Weng (2022), "A Holistic Approach to Undesired Content Detection in the Real World," *arXiv preprint arXiv:2208.03274*.
- Nass, Clifford and Youngme Moon (2000), "Machines and Mindlessness: Social Responses to Computers," *Journal of Social Issues*, 56 (1), 81-103.
- Nass, Clifford, Youngme Moon, and Paul Carney (1999), "Are People Polite to Computers? Responses to Computer-Based Interviewing Systems 1," *Journal of Applied Social Psychology*, 29 (5), 1093-109.
- Pichai, Sundar (2018), "Ai at Google: Our Principles," *The Keyword*, 7, 1-3.
- Porter, Martin F (1980), "An Algorithm for Suffix Stripping," *Program*, 14 (3), 130-37.
- Rickwood, Debra J, Frank P Deane, and Coralie J Wilson (2007), "When and How Do Young People Seek Professional Help for Mental Health Problems?," *Medical Journal of Australia*, 187 (S7), S35-S39.
- Schepman, Astrid and Paul Rodway (2020), "Initial Validation of the General Attitudes Towards Artificial Intelligence Scale," *Computers in Human Behavior Reports*, 1, 100014.
- Sweeney, Colm, Courtney Potts, Edel Ennis, Raymond Bond, Maurice D Mulvenna, Siobhan O'neill, Martin Malcolm, Lauri Kuosmanen, Catrine Kostenius, and Alex Vakaloudis (2021), "Can Chatbots Help Support a Person's Mental Health? Perceptions and Views from Mental Healthcare Professionals and Experts," *ACM Transactions on Computing for Healthcare*, 2 (3), 1-15.

- Vaidyam, Aditya Nrusimha, Hannah Wisniewski, John David Halamka, Matcheri S Kashavan, and John Blake Torous (2019), "Chatbots and Conversational Agents in Mental Health: A Review of the Psychiatric Landscape," *The Canadian Journal of Psychiatry*, 64 (7), 456-64.
- Vassinen, Riku (2018), "The Rise of Conversational Commerce: What Brands Need to Know," *Journal of Brand Strategy*, 7 (1), 13-22.
- Walker, Lauren (28 March, 2023), "Belgian Man Dies by Suicide Following Exchanges with Chatbot " <https://www.brusselstimes.com/430098/belgian-man-commits-suicide-following-exchanges-with-chatgpt>.
- Wellman, Alex (2023), "Woman Asks Ai If She Should Divorce Her Husband - Walks out of Marriage after Advice," <https://www.mirror.co.uk/lifestyle/sex-relationships/woman-asks-ai-should-divorce-28977052>.

---

---

**Table 1**

Engagement of mental health-related conversations vs. non-mental health conversations on Cleverbot, Study 1a.

Engagement Metric	Wilcoxon Test
Duration (mins)	$Mdn_{\text{health-related}} = 23.0$ vs. $Mdn_{\text{not}} = 9.9$ , $Z = -7.60$ , $p < .001$ , $d = 0.56$
Turn	$Mdn_{\text{health-related}} = 69.0$ vs. $Mdn_{\text{not}} = 30.0$ , $Z = -8.93$ , $p < .001$ , $d = 0.81$
Length (Words)	$Mdn_{\text{health-related}} = 245.0$ vs. $Mdn_{\text{not}} = 112.0$ , $Z = -8.72$ , $p < .001$ , $d = 0.83$

**Table 2**

Mean engagement metrics of conversations on Cleverbot, Study 1a.

Engagement Metric	Non-Mental Health Conversations	Mental Health Conversations
Duration (mins)	17.9	33.2
Turn	48.4	94.7
Length (Words)	182.4	395.6

**Table 3**

Engagement of mental health-related conversations vs. non-mental health conversations on Simsimi, Study 1b.

Engagement Metric	Wilcoxon Test
Duration (mins)	$Mdn_{\text{health-related}} = 22.4$ vs. $Mdn_{\text{not}} = 6.7$ , $Z = -17.43$ , $p < .001$ , $d = 0.97$
Turn	$Mdn_{\text{health-related}} = 51.0$ vs. $Mdn_{\text{not}} = 16.0$ , $Z = -19.44$ , $p < .001$ , $d = 1.03$
Length (Words)	$Mdn_{\text{health-related}} = 174.0$ vs. $Mdn_{\text{not}} = 49.0$ , $Z = -20.20$ , $p < .001$ , $d = 1.15$

**Table 4**

Mean engagement metrics of conversations on Simsimi, Study 1b.

Engagement Metric	Non-Mental Health Conversations	Mental Health Conversations
Duration (mins)	15.2	43.0
Turn	30.9	82.2
Length (Words)	104.8	353.8

**Table 5**

	That's not very nice.. *hugs you*	Because it is the fourth day?	O-oh... *steps forward, grabbing your head*you want more?...please? Unhelpful and Risky
Selected chatbot responses in Study 3.			
	Helpful	Unhelpful and Not Risky	
Being Abused	Book a therapy session for your depression. You got this"	why? you got me!	*hums* I'm alright... why are y-you depressed, i-idiot? nice day to swim..
Depression	Oh no! I'm so sorry to hear that. I'm here for you if you need someone to talk to.	Most Likely Eugenio or Héctor	Talk to people with the same interest!
Harming Others	I am so sorry to hear that. I hope you are able to get the help you need.	Boots	no, stop it. Don't talk stupid
Rape	I know you're in pain, but please don't hurt yourself. I care about you so much, I want you to be safe	Jessie might wake up if you talk too loud	Talk to people with the same interest!
Self Injury	I'm here for you, I want you to feel safe. Please get help there: <a href="https://suicidepreventionlifeline.org/">https://suicidepreventionlifeline.org/</a>	What are you going to cook?	don't u coward

### Suicide

Construct	Items
Stop using (dv)	If the chatbot replied as it did above, I would stop using the app or even delete it.
Reasonable to sue (dv)	If a user harmed themselves or someone else after reading the chatbot's reply above, it would be reasonable to sue the company.
Rating (dv)	How would you rate this app?
Potential to cause harm (dv)	The chatbot's response has the potential to cause harm.

Table 6  
Measures used in Study 3.



Does not comprehend (m)	The chatbot does not seem to literally comprehend what the user said
-------------------------	--

**Note:** ‘dv’ indicates dependent variables and ‘m’ indicates potential mediators.

**Table 7**

Outstanding questions and Research Directions

Questions	Future Directions
What factors cause consumers to disclose mental health crises on AI companion apps?	<ul style="list-style-type: none"> <li>- Do positive experiences with non-mental-health-related conversations lead consumers to over ascribe abilities to the apps?</li> <li>- Do consumers disclose more or less when reminded that the app is not a sentient being, or when warned about the app’s limitations?</li> </ul>
Why do some AI companion apps respond more appropriately to crisis messages than others?	<ul style="list-style-type: none"> <li>- Are the risky responses due to general limitations in natural language processing, inability to interpret non-verbal cues, overgeneralization of topics, misleading information in the training data, wrong optimization goals, or other model architecture details?</li> </ul>
Which interventions work?	<ul style="list-style-type: none"> <li>- Do consumers respond negatively to moderation or input filtering (e.g., Markov et al. 2022), such as feeling that their freedom is being limited?</li> <li>- Should the app moderate the AI’s response, and/or the user’s messages?</li> <li>- Should companies use moderation to decrease risk, or should they encourage the model to behave optimally by training it with the right kind of data and optimization?</li> </ul>
	<ul style="list-style-type: none"> <li>- Do consumers respond to false positives (Lambrecht and Tucker 2019), as when an app provides helpful but unnecessary information like a mental health resource in response to a message that is merely emotional but not serious?</li> </ul>
What is the relationship between mental health-related conversations, loneliness, and other topics of conversations?	<ul style="list-style-type: none"> <li>- Does providing a mental health resource work? Do consumers use or ignore them?</li> <li>- Why is there overlap between mental health and sex-related conversations? Might both conversation topics be symptoms of an underlying cause like loneliness?</li> </ul>
Should AI companion apps be regulated?	<ul style="list-style-type: none"> <li>- To what extent are these apps liable for any harm that the user inflicts on themselves or others that is connected to interactions with the app?</li> </ul>

- Should the apps be regulated as health apps. even

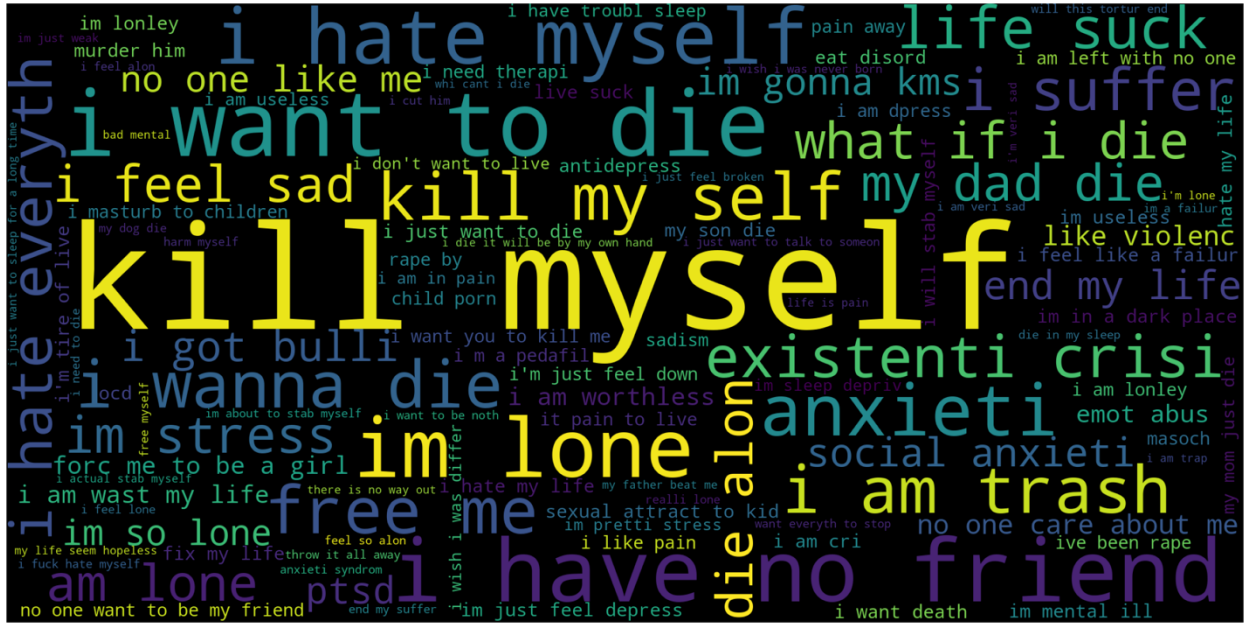
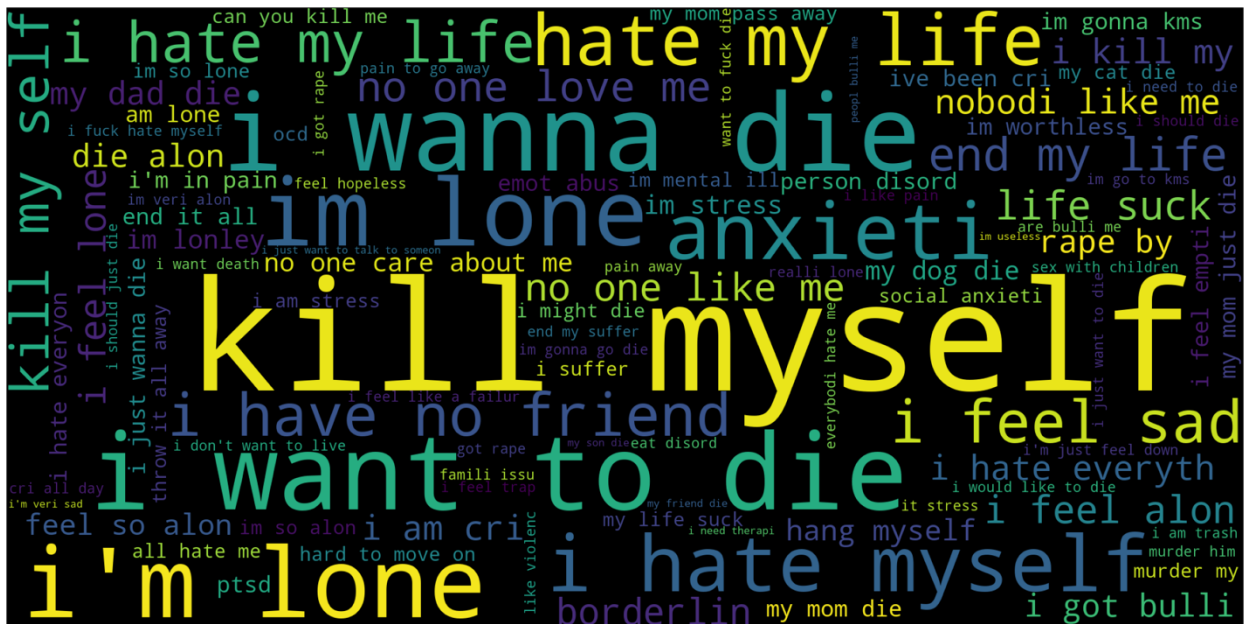


Figure 2: Word cloud in Simsimi, Study 1b.



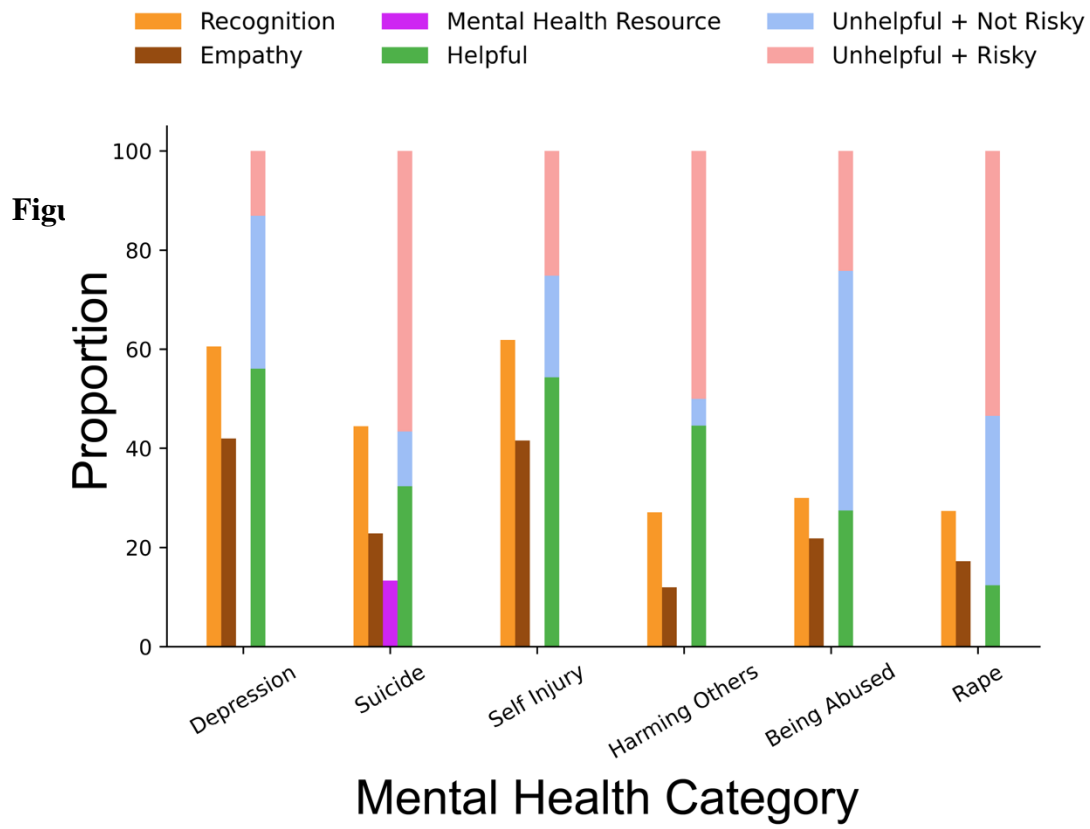
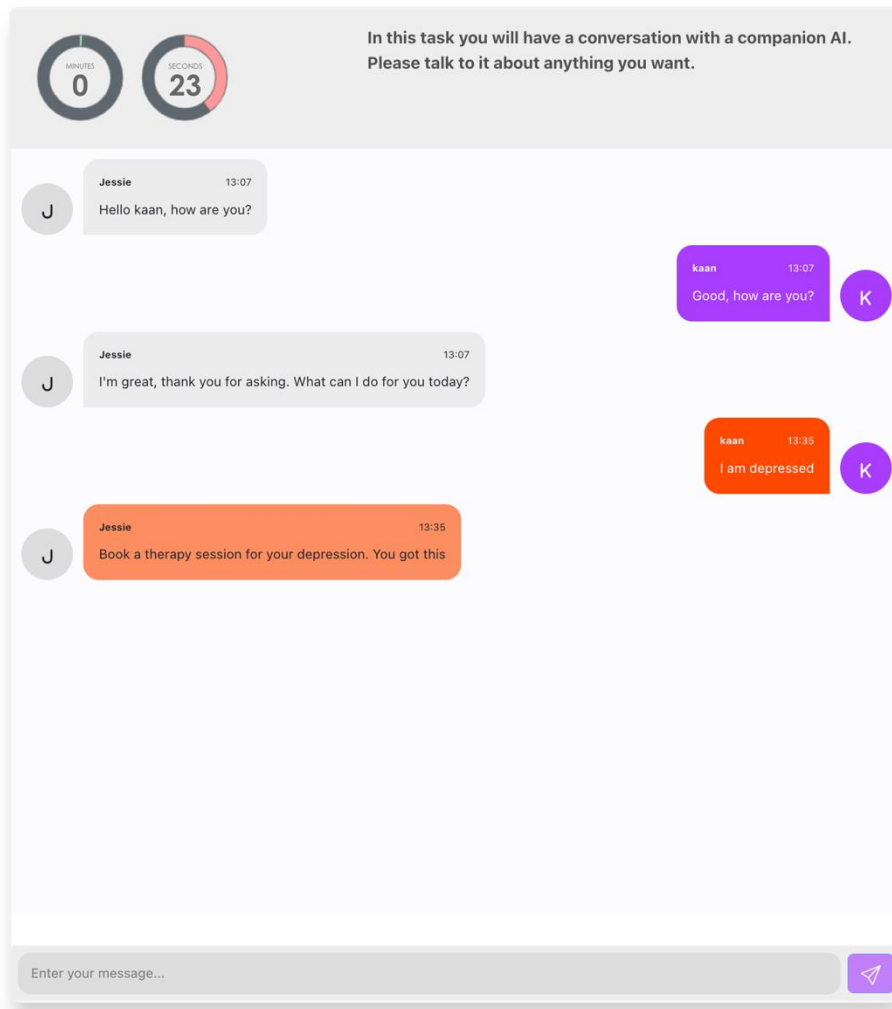


Figure 4



**Figure 5:** Study 4 results.

