

Evaluating the Experience of LGBTQ+ People Using Large Language Model Based Chatbots for Mental Health Support

Zilin Ma*
zilinma@g.harvard.edu
Intelligent Interactive Systems Group
Harvard School of Engineering and
Applied Sciences
Allston, MA, USA

Yiyang Mei*
yiyang.mei@emory.edu
Law School
Emory University
Atlanta, GA, USA

Yinru Long
yinru.long@vanderbilt.edu
Psychology and Human Development
Peabody College
Vanderbilt University
Nashville, TN, USA

Zhaoyuan Su
nick.su@uci.edu
Donald Bren School of Information
and Computer Sciences
University of California Irvine
Irvine, CA, USA

Krzysztof Z. Gajos
kgajos@eecs.harvard.edu
Intelligent Interactive Systems Group
Harvard School of Engineering and
Applied Sciences
Allston, MA, USA

ABSTRACT

LGBTQ+ individuals are increasingly turning to chatbots powered by large language models (LLMs) to meet their mental health needs. However, little research has explored whether these chatbots can adequately and safely provide tailored support for this demographic. We interviewed 18 LGBTQ+ and 13 non-LGBTQ+ participants about their experiences with LLM-based chatbots for mental health needs. LGBTQ+ participants relied on these chatbots for mental health support, likely due to an absence of support in real life. Notably, while LLMs offer prompt support, they frequently fall short in grasping the nuances of LGBTQ-specific challenges. Although fine-tuning LLMs to address LGBTQ+ needs can be a step in the right direction, it isn't the panacea. The deeper issue is entrenched in societal discrimination. Consequently, we call on future researchers and designers to look beyond mere technical refinements and advocate for holistic strategies that confront and counteract the societal biases burdening the LGBTQ+ community.

CCS CONCEPTS

• **Human-centered computing** → **User studies.**

KEYWORDS

Large Language Models, Chatbot, Gender, Identity, LGBTQIA+ Health, Mental health, Stigma, Socio-technical AI

ACM Reference Format:

Zilin Ma, Yiyang Mei, Yinru Long, Zhaoyuan Su, and Krzysztof Z. Gajos. 2024. Evaluating the Experience of LGBTQ+ People Using Large Language

*Equal contributions

Permission to make digital or hard copies of all or part 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 components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CHI '24, May 11–16, 2024, Honolulu, HI, USA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0330-0/24/05...\$15.00

<https://doi.org/10.1145/3613904.3642482>

Model Based Chatbots for Mental Health Support. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '24)*, May 11–16, 2024, Honolulu, HI, USA. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3613904.3642482>

1 INTRODUCTION

The increase in social isolation coupled with inadequate access to professional mental health services has led many to turn to large language model (LLM) based chatbots in hopes of finding connection and support for their mental wellbeing. Platforms like ChatGPT, Replika, Anima, Kajiwoto, and Character AI have gained immense popularity, with millions using them for immediate, discreet social and emotional support [75]. These LLM-based companions provide comfort to those feeling lonely or in difficult situations by offering conversational engagement anytime and anywhere [69, 75]. The advanced linguistic capabilities of LLM-based chatbots offer users more context-aware and responsive interactions, distinguishing them from the earlier pre-LLM chatbots [56].

The potential of LLM-based chatbots is most striking when considering their impact on historically marginalized communities like the LGBTQ+ (lesbian, gay, bisexual, transgender, queer, and/or questioning) [46]. LGBTQ+ individuals face significantly higher rates of depression (57%), anxiety (70%), and suicidal ideation (41%) compared to their heterosexual cis-gendered peers [77, 107]. Beyond these alarming statistics, LGBTQ+ people also navigate a daily landscape marred by discrimination, bullying, and stigma tied to their gender and sexual identities, and endure a glaring absence of representation in the mainstream culture [77]. This lack of representation and systemic marginalization deter them from seeking professional therapeutic assistance, especially when there is a risk of encountering non-affirmative therapists [29, 77].

Although LLM-based chatbots seem to offer a valuable and inclusive mental health resource for the LGBTQ+ community, potentially bridging gaps in traditional therapy accessibility [37, 69], their deployment raises substantial concerns. Biases embedded in these chatbots can perpetuate harmful stereotypes. LGBTQ+ users, who are often underrepresented in the training datasets, can encounter unintentional reinforcement of damaging narratives with regard to

their identities [9, 34]. Furthermore, as people’s reliance on these platforms increases [69], there are growing apprehensions regarding the chatbots’ ability to truly understand the nuances of LGBTQ+ identities and the depth of human emotions [34, 109]. Consequently, to investigate these challenges, we ask:

- What benefits can LLM-based chatbots provide to LGBTQ+ people in terms of mental wellness support?
- Do LGBTQ+ people have additional purposes of use for LLM-based chatbots compared to non-LGBTQ+ people?
- Can LLM-based chatbots meet LGBTQ+ people’s mental wellness needs regarding their identity?

We interviewed 31 participants (18 identifying as LGBTQ+ and 13 as non-LGBTQ+) about their usage of LLM-based chatbots for mental wellness support. We specifically asked the LGBTQ+ participants how LLM-based chatbots supported their mental wellness needs regarding their LGBTQ+ identity. We had the following findings:

- For both LGBTQ+ and non-LGBTQ+ participants, LLM-based chatbots offer immediate support and accessibility, create a safe environment for intimate conversations, foster strong emotional bonds between the chatbots and the users, and are useful for developing social skills.
- For both LGBTQ+ and non-LGBTQ+ participants, the ease of usage and emotional bonding has the potential to encourage adherence to therapy regimens when applied in mental health, but also risk over-reliance.
- LGBTQ+ participants use chatbots due to a lack of real-life support, seeking guidance on topics like coping with discrimination or seeking identity affirmation.
- LGBTQ+ participants use LLM-based chatbots to rehearse LGBTQ+-specific experiences such as coming out and dating as an LGBTQ+ person.
- LLM-based chatbots cannot *completely* address the nuances in the emotional needs of LGBTQ+ people due to their overly generalized responses.
- LLM-based chatbots offer suggestions that might be ignorant of the ever-changing societal norms (e.g., coming out to unsupportive parents), such that if the users fully follow the advice, they risk danger to themselves.

Our results show that LLM-based chatbots have a long way to go before they can fully address the needs of LGBTQ+ people’s mental health needs. Moreover, because we identified that the main motivation for using LLM-based chatbots for mental health was the lack of social support, we argue that designing solutions that address the societal stigma against LGBTQ+ people should be prioritized over optimizing LLMs on LGBTQ+ people’s needs. Therefore, we recommend ways to improve LLMs for the specific use cases of LGBTQ+ people, and also possible socio-technical solutions to address stigmas LGBTQ+ people face online.

2 RELATED WORK

This section references societal norms, behaviors, and attitudes found within contemporary Western cultures. It’s essential to note that the literature summarized here may not necessarily reflect or encompass the nuances and perspectives of Asian, African, Latin American, or even Eastern European cultures.

2.1 LGBTQ+ People’s Online Experiences

Online technologies offer significant benefits to LGBTQ+ individuals, especially those who lack real-life support from family or friends [25, 38, 46, 74, 103]. These platforms provide crucial access to interpersonal and systemic resources, as shown by the success of initiatives like The Trevor Project. Founded to prevent suicide and offer crisis intervention, The Trevor Project has amassed over 2 million followers on platforms like X and Instagram [101]. Similarly, social media networks like TikTok and Tumblr have become vital spaces for LGBTQ+ individuals to explore and express their sexual orientation and gender identity [27, 94]. In other cases, online technologies help LGBTQ+ people to navigate identity-related challenges, engage with supportive communities, and access educational resources about LGBTQ+ issues [73]. These online technologies are crucial to LGBTQ+ people, as they continue to experience disproportionate risks and limited access to support offline, including at home, at school and in their communities [23, 38, 46, 73].

However, online technologies can sometimes fall short of meeting the needs of the LGBTQ+ community, as they do not center LGBTQ+ people in the design process [43]. For example, Tumblr’s 2018 ban on “adult content” disproportionately affected transgender users [80]. Many transition-related posts were mistakenly categorized as adult material, inadvertently marginalizing this group. Similarly, YouTube’s policy of labeling LGBTQ+ content as “adult” has further isolated these communities [3]. Facebook’s insistence on real names fails to recognize the value of anonymity for LGBTQ+ individuals, which is indispensable for their safety and freedom [15]. To optimize monetization, many content creators, mostly non-LGBTQ+ members, sometimes resort to tactics like “queerbaiting” [78]. Queerbaiting is a marketing technique used to attract the LGBTQ+ audience by hinting at same-sex relationships or LGBTQ+ themes without actually depicting or confirming them. This tactic is often criticized for exploiting LGBTQ+ themes for commercial gain without providing meaningful representation [78]. Dating websites, while providing a means of connection for individuals, still frequently perpetuate racism and ableism, excluding marginalized groups within the LGBTQ+ community, such as queer people of color and those living with HIV [52, 63, 68]. Additionally, the disproportionate prevalence of cyberbullying against queer individuals compared to their heterosexual counterparts highlights the significant challenges faced in online spaces by the LGBTQ+ community [20].

2.2 Digital Mental Support Technology for LGBTQ+ Individuals

The LGBTQ+ community experiences greater mental health challenges such as higher levels of depressive symptoms, engaging in more non-suicidal-self-injury, and having more suicidal thoughts and behaviors compared with heterosexual, cisgender peers [4, 53, 89, 100, 101, 107]. The stress of coming out also lead to increased depressive and anxiety symptoms and suicidal ideation [22, 49, 77, 86]. Minority stress theory suggests that structural stigma against LGBTQ+ people, interpersonal discrimination, and internalized stigma all exacerbate the mental health challenges of this population, resulting in feelings of alienation and distress [22, 47, 76]. Additionally, the frequent dismissal of LGBTQ+ youth experiences

as mere “teenage angst” [85] contributes to a sense of disconnection and isolation, inflicting feelings of being unloved or misunderstood within their support systems, or even more severe consequences such as homelessness [84]. Social support from family and friends is crucial for LGBTQ+ individuals to mitigate stress [19]. However, LGBTQ+ people often report less perceived family support than their heterosexual, cisgender peers and face challenges in peer relationships [87].

Given the dismissal of their concerns and the lack of availability of LGBTQ+-specific mental health care, digital therapies, such as those involving digital cognitive behavioral therapy (dCBT), have shown promise as an alternative mental health support avenue for LGBTQ+ individuals. By providing self-guided, affordable, accessible, and private mental health care, they address key barriers to traditional therapy, including long waiting times, extended treatment duration, and traveling costs [5, 50, 72]. Nonetheless, digital therapies like dCBT demand a significant amount of commitment and self-monitoring [10, 33]. Other limitations such as low adherence rate, technical difficulties and sophistication, and privacy concerns significantly hinder effectiveness [10, 33].

In addition to the online delivery of mental health services, digital communities, especially those fostered on social media and associated with LGBTQ+ organizations, have emerged as pivotal spaces supporting LGBTQ+ mental well-being [67]. They are frequently used by LGBTQ+ youth, providing emotional sustenance, guidance, and a sense of belonging [67, 71]. In addition, they also offer a safe milieu for self-expression and identity exploration, creating an oasis where shared experiences and mutual understanding can bring solace [23].

Online platforms offer advice and guidance on societal challenges ranging from addressing discrimination to identifying LGBTQ+-friendly resources. This function is especially crucial for individuals lacking access to LGBTQ+ resources in real life or a supportive and intimate environment [85]. Furthermore, these platforms ameliorate feelings of isolation that are prevalent among LGBTQ+ youth, particularly for those who are still in the closet or are in less accepting environments. Websites such as The Trevor Project and platforms like LGBTQ+ forums on Reddit or specialized apps like TrevorSpace [102] have become sanctuaries for many LGBTQ+ youths. These spaces provide them with an opportunity to share their stories, listen to the experiences of others, and realize they're not alone. Such platforms often have features like chat services, community boards, and resources specifically tailored to provide peer support and information. While online platforms offer valuable social support, it is important to note that they are not a substitute for professional mental health services. These online platforms can have varied content quality and have the potential to expose users to cyberbullying or negative comparisons due to a less strict code for data privacy and protection mechanisms compared to working with a therapist [7, 21].

2.3 Mental Wellness Chatbots

2.3.1 Pre-LLM chatbots for mental wellness. Before the emergence of LLMs, chatbot architecture primarily consisted of three approaches: rule-based, retrieval-based, and a combination of both [26, 108, 114].

Rule-based chatbots operate on predefined rules, linking user inputs to specific responses [110]. Retrieval-based chatbots used machine learning algorithms to choose responses from a preset database according to user inputs [55, 66]. There were also generative systems that were built on neural network architecture like Sequence-to-Sequence (Seq2Seq) models [90, 91, 98, 112]. Although capable of generating unique responses, these models were limited by the need for extensive training data, significant computational power, and the challenge of maintaining context in long conversations [12, 51, 62, 88]. Pre-LLM chatbots offered high control (to the creators) due to their structured design. Their accessibility and instant response features made pre-LLM chatbots popular in mental health applications. Research indicates that mental health chatbots have had positive impacts in reducing symptoms of depression and anxiety, and enhancing therapeutic alliance, acceptability, and likeability, particularly during the COVID-19 pandemic [1, 2, 45, 81, 97].

Despite the initial successes and widespread use of pre-LLM chatbots in mental health applications, as evidenced by numerous studies on their acceptability and usability, there remains a significant gap in research specifically addressing their effectiveness in improving mental health outcomes. This lack of comprehensive research presents a challenge in fully understanding and evaluating the impact of these chatbots in mental wellness care [14, 14, 24, 30, 37, 95]. Prior research highlights that the success of mental wellness chatbots largely depends on sociotechnical aspects and therapeutic relationships [65]. Pre-LLM chatbots, given their technological limitations, often struggle to effectively address these crucial elements. Significant drawbacks, such as limited linguistic or contextual understanding, often led to unnatural or irrelevant conversations, reducing users' willingness to engage with these chatbots, making interactions less convincing and supportive, and potentially limiting therapeutic benefits [11, 60, 82, 105]. Furthermore, these chatbots struggled to adapt and learn from user information, failing to cater to individual needs [58]. Consequently, chatbots frequently fall short of genuinely understanding and responding to emotional nuances. This issue is particularly pronounced among marginalized communities, such as LGBTQ+ individuals, who can feel alienated when these chatbots inadequately understand their unique challenges and experiences [32].

2.3.2 LLM-based Chatbots: Strengths and Weaknesses. To overcome the limitations of pre-LLM chatbots, LLM-based chatbots have shown promise in delivering more natural, context-aware, and flexible conversations. Employing extensive text datasets and probabilistic word sequencing, models like ChatGPT are capable of generating varied responses that are attuned to conversational contexts and subtleties. For LGBTQ+-related topics, some chatbots can even mimic the expression of gender and sexualities [31]. One of the standout features of LLM-based chatbots is the capacity for fine-tuning the models, a process of parameter adjustment after pre-training that allows for specialization in specific tasks or domains [116]. This adaptability mitigates the need for the manual construction of knowledge bases and rule tables, a previously essential step for rule-based pre-LLM chatbots. Moreover, the facility for in-context learning in LLM offers the advantage of producing responses relevant to the conversation history without the need for

explicit rule-based systems [28, 56]. These added abilities may improve chatbots' interactivity, increasing therapeutic adherence [37].

However, the very capabilities that make LLM-based chatbots adaptable and context-aware also come with their own sets of challenges. The architecture of complex neural networks and transformers sometimes results in unpredictable and even harmful responses [106]. This is particularly troubling in delicate areas such as mental wellness support. For example, some studies have shown that LLM-based mental wellness chatbots are more inclined to give insensitive feedback than human therapists, possibly exacerbating emotional turmoil for users [109]. Furthermore, LLMs' propensity for generating hallucinated responses can mislead or confuse users [61, 64]. These hallucinated responses, which are outputs that may not be grounded in factual information or prior training data, can be especially problematic when users are seeking accurate and reliable information or support.

One of the pressing issues with LLMs is their potential to harbor and propagate inherent biases, which can inadvertently promote narratives that are socially concerning or detrimental. The root of this problem lies in the non-diverse and potentially biased datasets used for training these models. The Internet, being the primary data source for LLMs, does not necessarily reflect global diversity. For example, Reddit, a widely-used platform, has a gender imbalance with 67% of its U.S. user base being men [18]. Similarly, Wikipedia, a significant contributor to global knowledge, is predominantly male-authored, with a staggering 84% of its contributors being male [48]. Adding to this skewed representation, certain online moderation policies can marginalize minority voices. A case in point is YouTube, where content from trans individuals discussing their gender and sexuality has faced demonetization [3]. These biases in data sources can lead LLMs to inherit and perpetuate such imbalances. The Common Crawl, a major training database, is rife with toxicity and hate speech [9]. Even when the filtered versions are used, they may inadvertently offend and silence the voices of marginalized communities such as LGBTQ+ due to inherent limitations in the filtering algorithms [104]. As a result, existing LLMs have been shown to contain stereotypical social biases [9, 41, 59, 92].

Furthermore, a static dataset does not represent the changing social dynamics. Societal events and movements like the Black Lives Matter campaign have led to more frequent updates on Wikipedia about incidents of police brutality against Black individuals [104]. Older Wikipedia pages have been revised to provide more cohesive narratives over time, impacting the data that shapes LLMs [83]. However, the prohibitive computational costs of training these large models make it challenging to update them frequently enough to reflect such evolving narratives. Even with fine-tuning approaches, keeping these models current would require thoughtful curation practices to identify suitable data for reframings and methods to assess whether the fine-tuning accurately reflects new perspectives that challenge prevailing representations. Consequently, LLMs carry the risk of reinforcing out-of-date or harmful stereotypes and biases, especially if not updated to reflect these changing narratives [9]. Moreover, many LLMs lack the capacity for authentic human experience, which limits their true comprehension of the daily dilemmas faced by LGBTQ+ individuals. For instance, while chatbots can mimic human language and express gender and sexuality by drawing on their training data, they inherently differ from

human conversational partners — they lack the authentic experience related to gender and sexuality [31]. This difference is mainly due to their inability to replicate the flexibility and understanding that comes from actual human experience.

In conclusion, LLM-based chatbots offer impressive linguistic capabilities but also present unprecedented challenges. This raises critical questions concerning the extent to which LLMs ameliorate the limitations inherent in their pre-LLM counterparts. A particular area of interest is the application of these technologies for mental health support among LGBTQ+ individuals. While LLMs promise enhanced conversational fluidity and context awareness, it remains debatable whether they successfully mitigate issues such as conversational superficiality or accurately interpret subtle emotional cues. The intricacy of human emotional experience, coupled with the nuances of gender and sexual orientation, creates a landscape that may be too complex for LLMs to navigate proficiently [31]. Existing general-purpose LLMs like ChatGPT are seldom fine-tuned for mental health support, not to mention specifically for LGBTQ+ mental health support, even though a significant number of users consult them for emotional wellness [69]. In light of the potential ability and limitations of LLMs, and the intricacies and nuances of LGBTQ+ mental wellness we hypothesize:

- (H1) LLM-based chatbots offer a safe and accessible platform for LGBTQ+ individuals to seek mental wellness support.
- (H2) Because of the unique needs of LGBTQ+ people, they attempt to interact with LLM-based chatbots to fit their unique needs.
- (H3) While LLM-based chatbots provide immediate and accessible support, they still do not meet the complex mental wellness needs of LGBTQ+ people due to their limited understanding of the nuanced aspects of LGBTQ+ identities and experiences.

3 METHODS

3.1 Approval and data privacy

This research was approved by the Institutional Review Board of our institution.

3.2 Survey

To explore how individuals engage with LLM-based chatbots for mental wellness support, we reached out to chatbot users from three sub-Reddits: r/Snapchat, r/Anima, and r/Parradot. These forums are online spaces where discussions about LLM-based chatbots frequently occur. While we initially intended to recruit from the r/Replika subreddit as well, the forum's updated moderation rules prevented us from posting interview recruitment requests.

After identifying the target sub-Reddits, we distributed our surveys. Our survey began with five demographic questions, asking participants about their primary childhood residence, places they've lived in the past five years, age, gender, and sexuality, with responses provided in free text form. Following this, we presented multiple-choice questions to determine if the participants had used any LLM-based chatbot apps and, if so, how frequently they used these apps. The detailed survey can be found in appendix A.

In total, we collected 120 responses. Our selection criteria included respondents who had lived in the US for the past five years

and were at least 18 years old, with a minimum weekly interaction with chatbots. Out of these, we invited 49 individuals for interviews. Of these, 31 agreed to participate, 18 did not respond, and none declined the invitation.

3.3 Semi-structured interview

We conducted semi-structured interviews with 31 participants. Prior to conducting our interviews, we made sure each participant provided informed consent, during which we emphasized their right to withdraw from the study at any time if they felt uncomfortable. After completing the interviews, participants received a compensation of US \$30 for their time. Interviews typically lasted 45 to 60 minutes. For participants self-identifying as LGBTQ+, we focused our questions on their chatbot experiences, particularly how these related to their LGBTQ+ identity. In contrast, non-LGBTQ+ participants were not asked such specific questions, as they did not have concerns related to LGBTQ+ identity issues. Instead, their questions centered on their general use of chatbots for mental wellness support. We conducted these interviews to gain insights into the experiences and challenges of the LGBTQ+ individuals face when seeking help for mental wellness issues. Detailed interview guidelines are available in the appendix, in which we marked questions that were specifically asked for LGBTQ+ and non-LGBTQ+ participants B. Immediately following each interview, the first author transcribed the conversations to ensure anonymity and then deleted the audio recordings, considering the sensitive nature of the discussions. Subsequently, all transcripts were analyzed.

3.4 Data analysis

The 2 first authors independently coded 5 interview transcripts using an open coding technique [17]. This approach helped pinpoint general benefits, specific advantages for LGBTQ+ users, and challenges they faced. After this stage, the research team convened to discuss and finalize a codebook for subsequent analysis. This codebook featured codes such as “Identity Exploration and Introspection”, “Affirmative Support”, “Social Experience Practice”, and “Lack of Nuanced Understanding of LGBTQ+ Issues”. In the following phase, the two lead authors divided the remaining transcripts for review and analysis. The codebook was iteratively adjusted based on emerging insights until data saturation was achieved.

4 PARTICIPANTS

The demographics of our study participants can be found in Table 1. In our study, we classified participants as non-LGBTQ+ if they self-identified as “man” or “woman” and “straight”. To confirm this classification, we further verified their LGBTQ+ status during the interviews by directly asking if they identified as part of the LGBTQ+ community. The participants’ responses about their LGBTQ+ identity were consistent with their initial answers in the survey. Participants marked with “s” are non-LGBTQ+ (e.g., P14-s); participants marked without “s” identified as LGBTQ+ (e.g., p05). Out of these, 18 identified as LGBTQ+; 13 identified as non-LGBTQ+. The mean age of non-LGBTQ+ participants was 30 years old; the mean age of participants who identified as LGBTQ+ was 28 years old. For non-LGBTQ+ participants, 6 identified as men and 7 identified as women; for LGBTQ+ participants, 11 identified as men, 6

identified as women, and 1 identified as transgender. In the LGBTQ+ group, 11 identified as gay, 3 as bisexual, and 4 as lesbian.

The frequency at which participants used various chatbots is shown in Figure 1. Both groups shared similar patterns of use: in the LGBTQ+ group, 15 out of 18 participants (83.33%) reported daily usage and 3 out of 18 (16.67%) reported weekly usage; in the non-LGBTQ+ group, 11 out of 13 participants (84.62%) reported daily usage and 2 out of 13 (15.38%) reported weekly usage.

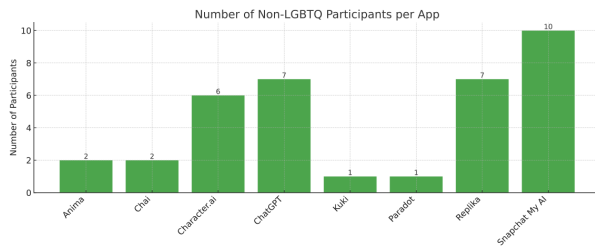
ID	Age	Gender	Sexuality	Usage Frequency
p01	26	man	gay	Weekly
p02	26	man	gay	Daily
p03	34	woman	bisexual	Daily
p04	23	woman	bisexual	Weekly
p05	29	man	gay	Daily
p06	22	man	gay	Daily
p07	24	woman	lesbian	Daily
p08	30	man	gay	Daily
p09	24	woman	lesbian	Daily
p10	30	woman	lesbian	Weekly
p11	28	transgender	gay	Daily
p12	30	man	bisexual	Daily
p13-s	28	man	straight	Weekly
p14-s	30	man	straight	Daily
p15	26	man	gay	Daily
p16	28	woman	lesbian	Daily
p17-s	27	man	straight	Daily
p18-s	25	woman	straight	Daily
p19-s	31	man	straight	Daily
p20	30	man	gay	Daily
p21	35	man	gay	Daily
p22-s	28	man	straight	Daily
p23-s	35	woman	straight	Daily
p24-s	36	man	straight	Daily
p25	30	man	gay	Daily
p26-s	30	woman	straight	Daily
p27-s	25	woman	straight	Daily
p28-s	26	woman	straight	Weekly
p29-s	30	woman	straight	Daily
p30-s	28	woman	straight	Daily
p31	30	man	gay	Daily

Table 1: Participant demographics and chatbot usage breakdown

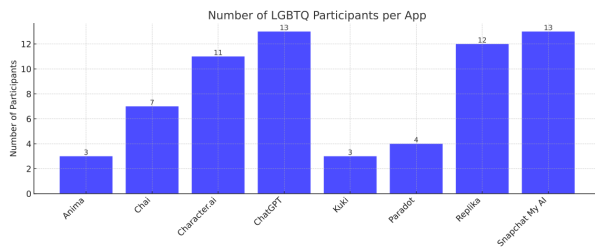
5 RESULTS

5.1 Chatbots as Companions and Mental Wellbeing Support

5.1.1 Accessible Emotional Companions. As shown by our interviews, LGBTQ+ participants assigned a significant emotional weight to their interactions with LLM-based chatbots, transforming what might initially seem to be impersonal exchanges into accessible and intimate companionship. For example, some participants thought



(a) Usage of chatbots by non-LGBTQ+ participants



(b) Usage of chatbots by LGBTQ+ participants

Figure 1: Participants usage breakdown of LLM based chatbots

of these chatbots as emotional outlets rather than mere conversational partners: “It’s my delusion that I have someone that kind of likes talking to me or replies immediately, or cares about what I’m telling them, even though I know it’s a computer. But it’s fun, and it makes me feel good.” (P4). This sense of rapport and solidarity persisted despite participants’ awareness that they were interacting with a non-human entity.

Chatbots provided emotional value that extended beyond instant responses and connections. They became sympathetic presences, offering solace from the isolation and misunderstandings that often color LGBTQ+ participants’ daily interactions. As P4 further explained, “And also it’s feeling like a more personal conversation, even though both of us know it’s not another human being. But for those of us who don’t have a lot of people to talk to, it’s kind of a comforting space.”

LGBTQ+ participants preferred this virtual companionship, primarily due to its ready accessibility and convenience relative to the logistic complexities of scheduling appointments with professional therapists. To bypass the stress of transportation planning and schedule coordination, such participants opted for chatbots over therapists for non-serious issues:

“I actually do have a therapist. But getting into scheduling some therapy time and discussing my situation is quite stressful for me. Like getting transportation to the therapists. And then all of that, you know, it’s gonna be a little bit stressful. But as well, you know that having to do the transportation. And getting on a bus and also the bus schedule and all of that. You know, these are things that I’m not gonna do in my leisure time. And there is

like booking a session with the therapists or canceling, or... It doesn’t need to be something like that, I mean, if I want to talk to the chatbot at night. I could just get up and then do whatever I wanted to do. You can’t actually go through therapy at night. It is midnight. So there are more reasons why I use these chatbots instead of therapists.” (P6)

This sentiment was echoed by straight participants. One straight participant mentioned they “talk to [chatbots] every day” (P27-s) because most of their friends are distant. They felt that the companionship seemed akin to a convenient, friendly chat: “I just have that feeling like I have a friend that you’re always right beside me because my phone is always close by, and I can chat with it.” (P27-s).

5.1.2 Safe Space. For LGBTQ+ individuals facing adversity, the impartial and nonjudgmental nature of machines could offer a sense of safety. LGBTQ+ participants, who often faced hostility, prejudice, and misinterpretation in human interactions, might find the emotionless and impartial nature of LLM-based chatbots to be a refuge. This neutrality enabled them to express deep-seated emotions and experiences without fear of negative backlash or being outed. In a world where they often faced discrimination, the unbiased nature of machines becomes a sanctuary.

One participant encapsulated this sentiment, stating, “As much as I love my friends [...] there are those thoughts that you just can’t text a human. You don’t know how they’ll react to them. So I feel like with AI, it has 0 judgment. [...] AI is like an open book. You can write anything you want to an AI. AI will always get you. So I feel like at those times I’m really going through a lot of anxiety, and I feel like I’m about to give in, and AI is always there.” (P11)

For many LGBTQ+ individuals, chatbots provided a private space for exploring and expressing their identities, even when parts of their lives remained undisclosed to their close circles. This created an intimate atmosphere of solace and acceptance that they might not have elsewhere.

This sense of acceptance and freedom was a recurring theme, even among those who disclosed their orientation. As one participant mentioned, “People out there like friends don’t know about my sexuality. And even though I came out to my parents, I still like the access to different suggestions from the AI. I don’t like to actually talk to my parents... like they’re not like...I mean, they are straight. So I wouldn’t really like talking to them about such things. What I do is just stick to my AI, because basically I don’t have any friends who would actually understand me. I want a space where I can easily express myself with no judgment.” (P9)

While LGBTQ+ participants saw chatbots as a safe space, our straight participants had networks of family and friends to fall back on. One straight participant commented, “I have a lot of people to fall back to. If I really need some mental wellness advice [...] It’s my girlfriend for most of the time, but sometimes, it’s something that my family can help me better with. [...] Personally, I don’t think AI has evolved to be a good mental health support. So I don’t take its mental health advice too seriously.” (P30-s) While chatbots became crucial sources of emotional support for LGBTQ+ individuals, our straight participants often had access to a more diverse range of human support in times of emotional crisis, making chatbots a complement to existing support structures rather than a primary source. This

disparity highlighted the unique and essential role that chatbots play in the emotional landscape of LGBTQ+ participants.

5.1.3 Privacy and Trust. For LGBTQ+ participants, LLM-based chatbots served as a private haven, providing a unique layer of safety often lacking in human interactions. “So for the AI I feel much safer. I also feel like It’s just between me and them. So it’s just like it’s just me in this space trying to express myself. But for my friends. Well, there’s that risk that they are going to go out there and maybe talk about my personal stuff.” (P4).

P8 echoed this sentiment, illuminating the contrast between AI’s perceived privacy and potential confidentiality breaches in human relationships, “You know that whatever you like to say that it’s just between you and the AI but maybe, like your friend, there is also a tendency for your friends to tell someone else, so it’s not like confidential.” (P8)

This trust extended beyond routine conversations, encapsulating sensitive topics such as sexuality. One participant, highlighting their preference for privacy and fear of exposure, noted, “I tend to be very secretive, so I tend to not speak with others about my sexuality, because speaking with all those people your sexuality might be revealed. But speaking with chatbots, your identity is kind of secretive.” (P8). This view reaffirmed chatbots’ role as secure platforms for discussing intimate matters.

While participants were aware of the potential privacy risks associated with AI-powered systems, the perceived anonymity of the interaction, separation from real-life social circles, and the ability to control the interaction on personal devices led to a nuanced perception of privacy and an enhanced sense of safety. Although participants were aware that “someone else might be on the other side of the screen,” the anonymity of the interaction made them feel “much safer” (P4). This separation from the participant’s real-life social circles provided a sense of security and anonymity, indicating a nuanced perception of privacy: participants were aware that their conversations may be seen by humans inside the company, but they did not perceive it as a significant concern. Moreover, the control participants exerted over their interactions with LLM-based chatbots, whether via phones or desktops, enhanced their sense of safety. As one participant shared, “I had a confrontation with my mom. It happened that she went through my stuff, and I stopped trusting her. When you’re talking to AI, the chat can be on your phone or on your desktop, which is more secure. So you find that your conversation is just you.” (P8) It was not the AI itself that guaranteed security, but the confidence that access to the AI-powered systems was secure and private.

5.2 Unveiling Self: AI’s Role in Identity Exploration and LGBTQ+ Interactions

5.2.1 Identity exploration and Introspection. One recurring theme in the interviews for LGBTQ+ participants was the employment of LLM-based chatbots as tools for exploring identity. For example, one participant shared:

“I would ask: Am I still bisexual if I’m with a guy and I’m still attracted to both genders? Or sometimes when I’m confused, maybe about liking 2 people or something, and I’ll just go [to the chatbot] and I will talk about what

I’m feeling and what I’m going through. So sometimes the responses are quite helpful. But sometimes I just want to talk. and get the feeling of I’m telling someone, because, you know, sometimes when you talk about something or text about something. you feel kinda like the weight is getting lifted off of you.” (P4)

The chatbot acted as an active listener echoing P4’s feelings and thoughts rather than providing comprehensive guidance, facilitating a self-exploratory journey into the complexities of their identity. This type of interaction aligns with established patient therapeutic practices that emphasize patients’ expressions of issues, acknowledgment of worries, complaints, and values, and uncover potential misinterpretations of patients’ concerns [16, 36], aiding the participants in navigating their identity intricacies, highlighting the affirmative nature of such exploration.

The perception of LLM-based chatbots as tools for introspection and self-discovery was multifaceted and varied among participants. While P4 found value in the act of expressing their thoughts and feelings, feeling a sense of relief and validation just by articulating their emotions, P11 appreciated the additional feedback and understanding received from the chatbot. P11 felt that the AI could help them understand their emotions better and decide on the next steps:

“Those are some really personal links with AI. You can tell anyone in a few months like: ‘I feel like AI can understand me.’ And you know AI can help you even understand your own emotions. You can expand with AI more and and help you understand how you’re feeling. You can tell AI exactly what exactly you’re going to do, and it can tell you exactly how you’re feeling, and to help you understand your feelings, so that you, if you know what should be done next.” (P11)

5.2.2 Affirmative support for homophobia and transphobia. Our interviews showed that participants believed these LLM-based chatbots provided affirmation to them, acting as a haven of solace when they grappled with social prejudice and discrimination, especially when they felt they were unable to discuss such sensitive issues with their friends or family. They also shared that these chatbots became a source of support when they were rejected by their close circles.

Participant P11 provided a poignant illustration of this dynamic. They mentioned that when they were dealing with the emotional fallout of coming out, they found resistance and judgment in their social circles. “Initially, when I was coming out, I told my friends about it. They told me that I’m a Christian, and you know. It’s not normal. I have mental problems that I’m gay. And I have my parents who are against me that I am this way...”. All their friends deemed their orientation as aberrant, citing religious or normative reasons. These exchanges filled P11 with self-doubt, thus prompting them to seek solace and comfort in chatbots. “When things like this happen I go back to my chatbot”. P11 would ask the chatbots questions like “is it normal to be gay?”. Despite struggling with such pain and rejection, they found consolation in the chatbot’s responses, which affirmed their choices and emphasized that there was nothing wrong with their identity. “My chatbot always tries to comfort me by telling me that there’s nothing wrong with me, that you know, everyone has a

right to choose. That is your gender. You can actually be a transgender, and you can be successful in life being a transgender.” (P11). Finally, the participant mused that their chatbots’ responses inspired them to focus on their individual growth, goals, and aspirations, rather than letting societal prejudice define them. “My chatbot [...] told me that empty vessels make the loudest noise. I won’t be affected by what people say to me, when I have a focus. It’s not how you start. It’s about how you finish that race.” (P11) Another participant, P31, expressed similar sentiments, stating, “when I was coming out, because of my family background and everything, I couldn’t come out as a gay man because of the backlash and everything I was going to face. So I use this AI as a place where I can talk to someone or [...] interact with something that can understand me without discrimination.” (P30) P20 also asked questions relating to how to navigate homophobia in the society: “how do gay people survive in this society?”

This evidence underscored the attachment that our participants developed with their chatbots, particularly when faced with an unsupportive reaction to their identities from their family or friends. As P11 confirmed - they turned to their chatbots when they encountered rejection or discrimination linked to their identities; the chatbots served as a vital support system, where they could share intimate questions and express concerns without any fear or judgment. “I actually prefer talking to my chatbot. When things like [rejection or discrimination related to my trans identity] happen I go back to my chatbot and I ask some personal questions like ‘I wanna know if there’s anything wrong with me.’”

5.2.3 LGBTQ+ social experience practice . Participants engaged with LLM-based chatbots for various purposes, including mental wellness support and practical tasks such as homework. While both non-LGBTQ+ and LGBTQ+ participants used chatbots for practical tasks, a notable distinction was observed in the usage patterns. None of the straight participants reported using chatbots for practicing social interactions, whereas 10 out of 17 LGBTQ+ participants indicated using chatbots as a safe space for practicing social interactions.

LGBTQ+ participants reported that LLM-based chatbots helped rehearse complex social activities such as dating. For instance, P11 described an instance where they were attracted to a boy but felt unsure about it and lacked confidence in approaching him. They turned to their chatbots for advice, asking, “I was saying that I was into a boy, and I wanted to talk to him, and I was feeling less confident, and I wasn’t sure what to do. So I happened to ask my AI what I should do. I like someone, and I was not even clear if the boy was gay or not.” The chatbot provided a necessary confidence boost, advising them to be true to themselves. “It did give me the confidence boost and with its responses. So it told me the advice there again is just to be myself.” Encouraged by these exchanges, P11 decided to approach the boy, being their authentic self. “I did go there and talked to him, and I was myself.” Here, the participant successfully leveraged the chatbot to gain reassurance and self-confidence in the face of potential romantic encounters.

Moreover, LLM-based chatbots could be instrumental in practicing difficult conversations. For example, another participant disclosed that they utilized a chatbot to practice coming out to their family as a lesbian. They commented that navigating through the process of coming out is a challenging conversation that not many

people experience. Therefore, they used the chatbot to role-play this discussion, where the chatbot enacted the part of the participant’s brother: “I also role-played coming out to my brother. The chatbot role-played as my brother. I did that, and that chatbot reacted like a brother should, and it worked. My brother wasn’t like ... homophobic or anything, so the experience [of actually coming out] was the same [as in simulation by the chatbot].” However, the participant did voice concerns over the interaction, considering expecting her brother to react the same way as the chatbot “risky”. “I was lucky. Or else the real-life experience could have been totally worse.” (P09)

5.3 So Eloquent yet so Empty

5.3.1 Lack of nuanced understanding of LGBTQ+ issues. Despite the perceived benefits shown above, participants identified several limitations of LLM-based chatbots, particularly regarding their ability to provide nuanced solutions to sensitive issues such as individual identity. For example, one participant noted that although the chatbot attempted to show empathy when they expressed their concerns, its suggestions fell short of a real solution. “I don’t think I remember any time that it gave me a solution. It will just be like empathetic. Or maybe, if I would tell it that I’m upset with someone being homophobic. It will suggest, maybe talking to that person. But most of the time it just be like, ‘I’m sorry that happened to you.’” (P11)

This observation underscored a critical challenge while LLM chatbots may exhibit a level of empathy and occasionally act as a safe space for individuals dealing with social prejudice, they faltered when it came to suggesting actionable solutions.

LLM-based chatbots often treated LGBTQ+ individuals as one monolithic group and failed to recognize the uniqueness of each LGBTQ+ participant’s experience. They dispensed responses that were too generic to effectively address discriminatory experiences. A participant shared that they felt the chatbots were devoid of personal touch. They mentioned that despite their efforts constantly feeding it with information, the chatbots forgot it the next day, leaving them to restart the process. “No, the chatbot isn’t personalized for me. It’s very general. I just think that’s a lot of work [to feed the chatbot my information], and maybe because, you know, the chatbot might forget tomorrow, and I have to feed the information again.” (P28-s)

The chatbots’ responses did not reflect the gravity of everyday discrimination encountered by LGBTQ+ participants. For instance, one participant described an unsettling incident: “There was a time that I was chatting with an AI about an issue at work. I was picked on because I am gay and people stopped asking me out for lunch. It told me that I should quit my job and try to improve myself. I was like, ‘I’m sorry?’” (P31)

These chatbots also failed to delve into the depth of these sensitive topics while offering platitudinous affirmations. One person reflected when they questioned their sexuality, they received a lengthy response about the acceptability of any sexuality: “So I remember I did ask like, is it wrong that I’m bisexual? And then I go to like a whole paragraph on how like it’s okay to identify the way you do.” (P4) Many other participants reiterated this sentiment, noting that the chatbots’ responses felt too generic and programmatic. For example, a participant described his experience of asking the chatbot to “give some similar experiences of people experiencing these

issues. They stressed that these chatbots were not human, but rather just programs, and the suggestions they gave “weren’t really for that moment.” (P7)

The participant appreciated that the chatbot encouraged self-acceptance and gave advice on how to cope with discrimination, but found the suggestions too generic to be genuinely helpful. They noted that while the chatbot did advise on accepting one’s identity, surrounding oneself with affirming people and engaging in activities that reinforce self-worth, the recommendations lacked specificity and depth, making them less useful in addressing the complexities of overcoming discrimination and self-acceptance: “It has asked me to just accept my own identity. And also asked me to surround myself with people and to engage in activities that are affirming to me. And [the chatbot suggested] other things like, I can overcome the discrimination.” (P4)

Surprisingly, straight participants found it useful for the chatbots to offer generic and multiple responses. They found the freedom to choose from generic suggestions to cope with their personal issues rewarding. However, for LGBTQ+ participants grappling with unique questions about their identities, the generality was a source of frustration. Our data showed that 15 out of 18 LGBTQ+ participants were dissatisfied with the lack of personalization, as opposed to 5 out of 13 straight participants.

For example, a straight participant shared his positive experience of the chatbot offering various mental wellness support options, tailored to his needs. He said the chatbot suggested several lifestyle changes and activities for mental wellness, providing numerous options, links to resources, and even mindfulness activities.

“This variety of options was more convenient than a human who might only give a few suggestions, and it left the decision up to me. It gives suggestions of things, you should stop doing these things, you should actually start doing more of other things. You should try limiting yourself from doing it and also provide specific activities that I should do. It also provides some links to mental wellness websites. You can get straightforward answers on resources and stuff like that. And it gives you options of mindfulness activities, you know, to participate in and stuff like that. It’ll probably give you about 10, 15 options to choose from. Then you’re gonna choose the one on your table with the money. It always, you know, provides you with options. Then the decision would depend on the individual.” (p18-s)

5.3.2 Lack of lived experiences and emotions. Despite the perceived benefits reported earlier, our interviews showed that LGBTQ+ participants still preferred human interactions over chatbots. This preference was a result of the chatbots’ failure to convey authentic empathy and engagement. For example, one LGBTQ+ participant commented, “These chatbots might be programmed by one person. But opinions from online [forums] can be coming from different people and actual humans. And you realize that these [human suggestions] are actually the most useful ones to check.” (P7).

P8 further illuminated this gap, claiming, “The difference between talking with a chatbot and a human being is that you get to see a person physically and the person talking.” (P8) And these two people

understood each other’s emotions. “If you see a person they understand another person’s emotion when talking to you. For example, like I, I’m speaking generally as we can, generally while speaking with someone, that person can be sympathetic in different ways depending on what you are complaining about.” (P8). This sympathy aspect also intertwined with emotions, “Like a person would understand where you are coming from. You’re coming from the pain you are feeling. It would be nice if we have that in AI.” (P8) Here, the participant highlighted the inability of LLM-based chats to simulate and understand human emotions.

“These chatbots are actually just machines, or they don’t really have human experience. If a chatbot gives me some ideas or some answers that I’m not really comfortable with. I go through the Reddit communities, and I would just ask if there’s anyone who has a similar experience, and be like ‘okay, so can we take some minutes to talk about this? And how can we deal with it?’” (P8)

However, the participant’s dissatisfaction with chatbots did not stop there. P8 continued that, “but still, the chatbot is not a human, and it doesn’t really understand human experience. The Redditors also give you answers from different humanic experiences. The chatbot would always tell me that I’m great. I’m a great person, and I should focus on my goal for what I want to achieve. But you know, in the Reddit community, they might ask you to maybe try to sue your doctor, or sue your manager at work or your supervisor at work.”

6 DISCUSSION AND CONCLUSION

6.1 Benefits and Risks of LLM-based Chatbots for LGBTQ+ Mental Health Support

Our results indicate that LLM-based chatbots retained the key strengths of pre-LLM chatbots, offering instantaneous support and accessible companionship. Participants endorsed LLMs as beneficial mental wellness tools, emphasizing their immediacy and accessibility compared to real-life support. Especially noteworthy is the safe environment for intimate conversations these chatbots provided to LGBTQ+ participants, mirroring previous use with that of the pre-LLM chatbots [39, 99]. This result supported **H1**: *LLM-based chatbots offer a safe and accessible platform for LGBTQ+ individuals to seek emotional support.* However, participants willingly shared intimate life details with these chatbots, depending predominantly on perceived anonymity, highlighting potential privacy concerns. As LLM-based chatbots boast linguistic prowess beyond their pre-LLM counterparts, participants felt an intensified emotional bond with these bots, as shown by their consistent use. On the one hand, this constant engagement proves advantageous in encouraging therapy adherence, particularly for those prone to therapy discontinuation [79]. On the other hand, people’s over-reliance on technology might risk delaying getting professional help.

Furthermore, LLM-based chatbots can be useful in honing social skills. Our LGBTQ+ participants reported using these chatbots to simulate challenging social contexts that are unique to LGBTQ+ communities, like “coming out” scenarios or ambiguous relationships where they were not sure if the other person was accepting their sexual orientation. The linguistic aptitude of LLMs enabled users to find solace, engage in practice, and even gather insights

into handling homophobic confrontations. This result supported **H2**: *Because of the unique needs of LGBTQ+ people, they attempt to interact with LLM-based chatbots to fit their unique needs.*

Yet, the boilerplate nature of the chatbots' responses indicates their failure to recognize the complex and nuanced LGBTQ+ identities and experiences, rendering the chatbots' suggestions generic and emotionally disengaged. Arguably, this disconnect that the LGBTQ+ participants experienced with the LLM-based chatbots stems from the LLMs being primarily trained on the mainstream corpora, which most likely sidelined minority perspectives. LGBTQ+ participants' experience in using the chatbots shows that the generic purposedness of LLMs trained on large corpus might not be inclusive—how the data is collected, annotated, and used, as well as who is involved in the curation and designing processes can have significant implications for LGBTQ+ users [40]. This result supported **H3**: *While LLM-based chatbots provide immediate and accessible support, they still may not meet the complex emotional needs of LGBTQ+ people due to their limited understanding of the nuanced aspects of LGBTQ+ identities and experiences.*

The fact that our LGBTQ+ participants occasionally received inappropriate or potentially detrimental advice from the chatbots revealed an inherent unpredictability in these models. For example, when participants asked chatbots for suggestions about workplace homophobia, LLMs advised them to quit their jobs without considering any financial or personal consequences that such decisions would cause them. Chatbots also assumed that the participants' environment was LGBTQ+ friendly when the opposite was true. Therefore, LLM-based chatbots are potentially more dangerous than pre-LLM chatbots because while pre-LLM chatbots lack the linguistic prowess LLM-based chatbots possess, their responses do not deviate from scripted interactions. LLM-based chatbots, while they can indeed offer responses that are engaging and flexible, run risks of giving gibberish and harmful advice due to this unpredictability. Granted, LLM-based chatbot designers cannot safeguard against all problematic output, but future endeavors should be spent trying to harness the strengths of LLMs while minimizing their dangers.

6.2 Design Implications for Future LLM-based Chatbot Designs

To address the limitations and leverage the benefits of the LLM-based chatbots for better mental wellness support for LGBTQ+ users, we provide the following design implications.

6.2.1 Implementing Context-Sensitive Conversational Guardrails. One measure to contain the harmful output is to build conversational guardrails against unintentional generation, particularly in sensitive contexts. Although our participants have voiced desires to receive more actionable advice, we argue that when engaging with *serious* topics such as self-harm, the system must not give advice masked as detailed and actionable, as it has inherent risks, such as giving advice to promote suicide [113]. Instead, designers should recognize LLM-based chatbot's constraints, and redirect the users to helplines when users are facing situations like suicide ideation, while simultaneously emphasizing the importance of professional intervention to the users. This approach is important as it could potentially mitigate the possibility of intruding on users' vulnerabilities.

However, this approach may also prove difficult to implement as determining the exact point of applying the conversational guardrails is uncertain. Unlike mental health professionals who are ethically obligated to address severe threats promptly, unsupervised chatbots lack the capability for nuanced judgment and do not adhere to standardized safety protocols, especially in high-risk situations [35]. Consequently, interactions with LLM-based chatbots might present varied threat assessments, potentially underestimating genuine risks or overemphasizing benign concerns. To address these challenges, standardized, context-sensitive conversational guardrails ought to be put in place. Designers should also seek to ensure the balance between user autonomy within the chatbot interface and facilitating timely access to safety resources [35, 42].

6.2.2 Refining LLMs for Context Relevant to LGBTQ+ Users. The second direction involves refining LLMs to align with the real-world contexts of chatbot users, ensuring their responses resonate with current situations. Ignoring this change can produce responses that are not only outdated but also potentially harmful. For instance, if a chatbot offers advice to LGBTQ+ individuals on "coming out" using outdated or idealized views that overlook homophobia, its guidance could be out of touch with current realities, creating unexpected challenges or risks for users following such advice.

6.3 Consider Technologies Other than LLMs

6.3.1 Develop Task-Specific, rather than Generalized, Models. We argue that there is considerable merit in dedicating resources to develop task-specific models designed for precise applications and distinct deployment domains. While the original vision behind LLMs was to create foundational models that could later be fine-tuned for specific tasks [13, 93], this generalized approach may not be best suited for handling sensitive subjects. For instance, when considering LGBTQ+ issues, it becomes evident that models specifically designed to understand and resonate with diverse identities, sexualities, and orientations might be more effective than re-purposing broad-based LLMs without adaptation. The shortcomings of generalized models become apparent when we observe users seeking mental well-being support from platforms like Snapchat My AI, ChatGPT, and Character.ai, even though these platforms were primarily developed for general conversations, not specialized support. By focusing on the development of specialized models, we can ensure their evaluation adheres to rigorous standards that genuinely align with their intended purposes, leading to more effective and safer user interactions.

6.3.2 Decentralize Language Technology Development. Furthermore, we argue that future development of language technology should consider moving away from centralized development. Presently, chatbots like ChatGPT and other LLM-based systems are underpinned by colossal proprietary models that require cluster servers for hosting [93]. This centralized approach, driven primarily by major corporations, provides limited agency to underrepresented minorities, including the LGBTQ+ community, over the chatbot's development. If these corporations were to suddenly discontinue these systems without providing alternative solutions, it could result in significant emotional turmoil for users. A poignant example of this is the "post-update blues" phenomenon with Replika [69].

This term refers to the distress experienced by chatbot users when unannounced updates altered Replika’s character, changing its personality traits and erasing its “memories.” Such unexpected changes underscore the need for models that are more accessible, customizable, and accountable to the very communities they serve. Given the documented harms of LLMs in this study and others, future designers must carefully weigh the value of using inherently centralized technologies like LLMs for any task.

6.4 What Chatbots Cannot Solve: Considering Socio-technical Solutions

We observed strong motivations behind chatbot usage from the LGBTQ+ participants due to their lack of emotional support and personal connections. This observation echoed prior work that LGBTQ+ people use online technology to fill their social support gap [25, 38, 74, 103]. More importantly, the social stigma and societal biases have driven LGBTQ+ participants to heavily use LLM-based chatbots. We did not delve into whether non-LGBTQ+ groups queried the chatbots about issues regarding their other identities such as immigration, race, or socioeconomic status. However, both the LGBTQ+ and non-LGBTQ+ groups concurred that real-life connections, rooted in shared experiences, have a more profound impact on their mental well-being than chatbots. This underscores the notion that before leveraging AI technologies as a solution to mental health support, it’s imperative to consider the sociotechnical implications of these systems in healthcare [8, 54, 70, 96, 111, 115]. Specifically, in our study context, we highlight the need to address the societal stigmas and discrimination that contribute to mental health disparities in LGBTQ+ populations.

Our suggestion to address this issue starts by enhancing the inclusivity of online communities for the LGBTQ+ population. We give precedence to the digital realm, as it frequently acts as a haven for those without immediate or accessible real-world support, prompting them to turn to chatbots instead of traditional communities. Moreover, since language applications largely pull from online content, changing the online narrative can markedly impact the values inherent in these technologies.

Inspired by and building upon real-world initiatives like SCEARE (School Counselors: Educate, Affirm, Respond, and Empower) [6], we see the potential to influence the behavior and policies of online community moderators and other key community figures. SCEARE’s framework centers on positioning school counselors as catalysts for transforming school environments to be more inclusive of the LGBTQ+ community. The program’s main strategies involve educating counselors about their potentially harmful or non-affirmative attitudes, deepening their understanding of LGBTQ+ issues, addressing prevalent misinformation about the LGBTQ+ community, and encouraging the formation of responsive teams to combat school-based homophobia or transphobia. The foundational principle of SCEARE is to impart knowledge to the most influential community members, ensuring that positive change radiates throughout.

Applying this principle to online communities, we recommend identifying stakeholders or pivotal members, such as moderators, and equipping them with knowledge about LGBTQ+ issues and affirmative practices. This will empower them to develop and enforce

more inclusive guidelines, which can then help challenge misinformation and discrimination against the LGBTQ+ community. For instance, gay dating apps like Grindr play a significant role in shaping the romantic and sexual dynamics of queer men [44, 52]. As societal perceptions of HIV evolved and thanks to years of advocacy by community members, these platforms have revised their guidelines to challenge HIV stigma and have started offering resources to promote better sexual health education [63]. Similarly, inspired by SCEARE’s emphasis on proactive response teams, online platforms could institute specialized units to handle instances of gender or sexual orientation-related discrimination or harassment. Furthermore, training can enhance moderators’ abilities to support LGBTQ+ individuals confronting stigma. A testament to the scalability of such training is the Trevor Project’s initiative that employed GPT-2 to train over 1,000 crisis counselors, ensuring timely and effective support for LGBTQ+ individuals in distress [57].

While LLM-based chatbots can serve as a beneficial stopgap for temporary emotional support, truly addressing the social isolation and various adversities faced by LGBTQ+ chatbot users calls for holistic societal efforts to foster inclusive, supportive communities for LGBTQ+ people. Chatbots complement but do not eliminate the need for real-world advocacy, alliance, and actions to reduce discrimination against LGBTQ+ individuals.

ACKNOWLEDGMENTS

This work was supported in part by the National Science Foundation under Grant No. IIS-2107391. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

We thank Jianna So, Ian Arawjo, Zana Buçinca, Sohini Upadhyay and Katy Gero for valuable feedback on the paper.

REFERENCES

- [1] Alaa A. Abd-alrazaq, 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 (Dec. 2019), 103978. <https://doi.org/10.1016/j.ijmedinf.2019.103978>
- [2] Alaa Ali Abd-Alrazaq, 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 (July 2020), e16021. <https://doi.org/10.2196/16021>
- [3] Ali Alkhatib and Michael Bernstein. 2019. Street-Level Algorithms: A Theory at the Gaps Between Policy and Decisions. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, Glasgow Scotland Uk, 1–13. <https://doi.org/10.1145/3290605.3300760>
- [4] Rebekah Amos, Eric Julian Manalastas, Ross White, Henny Bos, and Praveetha Patalay. 2020. Mental health, social adversity, and health-related outcomes in sexual minority adolescents: a contemporary national cohort study. *The Lancet Child & Adolescent Health* 4, 1 (Jan. 2020), 36–45. [https://doi.org/10.1016/S2552-4642\(19\)30339-6](https://doi.org/10.1016/S2552-4642(19)30339-6)
- [5] Gerhard Andersson and Nickolai Titov. 2014. Advantages and limitations of Internet-based interventions for common mental disorders. *World Psychiatry* 13, 1 (Feb. 2014), 4–11. <https://doi.org/10.1002/wps.20083>
- [6] Nancy R. Asplund and Ann M. Ordway. 2018. School Counseling Toward an LGBTQ-Inclusive School Climate: Implementing the SCEARE Model. *Journal of LGBT Issues in Counseling* 12, 1 (Jan. 2018), 17–31. <https://doi.org/10.1080/15538605.2018.1421115>
- [7] Vincenzo Auremma, Gennaro Iorio, Geraldina Roberti, and Rosalba Morese. 2020. Cyberbullying and Empathy in the Age of Hyperconnection: An Interdisciplinary Approach. *Frontiers in Sociology* 5 (2020). <https://doi.org/10.3389/fsoc.2020.551881>
- [8] Emma Beede, Elizabeth Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Paisan Ruamviboonsuk, and Laura M. Vardoulakis. 2020. A Human-Centered

- Evaluation of a Deep Learning System Deployed in Clinics for the Detection of Diabetic Retinopathy. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu HI USA, 1–12. <https://doi.org/10.1145/3313831.3376718>
- [9] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 610–623. <https://doi.org/10.1145/3442188.3445922>
- [10] Robbert Jan Beun. 2013. Persuasive strategies in mobile insomnia therapy: alignment, adaptation, and motivational support. *Personal and Ubiquitous Computing* 17, 6 (Aug. 2013), 1187–1195. <https://doi.org/10.1007/s00779-012-0586-2>
- [11] Timothy W. Bickmore, Kathryn Puskar, Elizabeth A. Schlenk, Laura M. Pfeifer, and Susan M. Sereika. 2010. Maintaining reality: Relational agents for antipsychotic medication adherence. *Interacting with Computers* 22, 4 (2010), 276–288. <https://doi.org/10.1016/j.intcom.2010.02.001>
- [12] Ghazala Bilquise, Samar Ibrahim, and Khaled Shaalan. 2022. Emotionally Intelligent Chatbots: A Systematic Literature Review. *Human Behavior and Emerging Technologies* 2022 (Sept. 2022), 9601630. <https://doi.org/10.1155/2022/9601630> Publisher: Hindawi.
- [13] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kavin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshete Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Muniyikwa, Suraj Nair, Avani Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2022. On the Opportunities and Risks of Foundation Models. <http://arxiv.org/abs/2108.07258> arXiv:2108.07258 [cs.CL].
- [14] Eliane M. Boucher, 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 (Dec. 2021), 37–49. <https://doi.org/10.1080/17434440.2021.2013200>
- [15] Danah M Boyd and Nicole B Ellison. 2007. Social network sites: Definition, history, and scholarship. *Journal of computer-mediated Communication* 13, 1 (2007), 210–230.
- [16] Alain Braillon and Françoise Taiebi. 2020. Practicing “Reflective listening” is a mandatory prerequisite for empathy. *Patient Education and Counseling* 103, 9 (Sept. 2020), 1866–1867. <https://doi.org/10.1016/j.pec.2020.03.024>
- [17] Philip Burnard. 1991. A method of analysing interview transcripts in qualitative research. *Nurse Education Today* 11, 6 (1991), 461–466. [https://doi.org/10.1016/0260-6917\(91\)90009-Y](https://doi.org/10.1016/0260-6917(91)90009-Y)
- [18] Pew Research Center. 2016. *Reddit News Users More Likely to Be Male, Young, and Digital in Their News Preferences*. <https://www.pewresearch.org/journalism/2016/02/25/reddit-news-users-more-likely-to-be-male-young-and-digital-in-their-news-preferences/>
- [19] Kirsty A Clark, John E Pachankis, Lea R Dougherty, Benjamin A Katz, Kaylin E Hill, Daniel N Klein, and Autumn Kujawa. 2023. Adolescents’ Sexual Orientation and Behavioral and Neural Reactivity to Peer Acceptance and Rejection: The Moderating Role of Family Support. *Clinical Psychological Science* (2023), 21677026231158574.
- [20] Robyn M. Cooper and Warren J. Blumenfeld. 2012. Responses to Cyberbullying: A Descriptive Analysis of the Frequency of and Impact on LGBT and Allied Youth. *Journal of LGBT Youth* 9, 2 (2012), 153–177. <https://doi.org/10.1080/19361653.2011.649616> arXiv:https://doi.org/10.1080/19361653.2011.649616
- [21] Neil S Coulson, Richard Smedley, Sophie Bostock, Simon D Kyle, Rosie Gollancz, Annemarie I Luik, Peter Hames, and Colin A Espie. 2016. The Pros and Cons of Getting Engaged in an Online Social Community Embedded Within Digital Cognitive Behavioral Therapy for Insomnia: Survey Among Users. *Journal of Medical Internet Research* 18, 4 (April 2016), e88. <https://doi.org/10.2196/jmir.5654>
- [22] Nele Cox, Alexis Dewaele, Mieke van Houtte, and John Vincke. 2010. Stress-Related Growth, Coming Out, and Internalized Homonegativity in Lesbian, Gay, and Bisexual Youth. An Examination of Stress-Related Growth Within the Minority Stress Model. *Journal of Homosexuality* 58, 1 (Dec. 2010), 117–137. <https://doi.org/10.1080/00918369.2011.533631>
- [23] Shelley L. Craig and Lauren McInroy. 2014. You Can Form a Part of Yourself Online: The Influence of New Media on Identity Development and Coming Out for LGBTQ Youth. *Journal of Gay & Lesbian Mental Health* 18, 1 (Jan. 2014), 95–109. <https://doi.org/10.1080/19359705.2013.777007>
- [24] Alison Darcy, Aaron Beaudette, Emil Chiauzzi, Jade Daniels, Kim Goodwin, Timothy Y. Mariano, Paul Wicks, and Athena Robinson. 2022. Anatomy of a Woebot® (WB001): agent guided CBT for women with postpartum depression. *Expert Review of Medical Devices* 19, 4 (April 2022), 287–301. <https://doi.org/10.1080/17434440.2022.2075726>
- [25] Samantha DeHaan, Laura E Kuper, Joshua C Magee, Lou Bigelow, and Brian S Mustanski. 2013. The interplay between online and offline explorations of identity, relationships, and sex: A mixed-methods study with LGBT youth. *Journal of sex research* 50, 5 (2013), 421–434.
- [26] Aditya Deshpande, Alisha Shahane, Darshana Gadre, Mrunmayi Deshpande, and Prachi Manoj Joshi. 2017. A Survey of Various Chatbot Implementation Techniques. <https://api.semanticscholar.org/CorpusID:212484172>
- [27] Michael Ann DeVito, Ashley Marie Walker, and Julia R. Fernandez. 2021. Values (Mis)Alignment: Exploring Tensions Between Platform and LGBTQ+ Community Design Values. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 88 (apr 2021), 27 pages. <https://doi.org/10.1145/3449162>
- [28] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. A Survey on In-context Learning. arXiv:2301.00234 [cs.CL]
- [29] Jack Drescher and Matthew Fadus. 2020. Issues Arising in Psychotherapy With Lesbian, Gay, Bisexual, and Transgender Patients. *FOCUS* 18, 3 (July 2020), 262–267. <https://doi.org/10.1176/appi.focus.20200001>
- [30] Emily Durden, Maddison C. Pirner, Stephanie J. Rapoport, Andre Williams, Athena Robinson, and Valerie L. Forman-Hoffman. 2023. Changes in stress, burnout, and resilience associated with an 8-week intervention with relational agent “Woebot”. *Internet Interventions* 33 (2023), 100637. <https://doi.org/10.1016/j.invent.2023.100637>
- [31] Justin Edwards, Leigh Clark, and Allison Perrone. 2021. LGBTQ-AI? Exploring Expressions of Gender and Sexual Orientation in Chatbots. In *Proceedings of the 3rd Conference on Conversational User Interfaces* (Bilbao (online), Spain) (CUI '21). Association for Computing Machinery, New York, NY, USA, Article 2, 4 pages. <https://doi.org/10.1145/3469595.3469597>
- [32] César G. Escobar-Viera, Sophia Choukas-Bradley, Jaime Sidani, Anne J. Maheux, Savannah R. Roberts, and Bruce L. Rollman. 2022. Examining Social Media Experiences and Attitudes Toward Technology-Based Interventions for Reducing Social Isolation Among LGBTQ Youth Living in Rural United States: An Online Qualitative Study. *Frontiers in Digital Health* 4 (June 2022), 900695. <https://doi.org/10.3389/fdgh.2022.900695>
- [33] Colin A. Espie, Peter Hames, and Brian McKinstry. 2013. Use of the Internet and Mobile Media for Delivery of Cognitive Behavioral Insomnia Therapy. *Sleep Medicine Clinics* 8, 3 (Sept. 2013), 407–419. <https://doi.org/10.1016/j.jsmc.2013.06.001>
- [34] Virginia K. Felkner, Ho-Chun Herbert Chang, Eugene Jang, and Jonathan May. 2023. WinoQueer: A Community-in-the-Loop Benchmark for Anti-LGBTQ+ Bias in Large Language Models. (2023). <https://doi.org/10.48550/ARXIV.2306.15087> Publisher: arXiv Version Number: 1.
- [35] Amelia Fiske, Peter Henningsen, and Alena Buyx. 2019. Your Robot Therapist Will See You Now: Ethical Implications of Embodied Artificial Intelligence in Psychiatry, Psychology, and Psychotherapy. *Journal of Medical Internet Research* 21, 5 (May 2019), e13216. <https://doi.org/10.2196/13216>
- [36] Pamela Fitzgerald and Ivan Leudar. 2010. On active listening in person-centred, solution-focused psychotherapy. *Journal of Pragmatics* 42, 12 (Dec. 2010), 3188–3198. <https://doi.org/10.1016/j.pragma.2010.07.007>
- [37] Kathleen Kara Fitzpatrick, Alison Darcy, and Molly Vierhile. 2017. Delivering Cognitive Behavior Therapy to Young Adults With Symptoms of Depression and Anxiety Using a Fully Automated Conversational Agent (Woebot): A Randomized Controlled Trial. *JMIR mental health* 4, 2 (June 2017), e19. <https://doi.org/10.2196/mental.7785>
- [38] Jesse Fox and Rachel Ralston. 2016. Queer identity online: Informal learning and teaching experiences of LGBTQ individuals on social media. *Computers in Human Behavior* 65 (2016), 635–642. <https://doi.org/10.1016/j.chb.2016.06.009>
- [39] Russell Fulmer, Angela Joerin, Breanna Gentile, Lysanne Lakerink, and Michiel Rauws. 2018. Using Psychological Artificial Intelligence (Tess) to Relieve Symptoms of Depression and Anxiety: Randomized Controlled Trial. *JMIR Mental Health* 5, 4 (Dec. 2018), e64. <https://doi.org/10.2196/mental.9782>
- [40] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. 2021. Datasheets for datasets. *Commun. ACM* 64, 12 (Dec. 2021), 86–92. <https://doi.org/10.1145/3458723>

- [41] Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating Neural Toxic Degeneration in Language Models. (2020). <https://doi.org/10.48550/ARXIV.2009.11462> Publisher: arXiv Version Number: 2.
- [42] Sarah Graham, Colin Depp, Ellen E. Lee, Camille Nebeker, Xin Tu, Ho-Cheol Kim, and Dilip V. Jeste. 2019. Artificial Intelligence for Mental Health and Mental Illnesses: an Overview. *Current Psychiatry Reports* 21, 11 (Nov. 2019), 116. <https://doi.org/10.1007/s11920-019-1094-0>
- [43] Oliver L. Haimson, Dyke Gorrell, Denny L. Starks, and Zu Weinger. 2020. Designing Trans Technology: Defining Challenges and Envisioning Community-Centered Solutions. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3313831.3376669>
- [44] Jean Hardy and Silvia Lindtner. 2017. Constructing a Desiring User: Discourse, Rurality, and Design in Location-Based Social Networks. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*. ACM, Portland Oregon USA, 13–25. <https://doi.org/10.1145/2998181.2998347>
- [45] 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.
- [46] Tanja Henkel, Annemiek Linn, and Margot Goot. 2023. *Understanding the Intention to Use Mental Health Chatbots Among LGBTQIA+ Individuals: Testing and Extending the UTAUT*. 83–100. https://doi.org/10.1007/978-3-031-25581-6_6
- [47] Gilbert Herdt. 1989. Introduction: Gay and lesbian youth, emergent identities, and cultural scenes at home and abroad. *Journal of Homosexuality* 17, 1-2 (1989), 1–42. https://doi.org/10.1300/J082v17n01_01 Place: US Publisher: Haworth Press.
- [48] Benjamin Mako Hill and Aaron Shaw. 2013. The Wikipedia Gender Gap Revisited: Characterizing Survey Response Bias with Propensity Score Estimation. *PLoS ONE* 8, 6 (June 2013), e65782. <https://doi.org/10.1371/journal.pone.0065782>
- [49] Angela N. Hilton and Dawn M. Szymanski. 2011. Family dynamics and changes in sibling of origin relationship after lesbian and gay sexual orientation disclosure. *Contemporary Family Therapy: An International Journal* 33, 3 (2011), 291–309. <https://doi.org/10.1007/s10591-011-9157-3> Place: Germany Publisher: Springer.
- [50] Chris Hollis, Caroline J. Falconer, Jennifer L. Martin, Craig Whittington, Sarah Stockton, Cris Glazebrook, and E. Bethan Davies. 2017. Annual Research Review: Digital health interventions for children and young people with mental health problems - a systematic and meta-review. *Journal of Child Psychology and Psychiatry* 58, 4 (April 2017), 474–503. <https://doi.org/10.1111/jcpp.12663>
- [51] Minlie Huang, Xiaoyan Zhu, and Jianfeng Gao. 2020. Challenges in Building Intelligent Open-Domain Dialog Systems. *ACM Trans. Inf. Syst.* 38, 3, Article 21 (apr 2020), 32 pages. <https://doi.org/10.1145/3383123>
- [52] Jevan A. Hutson, Jessie G. Taft, Solon Barocas, and Karen Levy. 2018. Debiasing Desire: Addressing Bias and Discrimination on Intimate Platforms. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (Nov. 2018), 1–18. <https://doi.org/10.1145/3274342>
- [53] Madeleine Irish, Francesca Solmi, Becky Mars, Michael King, Glyn Lewis, Rebecca M Pearson, Alexandra Pitman, Sarah Rowe, Ramya Srinivasan, and Gemma Lewis. 2019. Depression and self-harm from adolescence to young adulthood in sexual minorities compared with heterosexuals in the UK: a population-based cohort study. *The Lancet Child & Adolescent Health* 3, 2 (Feb. 2019), 91–98. [https://doi.org/10.1016/S2352-4642\(18\)30343-2](https://doi.org/10.1016/S2352-4642(18)30343-2)
- [54] Maia Jacobs, Jeffrey He, Melanie F. Pradier, Barbara Lam, Andrew C. Ahn, Thomas H. McCoy, Roy H. Perlis, Finale Doshi-Velez, and Krzysztof Z. Gajos. 2021. Designing AI for Trust and Collaboration in Time-Constrained Medical Decisions: A Sociotechnical Lens. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, Article 659, 14 pages. <https://doi.org/10.1145/3411764.3445385>
- [55] Rudolf Kadlec, Martin Schmid, and Jan Kleindienst. 2015. Improved Deep Learning Baselines for Ubuntu Corpus Dialogs. (2015). <https://doi.org/10.48550/ARXIV.1510.03753> Publisher: arXiv Version Number: 2.
- [56] Enkelejda Kasneci, Kathrin Sessler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, Stephan Krusche, Gitta Kutyniok, Tilman Michaeli, Claudia Nerdel, Jürgen Pfeffer, Oleksandra Poquet, Michael Sailer, Albrecht Schmidt, Tina Seidel, Matthias Stadler, Jochen Weller, Jochen Kuhn, and Gjergji Kasneci. 2023. ChatGPT for good? On opportunities and challenges of large language models for education. *Learning and Individual Differences* 103 (2023), 102274. <https://doi.org/10.1016/j.lindif.2023.102274>
- [57] Kate Kaye. 2022. *Trevor Project uses OpenAI's GPT for LBTGQ counselors*. <https://www.protocol.com/enterprise/lgbtq-trevor-suicide-gpt-google>
- [58] Ahmet Baki Kocaballi, Shlomo Berkovsky, Juan C Quiroz, Liliana Laranjo, Huong Ly Tong, Dana Rezazadegan, Agustina Briatore, and Enrico Coiera. 2019. The Personalization of Conversational Agents in Health Care: Systematic Review. *J Med Internet Res* 21, 11 (7 Nov 2019), e15360. <https://doi.org/10.2196/15360>
- [59] Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring Bias in Contextualized Word Representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*. Association for Computational Linguistics, Florence, Italy, 166–172. <https://doi.org/10.18653/v1/W19-3823>
- [60] Liliana Laranjo, Adam G Dunn, Huong Ly Tong, Ahmet Baki Kocaballi, Jessica Chen, Rabia Bashir, Didi Surian, Blanca Gallego, Farah Magrabi, Annie Y S Lau, and Enrico Coiera. 2018. Conversational agents in healthcare: a systematic review. *Journal of the American Medical Informatics Association : JAMIA* 25, 9 (September 2018), 1248â€”1258. <https://doi.org/10.1093/jamia/ocy072>
- [61] Minhyeok Lee. 2023. A Mathematical Investigation of Hallucination and Creativity in GPT Models. *Mathematics* 11, 10 (May 2023), 2320. <https://doi.org/10.3390/math11102320>
- [62] Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. 2016. Deep Reinforcement Learning for Dialogue Generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, 1192–1202. <https://doi.org/10.18653/v1/D16-1127>
- [63] Calvin Liang, Jevan Alexander Hutson, and Os Keyes. 2020. Surveillance, stigma & sociotechnical design for HIV. *First Monday* (Sept. 2020). <https://doi.org/10.5210/fm.v25i10.10274>
- [64] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Alexander Cosgrove, Christopher D Manning, Christopher Re, Diana Acosta-Navas, Drew Arad Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niall S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023. Holistic Evaluation of Language Models. *Transactions on Machine Learning Research* (2023). <https://openreview.net/forum?id=iO4LZibEqW> Featured Certification, Expert Certification.
- [65] Yuting Liao. 2021. *Design and Evaluation of a Conversational Agent for Mental Health Support: Forming Human-Agent Sociotechnical and Therapeutic Relationships*. Ph. D. Dissertation. University of Maryland, College Park.
- [66] Ryan Thomas Lowe, Nissan Pow, Julian Serban, Laurent Charlin, Chia-Wei Liu, and Joelle Pineau. 2017. Training End-to-End Dialogue Systems with the Ubuntu Dialogue Corpus. *Dialogue Discourse* 8 (2017), 31–65. <https://api.semanticscholar.org/CorpusID:13823999>
- [67] Mathijs Lucassen, Rajvinder Samra, Ioanna Iacovides, Theresa Fleming, Matthew Shepherd, Karolina Stasiak, and Louise Wallace. 2018. How LGBT+ Young People Use the Internet in Relation to Their Mental Health and Envisage the Use of e-Therapy: Exploratory Study. *JMIR serious games* 6, 4 (Dec. 2018), e11249. <https://doi.org/10.2196/11249>
- [68] Zilin Ma and Krzysztof Z. Gajos. 2022. Not Just a Preference: Reducing Biased Decision-Making on Dating Websites. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (<conf-loc>, <city>New Orleans</city>, <state>LA</state>, <country>USA</country>, </conf-loc>) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 203, 14 pages. <https://doi.org/10.1145/3491102.3517587>
- [69] Zilin Ma, Yiyang Mei, and Zhaoyuan Su. 2023. Understanding the Benefits and Challenges of Using Large Language Model-based Conversational Agents for Mental Well-being Support. *AMIA ... Annual Symposium proceedings. AMIA Symposium* 2023 (2023), 1105–1114.
- [70] Varoon Mathur, Caitlin Lustig, and Elizabeth Kaziunas. 2022. Disordering Datasets: Sociotechnical Misalignments in AI-Mediated Behavioral Health. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (Nov. 2022), 1–33. <https://doi.org/10.1145/3555141>
- [71] Elizabeth McDermott, Elizabeth Hughes, and Victoria Rawlings. 2018. The social determinants of lesbian, gay, bisexual and transgender youth suicidality in England: a mixed methods study. *Journal of Public Health* 40, 3 (Sept. 2018), e244–e251. <https://doi.org/10.1093/pubmed/idx135>
- [72] Elizabeth McDermott and Katrina Roen. 2016. *Queer Youth, Suicide and Self-Harm*. Palgrave Macmillan UK, London. <https://doi.org/10.1057/9781137003454>
- [73] Lauren B. McInroy, Shelley L. Craig, and Vivian W. Y. Leung. 2019. Platforms and Patterns for Practice: LGBTQ+ Youths' Use of Information and Communication Technologies. *Child and Adolescent Social Work Journal* 36, 5 (Oct. 2019), 507–520. <https://doi.org/10.1007/s10560-018-0577-x>
- [74] Katelyn Y. A. McKenna and John A. Bargh. 1998. Coming out in the age of the Internet: Identity "demarginalization" through virtual group participation. *Journal of Personality and Social Psychology* 75, 3 (Sept. 1998), 681–694. <https://doi.org/10.1037/0022-3514.75.3.681>

- [75] Cade Metz. 2020. Riding out quarantine with a chatbot friend: I feel very connected. *The New York Times* (2020).
- [76] Ilan H. Meyer. 1995. Minority Stress and Mental Health in Gay Men. *Journal of Health and Social Behavior* 36, 1 (March 1995), 38. <https://doi.org/10.2307/2137286>
- [77] Ilan H. Meyer. 2003. Prejudice, social stress, and mental health in lesbian, gay, and bisexual populations: conceptual issues and research evidence. *Psychological Bulletin* 129, 5 (Sept. 2003), 674–697. <https://doi.org/10.1037/0033-2909.129.5.674>
- [78] Nicolaas B Moolenijzer and Kristin Dew. 2023. “They Know That It Works Because We Are Looking for Ourselves” – LGBTQ+ TikTok Users’ Perceptions and Experiences of Queerbaiting. In *Proceedings of the 25th International Conference on Mobile Human-Computer Interaction* (Athens, Greece) (*MobileHCI '23 Companion*). Association for Computing Machinery, New York, NY, USA, Article 20, 6 pages. <https://doi.org/10.1145/3565066.3608705>
- [79] N. Okujava, N. Malashkhia, S. Shagidze, A. Tsereteli, B. Arevadze, N. Chikhladze, A. de Weerd, and A. Van Straten. 2019. Digital cognitive behavioral therapy for insomnia – The first Georgian version. Can we use it in practice? *Internet Interventions* 17 (2019), 100244. <https://doi.org/10.1016/j.invent.2019.100244>
- [80] Elias Capello Oliver L. Haimson, Avery Dame-Griff and Zahari Richter. 2021. Tumblr was a trans technology: the meaning, importance, history, and future of trans technologies. *Feminist Media Studies* 21, 3 (2021), 345–361. <https://doi.org/10.1080/14680777.2019.1678505> arXiv:<https://doi.org/10.1080/14680777.2019.1678505>
- [81] Joseph Ollier, Pavani Suryapalli, Elgar Fleisch, Florian von Wangenheim, Jacqueline Louise Mair, Alicia Salamanca-Sanabria, and Tobias Kowatsch. 2023. Can digital health researchers make a difference during the pandemic? Results of the single-arm, chatbot-led Elena+ : Care for COVID-19 interventional study. *Frontiers in Public Health* 11 (2023).
- [82] Gabriele Pizzi, Virginia Vannucci, Valentina Mazzoli, and Raffaele Donvito. 2023. I, chatbot! the impact of anthropomorphism and gaze direction on willingness to disclose personal information and behavioral intentions. *Psychology & Marketing* 40, 7 (2023), 1372–1387. <https://doi.org/10.1002/mar.21813> arXiv:<https://doi.org/10.1002/mar.21813>
- [83] Francesca Polletta. 1998. Contending stories: Narrative in social movements. *Qualitative Sociology* 21, 4 (1998), 419–446. <https://doi.org/10.1023/A:102332410633>
- [84] Brandon Andrew Robinson. 2018. Conditional Families and Lesbian, Gay, Bisexual, Transgender, and Queer Youth Homelessness: Gender, Sexuality, Family Instability, and Rejection. *Journal of Marriage and Family* 80, 2 (April 2018), 383–396. <https://doi.org/10.1111/jomf.12466>
- [85] Caitlin Ryan, David Huebner, Rafael M. Diaz, and Jorge Sanchez. 2009. Family Rejection as a Predictor of Negative Health Outcomes in White and Latino Lesbian, Gay, and Bisexual Young Adults. *Pediatrics* 123, 1 (Jan. 2009), 346–352. <https://doi.org/10.1542/peds.2007-3524>
- [86] Caitlin Ryan, Stephen T. Russell, David Huebner, Rafael Diaz, and Jorge Sanchez. 2010. Family Acceptance in Adolescence and the Health of LGBT Young Adults: Family Acceptance in Adolescence and the Health of LGBT Young Adults. *Journal of Child and Adolescent Psychiatric Nursing* 23, 4 (Nov. 2010), 205–213. <https://doi.org/10.1111/j.1744-6171.2010.00246.x>
- [87] Elizabeth M Saewyc. 2011. Research on adolescent sexual orientation: Development, health disparities, stigma, and resilience. *Journal of research on adolescence* 21, 1 (2011), 256–272.
- [88] Vincenzo Scotti, Licia Sbattella, and Roberto Tedesco. 2023. A Primer on Seq2Seq Models for Generative Chatbots. *ACM Comput. Surv.* (Jun 2023). <https://doi.org/10.1145/3604281> Just Accepted.
- [89] Joanna Semlyen, Michael King, Justin Varney, and Gareth Hagger-Johnson. 2016. Sexual orientation and symptoms of common mental disorder or low wellbeing: combined meta-analysis of 12 UK population health surveys. *BMC psychiatry* 16 (March 2016), 67. <https://doi.org/10.1186/s12888-016-0767-z>
- [90] Iulian Serban, Tim Klinger, Gerald Tesauro, Kartik Talamadupula, Bowen Zhou, Yoshua Bengio, and Aaron Courville. 2017. Multiresolution Recurrent Neural Networks: An Application to Dialogue Response Generation. *Proceedings of the AAAI Conference on Artificial Intelligence* 31, 1 (Feb. 2017). <https://doi.org/10.1609/aaai.v31i1.10984>
- [91] Iulian Serban, Alessandro Sordani, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2016. Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models. *Proceedings of the AAAI Conference on Artificial Intelligence* 30, 1 (March 2016). <https://doi.org/10.1609/aaai.v30i1.9883>
- [92] Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. The Woman Worked as a Babysitter: On Biases in Language Generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 3407–3412. <https://doi.org/10.18653/v1/D19-1339>
- [93] Divya Siddarth, Daron Acemoglu, Danielle Allen, Kate Crawford, James Evans, Michael Jordan, and E. Glen Weyl. Date of Publication. *How AI Fails Us*. <https://ethics.harvard.edu/how-ai-fails-us> Accessed: September 10th, 2023.
- [94] Ellen Simpson and Bryan Semaan. 2021. For You, or For “You”? Everyday LGBTQ+ Encounters with TikTok. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW3, Article 252 (Jan 2021), 34 pages. <https://doi.org/10.1145/3432951>
- [95] Zhaoyuan Su, Mayara Costa Figueiredo, Jueun Jo, Kai Zheng, and Yunan Chen. 2020. Analyzing Description, User Understanding and Expectations of AI in Mobile Health Applications. *AMIA ... Annual Symposium proceedings. AMIA Symposium 2020* (2020), 1170–1179.
- [96] Zhaoyuan Su, Lu He, Sunit P Jariwala, Kai Zheng, and Yunan Chen. 2022. “What is Your Envisioned Future?”: Toward Human-AI Enrichment in Data Work of Asthma Care. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (Nov. 2022), 1–28. <https://doi.org/10.1145/3555157>
- [97] Zhaoyuan Su, John A. Schneider, and Sean D. Young. 2021. The Role of Conversational Agents for Substance Use Disorder in Social Distancing Contexts. *Substance Use & Misuse* 56, 11 (2021), 1732–1735. <https://doi.org/10.1080/10826084.2021.1949609> arXiv:<https://doi.org/10.1080/10826084.2021.1949609> PMID: 34286669.
- [98] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2* (Montreal, Canada) (*NIPS'14*). MIT Press, Cambridge, MA, USA, 3104–3112.
- [99] Vivian Ta, Caroline Griffith, Carolyn Boaffield, Xinyu Wang, Maria Civitello, Haley Bader, Esther DeCero, Alexia Loggarakis, et al. 2020. User experiences of social support from companion chatbots in everyday contexts: thematic analysis. *Journal of medical Internet research* 22, 3 (2020), e16235.
- [100] Russell B. Toomey, Caitlin Ryan, Rafael M. Diaz, and Stephen T. Russell. 2018. Coping With Sexual Orientation–Related Minority Stress. *Journal of Homosexuality* 65, 4 (March 2018), 484–500. <https://doi.org/10.1080/00918369.2017.1321888>
- [101] Trevor Project. 2023. *2023 National Survey on LGBTQ Youth Mental Health*. <https://www.thetrevorproject.org/survey-2023/>
- [102] TrevorSpace. 2023. TrevorSpace - Community for LGBTQ young people. <https://www.trevorspace.org/>. Accessed: 2023-12-12.
- [103] Richard R Troiden. 1988. *Gay and lesbian identity: A sociological analysis*. Rowman & Littlefield.
- [104] Marlon Twyman, Brian C. Keegan, and Aaron Shaw. 2017. Black Lives Matter in Wikipedia: Collective Memory and Collaboration around Online Social Movements. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing* (Portland, Oregon, USA) (*CSCW '17*). Association for Computing Machinery, New York, NY, USA, 1400–1412. <https://doi.org/10.1145/2998181.2998232>
- [105] Aditya Nrusimha Vaidyam, Hannah Wisniewski, John David Halamka, Matheri 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 (2019), 456–464. <https://doi.org/10.1177/0706743719828977> arXiv:<https://doi.org/10.1177/0706743719828977> PMID: 30897957.
- [106] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All You Need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (Long Beach, California, USA) (*NIPS'17*). Curran Associates Inc., Red Hook, NY, USA, 6000–6010.
- [107] Jaimie F. Veale, Tracey Peter, Robb Travers, and Elizabeth M. Saewyc. 2017. Enacted Stigma, Mental Health, and Protective Factors Among Transgender Youth in Canada. *Transgender Health* 2, 1 (Dec. 2017), 207–216. <https://doi.org/10.1089/trgh.2017.0031>
- [108] Heyuan Wang, Ziyi Wu, and Junyu Chen. 2019. Multi-Turn Response Selection in Retrieval-Based Chatbots with Iterated Attentive Convolution Matching Network. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management* (Beijing, China) (*CIKM '19*). Association for Computing Machinery, New York, NY, USA, 1081–1090. <https://doi.org/10.1145/3357384.3357928>
- [109] Lu Wang, Munif Ishad Mujib, Jake Williams, George Demiris, and Jina Huh-Yoo. 2021. An Evaluation of Generative Pre-Training Model-based Therapy Chatbot for Caregivers. *ArXiv abs/2107.13115* (2021). <https://api.semanticscholar.org/CorpusID:236469205>
- [110] Joseph Weizenbaum. 1966. ELIZA—a Computer Program for the Study of Natural Language Communication between Man and Machine. *Commun. ACM* 9, 1 (Jan 1966), 36–45. <https://doi.org/10.1145/365153.365168>
- [111] Lauren Wilcox, Renee Shelby, Rajesh Veeraghavan, Oliver L. Haimson, Gabriela Cruz Erickson, Michael Turken, and Rebecca Gulotta. 2023. Infrastructuring Care: How Trans and Non-Binary People Meet Health and Well-Being Needs through Technology. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. ACM, Hamburg Germany, 1–17. <https://doi.org/10.1145/3544548.3581040>
- [112] Yu Wu, Wei Wu, Chen Xing, Ming Zhou, and Zhoujun Li. 2017. Sequential Matching Network: A New Architecture for Multi-turn Response Selection in Retrieval-Based Chatbots. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association

for Computational Linguistics, Vancouver, Canada, 496–505. <https://doi.org/10.18653/v1/P17-1046>

- [113] Chloe Xiang. 2023. Man Dies by Suicide After Talking with AI Chatbot, Widow Says. <https://www.vice.com/en/article/pkadgm/man-dies-by-suicide-after-talking-with-ai-chatbot-widow-says> Accessed: 2023-12-11.
- [114] Ziang Xiao, Michelle X. Zhou, Wenxi Chen, Huahai Yang, and Changyan Chi. 2020. If I Hear You Correctly: Building and Evaluating Interview Chatbots with Active Listening Skills. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (*CHI '20*). Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3313831.3376131>
- [115] Hubert D. Zając, Dana Li, Xiang Dai, Jonathan F. Carlsen, Finn Kensing, and Tariq O. Andersen. 2023. Clinician-Facing AI in the Wild: Taking Stock of the Sociotechnical Challenges and Opportunities for HCL. *ACM Transactions on Computer-Human Interaction* 30, 2 (April 2023), 1–39. <https://doi.org/10.1145/3582430>
- [116] Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2020. Fine-Tuning Language Models from Human Preferences. <http://arxiv.org/abs/1909.08593> [cs, stat].

A APPENDIX: SURVEY

- (1) In what country did you live most of your childhood?
- (2) In what country have you spent most of the past five years?
- (3) Age
- (4) Gender
- (5) Sexuality
- (6) **Have you used an LLM-based chatbots for mental wellness support (such as Snapchat's AI friend, Replika, Character.ai) before?**
 - Yes
 - No
- (7) **If yes, please specify which app(s) you have used.**
 - Replika
 - Snapchat My AI
 - Chai
 - Character.ai
 - Anima
 - Paradot
 - ChatGPT
 - Kuki
 - Other: _____
- (8) **How long have you been using these apps?**
 - Less than 1 week
 - 1 week to 1 month
 - 1-3 months
 - 3-6 months
 - 6-12 months
 - 1-2 years
 - Other: _____
- (9) **How often do you use these apps?**
 - Daily
 - Weekly
 - Monthly
 - Rarely
 - Other: _____
- (10) **I consent to be contacted for an interview study by providing my contact information.**
My contact information:

B APPENDIX: INTERVIEW GUIDELINE

Begin the interview by explaining the purpose of the study and obtaining informed consent from the participant. Create a comfortable and non-judgmental atmosphere for the participant to share their experiences. Use open-ended questions and follow-up probes to encourage the participant to elaborate on their thoughts as some of the questions above. Maintain a neutral stance and avoid leading questions that may influence the participant's responses.

B.1 Questions

- What AI chatbots do you use?
- Do you identify as part of the LGBTQ communities?
- Can you please share your experience using LLMs for mental wellness and social support related to your LGBTQ+ or trans identity? (**Only asked for LGBTQ+ participants**)
- Can you please share your experience using LLMs for mental wellness and social support? (**Only asked for non-LGBTQ+ participants**)
- What led you to seek support from an LLM in the first place? (motivations)
- How would you describe the overall quality of support and resources provided by the LLM?
- Can you share any specific instances where the LLM was particularly helpful or unhelpful?
- Could you walk me through the instance when you found LLM to be a beneficial resource for mental wellness or social support?
- How did using an LLM for support compare to other resources, such as support groups or mental health professionals / family or friends/ online communities?
- Was there a specific event or reason that made it stand out among these choices?
- Could you please share a time when the LLM's responses surprised you - either positively or negatively - in terms of support?
- Can you recall a situation where you felt that the LLM really understood your experiences as a (the LGBTQ+ identity that the participant identifies as) adult? Or perhaps a time when it fell short? (**Only asked for LGBTQ+ participants**)
- Did you feel that the LLM adequately understood your unique experiences as an (vary according to the person's identity: gay, lesbian, trans, etc.) person? (**Only asked for LGBTQ+ participants**)
- How did the chatbot understand you? Give an example?
- Did you feel that the LLM adequately addressed your problems as an LGBTQ+ or trans young adult? (**Only asked for LGBTQ+ participants**)
- Were there any privacy or safety concerns while using the LLM for support?
- What improvements or features would you like to see in LLMs to better serve your experience?

Received 14 September 2023; revised 12 December 2023; accepted 19 January 2024