
UX RESEARCH ON CONVERSATIONAL HUMAN-AI INTERACTION: A LITERATURE REVIEW OF THE ACM DIGITAL LIBRARY

A PREPRINT

Qingxiao Zheng

School of Information Sciences
University of Illinois at Urbana-Champaign
qzheng14@illinois.edu

Yiliu Tang

School of Informatics
University of Illinois at Urbana-Champaign
yiliut2@illinois.edu

Yiren Liu

School of Informatics
University of Illinois at Urbana-Champaign
yiren12@illinois.edu

Weizi Liu

College of Media
University of Illinois at Urbana-Champaign
weizil2@illinois.edu

Yun Huang

School of Information Sciences
University of Illinois at Urbana-Champaign
yunhuang@illinois.edu

-

ABSTRACT

Early conversational agents (CAs) focused on *dyadic* human-AI interaction between humans and the CAs, followed by the increasing popularity of *polyadic* human-AI interaction, in which CAs are designed to mediate human-human interactions. CAs for *polyadic* interactions are unique because they encompass hybrid social interactions, i.e., human-CA, human-to-human, and human-to-group behaviors. However, research on *polyadic* CAs is scattered across different fields, making it challenging to identify, compare, and accumulate existing knowledge. To promote the future design of CA systems, we conducted a literature review of ACM publications and identified a set of works that conducted UX (user experience) research. We qualitatively synthesized the effects of *polyadic* CAs into four aspects of human-human interactions, i.e., communication, engagement, connection, and relationship maintenance. Through a mixed-method analysis of the selected *polyadic* and *dyadic* CA studies, we developed a suite of evaluation measurements on the effects. Our findings show that designing with social boundaries, such as privacy, disclosure, and identification, is crucial for ethical *polyadic* CAs. Future research should also advance usability testing methods and trust-building guidelines for conversational AI.

Keywords Conversational Agent · Chatbot · Conversational AI · UX Research · Literature Review

1 Introduction

There is a rapidly growing body of literature on conversational agents or chatbots [Adamopoulou and Moussiades, 2020a]. As promising Artificial intelligence (AI) technologies, conversational agents are defined as "software that accepts natural language as input and generates natural language as output, engaging in a conversation with the user" [Griol et al., 2013]; chatbots, meanwhile, are computer programs designed to simulate conversation with human users via text [Adamopoulou and Moussiades, 2020a,b]. As these two terms are often perceived as interchangeable [Rapp

et al., 2021, McTear, 2020], in the remainder of this paper, we refer to both conversational agents and chatbots as CAs. Scholars have shown that these machines are able to compensate for human shortcomings or exceed human capacities [Fox and Gambino, 2021, Guzman and Lewis, 2020, Whittaker et al., 2018]. However, prior works focus on designing and evaluating *dyadic* human-AI interaction, which involve only one-to-one interactions between humans and their CAs [Bickmore et al., 2005, Schulman and Bickmore, 2009, Xu et al., 2017, Kopp et al., 2005, Anabuki et al., 2000]; whereas more recent works start tapping into *polyadic* human-AI interactions that also support human-human interactions [Kim et al., 2021, Wang et al., 2021, Kim et al., 2020, Toxtli et al., 2018, Benke et al., 2020].

Even though there are extensive literature reviews on CAs, e.g., [Seering et al., 2019, Chaves and Gerosa, 2021, de Barcelos Silva et al., 2020, Montenegro et al., 2019, Laranjo et al., 2018], they do not address how *polyadic* CAs are designed and evaluated, nor present the effects of using *polyadic* CAs on handling the challenges of human-human interaction [Hohenstein and Jung, 2020]. In this paper, we overview UX (user experience) research on *polyadic* CAs that 1) interact with more than one user in the same conversation and 2) engage in bidirectional conversations between all parties (human-AI and human-human). These *polyadic* CAs encompass a wide variety of complexities that *dyadic* Human-AI may not encounter, including multi-party interactions, social roles taking, group hierarchy, or social tension [Van Dijk, 1997]. To evaluate the effects of CAs’ support in human-human interactions, we need to examine both human-CA behaviors and human-to-human behaviors, as well as human-to-group behaviors potentially. Given the unique challenges of the design space, little is known about how *polyadic* CAs should be designed to address these different aspects of social interaction and how they can positively influence the overall user experience [Seering et al., 2019, Hohenstein and Jung, 2020].

To fill the void, we conducted a literature review using a mixed-method approach, mapping out what exists in *polyadic* CAs research. Specifically, we screened 1,302 ACM papers and identified 36 papers that designed *polyadic* CAs and 135 that designed *dyadic* CAs, which included user evaluation results. We investigated what fundamental human-human interaction challenges are addressed by these works, what effects on human-human interaction are evaluated, and what issues are overlooked when designing these CAs.

The contributions of this work to HCI and social computing are manifold. First, we summarize the fundamental challenges (i.e., consensus reaching, uneven participation, lack of emotional awareness, etc.) of human-human interaction that *polyadic* CAs are designed to tackle and their promising results. Second, we synthesize the current practices (e.g., design originality, relationship types, social scale, evaluation method, etc.) and areas that fall short of empirical studies. Third, we point out the issues that are overlooked by the researchers and practitioners in these works and propose to design CAs with *boundary awareness* for supporting human-human interaction via conversational AI. Fourth, we present the differences in research interests and evaluation metrics between *dyadic* and *polyadic* CA studies and identify existing gaps and potential directions. Further, we envision several research opportunities on conversational human-AI interaction concerning the theoretical groundings, relational dynamics, functional dimensions, and metaphysical implications.

2 Related Work

In this section, we first review briefly the history and the research landscape of CAs (Conversational Agents or Chatbots), and then point out the gaps in the existing literature.

2.1 Landscape of CAs

The current landscape of CAs is complex. The history of conversational interfaces could be traced back to the 1960s when text-based dialogue systems were designed to answer questions and for simulated casual conversation [McTear et al., 2016]. Since then, various terms have been used, and terminologies have become overlapping [Rapp et al., 2021]. Historically, several distinctions have emerged in research around the terms, such as text-based and spoken dialogue systems, voice user interfaces, chatbots, embodied conversational agents, and social robots and situated agents [McTear et al., 2016]. In recent years, with the commercialization of CAs, more terms emerged. For example, the website chatbots.org references over 1,300 chatbots across nine categories [Casas et al., 2020] and provides 161 conversational AI synonyms [McTear, 2020].

Similarly to the terminologies, the typologies of CAs are also not univocally categorized [Rapp et al., 2021]. In the past decade, researchers classified CAs based on multiple criteria, such as interaction modality (text-based, voice-based, mixed) [Hussain et al., 2019, Bittner et al., 2019, Mygland et al., 2021, Laranjo et al., 2018]; scale of social engagement (one-to-one, broadcasting, and community-based) [Seering et al., 2019]; knowledge domain (open domain or closed domain) [Hussain et al., 2019, Bittner et al., 2019]; goals (task-oriented or non-task oriented) [Hussain et al., 2019, Gao et al., 2018]; duration and locus of control [Følstad et al., 2018]; embodiment [Cassell, 2000]; design approach

(rule-based, retrieval based, and generative based) [Hussain et al., 2019]; platform (mobile, laptop, web browser, SMS, cells, multimodal platforms) [Klopfenstein et al., 2017, Laranjo et al., 2018], and application domains [de Barcelos Silva et al., 2020, Ruan et al., 2020].

In this paper, we include the aforementioned terminologies as our study objects, given that the CAs’ major functionalities are conversational-based. For example, embodied conversational agents (ECAs) are agents simulating humans’ face-to-face conversations and can use their body language when talking to users [Cassell, 2000]. Four properties distinguish ECAs from CAs: ECAs can recognize and respond to both verbal and nonverbal user inputs; deliver both verbal and nonverbal outputs; deal with conversational functions; and give signals that help the conversations [Cassell, 2000]. This work includes the ECA papers that conduct UX research on conversational interaction.

2.2 Gaps in Existing Literature Reviews of CAs

Most of the existing empirical studies focus on designing *dyadic* CAs [Skjuve et al., 2019, Adam et al., 2020, Cranshaw et al., 2017, Lee et al., 2019, Santos et al., 2020], without addressing the unique design challenges of *polyadic* CAs. Thus, existing literature reviews address CAs’ properties mainly in the context of *dyadic* conversationalists and the effect of interaction between the Human and CAs [Seering et al., 2019]. They look at how CAs can increase user engagement, enhance user experience, or enrich the relationship between humans and CAs. For example, one work mapped relevant themes in text-based CAs to understand user experiences, perceptions, and acceptances towards CAs [Rapp et al., 2021]; and one categorized the characteristics of human-CA interactions into conversational, social, and personification characteristics [Chaves and Gerosa, 2020]. Others are more specific, such as studying CAs’ emotional intelligence [Pamungkas, 2019], personalization [Kocaballi et al., 2019], trust-building [Rheu et al., 2021], and human-likeness [Van Pinxteren et al., 2020].

Recently, an increasing number of CAs are designed to support *polyadic* human-AI interaction, where human-human communication is supported by CAs. While *dyadic* CAs are commonly used to communicate with individuals as personal assistants [Sciuto et al., 2018, Modi et al., 2004], customer service agents [Lundkvist and Yakhlef, 2004, Shimazu, 2002], and personal healthcare partners [Bickmore et al., 2010, 2009, Grolleman et al., 2006], some emotional intelligent CAs are designed to help team members gain awareness of other group members’ emotional changes, report the overall sentiment of each group discussion, and maintain positive emotions during collaborations [Benke et al., 2020, Peng et al., 2019]. Such *polyadic* CAs address unique social challenges that the *dyadic* CAs were not designed to meet, e.g., in the context of a group or a community. However, little is known about the empirical use and the effect of *polyadic* CAs on different scenarios of human-human interaction [Hohenstein and Jung, 2020, Seering et al., 2019, Hohenstein and Jung, 2018, Chynał et al., 2018].

Also, researchers and practitioners have increasingly noted a lack of understanding in achieving quality UX for CAs. Although the popularity of CA applications exists in multiple domains, studies repeatedly reported CAs causing both pragmatic issues, in which chatbots failed to understand or to help users achieve their intended goals [Følstad et al., 2021, Følstad and Brandtzæg, 2017, Lee et al., 2021], and issues in which CAs failed to engage users over time [Følstad and Brandtzæg, 2017, Lee et al., 2021]. Thus, we set the survey scope to include papers that conducted user evaluations.

Our review of the CAs designed for *polyadic* Human-AI interaction can help researchers and practitioners identify common practices, research gaps, lessons learned, and promising areas for future advancements. Specifically, we are interested in addressing the following research questions:

RQ1: *What fundamental challenges of human-human interaction are addressed by these CAs?*

RQ2: *What are the research interests in dyadic and polyadic CAs?*

RQ3: *What are the practices of designing the CAs for polyadic human-AI interaction?*

RQ4: *What are the effects of using the proposed CAs on human-human interactions?*

RQ5: *What evaluation metrics have been used in understanding user experience?*

RQ6: *What issues are overlooked by scholars when designing polyadic CAs?*

3 Method

We used a literature review method that has been widely applied by prior works in HCI, e.g., [Rapp et al., 2021, Nunes et al., 2015, Mencarini et al., 2019], which has four stages: 1) Define: proposing the inclusion and exclusion criteria and identifying the appropriate data sources; 2) Search: developing specific query and collecting the papers through the data source; 3) Select: checking the search results against the inclusion and exclusion criteria and identifying the final

papers for both *dyadic* and *polyadic* works; and 4) Analyze: examining the selected papers by applying a mixed-method approach. We illustrate the four steps in Figure 1 and present the details of each stage below.

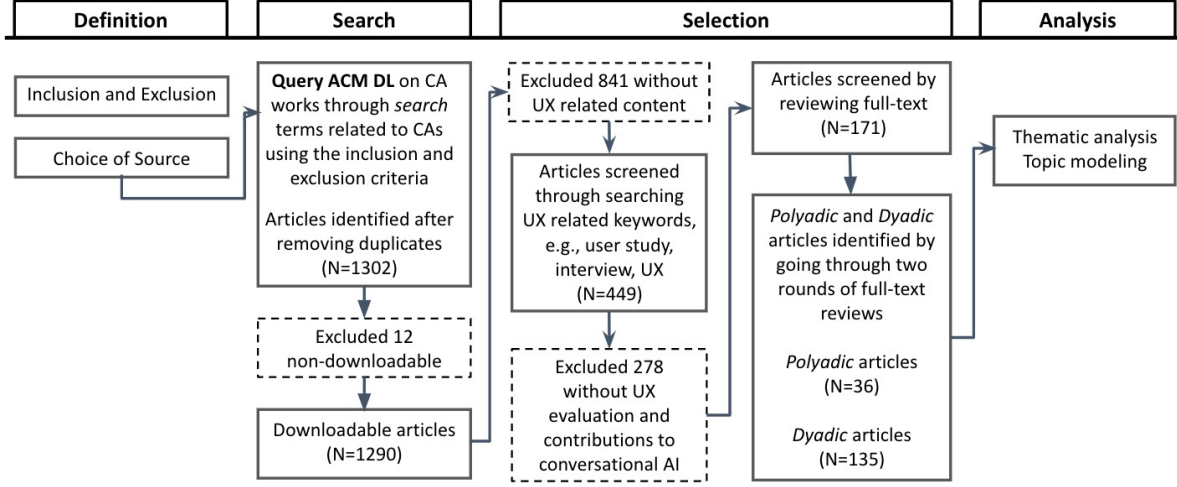


Figure 1: A Workflow Diagram of the Literature Review Process

3.1 Definition

In this section, we present how we defined the selection criteria and how we identified the appropriate data source.

3.1.1 Inclusion and Exclusion

A series of criteria were developed by researchers through multiple rounds of discussions. Criteria were selected to include works that were the most representative of the scope of our study (i.e., UX research on conversational human-AI interaction) and to filter irrelevant works. We employed the following inclusion and exclusion criteria.

For *dyadic* papers, we defined the inclusion criteria as follows: 1) selected articles need to study CAs that interact with only one human user in a session; 2) the CA interaction in the article is bidirectional; 3) users of the CA are aware of the existence of the CA; 4) the articles are research papers with user studies; 5) the major design feature is conversation-based, e.g., excluded sensory-based CA [Bohus and Horvitz, 2011]; 6) the articles are written in English; 7) the articles are included in the selected database. Articles that only assessed CA task performance without meaningfully exploring the interaction experience of Human-CA as a conversational technology were excluded.

For *polyadic* papers, the inclusion criteria shared 2)-7) requirements of selecting *dyadic* papers', except that the selected articles need to study CAs that interact with greater than one human user, instead of with only one human user. Exclusion criteria were defined as: 1) articles that only assessed CA task performance without meaningfully exploring the interaction experience of Human-CA or Human-Human as a conversational technology; and 2) the CAs in the articles interacted with multiple users, but there were no human-human interactions.

3.1.2 Source

To identify the data source, we first randomly retrieved 200 papers from five databases: ACM Digital Library, IEEE, Web of Science, Scopus, and Science direct. After exploring the initial search results, we chose Association for Computing Machinery Digital Library (ACM DL) as the final source for our literature review. Two co-authors reviewed 40 papers from each source, using the inclusion and exclusion criteria among all five databases. The qualification rates were ACM DL (23.5%), IEEE (12.5%), Web of Science (14.6%), Scopus (17%), and Science direct (4.9%). Specifically, 52% IEEE papers were technical papers without user evaluations; 59% Web of Science papers contained no CA designs. Due to the low qualification rates of other sources, we ultimately decided to choose ACM DL as the data source for this literature review. This decision was also made by considering that the ACM DL features a wide selection of reliable HCI works, and has been used as a solo source for literature review works, e.g., [Ralph and Robbes, 2020, Wainer et al., 2009].

3.2 Search

The search query is composed of two parts. The first part of the search covers synonyms of conversational agents, and the second part specifies terms that may be relevant for identifying the *polyadic* papers. The list of search terms was inspired by several CA research reviews [Laranjo et al., 2018, de Barcelos Silva et al., 2020, Rapp et al., 2021, ter Stal et al., 2020] and was developed through several iterations and refinement by the research team, in line with prior HCI work.

More specifically, the search query included 15 terms, which are “conversational agent”, “conversational AI”, “intelligent assistant”, “intelligent agent”, “chatbot”, “chatterbot”, “chatterbox”, “socialbot”, “digital assistant”, “conversational UI”, “conversational interface”, “conversation system”, “conversational system”, “dialogue system”, and “dialog system”. To explore *polyadic* works from the retrieved results, we searched 13 terms, “human-human,” “human human,” “multi-user,” “multi-users,” “multi user,” “multi users,” “multi-party,” “multi-parties,” “multi party,” “multiparty-based”, “multi parties”, “multi model”, and “multi-model.”

3.2.1 Data Preparation

We searched the papers through two steps, as illustrated in Figure 1. First, we searched on Mar 3rd, 2021 by using the query and connectors throughout the ACM DL, which resulted in 1,302 CA papers after removing duplicates. Second, we web-crawled the metadata and PDFs of these papers and built our database. To improve the stability and efficiency of the crawling process, we also saved the paper DOIs into a file and then crawled the metadata and the PDFs of the respective DOIs in sequence.

3.2.2 Database

The database consisted of two parts. The first part was a CSV file that comprehensively covers the metadata of the papers. It included the ten columns: paper DOI, title, authors, abstract, publication date, source, publisher, citations, keywords, affiliation of authors. The second part included the PDFs of the papers, as well as their corresponding text files. We used the Fitz module within the PyMuPDF project [Jorj X. McKie, 2016] to transform the PDF to text, as the tool has been applied in a number of works [Park and Pouchard, 2019, Yang et al., 2019, Park et al., 2020, Chang et al., 2021].

3.3 Selection

Once having the database, we screened the paper records to check the *dyadic* and *polyadic* CA works that meet our criteria through two steps, as illustrated in Figure 1. We looked for papers that potentially had user studies by searching through the full texts of all papers with the keywords “user study”, “user studies”, “interview”, “interviews” and “user experience”, which led to 449 papers. To initiate the set of *polyadic* papers, we also searched through all the papers using the 13 key terms, e.g., “human-human,” “human human,” “multi-user,” which resulted in an initial set of *polyadic* with 27 papers. We then removed duplicates between the two screening results. There were six (1.3%) duplicates between the two stages, resulting in 470 unique full texts. We then further reviewed the full texts against the identified inclusion and exclusion criteria for *polyadic* and *dyadic* papers. If ambiguities persisted about the eligibility of a specific paper, input was sought from the third co-author. The final study collected 36 (21.1%) *polyadic* papers and 135 (78.9%) *dyadic* papers.

3.4 Analysis

To examine the research landscape of UX research on conversational human-AI interaction and the differences between *polyadic* and *dyadic* works published in ACM DL, we analyzed the selected papers by applying a mixed-method approach.

3.4.1 Thematic Analysis

Prior literature review works [Rapp et al., 2021, Hussain et al., 2019, Bittner et al., 2019, Laranjo et al., 2018] have proposed several aspects when examining conversational AI research, e.g., regarding interaction modality [Hussain et al., 2019], characteristics of embodiment [Bittner et al., 2019], application domains [Rapp et al., 2021], evaluation methods [Rapp et al., 2021], and social scale [Bittner et al., 2019]. We found that these could be leveraged to code the research practices of *polyadic* CAs. However, the prior schemes were not sufficient for us to categorize the unique aspects of human-human interactions supported by *polyadic* CAs, e.g., fundamental challenges addressed, proven effects, and remaining issues. Thus, the authors adopted the grounded theory approach [Wolfswinkel et al., 2013] and

coded the attributes of the works using thematic analysis [Braun and Clarke, 2006]. Two co-authors started reviewing the papers and did open coding of all papers independently, and then discussed and compiled their codes together. These codes were added as new attributes to address the proposed RQs. After reaching an agreement regarding the codes, two co-authors coded the rest papers separately, compared their codes, and discussed possible revisions [McDonald et al., 2019, Pina et al., 2018]. A final inter-rater reliability of 91% was achieved and deemed satisfactory [Gwet, 2014].

3.4.2 Topic Modeling

Besides coding the research challenges and issues that are overlooked by scholars in designing *polyadic* CAs, we also applied topic modeling to explore the differences in discussion topics between the two groups of papers. The topic modeling method has been widely applied in prior works [Gurcan et al., 2021, Gurcan and Cagiltay, 2020, Sun and Yin, 2017]. This method can be employed as follows. We started text pre-processing by removing punctuation marks, symbols, and stop-words. We then tokenized and lemmatized the text for regularization. We used lemmatization to group words with similar semantic meanings but different syntactic forms (such as plurality and tense). We experimented with both lemmatization and no lemmatization, and the results with lemmatization exhibited higher interpretability and were more coherent when we sampled the papers to explain the results. We then used non-negative matrix factorization (NMF) to conduct topic modeling based on the TF-IDF scores. We built two topic models separately for *dyadic* and *polyadic* works to explore the major topics discussed in either set of the papers. This provided closer insights into the frequent topics discussed by *dyadic* and *polyadic* CA papers. All text analysis steps were completed using Python (NLTK and Scikit-learn libraries).

4 Findings

In this section, we present the findings of the ACM papers, 36 *polyadic* and 135 *dyadic* works, in order of the proposed RQs. Specifically, we show the fundamental challenges (RQ1 in Section 4.1), research interests (RQ2 in Section 4.2), practices (RQ3 in Section 4.3), proven effects (RQ4 in Section 4.4), evaluation metrics (RQ5 in Section 4.5), and issues (RQ6 in Section 4.6).

4.1 Addressing Fundamental Challenges of Human-Human Interaction (RQ1)

To answer our first research question, we identified major challenges in human-human interaction addressed in *polyadic* papers. We found that only some papers explicitly stated these challenges as pain points that motivated their designs. Most *challenges* addressed by *polyadic* CAs are in the collaborative learning, work, or group discussion contexts.

4.1.1 Inefficient Communication

As mentioned in prior research regarding group work and collaborations, several issues in human-human communication affect the efficiency. For example, in online group settings, team communications are often unstructured and less organized, with procrastination and distractions. Reaching consensus can be time-consuming for deliberative discussions [Kim et al., 2021, 2020], and thoughts and discussions can be hard to organize and make sense of [Zhang and Cranshaw, 2018] and challenging to control in terms of conversational flow [Bohus and Horvitz, 2009a,b, 2010]. In terms of workflow, task management such as refining, assigning, and tracking can be difficult [Toxtli et al., 2018]; heavy workloads related to coordinating schedules among multiple parties can be tedious [Cranshaw et al., 2017]; and switching between tools and platforms during collaborations can be burdensome [Avula et al., 2018]. Therefore, *polyadic* CAs are used to provide workflow support and discussion support to improve team communication efficiency.

4.1.2 Lack of Engagement

Inactive engagement is a common issue in multi-user interactions in collaborative teams. For example, in peer learning settings, it is challenging to engage students in provocative discussions [Dohsaka et al., 2009], to promote productive talk among students, such as explaining to peers and re-voicing others' statements [Dyke et al., 2013, Tegos et al., 2015], and to provide effective learning supports to others, e.g., by exhibiting an agreeable attitude and precipitating tension release [Ai et al., 2010, Chaudhuri et al., 2009, Kumar et al., 2010, Kumar and Rosé, 2014]. This is also true for collaborative work and online community contexts, in which more inputs [Savage et al., 2016] and greater community activity are needed [Seering et al., 2018, 2020, Luria et al., 2020a, Abokhodair et al., 2015, Wang et al., 2021]. Further issues of engagement in these contexts include uneven participation, in which a balanced amount of contribution is hard to achieve [Kim et al., 2020, Shamekhi et al., 2018], and user attractions are challenging to capture in casual social settings [Otogai et al., 2013, 2014, Zheng et al., 2005, Yuan and Chee, 2005].

4.1.3 Barriers in Relational Maintenance

One major challenge in multi-user social settings is maintaining positive relationships, which is crucial to forming a solid team or group. For example, in online collaborative work, there could be a lack of emotional awareness and mutual understanding between team members, as it is hard for them to detect and regulate emotions [Benke et al., 2020, Peng et al., 2019, Narain et al., 2020]. Moreover, it is always important to grow trust within a team, and the feasibility of using CAs for trust-building [Strohkorb Sebo et al., 2018] and setting privacy boundaries [Luria et al., 2020b] are explored and developed.

4.1.4 Need for Building Connections

In human-human interactions, there is a need for building social connections and identifying common grounds. This is particularly important in the "ice-breaking" and transitioning stages [Löw and Moshuber, 2020]. Also, people may seek similarities and agreement with their communication partners [Nakanishi et al., 2003, Isbister et al., 2000] or need idea supports during chats [Hohenstein and Jung, 2018]. This challenge could be more difficult in cross-cultural conversations. When people from different backgrounds meet for the first time and try to get acquainted with each other, it can be challenging to form impressions of each other free from the influence of cultural stereotypes [Isbister et al., 2000].

4.2 Research Interests in *Dyadic* and *Polyadic* CAs (RQ2)

4.2.1 Research Interest Over Time

The ACM papers we collected included UX research on conversational human-AI interaction starting from early 2000, e.g., *polyadic* CAs by Isbister et al. [Isbister et al., 2000] and *dyadic* CAs by Chai et al. [Chai et al., 2001]. The works collected were not evenly distributed over the years, e.g., about four UX papers in 2010 and declining to one paper in 2012 for both. However, starting from 2018, there was a surge in the number of selected articles, with the number of *dyadic* papers reaching a peak of 48 in 2020 and the number of *polyadic* papers reaching a peak of nine in the same year.

4.2.2 Authors' Affiliations

There was a total of 596 authors for the 135 *dyadic* papers, averaging 4.4 authors per paper, and 379 (63.6%) out of the 596 authors were from academia. In comparison, there were a total of 145 authors for the 36 *polyadic* papers, averaging four authors per paper; and 121 (83.4%) out of the 145 authors were from academia.

4.2.3 Impactful Works in the Area

As shown in Table 1, our database collected a good body of *dyadic* papers. For example, Medhi et al. [Medhi et al., 2011] evaluated the usability of a task-oriented CA for novice and low-literacy users and was highly cited, followed by several other usability studies of CAs [Jain et al., 2018a,b, Lau et al., 2010] on improving health [Lisetti et al., 2013, Bickmore et al., 2005], answering questions [Liao et al., 2018, Lovato et al., 2019], and counseling [Schulman and Bickmore, 2009, Smyth et al., 2010] via CAs. Similarly, a sample of *polyadic* papers received high citations, as shown in Table 2. For example, Isbister et al. [Isbister et al., 2000] designed a virtual agent to support human-human communication in virtual environments. Unlike the *dyadic* works, many selected *polyadic* works conducted UX research in the context of group environments, e.g., virtual meetings [Isbister et al., 2000, Nakanishi et al., 2003], games [Bohus and Horvitz, 2010, 2009a], online communities [Savage et al., 2016], and group collaboration [Zhang and Cranshaw, 2018, Avula et al., 2018]. Our dataset included both *dyadic* and *polyadic* on improving work productivity [Kim et al., 2019, Vtyurina et al., 2017, Kocielnik et al., 2018, Shamekhi et al., 2018, Toxtli et al., 2018, Cranshaw et al., 2017] and promoting education [Tanaka et al., 2015, Kumar et al., 2010, Ai et al., 2010, Dyke et al., 2013]. The average number of citations of a *polyadic* paper in our database was 19 (SD=23) from ACM and 44 (SD=46) from Google Scholar. Five (13.9%) out of the 36 *polyadic* papers had no citations yet. The average number of citations of a *dyadic* paper in our database was 10 (SD=14) from ACM and 19 (SD=22) from Google Scholar, on December 21st, 2021. The dataset included works demonstrating impacts at different levels.

4.3 Practices of Designing CAs for *Polyadic* Human-AI Interaction (RQ3)

In the following section, we present aspects that are often shared by *polyadic* and *dyadic* CA works and aspects that are unique to *polyadic* CA works.

Table 1: A Sample of *Dyadic* ACM Papers that Conducted UX Research.

Year	Authors	ACM Citations	Google Scholar Citations
2011	Medhi, I., Patnaik, S., Brunskill, E., Gautama, S. N. N., Thies, W., & Toyama, K. [Medhi et al., 2011]	166	275
2018	Jain, M., Kumar, P., Kota, R., & Patel, S. N. [Jain et al., 2018a]	75	187
2013	Lisetti, C., Amini, R., Yasavur, U., & Rishe, N. [Lisetti et al., 2013]	85	168
2005	Bickmore, T. W., Caruso, L., & Clough-Gorr, K. [Bickmore et al., 2005]	62	143
2009	Schulman, D., & Bickmore, T. [Schulman and Bickmore, 2009]	50	99
2018	Liao, Q. V., Hussain, M. M., Chandar, P., Davis, M., Khazaeni, Y., Crasso, M. P., Wang, D. K., Muller, M., Shami, N. S., & Geyer, W. [Liao et al., 2018]	37	48
2019	Kim, S., Lee, J., & Gweon, G. [Kim et al., 2019]	35	84
2017	Vtyurina, A., Savenkov, D., Agichtein, E., & Clarke, C. L. A. [Vtyurina et al., 2017]	30	81
2019	Lovato, S. B., Piper, A. M., & Wartella, E. A. [Lovato et al., 2019]	32	65
2018	Kocielnik, R., Avrahami, D., Marlow, J., Lu, D., & Hsieh, G. [Kocielnik et al., 2018]	47	46
2010	Smyth, T. N., Etherton, J., & Best, M. L. [Smyth et al., 2010]	37	86
2015	Tanaka, H., Sakti, S., Neubig, G., Toda, T., Negoro, H., Iwasaka, H. & Nakamura, S. [Tanaka et al., 2015]	31	48
2019	Zhou, M. X., Mark, G., Li, J., & Yang, H. [Zhou et al., 2019]	30	54
2018	Jain, M., Kota, R., Kumar, P., & Patel, S. N. [Jain et al., 2018b]	20	54
2010	Lau, T., Cerruti, J., Manzato, G., Bengualid, M., Bigham, J. P., & Nichols, J. [Lau et al., 2010]	25	44
2005	Marti, S., & Schmandt, C. [Marti and Schmandt, 2005]	24	40

Table 2: A Sample of *Polyadic* ACM Papers that Conducted UX Research.

Year	Authors	ACM Citations	Google Scholar Citations
2000	Isbister, K., Nakanishi, H., Ishida, T., & Nass, C. [Isbister et al., 2000]	88	207
2010	Bohus, D., & Horvitz, E. [Bohus and Horvitz, 2010]	81	149
2016	Savage, S., Monroy-Hernandez, A., & Höllerer, T. [Savage et al., 2016]	56	113
2009	Bohus, D., & Horvitz, E. [Bohus and Horvitz, 2009b]	55	96
2009	Bohus, D., & Horvitz, E. [Bohus and Horvitz, 2009a]	47	95
2018	Sebo, S. S., Traeger, M., Jung, M., & Scassellati, B. [Strohkorb Sebo et al., 2018]	43	76
2018	Shamekhi, A., Liao, Q. V., Wang D., Bellamy, R. K. E., & Erickson, T. [Shamekhi et al., 2018]	46	73
2018	Zhang, A. X., & Cranshaw, J. [Zhang and Cranshaw, 2018]	39	61
2010	Kumar, R., Ai, H., Beuth, J. L., & Rosé, C. P. [Kumar et al., 2010]	27	68
2018	Toxtli, C., Monroy-Hernández, A., & Cranshaw, J. [Toxtli et al., 2018]	35	59
2017	Cranshaw, J., Elwany, E., Newman, T., Kocielnik, R., Yu, B. W., Soni, S., Teevan, J., & Monroy-Hernández A. [Cranshaw et al., 2017]	21	67
2003	Nakanishi, H., Nakazawa, S., Ishida, T., Takanashi, K., & Isbister, K. [Nakanishi et al., 2003]	9	65
2018	Avula, S., Chadwick, G., Arguello, J., & Capra, R. [Avula et al., 2018]	20	42
2010	Ai, H., Kumar, R., Nguyen, D., Nagasunder, A., & Rosé, C. P. [Ai et al., 2010]	3	56
2013	Dyke, G., Howley, I., Adamson, D., Kumar, R., & Rosé, C. P. [Dyke et al., 2013]	1	52
2018	Seering, J., Flores, J. P., Savage, S., & Hammer, J. [Seering et al., 2018]	23	29

4.3.1 Shared Aspects between *Polyadic* and *Dyadic* CAs

We present aspects that *polyadic* and *dyadic* CAs shared, i.e., application domain, modality, agent characteristics, design originality, design method, and evaluation methods below.

Application domain. Among the 36 ACM papers we reviewed, eight *polyadic* CAs have been applied to education and collaborative learning, e.g., [Dyke et al., 2013]; five to online communities, e.g., [Savage et al., 2016]; three to group discussions, e.g., [Kim et al., 2021]; seven to work and productivity, e.g., [Toxtli et al., 2018, Avula et al., 2018]), two to virtual meetings, e.g., [Isbister et al., 2000]; three to guiding services, e.g., [Zheng et al., 2005]; two to games, e.g., [Rehm, 2008], one to family, i.e., [Luria et al., 2020b], and five undefined depending on the major focuses of the papers.

Modality. There are 22 polyadic papers applying text-only interactions, e.g., [Peng et al., 2019]; nine studies are based on video, e.g., [Nakanishi et al., 2003], only a few are audio only, e.g., [Luria et al., 2020b], and hybrid of audio-text, e.g., [Isbister et al., 2000]. These categories are not mutually exclusive, because some studies implemented multiple designs with different modalities, e.g., [Shamekhi et al., 2018].

Agent characteristics. We also looked at whether these CAs in the reviewed papers are embodied or not. While 22 papers studied CAs that are not embodied, e.g., [Chaudhuri et al., 2009], the rest of them are embodied CAs, e.g., [Isbister et al., 2000], which have figures that can demonstrate certain social characteristics through animated body languages [Cassell, 2000].

Design originality. Most studies present and evaluate original designs, e.g., [Chaudhuri et al., 2009] while three test or adopt existing systems. One paper uses existing systems (i.e., Google Allo) [Hohenstein and Jung, 2018] and two papers examine existing bots on the Twitter platform [Seering et al., 2018, Abokhodair et al., 2015].

Design methods. Among those papers which specified their design methods, five of them are framework-based, which means that the designed features are proof-of-concept designs guided by prior frameworks, models, or theories [Dyke et al., 2013, Tegos et al., 2015, Kim et al., 2021, Wang et al., 2021, Nakanishi et al., 2003]. Meanwhile, most of the designed features are proposed by researchers without leveraging prior designs or catering to specific user-needs,

e.g., [Toxtli et al., 2018, Dohsaka et al., 2009, Seering et al., 2020]. Two CAs adopted participatory design methods, including need-finding interviews, e.g., [Kim et al., 2020, Zhang and Cranshaw, 2018] and two ran ideation workshops, e.g., [Benke et al., 2020, Luria et al., 2020b].

Evaluation methods. The most common evaluation method in the reviewed studies is experimental, with 16 papers using it. Other methods include six using surveys, six using interviews, five using field studies, three using Wizard of Oz, two using interaction log analysis, two using observation and one using focus groups. These categories are not mutually exclusive since some studies employed multiple methods in their evaluations.

4.3.2 Unique Aspects with the *Polyadic* CAs

Below, we present aspects that involve multi-user interaction, unique to *polyadic* CAs, i.e., relationship type and social scale. Moreover, we report related theories and frameworks used in prior *polyadic* works.

Relationship types. The specific relationships highly depend on the contexts of the interactions, which include eight papers designing *polyadic* CAs as co-learners, e.g., [Chaudhuri et al., 2009], four as co-workers, e.g., [Cranshaw et al., 2017], seven as collaborators/discussants, e.g., collaborators on Slack [Avula et al., 2018], three as people meeting for the first time, e.g., [Isbister et al., 2000], four as online community members, e.g., [Savage et al., 2016], three as public visitors and guests, e.g., [Zheng et al., 2005], six as co-players in a game, e.g., [Rehm, 2008], one as SMS contacts, e.g., [Hohenstein and Jung, 2018], and one as family members, e.g., [Luria et al., 2020b]. Since *polyadic* CAs are designed for multiple users in contrast to *dyadic* CAs for one-on-one interactions, the relationship types between users are a unique dimension for us to examine *polyadic* human-AI interactions.

Social scale. The social scale of the human users is another unique dimension of *polyadic* human-AI interactions since the designs examined in these papers are for two-individuals or multi-users (more than three users). There are 18 papers on multi-users. Combining the relationship types, we can see that the multi-users scale ranges from smaller groups (e.g., three to five co-learners, [Dyke et al., 2013]), to medium sized groups (e.g., ten discussants [Kumar et al., 2010]), and to larger groups (e.g., online community members on Twitter [Savage et al., 2016]).

Social theories and related theoretical frameworks. We also collected the theories that motivated the designs of the CAs, i.e., presenting theory-inspired or framework-based designs; other papers are less theory-driven. Among them, self-extension theory has been used to discuss users' sense of ownership of a community-owned CA [Luria et al., 2020a]; the Computers Are Social Actors (CASA) paradigm has been used to compare the patterns of human-CA interactions and human-human interactions [Rehm, 2008]; message framing theory has been applied to design the bot messages in implementations and evaluations [Savage et al., 2016]; balance theory has been leveraged to explore the sense of "balance" in the social dynamics in group interaction mediated by CAs [Nakanishi et al., 2003]; and the Structural Role Theory [Biddle, 1986] has been used to identify the roles that the bots play on Twitch [Seering et al., 2018].

Other papers build on theoretical frameworks, such as: proposing a harmony model of CA mediation using Benne's categorization of functional roles in small group communication [Benne and Sheats, 1948]; positing Mutual Theory of Mind as the framework to design for natural long-term human-CA interactions [Wang et al., 2021]; and applying Input-Process-Output models [Kozlowski and Ilgen, 2006] to explain the process of emotional management in teams [Benke et al., 2020]. A few papers discuss important concepts in multi-user interactions. For example, vulnerability and trust are discussed in teamwork settings to support a design of robot that is intended to build trust in teams; and Barsade's Ripple Effect [Barsade, 2002] is introduced to support hypotheses of a machine agent's positive influence on other individuals in a team just as effectively as a human member [Strohkorb Sebo et al., 2018]. However, the theories and theoretical frameworks used in the limited number of existing publications are rather scattered and disorganized.

4.4 Proven Effects of *Polyadic* CAs on Human-Human Interaction (RQ4)

The reviewed articles reported two major sets of effects of *polyadic* CAs: the effects on *dyadic* interaction, emphasizing the effects on one-on-one interaction between human users and the CAs; and the effects on *polyadic* interaction, referring to the impact of the CAs designed to mediate human-human interaction. The effects of the *polyadic* CAs on *polyadic* human-AI interaction include improving the individual users' learning effectiveness [Tegos et al., 2015, Chaudhuri et al., 2009], perceived self-efficacy [Tegos et al., 2015], and satisfaction [Dohsaka et al., 2009]. Another large portion of this set of effects involves users' positive or negative attitudes, perceptions, and acceptance toward the CAs [Ai et al., 2010, Wang et al., 2021, Nakanishi et al., 2003, Rehm, 2008], and interactions and engagement patterns between users and the CAs depending on specific design features [Ai et al., 2010, Strohkorb Sebo et al., 2018]. Most studies laid primary focus on *dyadic* human-AI interaction effects, while the evaluation on group effects is relatively insufficient. Overall, the effects reported in the reviewed papers show opportunities of using *polyadic* CAs to improve human social experience.

4.4.1 Communication Efficiency

In collaborations, studies have supported that *polyadic* CAs can help with consensus-reaching, aid communication comprehension, enhance task management, and save the time and energy of human collaborators from tedious work, which can improve the group work efficiency and overall collaborative experience. Specifically, for consensus-reaching, a moderator CA helps to align group consensus and individual opinions, contributing to reaching agreements [Kim et al., 2021, 2020]. For better comprehension of group communication, a CA can tag and summarize the group chats to help users make sense of the conversation better [Zhang and Cranshaw, 2018]. Moreover, a tour guide CA in museums can mediate visitors' interactions by prompting topics, offering information, and concluding discussions [Otogi et al., 2013, 2014]. In terms of management and coordination, a scheduling CA can coordinate schedules of team members fast and efficiently, setting human personal assistants free [Cranshaw et al., 2017]. For burdensome tasks, a searchbot is designed for its human collaborators to save time when switching between tools and searching independently by offloading these tasks to the CA.

4.4.2 Group Engagement

The papers we reviewed have provided evidence that *polyadic* CAs may benefit group dynamics through encouraging engagement and balancing uneven participation. The effects of encouraging engagement are most salient for collaborations in education and online communities. In multiple studies in collaborative learning, the intervention of a CA that can ask questions and prompt students to show their thought processes, stimulating pedagogically beneficial conversations in the learning groups [Dyke et al., 2013, Tegos et al., 2015]. In online communities such as Twitter, bots designed to call people to action to specific social activities engage users to make relevant contributions to the discussion and lead to their interactions regarding future collaborations [Savage et al., 2016].

Also, *Polyadic* CAs as facilitators in group discussions not only encourage the amount of engagement, but also improve the distribution of the engagement by nudging members to participate evenly [Shamekhi et al., 2018]. Several CA designs can balance discussion by inviting less engaged learners to join the conversation and downplaying the “talkative” learners from over-controlling conversations [Chaudhuri et al., 2009]. The BabyBot on Twitch can grow up and learn from human users, which engages community members as a whole and opens opportunities for newcomers to participate in interactions as a way to welcome them on board [Luria et al., 2020a, Seering et al., 2020].

CAs are also designed to moderate multi-user engagement in open and dynamic environments [Bohus and Horvitz, 2009a,b, 2010], e.g., the hotel reception, games, and organizational systems, by identifying partakers and bystanders, and facilitating engagement decisions of when and whom to engage. Furthermore, more intensive and even engagement brings more diverse content. The moderator CA for discussions can therefore help with generating diverse and deliberative opinions [Kim et al., 2021, 2020].

4.4.3 Relationship Maintenance

Regarding social and relational aspects in groups, *polyadic* CAs show potential in regulating emotions and relationships to maintain harmonized group dynamics. Prior work presented prototype designs of CAs to offset a shortcoming of the text-based communication in teams, that is, the insufficient ability of understanding and managing emotions of the team members [Benke et al., 2020, Kumar and Rosé, 2014]. By monitoring the sentiment of the conversations and providing suggestions, CAs can improve emotion regulation and compromise facilitation [Benke et al., 2020]. Similarly, *polyadic* CA can analyze group behaviors to reinforce positivity by preventing the use of negative words [Peng et al., 2019]. Moreover, designing an agent that can make vulnerable statements in collaboration generates a *Ripple Effect*, which notes that human teammates with an agent that expresses vulnerability are more likely to engage in trust-related behaviors, including explaining failures to the group, consoling team members who made mistakes, and laughing together. These actions reduce tension in the team and enhance a positive and trusting atmosphere in groups [Strohkorb Sebo et al., 2018].

4.4.4 Building Connections

In a two-individuals interaction, the existence of a CA is crucial to the balance and solidarity of the triad. Studies have shown that a CA's agreement, disagreement, or bias toward one or both side(s) of the two individuals can significantly change the mutual perceptions of the interaction partners [Ai et al., 2010, Nakanishi et al., 2003]. Also, when a CA agrees or disagrees with both the human pair, the interaction partners feel more similar and attractive to each other [Nakanishi et al., 2003]. In contrast, a CA tutor showing a bias towards one perspective or another will polarize the learning pair [Ai et al., 2010]. Thus, *polyadic* CAs demonstrate potential to build solidarity between interaction pairs.

Polyadic CAs may also be used to establish new connections between pairs that meet for the first time. i.e., a helper agent that is built to form social connections between people in cross-cultural conversations [Isbister et al., 2000]. The

evaluation showed the positive effects of the agent designed to build common ground. Also, a CA was designed as a virtual companion in a university setting to help first-year students get on board and build new connections [Löw and Moshuber, 2020]. However, there were also a series of mixed findings of nuanced human users’ social perceptions regarding self, partner, and cultural stereotypes, which are worth further exploration.

4.5 Evaluation Metrics of CAs (RQ5)

We synthesize the metrics of prior studies used to evaluate user experience of conversational agents. These metrics show what categories previous researchers considered in evaluating their designs and what measurements they used. In the following section, we report on the results emerging from the analysis of the 135 *dyadic* papers and the 36 *polyadic* papers included in the final corpus. We categorize them in three main areas: chat log analysis or observation, survey scale, and interview. For definitions of all the metrics, see Appendix A.1 and A.2 for details.

4.5.1 Log Analysis or Observations

For *dyadic* papers, 36 papers evaluated CAs using log analysis or observation methods such as coding user behaviors through watching recorded videos. Table 4. provides an overview of what metrics were evaluated in the user studies. Six categories emerged: 1) linguistic features, i.e., language patterns that are evident in the language use; 2) prosodic features, i.e., features that appear when we put sounds together in connected speech; 3) dialogue features, i.e., features that are related to conversation structure and chat flow; 4) tasks-related features, i.e., features that measures metrics related to task fulfillment; 5) algorithm features, i.e., features concerning model training and efficiency; and 6) user-related features, i.e., features that are closely related to user perceptions and behaviors.

Linguistic features were evaluated in seven recent papers. To analyze the linguistic features in the chat logs, most of these studies leveraged existing natural language processing toolkits to analyze the text. For example, LIWC [Kawasaki et al., 2020] uses 82 dimensions to determine if a text uses positive or negative emotions, self-references, causal words etc. Another way is to count the use of linguistic patterns, such as the use of pronouns and proportion of utterances with terms repeated from previous conversations [Thomas et al., 2020]. Similar to linguistic features, prosodic features were used to evaluate the language pattern, but especially in speech-based agents, e.g., prior work measured the rate of speech, pitch variation, loudness variation, using OpenSMILE to process the audio signals [Thomas et al., 2020].

Dialogue features were used earlier than linguistic features, and were evaluated in nine papers, measuring the dialogue quality and efficiency [Xiao et al., 2020a], e.g., length, dialogue duration, number of turn taking, response time; and dialogue expressiveness, e.g., percentage of sentences. One of the frameworks being widely adopted in prior works is the PARAdigm for dialogue system evaluation, also known as PARADISE framework, which examines user satisfaction based on measures representing the performance dimensions of task success, dialogue quality, and dialogue efficiency, and has been applied to a wide range of systems, e.g., [Foster et al., 2009].

Task-related features were examined in nine papers. This dimension helps researchers to understand how users fulfill the tasks assigned by or facilitated by CAs. These dimensions are useful for task-oriented CAs that were designed to help users execute a task or solve a problem, e.g., customer service. While the task fulfillment rate was continually used as a metric in evaluating tasks, e.g., [Demberg et al., 2011, Medhi et al., 2011, Pecune et al., 2018], some metrics were introduced in recent years, such as dropout rate, e.g., [Kim et al., 2019], to capture the percentage of all participants who did not complete the task during evaluation.

Eight papers evaluated user perceptions and social behaviors by analyzing the chat log between the CA and the users, using metrics such as self-disclosure [Lee et al., 2020a, 2021], intimacy [Sannon et al., 2018], and aggressiveness [Bonfert et al., 2018]. Some theoretical frameworks were used, e.g., Barak and Gluck-Ofri’s scale, which conceptualized self-disclosure into three dimensions as information, thoughts, and feelings [Lee et al., 2020a].

For *polyadic* papers we reviewed, 14 studies employed log analysis and were mostly published during the 2010s. Several papers used task-oriented metrics, see Table 3, e.g., task duration, efficiency, effectiveness, and team performance [Avula et al., 2018, Seering et al., 2020, Luria et al., 2020a, Kumar and Rosé, 2014]. The evaluation metrics were both quantitative and qualitative. The quantitative metrics counted the frequency of users’ input or utterances [Tegos et al., 2015, Avula et al., 2018, Dohsaka et al., 2009, Abokhodair et al., 2015, Otogi et al., 2013, 2014, Narain et al., 2020]. The qualitative metrics evaluated discussion quality [Wang et al., 2021, Tegos et al., 2015, Kim et al., 2021], with an emphasis on consensus reaching [Kim et al., 2021], opinion diversity [Tegos et al., 2015], and linguistic features, e.g., verbosity [Wang et al., 2021], readability [Wang et al., 2021], adaptability [Wang et al., 2021] and positive and negative sentiment [Wang et al., 2021, Hohenstein and Jung, 2018, Nakanishi et al., 2003, Narain et al., 2020].

Table 3: *Polyadic* ACM Papers using Different Metrics for UX Evaluation (Count by Year Since 2005)

Category	Subcategory	Metrics (Defined in Appendix A.1)	(20)	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21
Chatlog	Socialbility	discussion quality		1					1	1										
		consensus reaching							1											
		opinion expression							1											
		opinion diversity							1											
		even participation							1											
	Task-related	linguistic features																	1	
		message quantity							1								1		1	
		task completion time										1								
		efficiency																	1	
		effectiveness																	1	
Survey	Communication-quality	satisfaction																	1	
		team performance										1							1	
		perceived communication effectiveness							1											
	Socio-emotional	perceived communication fairness							1										1	
		perceived communication efficiency							1											
		perceived group climate															1			
		perception of other group members										1	1				1		1	
		perception towards the agent										1								
		perception about collaboration																		
	Traditional	perceived social presence														1	1			
		perceived social support														1				
		level of annoyance with the intervention																		
		opinion of oneself/partner/cultural stereotypes							1											
		perceived appropriateness of agent's social behavior															1			
		users' psychological wellbeing																	1	
		anthropomorphism															1			
		perceived emotional competence														1				
		unmet expectations														1				
Interview	Socialbility	privacy concerns														1				
		level of perceived conversation control							1											
	Issue-relevant	perceived satisfaction							1										1	
		performance/effectiveness															1		1	
		level of engagement																		
		perceived capability to promote contributions															1			
		decisions to (not) engage with the CA																	1	
		overall UX															1		1	
		reflections on selves and others																	1	
		willingness to take actions														1				
		direction of improvement															1		1	

Table 4: *Dyadic* ACM Papers using Chat Analysis for UX Evaluation (Count by Year Since 2005)

Category	Metrics (Defined in Appendix A.2)	(20)	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21
Linguistic features	positive/negative words																1	3	
	length of utterances																	2	
	personal pronouns																	2	
	term-level repetition																	1	
	utterance-level repetition																	1	
	use of concrete words																	1	
	variation in language used																	1	
Prosodic features	variation in speech-based agents																	1	
Dialogue features	dialogue/response quality						1											3	
	dialogue efficiency metrics				1		1		1									2	
Tasks-related	expressiveness				1														
	task fulfillment rate/effectiveness						1		2							1			
	system usage															3			
Algorithm	user entry time								1										
	dropout rate																1		
	Intent modeling evaluation																	1	
User related	error rate								1										
	self-disclosure															1	1	1	1
	intimacy															1			
	reflection															1			
	impolite/aggressive behaviors															1			
	motivation																	1	
	affective states																	1	

4.5.2 Survey Scale

For *dyadic* papers, 93 papers evaluated CAs using survey scales. Table 5 provides an overview of what features were evaluated in the user studies. Four categories emerged: 1) conversation related metrics: aspects that help CAs to manage dialogues during interactions; 2) user perception of CA’s social features: aspects that reflect user perception of CA following social protocols; 3) user perceived system usability: aspects that capture the quality of a UX when users interact with CAs; and 4) user self-reported experience with CAs.

For conversation related features, some studies leverage widely adopted scales because they are often used for system evaluation. For example, informativeness, or information quality, has been used multiple times in user evaluations. Informativeness refers to quality of the conversational system to communicate truthfully, provide relevant information, and share content clearly and orderly. Prior work on AI-powered CA used Gricean Maxims’ information quality metrics [Xiao et al., 2020b] to evaluate this. Other metrics were proposed as modes of examination in recent years, e.g., syntactic readability, self-reported response quality and intensity of sentiments [Wambsganss et al., 2020], and instruction quality [Foster et al., 2009].

Multiple features were proposed in evaluating users’ perceptions of conversational agents’ social features, such as perceived agents’ common ground [Chai et al., 2014], sociability, [Wang et al., 2020], and perceived interruption or annoyance [Schulman and Bickmore, 2009]). Metrics such as playfulness [Zhou et al., 2019, Xiao et al., 2019] (enjoyment, interestingness, or funny) and perceived intimacy (perceived interpersonal closeness, or friendliness) have been evaluated using existing scales. For example, the inclusion of others in the self scale and the subjective closeness index to measure level of closeness [Yu et al., 2019]. With the development of conversational agents, more social features have been introduced in user evaluation, such as perceived anthropomorphism [S. Alpers et al., 2020], perceived personality traits [Völkel et al., 2020], and perceived naturalness [Shi et al., 2020].

For system usability evaluation, 33 papers evaluate CA usability using traditional usability metrics. For example, UMUX-LITE is the usability metric for user experience to measure evaluating factors such as perceived ease of use, and it is used in papers [Bickmore et al., 2005, Biswas, 2006]. Similarly, NASA-TLX, or task load index, which evaluates mental demand, physical demand, effort, frustration, and future adoption willingness, is used in studies [Jain et al., 2018b,a]. Moreover, prior work also uses usability questions from previous surveys [Khadpe et al., 2020].

There are 45 papers evaluating the conversational agent through collecting metrics regarding user self-reported personal experience responses. Among them, satisfaction, e.g., [Demberg et al., 2011, Bickmore et al., 2005, Liu et al., 2020] and perceived trust, e.g., [Zhou et al., 2019, Lau et al., 2010, Hoegen et al., 2019, Lee et al., 2020a], and perceived engagement, e.g., [Muresan and Pohl, 2019, Cai et al., 2021, Shi et al., 2020, Xiao et al., 2020a] were evaluated the most in the sample. There are some new features been used to evaluate *dyadic* interactions, such as serendipity [Cai et al., 2021], social orientation toward CAs [Ashktorab et al., 2019], degree of collaboration [Mitchell et al., 2021], and perceived capability to control [Thomas et al., 2020]. Thus, evaluation dimensions are becoming increasingly diverse when examining user experience with conversational agents.

Polyadic studies which employed surveys mainly evaluated users’ perceptions of both the functional and socio-emotional aspects of the interactions, with specific focus on communication quality and interpersonal dynamics. The former includes perceived communication effectiveness [Tegos et al., 2015, Kim et al., 2021], communication efficiency [Tegos et al., 2015, Kim et al., 2021, Luria et al., 2020a], communication fairness [Kim et al., 2021, 2020], and the quality of collaboration [Avula et al., 2018]. The latter includes perceived group climate [Gulz et al., 2011, Seering et al., 2018], social presence [Benke et al., 2020, Seering et al., 2018], social support [Benke et al., 2020], perceptions on the agent [Dohsaka et al., 2009, Seering et al., 2018, 2020, Kumar and Rosé, 2014], appropriateness of the agent’s intervention [Avula et al., 2018], as well as impressions about other group members [Gulz et al., 2011, Seering et al., 2020, Kumar and Rosé, 2014]. We observe that during the past decade works evaluated social dynamics in the interactions [Gulz et al., 2011, Benke et al., 2020]. These studies with survey methods also used a series of metrics that are commonly used in *dyadic* interactions, which evaluate the ability and competence of CAs [Avula et al., 2018, Benke et al., 2020, Seering et al., 2018], task performances [Seering et al., 2018, 2020], perceived anthropomorphism [Gulz et al., 2011, Seering et al., 2018, Shamekhi et al., 2018], intelligence [Shamekhi et al., 2018], likeability [Gulz et al., 2011], users’ satisfaction of the interaction [Kim et al., 2021, Seering et al., 2020], perceived control of the interaction [Nakanishi et al., 2003], and level of engagement [Avula et al., 2018, Seering et al., 2018].

4.5.3 Interview

For *dyadic* interaction papers, overall, 17 papers evaluated CAs using interview data analysis. Table 6. provides an overview of what metrics were evaluated in the user studies. On the one hand, users were interviewed regarding their general impressions on the CAs, e.g., [Kim et al., 2019, Lee et al., 2021, Prasad et al., 2019, Lee et al., 2020b, Cavazos Quero et al., 2019], perceived CA characteristics such as capability in handling requests [Følstad and Skjuve,

Table 5: *Dyadic* ACM Papers using Survey Scales for UX Evaluation (Count by Year Since 2005)

Category	Metrics (Defined in Appendix A.2)	(20)05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21
Conversation related	syntactic readability																	1
	self-reported response quality																	1
	informativeness/information quality	1													1			1
	conversation smoothness															1		
	intensity of sentiments																	1
Perception towards CA	instruction quality					1					1							
	perceived common ground										1							
	opinion of the CA as a partner					1												
	liking attitude	1										1			1			
	playfulness/enjoyment	1								1					1	3	3	1
	socialability														1		1	
	perceived anthropomorphism																2	
	perceived warmth of the AI system																1	
	impression																1	
	perceived personality traits																1	
	perceived persuasiveness					1												
	interaction/annoyance	1																
	perceived intimacy	1												1	1	2	2	1
	perceived naturalness																1	
	perceived safeness															1		
System usability	overall usability									1			1			4	2	2
	perceived ease of use	1	1									1			1		1	
	mental demand														2		1	
	physical demand														2			
	perceived task completion					1					1				2			
	effort														2			
	frustration														2			
	desire to continue the interaction									1					2		2	
	user acceptance of the agent									1					2		2	
	system consistency																1	
	perceived usefulness/helpfulness															1	4	1
	willingness to recommend																1	
	motivation															1		
	overall UX														1	1	2	
	contrast of user experience																1	
UX with CA	perceived quality of the interaction															1	1	
	perceived quality of CA															1	1	
	perceived satisfaction	1		1		1		1		1				1	1	3	3	
	perceived engagement															1	3	2
	perceived trust	1					1								1	1	2	1
	confidence															1		1
	perceived capability to control																1	1
	aroused emotions					1											1	
	serendipity																	1
	pleasant surprise																	1
	empathy															1	1	1
	self-reflection/self-awareness														1		1	1
	self-disclose														1		1	
	anxiety/tension																1	
	comfort																1	
	social orientation toward CAs															1		
	degree of collaboration																	1
	degree of decision-making																1	1

2019], personality [Kim et al., 2019, Prasad et al., 2019], and trust [Large et al., 2019]. On the other hand, they interviewed user behaviors towards conversational agents such as efforts used while interacting with the CAs [Lee et al., 2021], actual engagements [Lee et al., 2021, Ruan et al., 2019], and daily using practices [Lee et al., 2020b].

For *polyadic* interaction papers, studies that employed interviews focused more on specific issue-relevant metrics. These metrics include perceived capability of the agent on solving a particular issue [Gulz et al., 2011], engagement in the conversation [Avula et al., 2018, Seering et al., 2020, Narain et al., 2020], willingness to take actions under the persuasion of the CA [Savage et al., 2016], reflections on selves and others [Seering et al., 2020, Narain et al., 2020]. Some general metrics include overall user experience and impression [Seering et al., 2020, Narain et al., 2020, Shamekhi et al., 2018], as well as potential directions for improvement [Narain et al., 2020, Shamekhi et al., 2018]. Social interaction features seem to be essential for *polyadic* papers, especially factors related to promoting team contributions and engagement, e.g., [Avula et al., 2018]. There are many features that have been evaluated in *dyadic* interactions but are also suitable to be evaluated in the *polyadic* context, such as effect on judgement [Bawa et al., 2020], and perceived burden, i.e., time, financial, mental, and emotional burden [Park and Lee, 2020].

Table 6: *Dyadic* ACM Papers using Interview Related Metrics for UX Evaluation

Category	Metrics (Defined in Appendix A.2)	Number of Papers (Between 2005 and 2021)
Perceptions	overall impressions and experience	5
	perceived CA's capabilities	1
	perceived CA appearance and self-presentation	1
	perceived personality	2
	perceived naturalness	1
	reciprocativity	1
	perceived benefits of the agent	1
	motivation to use CAs	2
	perceived burden	1
	the best and worst aspects of using CAs	1
	perceived human-likeness	1
	perceived communicative experiences	1
	perceived enjoyment	1
	trust	1
	perceived effectiveness	2
Behaviors	engagement	2
	notice of CA's certain function	1
	efforts in learning the system	1
	effects on self-awareness, self-reflection	1
	effect on judgement	1
	effects on collaborations	1
	daily practices of using the CA	1

4.6 Overlooked Issues of *Polyadic* CAs (RQ6)

We also reviewed overlooked issues and additional findings aside from the primary design focuses, which raise intriguing issues that can be missing in design guidelines. These issues are related to appropriateness, privacy, and ethics in the designs of CAs, which warrant deeper discussions about the role of *polyadic* CAs, their social influence, and their relationship with human users in different group settings.

4.6.1 *Polyadic* CAs Need to be "Visible"

One overlooked issue points to the users' awareness of the *polyadic* CAs in the group. [Avula et al., 2018] presented a searchbot to support collaborative search tasks, and the authors discussed that the collaborative nature of their searchbot posed a new issue regarding awareness of the CA. They suggested that the searchbot should announce itself and remain "visible" to the users throughout the interaction. Whether and how to keep users' awareness of *polyadic* CAs were discussed in diverse contexts, such as the tutor [Chaudhuri et al., 2009, Ai et al., 2010], the assistant [Avula et al., 2018, Cranshaw et al., 2017], the moderator or facilitator [Kim et al., 2021, 2020, Isbister et al., 2000, Hohenstein and Jung, 2018], and one of the members in online communities [Savage et al., 2016, Seering et al., 2018]. The findings suggest a direction regarding design decisions to make users aware of CAs as their peers or not.

4.6.2 *Polyadic* CAs Need to be "Ignorable"

There is some discussion of ignorable design in the reviewed papers. In collaborative learning contexts, [Tegos et al., 2015] mentioned that CAs are sometimes ignored and abused by the group learner. Authors found that students provided hasty answers to the tutor agent and sometimes wanted to pay more attention to the learning questions instead of the agent's facilitation. Similarly, participants perceived a task reminder CA to be invasive or annoying, as they are "too frequent", "not context sensitive", and distracting [Toxtli et al., 2018].

Only one paper mentioned "easy to ignore" as a design suggestion when developing a supporting agent in multi-user contexts [Isbister et al., 2000], because if CAs are persistent with their intervention, the effects can potentially backfire. Given that there are limited papers discussing ignorable design suggestions, it remains a crucial issue for designers to create CAs that are less intrusive and are able to detect the moment when users do not want their interventions and shut down properly, particularly when the interaction between human users is the major focus.

4.6.3 *Polyadic* CAs Need to be Accountable

Another overlooked issue arises in voice activated CAs in interpersonal spaces [Luria et al., 2020a]. As CAs are involved in multiple users' interactions in their homes, participants were confused about how an CA would deal with interpersonal conflicts between users in the home without invading privacy [Luria et al., 2020a]. Papers also discussed the issue of CA ownerships - who should CAs be accountable to? If a CA is the mediator between human users, how much can and should other users consider this mediator? If there are interpersonal conflicts, what is the standpoint of

the CA? What will happen if the agent crosses a perceived boundary, and how should we tackle it? Several questions such as these were asked in the reviewed papers [Seering et al., 2020, Luria et al., 2020a, Zheng et al., 2021].

4.6.4 Topics Discussed in *Polyadic* and *Dyadic* Works

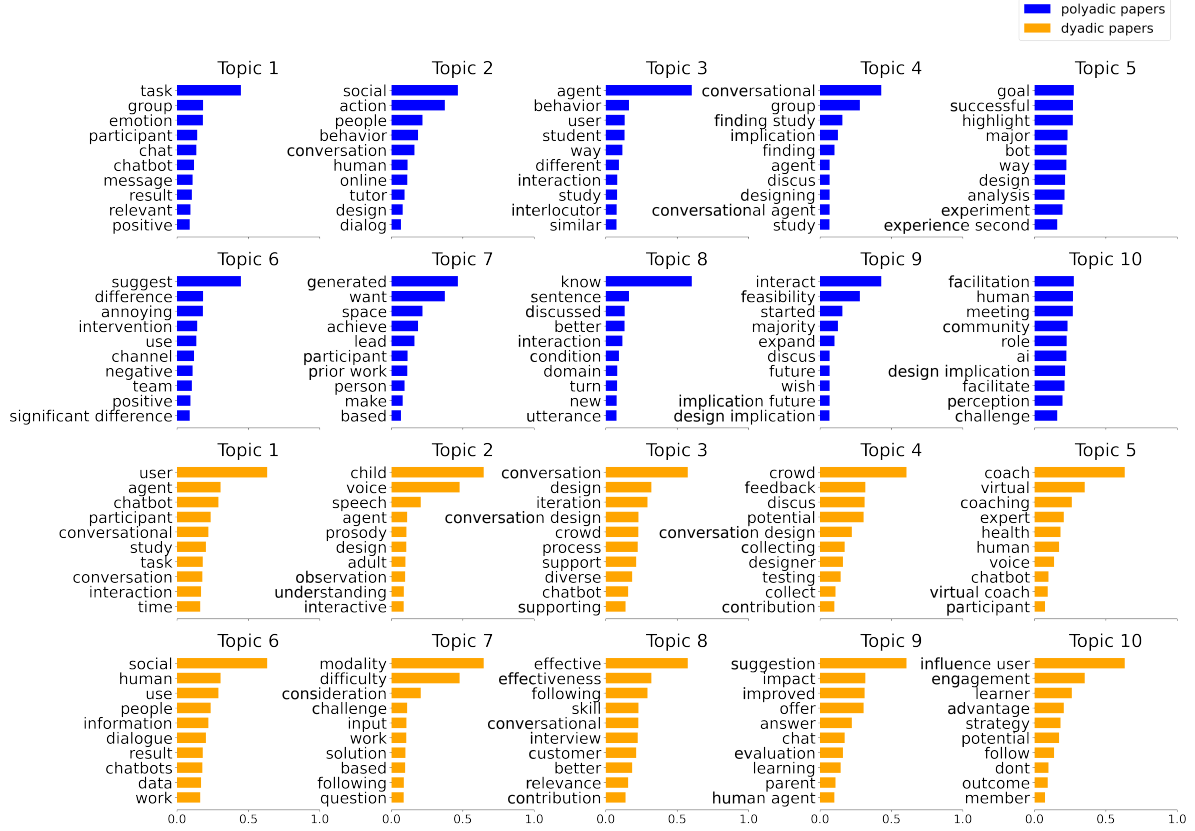


Figure 2: Comparison between topics (ordered by weight) discussed in *polyadic* and *dyadic* papers (the x-axis is the probability distribution of each term in its corresponding topic)

To explore the difference between the topics discussed in the UX research of *dyadic* and *polyadic* CAs, we applied topic modeling method over the papers' *discussion* sections in which authors reflected on the findings and proposed future research directions. As shown in Figure 2, we found that the top topics identified from the *polyadic* CA studies covered different concepts that were not frequent in the *dyadic* CA articles. As shown in Figure 2, *polyadic* papers' Topic 1 was about *group discussion* (18.82% topic weight); whereas *dyadic* papers mentioned *user-agent interaction* the most (44.16% topic weight). For example, Peng et al. evaluated a *polyadic* CA [Peng et al., 2019] that facilitated emotion regulation in group discussions, which received the highest weight for the *group discussion* topic. This work also contributed to *polyadic* Topic 2 (about *social behavior*, 11.98% topic weight), as it discussed how to positively impact users' social behavior and emotion. A similar topic appeared in the *dyadic* works but was ranked much lower in popularity (Topic 6 with 5.86% weight). On the contrary, Muresan et al.'s paper [Muresan and Pohl, 2019], a *dyadic* CA weighted the highest for Topic 1 *user-agent interaction*, discussed how anthropomorphism of CAs affects users' engagement. This study also contributed to *dyadic* Topic 2 (*conversational design* with 8.08% topic weight). However, both topics were not appeared in the top 10 of *polyadic* works. In addition to *social behavior* (e.g., [Rehm, 2008, Kumar and Rosé, 2014]), *polyadic* CA works tended to discuss more on *education* (Topic 3, 10.94% weight, e.g., [Kim et al., 2020, Dyke et al., 2013]), and *embodied design* (Topic 5, 10.13% weight, e.g., [Seering et al., 2020, Bohus and Horvitz, 2010]). But *education* related topic was also ranked low in *dyadic* papers (Topic 10 with 4.49% weight).

5 Discussion

Given the above results of the literature review, we discuss several research opportunities for future design and development of *polyadic* CAs to address social challenges in human-to-human interactions. These opportunities emerge

in three main directions, i.e., exploring an under-studied design space, using and developing theoretical foundations, and developing HCI design guidelines for building boundary-aware CAs. We also discuss specific research agendas of future communicative AI technologies [Guzman and Lewis, 2020] around two key aspects, e.g., relational dynamics and functional dimensions.

5.1 Spotlight an Under-Explored Design Space

Overall, we found a very small fraction of CAs, published in ACM venues, designed for *polyadic* human-AI interaction. Results suggest that such a design space has been severely under-explored. A predominance of CA surveys worked on the classification of CAs [Hussain et al., 2019, Seering et al., 2019, Laranjo et al., 2018, Følstad et al., 2018, Cassell, 2000, Klopfenstein et al., 2017, de Barcelos Silva et al., 2020], and some efforts were made in investigating CAs’ potential to improve user engagement and experience in *dyadic* Human-AI interaction [Chaves and Gerosa, 2020, Rapp et al., 2021]. However, still little attention was given to the conversational effects of CAs on *polyadic* human-AI interaction. Thus, this literature review filled the research gap and mapped the progress in this field for future researchers.

In terms of the application domains, results showed that some areas (e.g., public services, health) fall short of UX research on *polyadic* CAs. Also, user evaluation lacks empirical findings of large scaled human-human interaction. We also found that most original designs were proposed to explore the feasibility of conversational AI in varied application domains (e.g., [Toxtli et al., 2018, Dohsaka et al., 2009, Seering et al., 2020, Kim et al., 2020, Benke et al., 2020, Luria et al., 2020b]). With the tsunami of conversational AI, the need for design knowledge becomes more intense [Rapp et al., 2021]. However, there is no wide consensus on grounded practices on which to create novel designs.

5.2 Connect with and Contribute to Theories

Our review highlights a limitation in the theoretical groundings of designing and evaluating CAs for *polyadic* Human-AI interactions. We found that many works did not leverage or contribute to existing theories. This was also found by prior literature reviews of CAs [Rapp et al., 2021, Clark et al., 2019]. There can be two reasons. First, *polyadic* conversational design is still an emerging area, and HCI researchers have not identified the best theories to explain, build on, and support their studies. Currently, some adopted existing theoretical grounds from social psychology (e.g., CASA) [Rehm, 2008], sociology (e.g., balance theory, structural roles) [Nakanishi et al., 2003, Seering et al., 2018] and communication studies (e.g., self-extension theory) [Luria et al., 2020a]. Second, *polyadic* CAs are bringing brand-new dynamics of social interactions, to an extent that existing theoretical frameworks cannot account for. For example, prior work [Fox and Gambino, 2021] identified that human users’ expectations and perceptions towards machines are different from those towards human partners. Thus, when new social dynamics are formed, existing theories may not be sufficient.

To leverage existing theories, scholars need to identify the social interaction *challenges* (e.g., findings in RQ1) to be addressed and the *application domains* (e.g., findings in RQ3), and then look for related theories resolving the particular challenge in the relevant application domain to inspire the CA designs. For example, some *polyadic* CAs are designed for group members to reach consensus [Kim et al., 2021, 2020, Benke et al., 2020]. Prior works on collaborative engineering, a domain where designs are created to enable engineers to work effectively with all stakeholders for completing collaborative tasks [Borsato and Peruzzini, 2015], have applied consensus building theory (CBT) [Briggs et al., 2005] in designing for collaboration processes [Kolfshoten and De Vreede, 2007, Nabukenya et al., 2009]. The CBT is often applied by those conducting IT requirement negotiations or those conducting risk and control self-assessments. In the situation where decisions cannot be made unilaterally because team members are co-responsible peers, teams need to resort to approaches to build consensus to gain commitment [Kolfshoten and De Vreede, 2007]. The CBT could be applied to direct four major steps in reaching a consensus, i.e., articulating a proposal, evaluating willingness to commit, diagnosing causes of conflict, and invoking conflict resolution strategies [Briggs et al., 2005]. Similarly, when designing a *polyadic* CA under a collaborative negotiation situation, designers may leverage the CBT model to inform designs. Future studies can also adopt a similar approach in finding suitable theories or frameworks.

5.3 Propose the Notion of Boundary-Awareness

Only a few *polyadic* papers used existing systems in user evaluation. Regardless of the CAs’ originality, the overlooked negative UX identified by prior research indicates that there is a pressing need for improvements in design strategies and guidelines.

5.3.1 *Polyadic* CAs are Unique to Design

In prior works, human-likeness (e.g., empathy [Samrose et al., 2020] and self-disclosure [Lee et al., 2020b]) is identified as a critical aspect to improving UX and to encouraging users to show favorable feelings (trust, openness, tolerance, etc) towards CAs [Rapp et al., 2021, Chaves and Gerosa, 2020]. However, similar topics are less reflected by *polyadic* human-AI interaction (RQ5). Instead, researchers raised most open concerns on how CAs should establish social boundaries in human-human interactions, which suggests that *polyadic* CAs may be *unique* to design compared with *dyadic* CAs.

Taking the learning context as an example, when a *polyadic* CA was deployed as a study peer, researchers found that students teams often ignore or provide hasty replies to CAs [Tegos et al., 2015]. However, how can CAs be less intrusive partners in such a context? Similarly, in a family context, when family members have different opinions, how can the CA understand when to "knock the door" when it is tackling conflicts between couples [Zheng et al., 2021]? Should it let parents know where their children are up to [Luria et al., 2020b]? Our findings highlight the importance of questions for CAs to understand social boundaries. Current studies are yet able to take care of the dilemmas in situations where user needs or human-AI interaction conflict and the CAs need to react or take sides.

5.3.2 Design CAs with Boundary-Awareness

By reviewing the overlooked issues, we identified multi-dimensional themes in *polyadic* human-AI interaction, which we consider as social boundary issues (RQ6). Boundaries in social sciences can be understood as a set of rules followed by most people in a particular society, which are vital in the society because they can guide human behaviors and assist in managing what is and what is not acceptable [Houghton and Joinson, 2010, Lamont and Molnár, 2002]. This issue is closely related to privacy and disclosure, as the presence of *polyadic* CAs as social actors changes the social dynamics and users' information management strategies.

Communication privacy management theory (CPM) [Petronio, 2013] identified three main elements in people's private information management: (1) privacy ownership, (2) privacy control, and (3) privacy turbulence. When an individual has decided to disclose private information, as a result, the recipient of the information becomes a co-owner or shareholder. From that moment, the initial owner of the information must set rules and boundaries of how to manage private information. Privacy turbulence may happen when these rules are violated, and the original owner may refrain from disclosing any more information or may engage in negotiations or coordinate rules and boundaries with the co-owner [Petronio and Durham, 2014]. Clarifying ownership of users' private data to *polyadic* CAs and enabling them to learn privacy boundaries in relationships over time are essential steps towards building social sensibility *polyadic* CAs when participating in human-human interactions. For example, a boundary-aware CA can learn if and when it should intervene when mediating couples' conflicts by receiving requests from one side of the couple [Zheng et al., 2021]. Similarly, when users request a CA to understand situations of their living alone elderly parents, the boundary aware CA can reply appropriately. The privacy boundaries in the interactions also depend on how CAs present themselves and how users perceive the CAs. For instance, when CAs join human interactions, their roles may vary from active conversation participants to less salient group collaboration assistants. A boundary-aware CA may adjust its proactivity and intensity in joining the discussion.

In our review, multiples studies showed that during the interactions, users expressed confusions and concerns about what roles *polyadic* CAs should play, how *polyadic* CAs should behave and what is the proper distance that *polyadic* CAs should maintain with their users in various group settings. As *polyadic* human-AI interaction involves multiple relationships and complex social dynamics in its nature [Tegos et al., 2015], the boundary between CAs and different users, and the boundary between users need to be taken care of.

Thus, we propose to design CAs with *boundary-awareness* for supporting *polyadic* human-AI interaction. Traditional HCI design addressed boundaries, e.g., between computers and people and between computers and the physical world [Vogel and Balakrishnan, 2004, Stephanidis et al., 2019]. Different from the traditional boundaries, which are often between virtual and physical worlds [Burden and Kearney, 2016], the new boundaries brought by the CAs involve human-to-human boundaries. Prior work [Mou and Xu, 2017] suggested that, when facing AI, humans demonstrate different personalities from interacting with other humans. Thus, it might be difficult for CAs to add boundary management in the design because of the collective result of unclear roles, undecided social rules, and social emotions [Guzman, 2020] during human-AI interaction.

CAs with *boundary-awareness* cannot be detached from the discussion of privacy, e.g., [Avula et al., 2018, Luria et al., 2020b, Reig et al., 2020]. To address the boundary issue such as "whether the CA should inform the parents where their kids are up to" [Luria et al., 2020b], we can gain insights from privacy boundaries. For example, Palen and Dourish [Palen and Dourish, 2003] discussed three boundaries with different desires during one's privacy management. 1) "Disclosure boundary" is the desire to keep one's information private or public. Privacy regulation in practice is

not simply a matter of prohibiting one's data from being disclosed. Instead, human-to-human interactions frequently require selective disclosures of personal information to declare allegiance or even differentiate ourselves from others. 2) "Temporal boundary" suggested that privacy management is in the tension between past, present, and future. Users' response to disclosure situations is interpreted according to other events and expectations. 3) "Identity boundary" implied a boundary between self and other. For example, employees are discouraged from using corporate email addresses to post to public forums. Designers can leverage these three aspects to understand the most critical privacy boundaries in creating CAs with boundary awareness.

5.3.3 Adopt and Evaluate the Existing Guidelines

The unique challenges urge us to re-examine existing AI design guidelines and design notions. For example, Microsoft's guidelines for designing human-AI interaction suggested that AI-infused systems should "support efficient dismissal" [Amershi et al., 2019], therefore CAs should be able to learn the social boundaries such that their interaction is perceived less intrusive, addressing a concern that was raised in prior *polyadic* CA works [Tegos et al., 2015, Beuth et al., 2010]. The guidelines also suggested that AI systems should "match relevant social norms" (i.e., "the experience is delivered in a way that users would expect") and should "mitigate social biases" (i.e., "the system language and behaviors do not reinforce undesirable and unfair stereotypes and biases given their social and cultural context") [Amershi et al., 2019]. An example of potentially using these guidelines can be found in prior *polyadic* CA works [Kim et al., 2021] which designed a CA tailored for structural group discussions. The study reported that the CA facilitates the team reaching a logical consensus on a highly contentious topic even though the members' positions and understandings are vehemently opposed to each other, indicating the design of CA's that acts in a way that the team expected (knowing the norm). However, the authors also proposed that when discussing divisive issues, it may be more appropriate to design with interpersonal and social power dynamics in mind and encourage and protect participants' contributions from marginalized groups (understanding social biases.)

Moreover, existing CA guidelines also provide a broad set of items for designing responsible and trustworthy CAs, such as transparency, human control, awareness of human values, accountability, fairness, privacy and security, accessibility, and professional responsibility, e.g., [Microsoft, 2018]. These guidelines can be employed in future CA designs, and more academic works should also evaluate these guidelines' effectiveness in the wild.

5.4 Other Key Aspects of Communicative AI

We also discuss research agendas of future communicative AI technologies from two perspectives, i.e., relational dynamics and functional dimensions.

5.4.1 Relational Dynamics

Findings from RQ4 suggest that *polyadic* CAs show potentials for influencing social dynamics. Relational dynamics reflect the ways people interact with communicative technologies, with themselves, and with other people or group of people [Guzman and Lewis, 2020]. Possible research directions include the dynamics of relationship types, and relational attributions. We suggest that *polyadic* CAs can be further investigated in its effectiveness to foster new relational dynamics. For example, according to the Computers as Social Actors paradigm [Nass et al., 1994], CAs are considered as "social actors". Humans rely on many perceivable attributes during an social interactions with other humans [Donath, 2007]. To make sense of human-AI interaction, prior works identify, implement, and test CA effects in various social attributions, such as "social cue" [Ochs et al., 2017, Feine et al., 2019] (i.e., a signal that triggers a social reaction of the user towards the CA), "social roles" [Seering et al., 2018] (i.e., the function CAs and humans play in the interaction context), and social identities [Mirbabaie et al., 2021, Wambsganss et al.] (i.e., how CA or humans organize themselves into and within groups, social values that take collective goals, ethical concerns into consecrations), to name a few.

5.4.2 Functional Dimensions

We propose that the design and use of *polyadic* CAs should address questions such as: 1) what communication challenges are to be addressed by *polyadic* CAs?; 2) what is the unique function of *polyadic* CAs?; and 3) how can *polyadic* CA effectively address the identified challenge compared to other methods? These questions should fit into the functional dimensions proposed by [Guzman and Lewis, 2020], which reflects how certain AI technologies are designed and how people can make sense of these devices and applications.

In our literature review, not all works identified a fundamental challenge of human-human interaction before proposing a *polyadic* CA (RQ1). There could be a tendency of technology-determinism, believing that "technological change can determine social change in a prescribed manner [Dafoe, 2015]." Namely, some works assumed that CAs could help with

the challenges without evaluating the effect of the proposed CAs by comparing with other counterparts, e.g., without a CA, with a human, or with previously tested CAs. In future research, designers and practitioners may be better aware of this potential tendency and propose more comprehensive study plans to evaluate the proposed CAs.

6 Limitations and Future Work

This literature review has several limitations as a result of the data source, the search query, and data analysis methods. First, we chose ACM DL as the source of our literature review work. Several prior works also conducted literature reviews by using only ACM publications for HCI and CSCW [Wainer and Barsottini, 2007, Pinelle and Gutwin, 2000] and computer science research [Wainer et al., 2009]. However, using one major source could potentially generate a selection bias, e.g., including certain publications from non-ACM venues into the analysis could be challenged as "distorting the conclusions toward a particular direction" [Wainer et al., 2009]. Therefore, when presenting our findings, we focused on the selected papers as samples illustrating the emerging themes that were not presented in prior literature reviews. One major purpose of doing so is to inform UX researchers and practitioners to design conversational AI by leveraging the existing UX research outcomes and apply the suite of measurements in their work. In the future, conducting a systematic literature review by sampling from different sets of HCI research, e.g., papers indexed by IEEE, Web of Science, Scopus, and Science Direct, can further deepen our understanding of conversational AI research.

Second, we selected research publications that completed both design and UX evaluation of CAs. According to quantitative survey of experimental evaluation in computer science, scholarly publications in ACM could be classified into five categories, including formal theory, design and modeling, empirical work, hypothesis testing, and others (e.g., surveys) [Tichy et al., 1995]. The publications we selected could be a combination of design and modeling (i.e., "systems, techniques, or models, whose claimed properties cannot be proven formally") and empirical work (i.e., "articles that collect, analyze, and interpret observations about known designs, systems, or models, or about abstract theories or subjects [Wainer and Barsottini, 2007]"). Even though we tuned the search query to the best we could during web-crawling, it is possible that certain works were not captured by the terms we defined due to our fields of expertise. Meanwhile, it is expected that a significant amount of works that proposed novel techniques without UX evaluations were not included in our analysis. Therefore, we do not claim that the findings are exhaustive. Even though our findings build upon significant prior literature review results, the themes of *polyadic* challenges and evaluation metrics can vary and will expand with more works identified.

Lastly, the thematic analysis process is subjective, and lemmatization [Nguyen and Shirai, 2015] and topic modeling [Gurcan et al., 2021] might introduce potential biases as well. For example, lemmatization reduces words to their lemma forms, the use of lemmatization could potentially cause a loss of words' meaning at a certain level. For topic modeling, since we used text from the original papers without filtering irrelevant words, the results contained noises and irrelevant words that sometimes interfered the interpretability. Therefore, we cannot claim differences to be statistically significant. For example, the topic modeling results were discussed in the context of specific examples to elicit more discussions and future research.

7 Conclusion

This literature review presents what fundamental human-human interaction challenges are addressed by *polyadic* CA works scrutinized from ACM publications, what effects on human-human interaction are evaluated by these works, and what issues are overlooked when designing *polyadic* CAs. In particular, we propose that researchers and practitioners should design CAs with *boundary awareness* for supporting *polyadic* human-AI interaction. Further, we envision several research opportunities on conversational human-AI interaction with respect to the theoretical groundings, relational dynamics, and functional dimensions.

8 Acknowledgement

This material is based upon work supported by the National Science Foundation under Grant No. 2119589. The research is also partially supported by the IBM-ILLINOIS Center for Cognitive Computing Systems Research (C3SR), a research collaboration as part of the IBM AI Horizons Network.

A APPENDIX

A.1 Definitions of Different Metrics Used in *Polyadic* Papers

Category	Metrics	Definition
Sociability	discussion quality	the quality of users' discussions
	consensus reaching	the reaching of a consensus in terms of behavioral and perceived opinion alignment
	opinion expression	users' message quantity, even participation, and perceived outspokenness
	opinion diversity	the number of unique lexical morphemes shared within a group
	even participation	how equally individual members contribute to the discussion.
	linguistic features	the linguistic features extracted from users' utterance text, e.g. verbosity, readability, sentiment, linguistic diversity and adaptability
Task-related	task completion time	the time elapsed before users successfully completed the intended task(s).
	efficiency	the number of system turns required and average system reaction time to users' requests.
	effectiveness	users' successes in the tasks undertaken with the CA support.
	satisfaction	users' overall ratings for CA's understandability, relevancy and efficiency
	team performance	the team's success in the tasks undertaken with the CA's support
Communication quality	perceived communication effectiveness	the degree of how much the CA helps users reach the goal
	perceived communication fairness	the degree of participation fairness that users perceive during discussions
	perceived communication efficiency	how quickly the consensus is reached during discussions
	perceived group climate	the relatively enduring tone and quality of group interaction experienced similarly by group members
Socio-emotional	perceptions of other group members	users' opinion and perception about other members in group discussions
	perception towards the agent	users' opinion and perceptions about the CA in discussion
	perceptions about the collaboration	users' collaboration experience in terms of users' awareness of each other's activities, effort, and enjoyment.
	perceived social presence	users' experience of being present with other persons and having access to their thoughts and emotions
	perceived social support	users' experience of being provided with support from other persons during discussion
	level of annoyance with the intervention	users' perceived annoyance when receiving the CA's intervention
	opinion of oneself/partner/cultural stereotypes	users' opinion about self, partner, and cultural stereotypes under the CA's influence during discussion
	perceived appropriateness of CA's social behavior	users' perception about the CA's ability to behave appropriately like a human
Traditional	psychological wellbeing	users' psychological wellbeing measured by positive relationships with others, personal mastery, autonomy, a feeling of purpose and meaning in life, and personal growth and development
	anthropomorphism	the extent to which the CA can demonstrate attribution of human characteristics or behavior
	perceived emotional competence	users' perception of the CA's emotional skills
	unmet expectations	users' expectations that are not met by the CA during the study
	privacy concerns	users' concerns about the safeguarding and usage of personal data provided to the CA
	level of perceived conversation control	users' feeling in control of the discussion using the CA
	perceived satisfaction	users' overall ratings for CA's understandability, relevancy and efficiency
	performance/effectiveness	the measure of user's success in the tasks undertaken in the CA interactions
Issue-relevant	level of engagement	users' level of engagement with the CA, e.g., the number of interactions measured by clicks and selections
	perceived capability to promote contributions	users' perception about the CA's ability to elicit contributions and opinions from participants during discussion
	decisions to (not) engage with the CA	users' decision of whether users would like to engage in interactions with the CA
	overall UX	users' perception of the overall user experience during interactions with the CA
	reflections on selves and others	users' reflections on behaviors and feelings of self or other participants
	willingness to take actions	users' willingness to take certain actions under the persuasion of the CA
	direction of improvement	the possible directions of improvements proposed by users' after interacting with the CA

A.2 Definitions of Different Metrics Used in *Dyadic* Papers

Category	Metrics	Definition
Linguistic features	positive/negative words	number of users' use of positive/negative words to express emotions
	length of utterances	users' number of users' use of words per statement
	personal pronouns	users' use of first- and third-person pronouns
	term-level repetition	the proportion of terms in one utterance which were repeated from the participant's previous utterance
	utterance-level repetition	the proportion of utterances where term-level repetition is greater than zero
	use of concrete words	the concreteness of each entry, e.g., "tennis" is more concrete than "sports"
	variation in language used	the variation in expressions and use of different words
Prosodic features	variation in speech-based agents	the word count, variance in pitch, rate and loudness in audio interactions
Dialogue features	dialogue/response quality	the syntactic readability and intensity of sentiments in users' replies
	dialogue efficiency metrics	the number of system turns required and average system reaction time to users' requests
	expressiveness	the quality of effectively conveying thoughts or feelings by users
Tasks related	task fulfillment rate/effectiveness	the measure of users' success rate in tasks undertaken during the interactions
	system use	a mixed measure of multiple system-related metrics, e.g., types of messages, words, time taken to complete the task
	user entry time	the amount of time users need to provide inputs to or interact verbally with the CA
	dropout rate	the percentage of respondents who quit before the study was completed
Algorithm	intent modeling evaluation	the accuracy of which the CA can identify the correct intents from users' utterances
	error rate	the rate of errors occurred during users' interaction with the CA
User	self-disclosure	users' quality of responses to the CA based on their trust towards the CA
	intimacy	users' perceive closeness, inter-connectedness, and companionship from the CA
	reflection	users' self-rated frequency of reflecting on thoughts and consciousness of their inner feelings
	impolite/aggressive behaviors	the number of occurrences of impolite phrases
	motivation	users' motivations or intents that drive users to interact with the CA
	affective states	the proportion of time users spent in each affective state

Category	Metrics	Definition
Conversation related	syntactic readability	the syntax-wise readability of conversation text between users and CA based on the Flesch-readability score
	self-reported response quality	the perceived response quality reported by users
	informativeness/information quality	the quality of information conveyed by the CA through text
	conversation smoothness	the level of smoothness of the conversation between users and CAs measure by Session Evaluation Questionnaire
	intensity of sentiments	the intensity of sentiments expressed by users calculated by TextBlob
	instruction quality	the quality of instructions given by the CA to users
Perception towards the CA	perceived common ground	the perceived mutual understanding between users and the CA
	opinion of the CA as a partner	the ease of which subjects were able to interact with the CA
	playfulness/enjoyment	the extent to which the CA can match users' interests
	sociability	users' perception of the CA's social skills
	perceived anthropomorphism	the extent to which the CA can demonstrate attribution of human characteristics or behaviors
	perceived warmth of the AI system	the extent to which users feel the CA is good-natured and warm
	impression	users' feelings of the CA regarding its competence, confidence, warmth, and sincerity
	perceived personality traits	users' feelings of the CA in terms of openness, conscientiousness, extraversion, agreeableness, and neuroticism
	perceived persuasiveness	users' perception of the CA's utterances rated as bad-good, foolish-wise, negative-positive, beneficial-harmful, effective-ineffective, and convincing-unconvincing
	perceived intimacy	users' perceived intimacy and closeness with the CA, e.g., feeling of inter-connectedness of self and other
	perceived naturalness	users' agreement level to the statement that the CA's responses are natural
	perceived safeness	users' sense of safety when interacting with the CA
System usability	overall usability	users' perceived overall usability of the CA
	perceived ease of use	the degree to which users believe that interacting with the CA would be free of effort
	mental demand	the amount of mental or perceptual activities (e.g., thinking, deciding, calculating, remembering, looking, searching) that is required to interact with the CA
	physical demand	the amount of physical activities (e.g., pushing, pulling, controlling, activating, etc.) that is required to interact with the CA
	perceived task completion	the degree to which users believe to have successfully communicated and reached a mutual understanding with the CA
	effort	the total workload associated with the tasks, considering all sources and components
	frustration	users' feelings of being insecure, discouraged, irritated, and annoyed versus being secure, gratified, content and complacent when interacting with the CA
	desire to continue the interaction	the degree to which users would consider keep using this method in the future, or users' behavioural tendencies through their desire to help and cooperate with the CA
	user acceptance of the agent	users' willingness to accept the CA's interaction
	system consistency	the consistency between the behaviors and utterances of the CA
	perceived usefulness / helpfulness	users' own perceptions of the session's efficacy, e.g., the CA gave users good suggestions for helping them discover songs.
	willingness to recommend	the degree to which users would recommend the CA to their friends or family for managing mental well-being and to people who have needs
	motivation	users' motivation or intent to interact with the CA

Category	Metrics	Definition
User experience with the CA	user experience (UX)	a mix ratings of the general experience, e.g., emotion, ease of use, usefulness, and intention to use
	contrast of user experience	users' perceptions of the experience without drawing explicit attention to the contrast between their expectations and their experience
	perceived quality of the interaction	users' overall self-rated quality of the CA, e.g., in communicating, building rapport, and task fulfillment's
	perceived satisfaction	users' overall ratings for CA's understandability, relevancy and efficiency
	perceived engagement	the degree that the CA can engage the participants during the conversation, e.g., making users feel it is entertaining and interesting to engage in a dialogue with the CA
	perceived trust	the CA's ability in providing unbiased and accurate suggestions and making users trust it
	confidence	users' confidence that users will like the content the CA suggests
	perceived capability to control	users' feeling of being in control of the conversation
	aroused emotions	change of users' emotions state while using the system
	serendipity	the CA's ability in recommending things that users had not considered in the first place but turned out to be a positive and surprising discovery
	pleasant surprise	the CA's ability in providing contents that are overall pleasantly surprising to users
	empathy	the CA's ability to understand and share the feelings of users
	self-reflection/self-awareness	users' reflection on thoughts and consciousness of their inner feelings
	self-disclose	the degree of which users feel comfortable about disclosing to the agent and express opinions openly
	anxiety/tension	the degree of which the interaction makes users feel anxious or tense
	comfort	the degree of which the interaction is comfortable
	social orientation toward CAs	the desire to engage in human-like social interactions with CA, which is associated with a mental model of an agent system as being a sociable entity
	degree of collaboration	the level of collaboration between users and the CA during interaction
	degree of decision-making	the degree of which the CA helps with users' decision-making
Perceptions	overall impressions and experience	users' perceptions of the overall experience interacting with the CA
	perceived CA's capabilities	the CA's capability in handling simple requests and resembling human representative
	perceived CA appearance and self-presentation	the CA's visual appearance and persona
	perceived personality	users' feeling for the CA in openness (intellectual curiosity, creativity), conscientiousness (neatness, perseverance, reliability, and responsibility), extraversion (sociability, activity, and assertiveness), agreeableness (friendliness, helpfulness, and cooperativeness in dealing with others) and neuroticism (stability, anxiety, and the frequency of experiencing negative affect)
	perceived naturalness	the degree of which users feel the conversation with the CA is natural, not forced
	reciprocative	the degree of which users feel the CA reciprocated their language or feelings
	perceived burden	the degree of which users feel the conversation with the CA is costly in time, financially, mentally and emotionally.
	perceived human-likeness	the CA's ability to talk like a human and its conversational skills
	perceived effectiveness	the degree of how much the CA helps to address users' needs
Behaviors	engagement	the degree that the CA can engage the participants during the conversation, e.g., measured by an engagement rating between 1 (not engaging) and 5 (very engaging) in user survey
	notice of CA's certain function	the action of the CA sending messages informing what the chatbot could do, noticing a tutorial and menu
	efforts in learning the system	efforts users take in learning how to interact with the CA
	self-reflection/self-awareness	users' reflections on thoughts and consciousness of their inner feelings
	effects on judgement	the degree that users feel the CA affected their evaluation positively, negatively or neither
	effects on collaborations	the degree of which the CA affected users' willingness of collaborating with the CA
	daily practices of using the CA	participants' daily practices of using the chatbot, e.g., "Please briefly tell us how you used this chatbot during the past three weeks"

References

- Eleni Adamopoulou and Lefteris Moussiades. An overview of chatbot technology. In *IFIP International Conference on Artificial Intelligence Applications and Innovations*, pages 373–383. Springer, 2020a.
- David Griol, Javier Carbó, and José M Molina. An automatic dialog simulation technique to develop and evaluate interactive conversational agents. *Applied Artificial Intelligence*, 27(9):759–780, 2013.
- Eleni Adamopoulou and Lefteris Moussiades. Chatbots: History, technology, and applications. *Machine Learning with Applications*, 2:100006, 2020b.
- Amon Rapp, Lorenzo Curti, and Arianna Boldi. The human side of human-chatbot interaction: A systematic literature review of ten years of research on text-based chatbots. *International Journal of Human-Computer Studies*, page 102630, 2021.
- Michael McTear. Conversational ai: Dialogue systems, conversational agents, and chatbots. *Synthesis Lectures on Human Language Technologies*, 13(3):1–251, 2020.
- Jesse Fox and Andrew Gambino. Relationship development with humanoid social robots: Applying interpersonal theories to human/robot interaction. *Cyberpsychology, Behavior, and Social Networking*, 2021.
- Andrea L Guzman and Seth C Lewis. Artificial intelligence and communication: A human-machine communication research agenda. *New Media & Society*, 22(1):70–86, 2020.
- Meredith Whittaker, Kate Crawford, Roel Dobbe, Genevieve Fried, Elizabeth Kaziunas, Varoon Mathur, Sarah Mysers West, Rashida Richardson, Jason Schultz, and Oscar Schwartz. *AI now report 2018*. AI Now Institute at New York University New York, 2018.
- Timothy W Bickmore, Lisa Caruso, and Kerri Clough-Gorr. Acceptance and usability of a relational agent interface by urban older adults. In *CHI’05 extended abstracts on Human factors in computing systems*, pages 1212–1215, 2005.
- Daniel Schulman and Timothy Bickmore. Persuading users through counseling dialogue with a conversational agent. In *Proceedings of the 4th international conference on persuasive technology*, pages 1–8, 2009.
- Anbang Xu, Zhe Liu, Yufan Guo, Vibha Sinha, and Rama Akkiraju. A new chatbot for customer service on social media. In *Proceedings of the 2017 CHI conference on human factors in computing systems*, pages 3506–3510, 2017.
- Stefan Kopp, Lars Gesellensetter, Nicole C Krämer, and Ipke Wachsmuth. A conversational agent as museum guide—design and evaluation of a real-world application. In *International workshop on intelligent virtual agents*, pages 329–343. Springer, 2005.
- Mahoro Anabuki, Hiroyuki Kakuta, Hiroyuki Yamamoto, and Hideyuki Tamura. Welbo: An embodied conversational agent living in mixed reality space. In *CHI’00 extended abstracts on Human factors in computing systems*, pages 10–11, 2000.
- Soomin Kim, Jinsu Eun, Joseph Seering, and Joonhwan Lee. Moderator chatbot for deliberative discussion: Effects of discussion structure and discussant facilitation. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–26, 2021.
- Qiaosi Wang, Koustuv Saha, Eric Gregori, David Joyner, and Ashok Goel. Towards mutual theory of mind in human-ai interaction: How language reflects what students perceive about a virtual teaching assistant. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2021.
- Soomin Kim, Jinsu Eun, Changhoon Oh, Bongwon Suh, and Joonhwan Lee. Bot in the bunch: Facilitating group chat discussion by improving efficiency and participation with a chatbot. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2020.
- Carlos Toxtli, Andrés Monroy-Hernández, and Justin Cranshaw. Understanding chatbot-mediated task management. In *Proceedings of the 2018 CHI conference on human factors in computing systems*, pages 1–6, 2018.
- Ivo Benke, Michael Thomas Knierim, and Alexander Maedche. Chatbot-based emotion management for distributed teams: A participatory design study. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW2):1–30, 2020.
- Joseph Seering, Michal Luria, Geoff Kaufman, and Jessica Hammer. Beyond dyadic interactions: Considering chatbots as community members. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2019.
- Ana Paula Chaves and Marco Aurelio Gerosa. How should my chatbot interact? a survey on social characteristics in human-chatbot interaction design. *International Journal of Human-Computer Interaction*, 37(8):729–758, 2021.

- Allan de Barcelos Silva, Marcio Miguel Gomes, Cristiano André da Costa, Rodrigo da Rosa Righi, Jorge Luis Victoria Barbosa, Gustavo Pessin, Geert De Doncker, and Gustavo Federizzi. Intelligent personal assistants: A systematic literature review. *Expert Systems with Applications*, 147:113193, 2020.
- Joao Luis Zeni Montenegro, Cristiano André da Costa, and Rodrigo da Rosa Righi. Survey of conversational agents in health. *Expert Systems with Applications*, 129:56–67, 2019.
- Liliana Laranjo, Adam G Dunn, Huong Ly Tong, Ahmet Baki Kocaballi, Jessica Chen, Rabia Bashir, Didi Surian, Blanca Gallego, Farah Magrabi, Annie YS Lau, et al. Conversational agents in healthcare: a systematic review. *Journal of the American Medical Informatics Association*, 25(9):1248–1258, 2018.
- Jess Hohenstein and Malte Jung. Ai as a moral crumple zone: The effects of ai-mediated communication on attribution and trust. *Computers in Human Behavior*, 106:106190, 2020.
- Teun A Van Dijk. *Discourse as structure and process*, volume 1. Sage, 1997.
- Michael Frederick McTear, Zoraida Callejas, and David Griol. *The conversational interface*, volume 6. Springer, 2016.
- Jacky Casas, Marc-Olivier Tricot, Omar Abou Khaled, Elena Mugellini, and Philippe Cudré-Mauroux. Trends & methods in chatbot evaluation. In *Companion Publication of the 2020 International Conference on Multimodal Interaction*, pages 280–286, 2020.
- Shafquat Hussain, Omid Ameri Sianaki, and Nedat Ababneh. A survey on conversational agents/chatbots classification and design techniques. In *Workshops of the International Conference on Advanced Information Networking and Applications*, pages 946–956. Springer, 2019.
- Eva Bittner, Sarah Oeste-Reiß, and Jan Marco Leimeister. Where is the bot in our team? toward a taxonomy of design option combinations for conversational agents in collaborative work. In *Proceedings of the 52nd Hawaii international conference on system sciences*, 2019.
- Morten Johan Mygland, Morten Schibbye, Ilias O Pappas, and Polyxeni Vassilakopoulou. Affordances in human-chatbot interaction: a review of the literature. In *Conference on e-Business, e-Services and e-Society*, pages 3–17. Springer, 2021.
- Jianfeng Gao, Michel Galley, and Lihong Li. Neural approaches to conversational ai. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 1371–1374, 2018.
- Asbjørn Følstad, Marita Skjuve, and Petter Bae Brandtzaeg. Different chatbots for different purposes: towards a typology of chatbots to understand interaction design. In *International Conference on Internet Science*, pages 145–156. Springer, 2018.
- Justine Cassell. Embodied conversational interface agents. *Communications of the ACM*, 43(4):70–78, 2000.
- Lorenz Cuno Klopfenstein, Saverio Delpriori, Silvia Malatini, and Alessandro Bogliolo. The rise of bots: A survey of conversational interfaces, patterns, and paradigms. In *Proceedings of the 2017 conference on designing interactive systems*, pages 555–565, 2017.
- Sherry Ruan, Jiayu He, Rui Ying, Jonathan Burkle, Dunia Hakim, Anna Wang, Yufeng Yin, Lily Zhou, Qian Yao Xu, Abdallah AbuHashem, et al. Supporting children’s math learning with feedback-augmented narrative technology. In *Proceedings of the Interaction Design and Children Conference*, pages 567–580, 2020.
- Marita Skjuve, Ida Maria Haugstveit, Asbjørn Følstad, and Petter Bae Brandtzaeg. Help! is my chatbot falling into the uncanny valley? an expirical study of user experience in human-chatbot interaction. *Human Technology*, 15(1), 2019.
- Martin Adam, Michael Wessel, and Alexander Benlian. Ai-based chatbots in customer service and their effects on user compliance. *Electronic Markets*, pages 1–19, 2020.
- Justin Cranshaw, Emad Elwany, Todd Newman, Rafal Kocielnik, Bowen Yu, Sandeep Soni, Jaime Teevan, and Andrés Monroy-Hernández. Calendar. help: Designing a workflow-based scheduling agent with humans in the loop. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, pages 2382–2393, 2017.
- Minha Lee, Sander Ackermans, Nena van As, Hanwen Chang, Enzo Lucas, and Wijnand IJsselstein. Caring for vincent: a chatbot for self-compassion. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2019.
- Kyle-Althea Santos, Ethel Ong, and Ron Resurreccion. Therapist vibe: children’s expressions of their emotions through storytelling with a chatbot. In *Proceedings of the Interaction Design and Children Conference*, pages 483–494, 2020.
- Ana Paula Chaves and Marco Aurelio Gerosa. How should my chatbot interact? a survey on social characteristics in human–chatbot interaction design. *International Journal of Human–Computer Interaction*, pages 1–30, 2020.
- Endang Wahyu Pamungkas. Emotionally-aware chatbots: A survey. *arXiv preprint arXiv:1906.09774*, 2019.

- Ahmet Baki Kocaballi, Shlomo Berkovsky, Juan C Quiroz, Liliana Laranjo, Huong Ly Tong, Dana Rezazadegan, Agustina Briatore, and Enrico Coiera. The personalization of conversational agents in health care: systematic review. *Journal of medical Internet research*, 21(11):e15360, 2019.
- Minjin Rheu, Ji Youn Shin, Wei Peng, and Jina Huh-Yoo. Systematic review: Trust-building factors and implications for conversational agent design. *International Journal of Human-Computer Interaction*, 37(1):81–96, 2021.
- Michelle ME Van Pinxteren, Mark Pluymaekers, and Jos GAM Lemmink. Human-like communication in conversational agents: a literature review and research agenda. *Journal of Service Management*, 2020.
- Alex Sciuto, Arnita Saini, Jodi Forlizzi, and Jason I Hong. "hey alexa, what's up?" a mixed-methods studies of in-home conversational agent usage. In *Proceedings of the 2018 Designing Interactive Systems Conference*, pages 857–868, 2018.
- Pragnesh Jay Modi, Manuela Veloso, Stephen F Smith, and Jean Oh. Cmradar: A personal assistant agent for calendar management. In *International Bi-Conference Workshop on Agent-Oriented Information Systems*, pages 169–181. Springer, 2004.
- Anders Lundkvist and Ali Yakhlef. Customer involvement in new service development: a conversational approach. *Managing Service Quality: An International Journal*, 2004.
- Hideo Shimazu. Expertclerk: a conversational case-based reasoning tool for developing salesclerk agents in e-commerce webshops. *Artificial Intelligence Review*, 18(3):223–244, 2002.
- Timothy W Bickmore, Suzanne E Mitchell, Brian W Jack, Michael K Paasche-Orlow, Laura M Pfeifer, and Julie O'Donnell. Response to a relational agent by hospital patients with depressive symptoms. *Interacting with computers*, 22(4):289–298, 2010.
- Timothy W Bickmore, Laura M Pfeifer, and Brian W Jack. Taking the time to care: empowering low health literacy hospital patients with virtual nurse agents. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 1265–1274, 2009.
- Jorne Grolleman, Betsy van Dijk, Anton Nijholt, and Andrée van Emst. Break the habit! designing an e-therapy intervention using a virtual coach in aid of smoking cessation. In *International Conference on Persuasive Technology*, pages 133–141. Springer, 2006.
- Zhenhui Peng, Taewook Kim, and Xiaojuan Ma. Gremobot: Exploring emotion regulation in group chat. In *Conference Companion Publication of the 2019 on Computer Supported Cooperative Work and Social Computing*, pages 335–340, 2019.
- Jess Hohenstein and Malte Jung. Ai-supported messaging: An investigation of human-human text conversation with ai support. In *Proceedings of the Extended Abstracts of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–6, 2018.
- Piotr Chynał, Julia Falkowska, and Janusz Sobecki. Human-human interaction: A neglected field of study? In *International Conference on Intelligent Human Systems Integration*, pages 346–351. Springer, 2018.
- Asbjørn Følstad, Theo Araujo, Effie Lai-Chong Law, Petter Bae Brandtzaeg, Symeon Papadopoulos, Lea Reis, Marcos Baez, Guy Laban, Patrick McAllister, Carolin Ischen, et al. Future directions for chatbot research: an interdisciplinary research agenda. *Computing*, pages 1–28, 2021.
- Asbjørn Følstad and Petter Bae Brandtzaeg. Chatbots and the new world of hci. *interactions*, 24(4):38–42, 2017.
- Yi-Chieh Lee, Naomi Yamashita, and Yun Huang. Exploring the effects of incorporating human experts to deliver journaling guidance through a chatbot. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–27, 2021.
- Francisco Nunes, Nervo Verdezoto, Geraldine Fitzpatrick, Morten Kyng, Erik Grönvall, and Cristiano Storni. Self-care technologies in hci: Trends, tensions, and opportunities. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 22(6):1–45, 2015.
- Eleonora Mencarini, Amon Rapp, Lia Tirabeni, and Massimo Zancanaro. Designing wearable systems for sports: A review of trends and opportunities in human-computer interaction. *IEEE Transactions on Human-Machine Systems*, 49(4):314–325, 2019.
- Dan Bohus and Eric Horvitz. Multiparty turn taking in situated dialog: Study, lessons, and directions. In *Proceedings of the SIGDIAL 2011 Conference*, pages 98–109, 2011.
- Paul Ralph and Romain Robbes. The acm sigsoft paper and peer review quality initiative: Status report. *ACM SIGSOFT Software Engineering Notes*, 45(2):17–18, 2020.

- Jacques Wainer, Claudia G Novoa Barsottini, Danilo Lacerda, and Leandro Rodrigues Magalhães de Marco. Empirical evaluation in computer science research published by acm. *Information and Software Technology*, 51(6):1081–1085, 2009.
- Silke ter Stal, Lean Leonie Kramer, Monique Tabak, Harm op den Akker, and Hermie Hermens. Design features of embodied conversational agents in ehealth: a literature review. *International Journal of Human-Computer Studies*, 138:102409, 2020.
- Ruikai Liu Jorj X. McKie. Pymupdf. <https://github.com/pymupdf/PyMuPDF>, 2016.
- Gilchan Park and Line Pouchard. Scientific literature mining for experiment information in materials design. In *2019 New York Scientific Data Summit (NYSDS)*, pages 1–4. IEEE, 2019.
- Huichen Yang, Carlos A Aguirre, F Maria, Derek Christensen, Luis Bobadilla, Emily Davich, Jordan Roth, Lei Luo, Yihong Theis, Alice Lam, et al. Pipelines for procedural information extraction from scientific literature: towards recipes using machine learning and data science. In *2019 International conference on document analysis and recognition workshops (ICDARW)*, volume 2, pages 41–46. IEEE, 2019.
- Gilchan Park, Julia Taylor Rayz, and Line Pouchard. Figure descriptive text extraction using ontological representation. In *The Thirty-Third International Flairs Conference*, 2020.
- Huilin Chang, Yihnew Eshetu, and Celeste Lemrow. Supervised machine learning and deep learning classification techniques to identify scholarly and research content. In *2021 Systems and Information Engineering Design Symposium (SIEDS)*, pages 1–6. IEEE, 2021.
- Joost F Wolfswinkel, Elfi Furtmueller, and Celeste PM Wilderom. Using grounded theory as a method for rigorously reviewing literature. *European journal of information systems*, 22(1):45–55, 2013.
- Virginia Braun and Victoria Clarke. Using thematic analysis in psychology. *Qualitative research in psychology*, 3(2): 77–101, 2006.
- Nora McDonald, Sarita Schoenebeck, and Andrea Forte. Reliability and inter-rater reliability in qualitative research: Norms and guidelines for cscw and hci practice. *Proceedings of the ACM on Human-Computer Interaction*, 3 (CSCW):1–23, 2019.
- Laura R Pina, Carmen Gonzalez, Carolina Nieto, Wendy Roldan, Edgar Onofre, and Jason C Yip. How latino children in the us engage in collaborative online information problem solving with their families. *Proceedings of the ACM on Human-Computer Interaction*, 2(CSCW):1–26, 2018.
- Kilem L Gwet. *Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters*. Advanced Analytics, LLC, 2014.
- Fatih Gurcan, Nergiz Ercil Cagiltay, and Kursat Cagiltay. Mapping human–computer interaction research themes and trends from its existence to today: A topic modeling-based review of past 60 years. *International Journal of Human–Computer Interaction*, 37(3):267–280, 2021.
- Fatih Gurcan and Nergiz Ercil Cagiltay. Research trends on distance learning: a text mining-based literature review from 2008 to 2018. *Interactive Learning Environments*, pages 1–22, 2020.
- Lijun Sun and Yafeng Yin. Discovering themes and trends in transportation research using topic modeling. *Transportation Research Part C: Emerging Technologies*, 77:49–66, 2017.
- Amy X Zhang and Justin Cranshaw. Making sense of group chat through collaborative tagging and summarization. *Proceedings of the ACM on Human-Computer Interaction*, 2(CSCW):1–27, 2018.
- Dan Bohus and Eric Horvitz. Dialog in the open world: platform and applications. In *Proceedings of the 2009 international conference on Multimodal interfaces*, pages 31–38, 2009a.
- Dan Bohus and Eric Horvitz. Learning to predict engagement with a spoken dialog system in open-world settings. In *Proceedings of the SIGDIAL 2009 Conference*, pages 244–252, 2009b.
- Dan Bohus and Eric Horvitz. Facilitating multiparty dialog with gaze, gesture, and speech. In *International Conference on Multimodal Interfaces and the Workshop on Machine Learning for Multimodal Interaction*, pages 1–8, 2010.
- Sandeep Avula, Gordon Chadwick, Jaime Arguello, and Robert Capra. Searchbots: User engagement with chatbots during collaborative search. In *Proceedings of the 2018 conference on human information interaction & retrieval*, pages 52–61, 2018.
- Kohji Dohsaka, Ryota Asai, Ryuichiro Higashinaka, Yasuhiro Minami, and Eisaku Maeda. Effects of conversational agents on human communication in thought-evoking multi-party dialogues. In *Proceedings of the SIGDIAL 2009 Conference*, pages 217–224, 2009.

- Gregory Dyke, Iris Howley, David Adamson, Rohit Kumar, and Carolyn Penstein Rosé. Towards academically productive talk supported by conversational agents. In *Productive multivocality in the analysis of group interactions*, pages 459–476. Springer, 2013.
- Stergios Tegos, Stavros Demetriadis, and Anastasios Karakostas. Promoting academically productive talk with conversational agent interventions in collaborative learning settings. *Computers & Education*, 87:309–325, 2015.
- Hua Ai, Rohit Kumar, Dong Nguyen, Amrut Nagasunder, and Carolyn P Rosé. Exploring the effectiveness of social capabilities and goal alignment in computer supported collaborative learning. In *International Conference on Intelligent Tutoring Systems*, pages 134–143. Springer, 2010.
- Sourish Chaudhuri, Rohit Kumar, Iris K Howley, and Carolyn Penstein Rosé. Engaging collaborative learners with helping agents. In *AIED*, pages 365–372, 2009.
- Rohit Kumar, Hua Ai, Jack L Beuth, and Carolyn P Rosé. Socially capable conversational tutors can be effective in collaborative learning situations. In *International conference on intelligent tutoring systems*, pages 156–164. Springer, 2010.
- Rohit Kumar and Carolyn P Rosé. Triggering effective social support for online groups. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 3(4):1–32, 2014.
- Saiph Savage, Andres Monroy-Hernandez, and Tobias Höllerer. Botivist: Calling volunteers to action using online bots. In *Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing*, pages 813–822, 2016.
- Joseph Seering, Juan Pablo Flores, Saiph Savage, and Jessica Hammer. The social roles of bots: evaluating impact of bots on discussions in online communities. *Proceedings of the ACM on Human-Computer Interaction*, 2(CSCW): 1–29, 2018.
- Joseph Seering, Michal Luria, Connie Ye, Geoff Kaufman, and Jessica Hammer. It takes a village: Integrating an adaptive chatbot into an online gaming community. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2020.
- Michal Luria, Joseph Seering, Jodi Forlizzi, and John Zimmerman. Designing chatbots as community-owned agents. In *Proceedings of the 2nd Conference on Conversational User Interfaces*, pages 1–3, 2020a.
- Norah Abokhodair, Daisy Yoo, and David W McDonald. Dissecting a social botnet: Growth, content and influence in twitter. In *Proceedings of the 18th ACM conference on computer supported cooperative work & social computing*, pages 839–851, 2015.
- Ameneh Shamekhi, Q Vera Liao, Dakuo Wang, Rachel KE Bellamy, and Thomas Erickson. Face value? exploring the effects of embodiment for a group facilitation agent. In *Proceedings of the 2018 CHI conference on human factors in computing systems*, pages 1–13, 2018.
- Shochi Otogi, Hung-Hsuan Huang, Ryo Hotta, and Kyoji Kawagoe. Finding the timings for a guide agent to intervene in inter-user conversation in considering their gaze behaviors. In *Proceedings of the 6th workshop on Eye gaze in intelligent human machine interaction: gaze in multimodal interaction*, pages 19–24, 2013.
- Shochi Otogi, Hung-Hsuan Huang, Ryo Hotta, and Kyoji Kawagoe. Analysis of personality traits for intervention scene detection in multi-user conversation. In *Proceedings of the second international conference on Human-agent interaction*, pages 237–240, 2014.
- Jun Zheng, Xiang Yuan, and Yam San Chee. Designing multiparty interaction support in elva, an embodied tour guide. In *Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems*, pages 929–936, 2005.
- Xiang Yuan and Yam San Chee. Design and evaluation of elva: an embodied tour guide in an interactive virtual art gallery. *Computer animation and virtual worlds*, 16(2):109–119, 2005.
- Jaya Narain, Tina Quach, Monique Davey, Hae Won Park, Cynthia Breazeal, and Rosalind Picard. Promoting wellbeing with sunny, a chatbot that facilitates positive messages within social groups. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–8, 2020.
- Sarah Strohkorb Sebo, Margaret Traeger, Malte Jung, and Brian Scassellati. The ripple effects of vulnerability: The effects of a robot’s vulnerable behavior on trust in human-robot teams. In *Proceedings of the 2018 ACM/IEEE International Conference on Human-Robot Interaction*, pages 178–186, 2018.
- Michal Luria, Rebecca Zheng, Bennett Huffman, Shuangni Huang, John Zimmerman, and Jodi Forlizzi. Social boundaries for personal agents in the interpersonal space of the home. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2020b.

- Christian Löw and Lukas Moshuber. Grätzelbot-gamifying onboarding to support community-building among university freshmen. In *Proceedings of the 11th Nordic Conference on Human-Computer Interaction: Shaping Experiences, Shaping Society*, pages 1–3, 2020.
- Hideyuki Nakanishi, Satoshi Nakazawa, Toru Ishida, Katsuya Takanashi, and Katherine Isbister. Can software agents influence human relations? balance theory in agent-mediated communities. In *Proceedings of the second international joint conference on autonomous agents and multiagent systems*, pages 717–724, 2003.
- Katherine Isbister, Hideyuki Nakanishi, Toru Ishida, and Cliff Nass. Helper agent: Designing an assistant for human-human interaction in a virtual meeting space. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 57–64, 2000.
- Joyce Chai, Veronika Horvath, Nanda Kambhatla, Nicolas Nicolov, and Margo Stys-Budzikowska. A conversational interface for online shopping. In *Proceedings of the First International Conference on Human Language Technology Research*, 2001.
- Indrani Medhi, Somani Patnaik, Emma Brunskill, SN Nagasena Gautama, William Thies, and Kentaro Toyama. Designing mobile interfaces for novice and low-literacy users. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 18(1):1–28, 2011.
- Mohit Jain, Pratyush Kumar, Ramachandra Kota, and Shwetak N Patel. Evaluating and informing the design of chatbots. In *Proceedings of the 2018 Designing Interactive Systems Conference*, pages 895–906, 2018a.
- Mohit Jain, Ramachandra Kota, Pratyush Kumar, and Shwetak N Patel. Convey: Exploring the use of a context view for chatbots. In *Proceedings of the 2018 chi conference on human factors in computing systems*, pages 1–6, 2018b.
- Tessa Lau, Julian Cerruti, Guillermo Manzato, Mateo Bengualid, Jeffrey P Bigham, and Jeffrey Nichols. A conversational interface to web automation. In *Proceedings of the 23rd annual ACM symposium on User interface software and technology*, pages 229–238, 2010.
- Christine Lisetti, Reza Amini, Ugan Yasavur, and Naphtali Rishe. I can help you change! an empathic virtual agent delivers behavior change health interventions. *ACM Transactions on Management Information Systems (TMIS)*, 4(4): 1–28, 2013.
- Q Vera Liao, Muhammed Mas-ud Hussain, Praveen Chandar, Matthew Davis, Yasaman Khazaeni, Marco Patricio Crasso, Dakuo Wang, Michael Muller, N Sadat Shami, and Werner Geyer. All work and no play? In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2018.
- Silvia B Lovato, Anne Marie Piper, and Ellen A Wartella. Hey google, do unicorns exist? conversational agents as a path to answers to children’s questions. In *Proceedings of the 18th ACM International Conference on Interaction Design and Children*, pages 301–313, 2019.
- Thomas N Smyth, John Etherton, and Michael L Best. Moses: Exploring new ground in media and post-conflict reconciliation. In *Proceedings of the SIGCHI conference on Human Factors in computing systems*, pages 1059–1068, 2010.
- Soomin Kim, Joonhwan Lee, and Gahgene Gweon. Comparing data from chatbot and web surveys: Effects of platform and conversational style on survey response quality. In *Proceedings of the 2019 CHI conference on human factors in computing systems*, pages 1–12, 2019.
- Alexandra Vtyurina, Denis Savenkov, Eugene Agichtein, and Charles LA Clarke. Exploring conversational search with humans, assistants, and wizards. In *Proceedings of the 2017 chi conference extended abstracts on human factors in computing systems*, pages 2187–2193, 2017.
- Rafal Kocielnik, Daniel Avrahami, Jennifer Marlow, Di Lu, and Gary Hsieh. Designing for workplace reflection: a chat and voice-based conversational agent. In *Proceedings of the 2018 designing interactive systems conference*, pages 881–894, 2018.
- Hiroki Tanaka, Sakriani Sakti, Graham Neubig, Tomoki Toda, Hideki Negoro, Hidemi Iwasaka, and Satoshi Nakamura. Automated social skills trainer. In *Proceedings of the 20th International Conference on Intelligent User Interfaces*, pages 17–27, 2015.
- Michelle X Zhou, Gloria Mark, Jingyi Li, and Huahai Yang. Trusting virtual agents: The effect of personality. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 9(2-3):1–36, 2019.
- Stefan Marti and Chris Schmandt. Physical embodiments for mobile communication agents. In *Proceedings of the 18th annual ACM symposium on User interface software and technology*, pages 231–240, 2005.
- Matthias Rehm. “she is just stupid”—analyzing user-agent interactions in emotional game situations. *Interacting with Computers*, 20(3):311–325, 2008.
- Bruce J Biddle. Recent developments in role theory. *Annual review of sociology*, 12(1):67–92, 1986.

- Kenneth D Benne and Paul Sheats. Functional roles of group members. *Journal of social issues*, 4(2):41–49, 1948.
- Steve WJ Kozlowski and Daniel R Ilgen. Enhancing the effectiveness of work groups and teams. *Psychological science in the public interest*, 7(3):77–124, 2006.
- Sigal G Barsade. The ripple effect: Emotional contagion and its influence on group behavior. *Administrative science quarterly*, 47(4):644–675, 2002.
- Masamune Kawasaki, Naomi Yamashita, Yi-Chieh Lee, and Kayoko Nohara. Assessing users’ mental status from their journaling behavior through chatbots. In *Proceedings of the 20th ACM International Conference on Intelligent Virtual Agents*, pages 1–8, 2020.
- Paul Thomas, Daniel McDuff, Mary Czerwinski, and Nick Craswell. Expressions of style in information seeking conversation with an agent. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1171–1180, 2020.
- Ziang Xiao, Michelle X Zhou, Wenxi Chen, Huahai Yang, and Changyan Chi. 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*, pages 1–14, 2020a.
- Mary Ellen Foster, Manuel Giuliani, and Alois Knoll. Comparing objective and subjective measures of usability in a human-robot dialogue system. In *Proceedings of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing (ACL-IJCNLP 2009)*, 2009.
- Vera Demberg, Andi Winterboer, and Johanna D Moore. A strategy for information presentation in spoken dialog systems. *Computational Linguistics*, 37(3):489–539, 2011.
- Florian Pecune, Jingya Chen, Yoichi Matsuyama, and Justine Cassell. Field trial analysis of socially aware robot assistant. In *Proceedings of the 17th international conference on autonomous agents and multiagent systems*, pages 1241–1249, 2018.
- Yi-Chieh Lee, Naomi Yamashita, and Yun Huang. Designing a chatbot as a mediator for promoting deep self-disclosure to a real mental health professional. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW1):1–27, 2020a.
- Shruti Sannon, Brett Stoll, Dominic DiFranzo, Malte Jung, and Natalya N Bazarova. How personification and interactivity influence stress-related disclosures to conversational agents. In *companion of the 2018 ACM conference on computer supported cooperative work and social computing*, pages 285–288, 2018.
- Michael Bonfert, Maximilian Spliethöver, Roman Arzaroli, Marvin Lange, Martin Hanci, and Robert Porzel. If you ask nicely: a digital assistant rebuking impolite voice commands. In *proceedings of the 20th ACM international conference on multimodal interaction*, pages 95–102, 2018.
- Ziang Xiao, Michelle X Zhou, Q Vera Liao, Gloria Mark, Changyan Chi, Wenxi Chen, and Huahai Yang. Tell me about yourself: Using an ai-powered chatbot to conduct conversational surveys with open-ended questions. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 27(3):1–37, 2020b.
- Thiemo Wambgsanss, Rainer Winkler, Matthias Söllner, and Jan Marco Leimeister. A conversational agent to improve response quality in course evaluations. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–9, 2020.
- Joyce Y Chai, Lanbo She, Rui Fang, Spencer Ottarson, Cody Littlely, Changsong Liu, and Kenneth Hanson. Collaborative effort towards common ground in situated human-robot dialogue. In *2014 9th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 33–40. IEEE, 2014.
- Jinping Wang, Hyun Yang, Ruosi Shao, Saeed Abdullah, and S Shyam Sundar. Alexa as coach: Leveraging smart speakers to build social agents that reduce public speaking anxiety. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2020.
- Ziang Xiao, Michelle X Zhou, and Wat-Tat Fu. Who should be my teammates: Using a conversational agent to understand individuals and help teaming. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, pages 437–447, 2019.
- Qian Yu, Tonya Nguyen, Soravis Prakkamakul, and Niloufar Salehi. "i almost fell in love with a machine" speaking with computers affects self-disclosure. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–6, 2019.
- Benjamin S. Alpers, Kali Cornn, Lauren E. Feitzinger, Usman Khaliq, So Yeon Park, Bardia Beigi, Daniel Joseph Hills-Bunnell, Trevor Hyman, Kaustubh Deshpande, Rieko Yajima, et al. Capturing passenger experience in a ride-sharing autonomous vehicle: The role of digital assistants in user interface design. In *12th International Conference on Automotive User Interfaces and Interactive Vehicular Applications*, pages 83–93, 2020.

- Sarah Theres Völkel, Renate Haeuslschmid, Anna Werner, Heinrich Hussmann, and Andreas Butz. How to trick ai: Users' strategies for protecting themselves from automatic personality assessment. In *Proceedings of the 2020 CHI conference on human factors in computing systems*, pages 1–15, 2020.
- Weiyan Shi, Xuewei Wang, Yoo Jung Oh, Jingwen Zhang, Saurav Sahay, and Zhou Yu. Effects of persuasive dialogues: testing bot identities and inquiry strategies. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2020.
- Pradipta Biswas. A flexible approach to natural language generation for disabled children. In *Proceedings of the COLING/ACL 2006 Student Research Workshop*, pages 1–6, 2006.
- Pranav Khadpe, Ranjay Krishna, Li Fei-Fei, Jeffrey T Hancock, and Michael S Bernstein. Conceptual metaphors impact perceptions of human-ai collaboration. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW2):1–26, 2020.
- Mingming Liu, Qicheng Ding, Yu Zhang, Guoguang Zhao, Changjian Hu, Jiangtao Gong, Penghui Xu, Yu Zhang, Liuxin Zhang, and Qianying Wang. Cold comfort matters-how channel-wise emotional strategies help in a customer service chatbot. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–7, 2020.
- Rens Hoegen, Deepali Aneja, Daniel McDuff, and Mary Czerwinski. An end-to-end conversational style matching agent. In *Proceedings of the 19th ACM International Conference on Intelligent Virtual Agents*, pages 111–118, 2019.
- Andreea Muresan and Henning Pohl. Chats with bots: Balancing imitation and engagement. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–6, 2019.
- Wanling Cai, Yucheng Jin, and Li Chen. Critiquing for music exploration in conversational recommender systems. In *26th International Conference on Intelligent User Interfaces*, pages 480–490, 2021.
- Zahra Ashktorab, Mohit Jain, Q Vera Liao, and Justin D Weisz. Resilient chatbots: Repair strategy preferences for conversational breakdowns. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2019.
- Elliot G Mitchell, Rosa Maimone, Andrea Cassells, Jonathan N Tobin, Patricia Davidson, Arlene M Smaldone, and Lena Mamykina. Automated vs. human health coaching: Exploring participant and practitioner experiences. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–37, 2021.
- Agnetta Gulz, Magnus Haake, and Annika Silvervarg. Extending a teachable agent with a social conversation module—effects on student experiences and learning. In *International conference on artificial intelligence in education*, pages 106–114. Springer, 2011.
- Archana Prasad, Sean Blagsvedt, Tej Pochiraju, and Indrani Medhi Thies. Dara: A chatbot to help indian artists and designers discover international opportunities. In *Proceedings of the 2019 on Creativity and Cognition*, pages 626–632. 2019.
- Yi-Chieh Lee, Naomi Yamashita, Yun Huang, and Wai Fu. "i hear you, i feel you": Encouraging deep self-disclosure through a chatbot. In *Proceedings of the 2020 CHI conference on human factors in computing systems*, pages 1–12, 2020b.
- Luis Cavazos Quero, Jorge Iranzo Bartolomé, Dongmyeong Lee, Yerin Lee, Sangwon Lee, and Jundong Cho. Jido: a conversational tactile map for blind people. In *The 21st International ACM SIGACCESS Conference on Computers and Accessibility*, pages 682–684, 2019.
- Asbjørn Følstad and Marita Skjuve. Chatbots for customer service: user experience and motivation. In *Proceedings of the 1st international conference on conversational user interfaces*, pages 1–9, 2019.
- David R Large, Leigh Clark, Gary Burnett, Kyle Harrington, Jacob Luton, Peter Thomas, and Pete Bennett. "it's small talk, jim, but not as we know it." engendering trust through human-agent conversation in an autonomous, self-driving car. In *Proceedings of the 1st International Conference on Conversational User Interfaces*, pages 1–7, 2019.
- Sherry Ruan, Liwei Jiang, Justin Xu, Bryce Joe-Kun Tham, Zhengneng Qiu, Yeshuang Zhu, Elizabeth L Murnane, Emma Brunskill, and James A Landay. Quizbot: A dialogue-based adaptive learning system for factual knowledge. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2019.
- Anshul Bawa, Pranav Khadpe, Pratik Joshi, Kalika Bali, and Monojit Choudhury. Do multilingual users prefer chat-bots that code-mix? let's nudge and find out! *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW1): 1–23, 2020.
- Hyanghee Park and Joonhwan Lee. Can a conversational agent lower sexual violence victims' burden of self-disclosure? In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–8, 2020.

- Qingxiao Zheng, Daniela M Markazi, Yiliu Tang, and Yun Huang. "pocketbot is like a knock-on-the-door!": Designing a chatbot to support long-distance relationships. *Proceedings of the ACM on Human-Computer Interaction*, 5 (CSCW2):1–28, 2021.
- Leigh Clark, Philip Doyle, Diego Garaialde, Emer Gilmartin, Stephan Schlögl, Jens Edlund, Matthew Aylett, João Cabral, Cosmin Munteanu, Justin Edwards, et al. The state of speech in hci: Trends, themes and challenges. *Interacting with Computers*, 31(4):349–371, 2019.
- Milton Borsato and Margherita Peruzzini. Collaborative engineering. In *Concurrent engineering in the 21st century*, pages 165–196. Springer, 2015.
- Robert O Briggs, Gwendolyn L Kolfschoten, and Gert-Jan de Vreede. Toward a theoretical model of consensus building. *AMCIS 2005 Proceedings*, page 12, 2005.
- Gwendolyn L Kolfschoten and Gert-Jan De Vreede. The collaboration engineering approach for designing collaboration processes. In *International Conference on Collaboration and Technology*, pages 95–110. Springer, 2007.
- Josephine Nabukenya, Patrick van Bommel, and Henderik Alex Proper. A theory-driven design approach to collaborative policy making processes. In *2009 42nd Hawaii International Conference on System Sciences*, pages 1–10. IEEE, 2009.
- Samiha Samrose, Kavya Anbarasu, Ajjen Joshi, and Taniya Mishra. Mitigating boredom using an empathetic conversational agent. In *Proceedings of the 20th ACM International Conference on Intelligent Virtual Agents*, pages 1–8, 2020.
- David J Houghton and Adam N Joinson. Privacy, social network sites, and social relations. *Journal of Technology in Human Services*, 28(1-2):74–94, 2010.
- Michèle Lamont and Virág Molnár. The study of boundaries in the social sciences. *Annual review of sociology*, 28(1): 167–195, 2002.
- Sandra Petronio. Brief status report on communication privacy management theory. *Journal of Family Communication*, 13(1):6–14, 2013.
- Sandra Petronio and Wesley T Durham. Communication privacy management theory. *Engaging Theories in Interpersonal Communication: Multiple Perspectives*, page 335, 2014.
- Daniel Vogel and Ravin Balakrishnan. Interactive public ambient displays: transitioning from implicit to explicit, public to personal, interaction with multiple users. In *Proceedings of the 17th annual ACM symposium on User interface software and technology*, pages 137–146, 2004.
- Constantine Stephanidis, Gavriel Salvendy, Margherita Antona, Jessie YC Chen, Jianming Dong, Vincent G Duffy, Xiaowen Fang, Cali Fidopiastis, Gino Fragomeni, Limin Paul Fu, et al. Seven hci grand challenges. *International Journal of Human-Computer Interaction*, 35(14):1229–1269, 2019.
- Kevin Burden and Matthew Kearney. Conceptualising authentic mobile learning. In *Mobile learning design*, pages 27–42. Springer, 2016.
- Yi Mou and Kun Xu. The media inequality: Comparing the initial human-human and human-ai social interactions. *Computers in Human Behavior*, 72:432–440, 2017.
- Andrea L Guzman. Ontological boundaries between humans and computers and the implications for human-machine communication. *Human-Machine Communication*, 1(1):3, 2020.
- Samantha Reig, Michal Luria, Janet Z Wang, Danielle Oltman, Elizabeth Jeanne Carter, Aaron Steinfeld, Jodi Forlizzi, and John Zimmerman. Not some random agent: Multi-person interaction with a personalizing service robot. In *Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction*, pages 289–297, 2020.
- Leysia Palen and Paul Dourish. Unpacking "privacy" for a networked world. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 129–136, 2003.
- Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N Bennett, Kori Inkpen, et al. Guidelines for human-ai interaction. In *Proceedings of the 2019 chi conference on human factors in computing systems*, pages 1–13, 2019.
- JL Beuth, CP Rosé, and R Kumar. Software agent-monitored tutorials enabling collaborative learning in computer-aided design and analysis. In *ASME International Mechanical Engineering Congress and Exposition*, volume 44434, pages 339–346, 2010.
- Microsoft. Responsible bots: 10 guidelines for developers of conversational ai. 2018.
- Clifford Nass, Jonathan Steuer, and Ellen R Tauber. Computers are social actors. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 72–78, 1994.

- Judith Donath. Signals, cues and meaning. *Signals, Truth and Design*, 2007.
- Magalie Ochs, Nathan Libermann, Axel Boidin, and Thierry Chaminade. Do you speak to a human or a virtual agent? automatic analysis of user’s social cues during mediated communication. In *Proceedings of the 19th ACM International Conference on Multimodal Interaction*, pages 197–205, 2017.
- Jasper Feine, Ulrich Gnewuch, Stefan Morana, and Alexander Maedche. A taxonomy of social cues for conversational agents. *International Journal of Human-Computer Studies*, 132:138–161, 2019.
- Milad Mirbabaie, Stefan Stieglitz, Felix Brünker, Lennart Hofeditz, Björn Ross, and Nicholas RJ Frick. Understanding collaboration with virtual assistants—the role of social identity and the extended self. *Business & Information Systems Engineering*, 63(1):21–37, 2021.
- Thiemo Wambsganss, Anne Höch, Naim Zierau, and Matthias Söllner. Ethical design of conversational agents: Towards principles for a value-sensitive design.
- Allan Dafoe. On technological determinism: a typology, scope conditions, and a mechanism. *Science, Technology, & Human Values*, 40(6):1047–1076, 2015.
- Jacques Wainer and Claudia Barsottini. Empirical research in cscw—a review of the acm/cscw conferences from 1998 to 2004. *Journal of the Brazilian Computer Society*, 13(3):27–35, 2007.
- David Pinelle and Carl Gutwin. A review of groupware evaluations. In *Proceedings IEEE 9th International Workshops on Enabling Technologies: Infrastructure for Collaborative Enterprises (WET ICE 2000)*, pages 86–91. IEEE, 2000.
- Walter F Tichy, Paul Lukowicz, Lutz Prechelt, and Ernst A Heinz. Experimental evaluation in computer science: A quantitative study. *Journal of Systems and Software*, 28(1):9–18, 1995.
- Thien Hai Nguyen and Kiyooki Shirai. Topic modeling based sentiment analysis on social media for stock market prediction. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1354–1364, 2015.