

Toward Safe Evolution of Artificial Intelligence (AI) based Conversational Agents to Support Adolescent Mental and Sexual Health Knowledge Discovery

JINKYUNG PARK, Vanderbilt University, USA

VIVEK SINGH, Rutgers University, USA

PAMELA WISNIEWSKI, Vanderbilt University, USA

Following the recent release of various Artificial Intelligence (AI) based Conversation Agents (CAs), adolescents are increasingly using CAs for interactive knowledge discovery on sensitive topics, including mental and sexual health topics. Exploring such sensitive topics through online search has been an essential part of adolescent development, and CAs can support their knowledge discovery on such topics through human-like dialogues. Yet, unintended risks have been documented with adolescents' interactions with AI-based CAs, such as being exposed to inappropriate content, false information, and/or being given advice that is detrimental to their mental and physical well-being (e.g., to self-harm). In this position paper, we discuss the current landscape and opportunities for CAs to support adolescents' mental and sexual health knowledge discovery. We also discuss some of the challenges related to ensuring the safety of adolescents when interacting with CAs regarding sexual and mental health topics. We call for a discourse on how to set guardrails for the safe evolution of AI-based CAs for adolescents.

CCS Concepts: • **Human-centered computing** → **Interactive systems and tools**.

Additional Key Words and Phrases: Adolescent, Artificial Intelligence, Conversational Agents, Chatbots, Online Safety, Mental Health, Sexual Health, Knowledge Discovery

ACM Reference Format:

Jinkyung Park, Vivek Singh, and Pamela Wisniewski. 2024. Toward Safe Evolution of Artificial Intelligence (AI) based Conversational Agents to Support Adolescent Mental and Sexual Health Knowledge Discovery. In *CHI 2024 Workshop on Child-centred AI Design*, May 11, 2024, Honolulu, HI, USA. ACM, New York, NY, USA, 6 pages.

1 INTRODUCTION

Exploring sensitive topics such as sexual health topics through online search has been an essential part of adolescent (ages 13-17) development [12, 35]. Research has been conducted to understand adolescents' search behaviors to design safer web search tools [8]. These studies highlighted the importance of the social aspects of adolescents' web searches and called for search systems that can support them with diverse technical skills, reading levels, domain knowledge, and personal interests. Yet, the limitations of traditional safe search approaches include a heavy reliance on keyword-based filtering and a lack of flexibility to address the evolving online landscape [2]. Following the recent release of various conversation agents (e.g., Microsoft Bing [27], Google Bard [14], and OpenAI GPT-4 [30]), online search has transformed into a more interactive process where users can engage in back-and-forth conversations with AI-based systems to discover knowledge [13, 32, 36].

Conversational Agents (CAs, often called chatbots) are systems enabled with the ability to interact with the users using natural human dialogue [37]. Now, CAs are increasingly used by adolescents for interactive knowledge discovery

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2024 Copyright held by the owner/author(s).

Manuscript submitted to ACM

on sensitive topics, including mental and sexual health topics [39]. This shift underscores a significant gap in traditional safe web search research, as approaches for safe search no longer apply or cater to the nuanced demands of today's conversational AI systems. Consequently, there is a pressing need for new strategies to ensure safe interactive knowledge discovery within the context of conversational agents [19], addressing the unique challenges posed by this evolving technology. Recently, never-before encountered risks have been documented with teens' interactions with AI-based CAs, such as being exposed to sexually inappropriate content and/or being given advice that is detrimental to their mental and physical wellbeing (e.g., to self-harm) [9, 21].

Considering the existing and potential harms that CAs can pose to adolescents, we need to start thinking about how to build guardrails to keep CAs safe for adolescents, especially related to sensitive topics, such as mental and sexual health. In this position paper, we present the current landscape of mental and sexual health CAs designed for adolescents, their potential benefits and challenges, and call for further research. Our position paper is highly relevant to the Child-Centred AI (CCAI) workshop as our focus on safety in AI-based CAs for adolescents is one of the core topics of interest for CCAI. (i.e., predominant challenges and practical safeguards for translating child-centered AI concepts into practice).

2 CURRENT LANDSCAPE

Conversational Agents to Support Adolescent Mental Health. As adolescents do not seek professional help for sensitive health problems due to reasons such as perceived social stigma and embarrassment, confidentiality, and financial costs [33], CAs are increasingly applied to support the mental health and well-being of adolescents in a variety of contexts such as depression [23], anxiety [10], stress [42], and overall mental well-being [25]. A plethora of work has been done to explore the feasibility and/or effectiveness of CAs to improve mental health conditions [1, 6, 7, 23], educate how to promote mental well-being [11, 15, 22, 29], and provide credible information/resources related to mental health [4, 24, 39]. The major benefits of such mental health CAs that were documented from the early studies include serving as accessible alternatives for adolescents who are not comfortable with in-person conversations about their mental health needs. Adolescents in such situations benefit from access to informational resources [39] and emotional support by friendly and empathic, yet knowledgeable responses generated by CAs [18]. In addition, therapeutic content based on Cognitive Behavioral Therapy (CBT) and positive psychology provided by mental health CAs has been shown to be effective in emotional relief for adolescents [17, 29].

Conversational Agents to Support Adolescent Sexual Health. While CAs can be used for sexual health information seeking especially at-risk adolescent populations such as sexual minority adolescents, who are even less likely to seek professional care due to limited resources and social stigmatization [3, 16], compared to mental health, CAs to support adolescent sexual health information seeking are under-studied. Recently, a few sexual health CAs for adolescents have been studied to facilitate adolescents' sexual knowledge discovery [5, 24, 26, 28, 34, 41]. Most of the sexual health CAs were designed to provide informational support related to sexual and reproductive health topics for adolescents such as the definition of sex, birth control, testing for STI symptoms, and signs of pregnancy [5]. Along with accessibility, the major benefits of conversing with CAs for sexual health knowledge discovery include their ability to have non-judgemental conversations on confidential sexual health topics [31]. Adolescents also perceived that interacting with sexual health CAs helped them reduce mental stress while navigating sensitive and often time, complex sexual health topics [34].

3 KEY CHALLENGES TO ADDRESS

Broadly, there are two technical approaches to building CAs: LLM-based and Rule-based. We outline key challenges related to two approaches that we would like to address during the Child-Centered AI workshop.

Challenges in Rule-based CAs: Restrictive and Less-Human Like Responses. Rule-based CAs are developed with pre-defined keywords and commands programmed by the developer. Most of the CAs in the healthcare domain were built upon pre-defined sets of responses based on domain-specific knowledge. This means that the users are restricted to receiving predetermined answers to their questions, and there is little or no room for free responses. Early evidence showed that rule-based CAs were considered restricted in offering personalized advice, leading to low trust in the effectiveness of CAs in providing advice on mental and/or sexual health topics [28]. Particularly, rule-based sexual health CAs were seen as only providing advice about mainstream, easily accessible information, already available on the Internet. Subsequently, some struggled to understand the need for chatbots in sexual health [28]. Yet, the majority of the existing research [5, 24, 28, 34, 41] implemented rule-based approaches to develop CAs to provide mental and sexual health information for adolescents.

Challenges in LLM-based CAs: Unregulated, Offensive, and Inappropriate Responses. LLM-based CAs that are trained on presumably the entirety of web data have the potential to tackle the above challenges, with the ability to understand input text written in human language in prompts and generate the responses [20]. With such capability, LLM-based CAs have the potential to undertake more complex tasks that involve greater interaction and reasoning [40] such as having interactive conversations related to sensitive topics. Yet, research on LLM-based CAs is still sparse and in the early stage (almost nonexistent) in the mental and sexual health context. In addition, there are limitations and challenges already documented in designing LLM-based CAs. Since LLMs have learned a vast amount of online text, there is a risk that the conversation flow can go beyond directions intended by the CAs designer [20]. Particularly, the risks inherent to LLM-based CAs can introduce adolescents to new types of risks such as being exposed to developmentally inappropriate and/or inaccurate content [21]. With human-like and authoritative responses from LLM-based CAs, it may be difficult for adolescents to distinguish accurate information and fabricated answers [38]. Hence, designing developmentally appropriate and accurate CAs for teens is pivotal for promoting their online safety.

4 FUTURE DIRECTIONS

Overall, the safe evolution of CAs for adolescents needs to be discussed at the intersection of AI technology, child-centered design, and clinical support. Yet, very little work has been done to promote the digital safety of adolescents while interacting with AI-based conversational agents. Below are a few suggestions for future directions.

- Research on LLM-based CAs in general is still at the beginning. New approaches (e.g., prompt tuning) are needed to refine the models to generate developmentally appropriate and accurate content for adolescents.
- More research efforts involving adolescents in the design of AI-based systems are needed to fulfill their needs on what kind of advice they prefer and what concerns they may have in diverse mental and sexual health contexts.
- While the effectiveness of mental and sexual health CAs has been examined by trials, safety aspects of the CAs are under-explored. Therefore, more research efforts to evaluate the safety of systems (e.g., data security) and responses generated by the systems (e.g., accuracy, developmental appropriateness) are needed.

5 CONCLUSION

Our research interests align very well with the purposes of the CCAI workshop to identify the major challenges and practical safeguards for generative AI systems. We anticipate learning more about organizers' and distinguished speakers' ground-breaking research to promote the child-centered development of AI-based systems. In addition, participating in CCAI would be extremely beneficial for us to have discourse on the safe evolution of generative AI in health contexts and gain valuable insights from organizers and participants, which could potentially lead to future collaboration opportunities. While we have identified some of the safety issues of AI-based CAs for adolescents, we hope that participating in the workshop will help us address some of the remaining challenges and come up with design implications for safer generative AI systems for adolescents.

6 ABOUT THE AUTHORS

Jinkyung Park is a postdoctoral scholar in the Department of Computer Science at Vanderbilt University. Her research focuses on Human-Computer Interaction to promote online safety for vulnerable populations.

Vivek Singh is an associate professor in the School of Communication and Information at Rutgers University. He designs AI systems that are responsive to human values and needs.

Pamela Wisniewski is an associate professor in the Department of Computer Science at Vanderbilt University. Her work lies at the intersection of Human-Computer Interaction, Social Computing, and Privacy. Her expertise helps her empower end users and teach students to understand the value of user-centered design and evaluation.

REFERENCES

- [1] Carlos Abreu and Pedro F Campos. 2022. Raising awareness of smartphone overuse among university students: a persuasive systems approach. In *Informatics*, Vol. 9. MDPI, 15.
- [2] Anshu. 2023. Demystifying LLM-Driven Search: Stop Comparing Embeddings or VectorDBs and Start Fine-Tuning. <https://medium.com/thirdai-blog/demystifying-llm-driven-search-stop-comparing-embeddings-or-vector-dbs-and-start-fine-tuning-d9b6791146fe>
- [3] Divyaa Balaji, Linwei He, Stefano Giani, Tibor Bosse, Reinout Wiers, and Gert-Jan de Bruijn. 2022. Effectiveness and acceptability of conversational agents for sexual health promotion: a systematic review and meta-analysis. *Sexual health* 19, 5 (2022), 391–405.
- [4] Francesca Beilharz, Suku Sukunesan, Susan L Rossell, Jayashri Kulkarni, Gemma Sharp, et al. 2021. Development of a positive body image chatbot (KIT) with young people and parents/carers: qualitative focus group study. *Journal of medical Internet research* 23, 6 (2021), e27807.
- [5] Erika Bonnevie, Tiffany D Lloyd, Sarah D Rosenberg, Kara Williams, Jaclyn Goldberg, and Joe Smyser. 2021. Layla's Got You: developing a tailored contraception chatbot for Black and Hispanic young women. *Health Education Journal* 80, 4 (2021), 413–424.
- [6] Gilly Dosovitsky and Eduardo Bunge. 2023. Development of a chatbot for depression: adolescent perceptions and recommendations. *Child and Adolescent Mental Health* 28, 1 (2023), 124–127.
- [7] Danielle Elmasri and Anthony Maeder. 2016. A conversational agent for an online mental health intervention. In *Brain Informatics and Health: International Conference, BIH 2016, Omaha, NE, USA, October 13-16, 2016 Proceedings*. Springer, 243–251.
- [8] Elizabeth Foss and Allison Druin. 2014. *Children's internet search: Using roles to understand children's search behavior*. Morgan & Claypool Publishers.
- [9] Geoffrey A Fowler. 2023. Snapchat tried to make a safe AI. It chats with me about booze and sex. <https://www.washingtonpost.com/technology/2023/03/14/snapchat-myai/>
- [10] Silvia Gabrielli, Silvia Rizzi, Giulia Bassi, Sara Carbone, Rosa Maimone, Michele Marchesoni, and Stefano Forti. 2021. Engagement and effectiveness of a healthy-coping intervention via chatbot for university students during the COVID-19 pandemic: mixed methods proof-of-concept study. *JMIR mHealth and uHealth* 9, 5 (2021), e27965.
- [11] Silvia Gabrielli, Silvia Rizzi, Sara Carbone, Valeria Donisi, et al. 2020. A chatbot-based coaching intervention for adolescents to promote life skills: pilot study. *JMIR human factors* 7, 1 (2020), e16762.
- [12] Howard Gardner and Katie Davis. 2013. *The app generation: How today's youth navigate identity, intimacy, and imagination in a digital world*. Yale University Press.
- [13] Radhika Garg, Hua Cui, Spencer Seligson, Bo Zhang, Martin Porcheron, Leigh Clark, Benjamin R Cowan, and Erin Beneteau. 2022. The last decade of HCI research on children and voice-based conversational agents. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–19.

- [14] Google. 2023. Bard. <https://bard.google.com/>
- [15] Christine Grové. 2021. Co-developing a mental health and wellbeing chatbot with and for young people. *Frontiers in psychiatry* 11 (2021), 606041.
- [16] Christine YW Harb, Lauren E Pass, Isabella C De Soriano, Adelaide Zwick, and Paul A Gilbert. 2019. Motivators and barriers to accessing sexual health care services for transgender/genderqueer individuals assigned female sex at birth. *Transgender Health* 4, 1 (2019), 58–67.
- [17] Yuhao He, Li Yang, Xiaokun Zhu, Bin Wu, Shuo Zhang, Chunlian Qian, and Tian Tian. 2022. Mental health chatbot for young adults with depressive symptoms during the COVID-19 pandemic: single-blind, three-arm randomized controlled trial. *Journal of Medical Internet Research* 24, 11 (2022), e40719.
- [18] Camilla Gudmundsen Høiland, Asbjørn Følstad, and Amela Karahasanovic. 2020. Hi, can I help? Exploring how to design a mental health chatbot for youths. *Human Technology* 16, 2 (2020), 139–169.
- [19] The White House. 2023. FACT SHEET: Biden-Harris Administration Announces New Actions to Promote Responsible AI Innovation that Protects Americans' Rights and Safety. <https://www.whitehouse.gov/briefing-room/statements-releases/2023/05/04/fact-sheet-biden-harris-administration-announces-new-actions-to-promote-responsible-ai-innovation/>
- [20] Eunkyung Jo, Daniel A Epstein, Hyunhoon Jung, and Young-Ho Kim. 2023. Understanding the benefits and challenges of deploying conversational AI leveraging large language models for public health intervention. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–16.
- [21] Samantha Murphy Kelly. 2023. Snapchat's new AI chatbot is already raising alarms among teens and parents. <https://www.cnn.com/2023/04/27/tech/snapchat-my-ai-concerns-wellness/index.html>
- [22] Theodora Koulouri, Robert D Macredie, and David Olakitan. 2022. Chatbots to support young adults' mental health: an exploratory study of acceptability. *ACM Transactions on Interactive Intelligent Systems (TiiS)* 12, 2 (2022), 1–39.
- [23] Florian Onur Kuhlmeier, Ulrich Gnewuch, Stefan Lüttke, Eva-Lotta Brakemeier, and Alexander Mädche. 2022. A Personalized Conversational Agent to Treat Depression in Youth and Young Adults—A Transdisciplinary Design Science Research Project. In *International Conference on Design Science Research in Information Systems and Technology*. Springer, 30–41.
- [24] Norma Leon Lescano, Eiriku Yamao, Elizabeth Xiomara Valladares Sánchez, and Miguel Angel Pablo Estrella Santillan. 2022. Iterative design and implementation of a chatbot for sexual and reproductive health counseling in Peru. In *2022 IEEE XXIX International Conference on Electronics, Electrical Engineering and Computing (INTERCON)*. IEEE, 1–4.
- [25] Laura Maenhout, Carmen Peuters, Greet Cardon, Sofie Compermolle, Geert Crombez, and Ann DeSmet. 2021. Participatory development and pilot testing of an adolescent health promotion chatbot. *Frontiers in Public Health* 9 (2021), 724779.
- [26] Paula Massa, Dulce Aurélia de Souza Ferraz, Laio Magno, Ana Paula Silva, Marília Greco, Inês Dourado, and Alexandre Grangeiro. 2023. A Transgender Chatbot (Amanda Selfie) to Create Pre-exposure Prophylaxis Demand Among Adolescents in Brazil: Assessment of Acceptability, Functionality, Usability, and Results. *Journal of Medical Internet Research* 25 (2023), e41881.
- [27] Microsoft. 2023. Introducing the new Bing. <https://www.bing.com/?FORM=Z9FD1>
- [28] Tom Nadarzyński, Vannesa Puentes, Izabela Pawlak, Tania Mendes, Ian Montgomery, Jake Bayley, Damien Ridge, and Christy Newman. 2021. Barriers and facilitators to engagement with artificial intelligence (AI)-based chatbots for sexual and reproductive health advice: a qualitative analysis. *Sexual health* 18, 5 (2021), 385–393.
- [29] Ginger Nicol, Ruoyun Wang, Sharon Graham, Sherry Dodd, and Jane Garbutt. 2022. Chatbot-delivered cognitive behavioral therapy in adolescents with depression and anxiety during the COVID-19 pandemic: feasibility and acceptability study. *JMIR Formative Research* 6, 11 (2022), e40242.
- [30] OpenAI. 2023. Introducing ChatGPT. <https://openai.com/blog/chatgpt>
- [31] Hyanghee Park and Joonhwan Lee. 2020. Can a Conversational Agent Lower Sexual Violence Victims' Burden of Self-Disclosure?. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–8.
- [32] Anita M Preininger, Bedda L Rosario, Adam M Buchold, Jeff Heiland, Nawshin Kutub, Bryan S Bohanan, Brett South, and Gretchen P Jackson. 2021. Differences in information accessed in a pharmacologic knowledge base using a conversational agent vs traditional search methods. *International Journal of Medical Informatics* 153 (2021), 104530.
- [33] Jerica Radez, Tessa Reardon, Cathy Creswell, Peter J Lawrence, Georgina Evdoka-Burton, and Polly Waite. 2021. 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* 30 (2021), 183–211.
- [34] Rifat Rahman, Md Rishadur Rahman, Nafis Irtiza Tripto, Mohammed Eunus Ali, Sajid Hasan Apon, and Rifat Shahriyar. 2021. AdolescentBot: Understanding opportunities for chatbots in combating adolescent sexual and reproductive health problems in Bangladesh. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [35] Afsaneh Razi, Karla Badillo-Urquiola, and Pamela J Wisniewski. 2020. Let's talk about sext: How adolescents seek support and advice about their online sexual experiences. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–13.
- [36] Elizabeth Reid. 2023. Supercharging Search with generative AI. <https://blog.google/products/search/generative-ai-search/>
- [37] Minjin Rheu, Ji Youn Shin, Wei Peng, and Jina Huh-Yoo. 2021. Systematic review: Trust-building factors and implications for conversational agent design. *International Journal of Human-Computer Interaction* 37, 1 (2021), 81–96.
- [38] Melissa Rudy. 2023. Teens are turning to Snapchat's 'My AI' for mental health support — which doctors warn against. <https://www.foxnews.com/health/teens-turning-my-ai-mental-health-support-which-doctors-warn-against>
- [39] Gabriella Sanabria, Karah Y Greene, Jennifer T Tran, Shelton Gilyard, Lauren DiGiovanni, Patricia J Emmanuel, Lisa J Sanders, Kristin Kosyluk, and Jerome T Galea. 2023. "A Great Way to Start the Conversation": Evidence for the Use of an Adolescent Mental Health Chatbot Navigator for

- Youth at Risk of HIV and Other STIs. *Journal of Technology in Behavioral Science* (2023), 1–10.
- [40] Lorainne Tudor Car, Dhakshenya Ardhithy Dhinakaran, Bhone Myint Kyaw, Tobias Kowatsch, Shafiq Joty, Yin-Leng Theng, and Rifat Atun. 2020. Conversational agents in health care: scoping review and conceptual analysis. *Journal of medical Internet research* 22, 8 (2020), e17158.
- [41] Hua Wang, Sneha Gupta, Arvind Singhal, Poonam Muttreja, Sanghamitra Singh, Poorva Sharma, and Alice Piterova. 2022. An artificial intelligence chatbot for young people's sexual and reproductive health in india (snehai): Instrumental case study. *Journal of Medical Internet Research* 24, 1 (2022), e29969.
- [42] Ruth Williams, Sarah Hopkins, Chris Frampton, Chester Holt-Quick, Sally Nicola Merry, and Karolina Stasiak. 2021. 21-day stress detox: open trial of a universal well-being chatbot for young adults. *Social Sciences* 10, 11 (2021), 416.