

The AI Doctor Is In: A Survey of Task-Oriented Dialogue Systems for Healthcare Applications

Anonymous ACL submission

Abstract

Task-oriented dialogue systems in healthcare are increasingly common and have been characterized by diverse architectures and objectives. Although they have been surveyed in the medical community from a non-technical perspective, a systematic review from a rigorous computational perspective remains noticeably absent. This has resulted in limited knowledge of important implementation and replicability details, slowing the pace of innovation. To fill this gap, we investigated an initial pool of 4070 papers from well-known computer science, natural language processing, and artificial intelligence venues, identifying 70 papers discussing the system-level implementation of task-oriented dialogue systems for healthcare applications. We comprehensively reviewed these papers, and present our key findings including identified gaps and corresponding recommendations.

1 Introduction

Dialogue systems¹ have a daily presence in many individuals' lives, acting as virtual assistants (Hoy, 2018), customer service agents (Xu et al., 2017), or even companions (Zhou et al., 2020). While some systems are designed to conduct unstructured conversations in open domains (*chatbots*), others (*task-oriented dialogue systems*) help users to complete tasks in a specific domain (Jurafsky and Martin, 2009; Qin et al., 2019). Task-oriented dialogue systems can potentially play an important role in health and medical care (Laranjo et al., 2018), and they have been adopted by growing numbers of patients, caregivers, and clinicians (Kearns et al., 2019). Nonetheless, there remains a translational gap (Newman-Griffis et al., 2021) between cutting-edge, foundational work in dialogue systems and

¹We follow an inclusive definition of *dialogue systems*, encompassing any intelligent systems designed to converse with humans via natural language.

prototypical or deployed dialogue agents in healthcare settings. This limits the proliferation of scientific progress to real-world systems, constraining the potential benefits of fundamental research.

We move towards closing this gap by conducting a comprehensive, scientifically rigorous analysis of task-oriented healthcare dialogue systems. Our underlying objectives are to (a) explore how these systems have been employed to date, and (b) map out their characteristics, shortcomings, and subsequent opportunities for follow-up work. Importantly, we seek to address the limitations of prior systematic reviews by extensively investigating the included systems from a computational perspective. Our primary contributions are as follows. (1) We systematically search through 4070 papers from well-known technical venues and identify 70 papers fitting our inclusion criteria.² (2) We analyze these systems based on many factors, including system objective, language, architecture, modality, device type, and evaluation paradigm, among others. (3) We identify common limitations across systems, including an incomplete exploration of architecture, replicability concerns, ethical and privacy issues, and minimal investigation of usability or engagement. We offer practical suggestions for addressing these as an on-ramp for future work.

In the long term, we hope that the gaps and opportunities identified in this survey can stimulate more rapid advancements in the design of task-oriented healthcare dialogue systems. We also hope that the survey provides a useful starting point and synthesis of prior work for NLP researchers and practitioners entering this critical yet surprisingly under-studied application domain.

2 Related Work

Dialogue systems in healthcare have been the focus of several recent surveys conducted by the medical

²A full listing of these papers is provided in the appendix.

and clinical communities (Vaidyam et al., 2019; Laranjo et al., 2018; Kearns et al., 2019). These surveys have investigated the real-world utilization of deployed systems, rather than examining their design and implementation from a technical perspective. In contrast, studies examining these systems through the lens of AI and NLP research and practice have been limited. Zhang et al. (2020) and Chen et al. (2017) presented surveys of recent advances in general-domain task-oriented dialogue systems. Although they provide an excellent holistic portrait of the subfield, they do not delve into aspects of particular interest in healthcare settings (e.g., system objectives doubling as clinical goals), limiting their usefulness for this audience.

Vaidyam et al. (2019), Laranjo et al. (2018), and Kearns et al. (2019) conducted systematic reviews of dialogue systems deployed in mental health (Vaidyam et al., 2019) or general healthcare (Laranjo et al., 2018; Kearns et al., 2019) settings. Vaidyam et al. (2019) examined 10 articles, and Laranjo et al. (2018) and Kearns et al. (2019) examined 17 and 46 articles, respectively. All surveys were written for a medical audience and focused on healthcare issues and impact, covering few articles from AI, NLP, or general computer science venues.

Montenegro et al. (2019) and Tudor Car et al. (2020) recently reviewed 40 and 47 articles, respectively, covering conversational agents in the healthcare domain. These two surveys are the closest to ours, but differ in important ways. First, our focus is on a specific class of conversational agents: task-oriented dialogue systems. The surveys by Montenegro et al. (2019) and Tudor Car et al. (2020) used a wider search breaching their ability to provide extensive technical depth. We also reviewed more papers (70 articles), which were then screened using a more thorough taxonomy as part of the analysis. Some aspects that we considered that differ from these prior surveys include the overall dialogue system architecture, the dialogue management architecture, the system evaluation methods, and the dataset(s) used when developing and/or evaluating the system.

3 Search Criteria and Screening

We designed search criteria in concert with our goal of filling a translational information gap between fundamental dialogue systems research and applied systems in the healthcare domain. To do so, we retrieved articles from well-respected computer sci-

Screening Process	ACM	IEEE	ACL	AAAI	Total
Initial Search	1050	1400	1020	600	4070
Title Screening	151	273	106	55	585
Abstract Screening	32	45	26	8	110
Final Screening	21	31	16	2	70

Table 1: The number of papers included from each database in each step of the paper screening process.

ence, AI, and NLP databases and screened them for focus on task-oriented dialogue systems designed for healthcare settings. Our target databases were: (1) ACM,³ (2) IEEE,⁴ (3) the ACL Anthology,⁵ and (4) the AAAI Digital Library.⁶ ACM and IEEE are large databases of papers from prestigious conferences and journals across many CS fields, including but not limited to robotics, human-computer interaction, data mining, and multimedia systems. The ACL Anthology is the premier database of publications within NLP, hosting papers from major conferences and topic-specific venues (e.g., *SIGDIAL*, organized by the Special Interest Group on Discourse and Dialogue). The AAAI Digital Library hosts papers not only from the *AAAI Conference on Artificial Intelligence*, but also from other AI conferences, *AI Magazine*, and the *Journal of Artificial Intelligence Research*. We applied the following inclusion criteria when identifying papers:

- The main focus must be on the technical design or implementation of a task-oriented dialogue system.
- The system must be designed for health-related applications.
- The article must *not* be dedicated to one specific module of the system’s architecture (e.g., the natural language understanding component of a health-related dialogue system).

Although a narrower scope—e.g., developing improved methods for slot-filling—is common when publishing in the dialogue systems community,

³<https://dl.acm.org/>

⁴<https://ieeexplore.ieee.org/>

⁵<https://www.aclweb.org/anthology/>

⁶<https://aaai.org/Library/library.php>

these papers tend to place more emphasis on technical design irrespective of application context, offering less coverage of the system-level characteristics that are the target of this survey. We followed four steps in our screening process. First (*Initial Search*), we applied a predefined search query to the databases to populate our initial list of papers. To generate the query, we used the keywords “task-oriented,” “dialogue system,” “conversational agent,” “health,” and “healthcare,” and synonyms and abbreviations of these keywords. We short-listed papers using these keywords individually as well as in combination with one another.

Next (*Title Screening*), we performed a preliminary screening through the initial list of papers by reading the titles, keeping those that satisfied the inclusion criteria. Then (*Abstract Screening*), we went through the list of papers remaining after the title screening and read the abstracts, keeping those that satisfied the inclusion criteria. Lastly (*Final Screening*), we read the body of the papers remaining after the abstract screening and kept those that satisfied the inclusion criteria.

These funnel filtering processes were conducted by a computer science graduate student (a fluent L2 English speaker) using predefined search and screening guidelines. Questions or uncertainties regarding a paper’s compliance with inclusion criteria were forwarded along to the senior project lead (a computer science professor and fluent L1 English speaker with expertise in NLP) and final consensus was reached via discussion among the two parties. We detail the number of papers remaining after each screening step in Table 1.

In total, 70 papers (21 from ACM, 31 from IEEE, 16 from ACL, and 2 from AAAI⁷) satisfied the inclusion criteria. We survey papers meeting our inclusion criteria according to a wide range of parameters, and present our findings in the following subsections, grouped into thematic categories: ontology (§4), system architecture (§5), system design (§6), dataset (§7), and system evaluation (§8).

4 Ontology

We map each paper to its domain of research (§4.1), system objective (§4.2), target audience (§4.3), and language (§4.4), and present our findings.

⁷Papers about task-oriented dialogue systems published at AAAI often focus on one specific component of the system from a technical perspective, rather than proposing a conversational agent as a whole. Therefore, only two papers from the AAAI Digital Library satisfied the inclusion criteria.

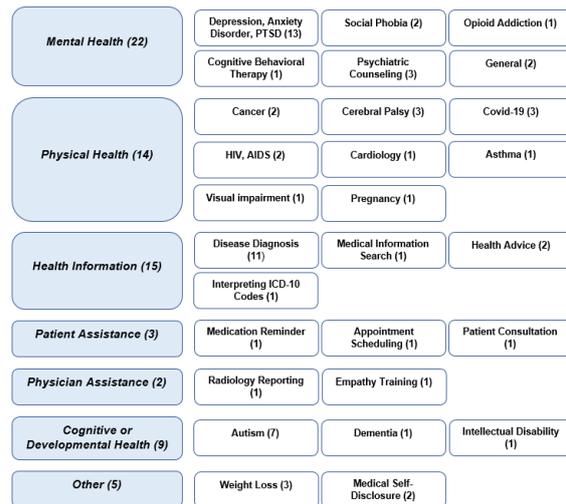


Figure 1: Research domains and corresponding subcategories for the included papers. Parentheses indicate the number of papers belonging to the (sub)category.

4.1 Domain of Research

Task-oriented dialogue systems can potentially impact many facets of healthcare in society (Bickmore and Giorgino, 2004). We define a *domain of research* as the healthcare area in which the system operates. We identify both broad domains and more specific subcategories thereof based on the systems surveyed, outlined in Figure 1. Broad domain categories include *mental health*, *physical health*, *health information*, *patient assistance*, *physician assistance*, *cognitive or developmental health*, and *other* (comprising subcategories not easily classifiable to one of the broader domains).

Systems in the *mental health* domain supported individuals with mental or psychological health conditions, and systems in the *cognitive or developmental health* domain were a close analogue for individuals with conditions impacting memory, executive, or other cognitive function. Systems in the *physical health* domain were targeted towards individuals with specific physical health concerns, including infectious (e.g., Covid-19), non-infectious (e.g., cancer), and temporary (e.g., pregnancy) conditions. Systems providing *health information* performed general-purpose actions such as offering advice or suggesting disease diagnoses. Finally, systems performing *patient assistance* or *physician assistance* supported specific patient- or physician-focused healthcare tasks. Dialogue systems designed for *mental health*, *physical health*, and *health information* were the most prevalent, covering 51 of the 70 included papers.

System Objective	# Papers
Diagnosis	7
Monitoring	8
Intervention	13
Counseling	5
Assistance	12
Multi-Objective	25

Table 2: Distribution of system objectives across the surveyed papers. Additional details regarding *multi-objective* papers are provided in the appendix.

4.2 System Objective

Task-oriented dialogue systems define value relative to the goals of a target task. We define the *system objective* as the healthcare task for which a system is designed. Some system objectives may be closely aligned with a single domain, whereas others may occur in many different domains (e.g., *monitoring* mental, physical, or cognitive conditions). Thus, although the domain of research and system objective may frequently correlate, there is not by necessity a direct association.

Included systems were categorized as being designed to: *diagnose* a health condition (e.g., by predicting whether the user suffers from cognitive decline); *monitor* user states (e.g., by tracking their diets or periodically checking their mood); *intervene* by addressing users’ health concerns or improving their states (e.g., by teaching children how to map facial expressions to emotions); *counsel* users without providing any direct intervention (e.g., by listening to users’ concerns and empathizing with them); or *assist* users by providing information or guidance (e.g., by answering questions from users who are filling out forms). Many systems were also categorized as *multi-objective*, meaning that they were designed for more than one of those goals.

Table 2 shows the number of systems having each objective. Many systems (25/70) were designed for more than one target objective. Among *multi-objective* systems, those that were designed for both diagnosis and assistance had the highest frequency (7/25); we provide additional details regarding these systems in Table 8 of the appendix.

Separately, we also considered the role of *engagement* as an objective of each system. We define this as a goal of engaging target users in interaction, irrespective of underlying health goals. Engagement may be of particular interest in health-

Target Audience	# Papers
Patients	59
Caregivers	3
Patients & Caregivers	2
Clinicians	11

Table 3: Distribution of the target audiences of the systems described in the surveyed papers.

care settings since it can be critical in encouraging adoption or adherence with respect to healthcare outcomes (Montenegro et al., 2019). Surprisingly, almost 60% of the papers (41 of the 70 surveyed) did not mention any goals pertaining to engaging users in more interactions.

4.3 Target Audience

The final consumers of healthcare systems often fall into three groups: *patients*, *caregivers*, and *clinicians*. Table 3 shows the number of systems surveyed that focus on each category. We find that out of 70 task-oriented dialogue systems, 59 are designed specifically for patients.

4.4 Language

Most general-domain dialogue systems research has been conducted in English and other high-resource languages (Artetxe et al., 2020). Expanding language diversity may extend the benefits of health-related dialogue systems more globally. As shown in Figure 2, among the systems included in our review a majority (56%) are designed for English speakers. Encouragingly, several of the included systems did focus on lower-resource languages, including Telugu (Duggenpudi et al., 2019), Bengali (Rahman et al., 2019), and Setswana (Grover et al., 2009).

5 System Architecture

We investigate both the general architecture of the system (§5.1), and if applicable, the dialogue management architecture specifically (§5.2).

5.1 General Architecture

Task-oriented dialogue systems are generally designed using *pipeline* or *end-to-end* architectures. Pipeline architectures typically consist of separate components for natural language understanding, dialogue state tracking, dialogue policy, and natural language generation. The ensemble of the dialogue state tracker and dialogue policy is the *dialogue*

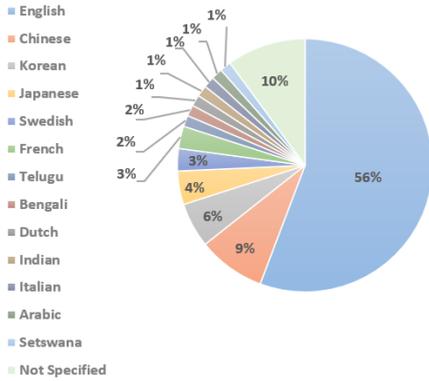


Figure 2: Language diversity across the surveyed systems. A small percentage (10%) of papers do not specify the system’s language.

System Architecture	# Papers
Pipeline	58
End-to-End	2
Not Specified	10

Table 4: Distribution of papers describing systems with pipeline or end-to-end architectures, or that do not specify the architecture.

manager (Chen et al., 2017). End-to-end architectures train a single model to produce output for a given input, often interacting with structured external databases and requiring extensive training data (Chen et al., 2017). As shown in Table 4, only 2.85% of papers (2 of the 70 surveyed) implemented an end-to-end system; this is unsurprising given the limited training data available in most healthcare domains. We also found that 14% (10 papers) did not directly specify the architecture of their developed system.

5.2 Dialogue Management Architecture

Unlike other pipeline components that impact user experience and engagement but not fundamental decision-making, the dialogue manager is central to overall functionality (Zhao et al., 2019); thus, we afford it special attention. In *rule-based* approaches, the system interacts with users based on a predefined set of rules, with success conditioned upon coverage of all relevant cases (Siangchin and Samanchuen, 2019). *Intent-based* approaches seek to extract the user’s intention from the dialogue, and then perform the relevant action (Jurafsky and Martin, 2009). In *hybrid* dialogue management architectures, the system leverages a combination of rule-based and intent-based approaches, and fi-

Dialogue Management Architecture	# Papers
Rule-based	17
Intent-based	20
Hybrid Architecture	21
Corpus-based	0

Table 5: Distribution of dialogue management architectures across the surveyed papers. This table does not include papers describing end-to-end architectures ($n = 2$) or for which system architecture was not specified ($n = 10$).

nally *corpus-based* approaches mine the dialogues of human-human conversations and produce responses using retrieval methods or generative methods (Jurafsky and Martin, 2009). As shown in Table 5, among papers reporting on dialogue management architecture, we observe a fairly even mix of rule-based, intent-based, and hybrid architectures.

6 System Design

6.1 Modality

Modality, the channel through which information is exchanged between a computer and a human (Karray et al., 2008), can play an important role in dialogue quality and user satisfaction (Bilici et al., 2000). *Unimodal* systems use a single modality for information exchange, whereas *multimodal* systems use multiple modalities (Karray et al., 2008). Systems reviewed in this survey operated using one or more of several modalities. In *text-based* or *spoken* interaction, users interact with the system by typing or speaking, respectively. In interaction via *graphical user interface (GUI)*, users interact with the system through the use of visual elements.

In general, multimodal dialogue systems can be flexible and robust, but especially challenging to implement in the medical domain (Sonntag et al., 2009). We find that 49 papers describe unimodal systems and 21 describe multimodal systems. Table 6 provides more details regarding their distribution across modalities.

6.2 Device

Dialogue systems may facilitate interaction using a variety of devices (Arora et al., 2013), ranging from telephones (Garvey and Sankaranarayanan, 2012) to computers (McTear, 2010) to any other technology that allows interaction (e.g., VR-based avatars (Brinkman et al., 2012b; McTear, 2010)). We categorized the included systems as *mobile*, *telephone*,

Unimodal		Multimodal	
Category	# Papers	Category	# Papers
Text	23	Spoken + Text	14
Spoken	25	Spoken + GUI	4
GUI	1	Text + GUI	3

Table 6: Distribution of modality type across the unimodal (49 total, left) and multimodal (21 total, right) systems surveyed.

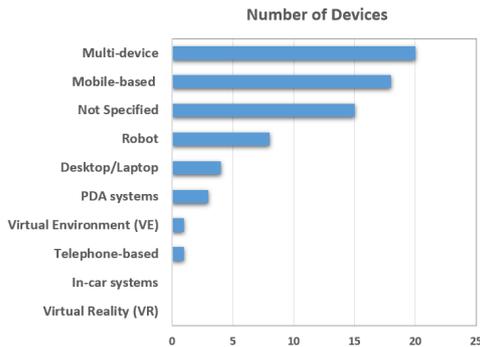


Figure 3: Distribution of device type across the surveyed papers.

desktop/laptop, in-car, PDA, robot, virtual environment, or virtual reality (including virtual agents and avatars) systems, considering systems as *multi-device* if they leveraged multiple devices for interaction. As shown in Figure 3, we found that multi-device and mobile-based dialogue systems were most popular. Table 9 in the appendix provides additional details regarding multi-device systems.

7 Dataset

Data is crucial for effective system development (Serban et al., 2015), but many datasets for training dialogue systems are smaller than those used for other NLP tasks (Lowe et al., 2017). This is even more pronounced in the healthcare domain, in part due to the risk of data misuse by others or the lack of data sharing incentives (Lee and Yoon, 2017).

We reviewed each paper for information regarding the data used during system development, focusing on dataset size, availability, and privacy-preserving measures. Only 20 papers provide details about the data used (two papers provided a link to the dataset, and the remaining 18 discussed the dataset size). Unfortunately, the remaining papers did not provide rationale for their lack of data or other replicability information. Our assumption is that often the data contained sensitive information,

Evaluation Type	# Papers
Human Evaluation	28
Automated Evaluation	7
Human & Automated Evaluation	9
Not Specified	26

Table 7: Distribution of evaluation methods across the surveyed papers.

preventing authors from releasing specific details, but only 19 of the 70 included papers provided information about data-related privacy or ethical considerations. Only 10 mentioned Institutional Review Board (IRB) approval for their dataset and/or task, despite IRB (or equivalent) review being a crucial step towards ensuring that research is conducted ethically and in such a way that protects human subjects to the extent possible (Amdur and Biddle, 1997).

8 System Evaluation

We examined the means through which systems were evaluated both qualitatively and quantitatively (Deriu et al., 2019; Hastie, 2012). We defined *human evaluation*, often implemented in prior work through questionnaires (Wang et al., 2020; Grover et al., 2009; Holmes et al., 2019) or direct feedback from real-world users (Deriu et al., 2019), as an evaluation that relies on subjective, first-hand, human user experience. In contrast, *automated evaluation* provides an objective, quantitative measurement of one or more dimensions of the system from a mathematical perspective (Finch and Choi, 2020). Some metrics used for automated evaluation of the reviewed systems include measures of task performance (Ali et al., 2020) and completion rates (Holmes et al., 2019), response correctness (Rosruen and Samanchuen, 2018), and response time (Grover et al., 2009).

In Table 7, we observe that nearly half of the papers conducted human evaluations; however, a large percentage (37%) also did not discuss evaluation at all. We further analyzed papers conducting human evaluations and found that they included an average of 26 (mode = 12) participants. More details regarding the human and automated evaluations are provided in Tables 10, 11, and 12 of the appendix. In a follow-up analysis of *system usability*, defined as the degree to which users are able to engage with a system safely, effectively, efficiently,

and enjoyably (Lee et al., 2019), we observed that 33 papers explicitly evaluated the usability of their system.

9 Discussion

We identify common limitations across many surveyed systems, accompanied by recommendations for addressing them in future work.

9.1 Incomplete Exploration of System Design

We observed little system-level architectural diversity across the surveyed systems, with most (83%) having a pipeline architecture. This architectural homogeneity limits our understanding of good design practice within this domain. Recent studies demonstrate that end-to-end architectures for task-oriented dialogue systems could compete with pipeline architectures given sufficient high-quality data (Hosseini-Asl et al., 2020; Ham et al., 2020; Bordes et al., 2017; Wen et al., 2016). However, the external knowledge sources often leveraged in end-to-end systems are notoriously complex in many healthcare sub-domains (Campillos-Llanos et al., 2020). Additionally, for healthcare applications interpretability is highly desired (Ham et al., 2020), but explanations are often obfuscated in end-to-end systems (Ham et al., 2020; Wen et al., 2016). Finally, users of these systems may seek guidance on sensitive topics, which can exacerbate privacy concerns (Xu et al., 2021). Any system trained on large, weakly curated datasets may also learn unpleasant behaviors and amplify biases in the training data, in turn producing harmful consequences (Dinan et al., 2021; Bender et al., 2021). We recommend further experimentation with architectural design, in parallel with work towards developing high-quality healthcare dialogue datasets, which to date remain scarce (Farzana et al., 2020).

We noticed that a considerable number of the systems (33%) allowed only text-based interaction. However, it is well-established that individuals from certain demographic groups are more comfortable conversing with dialogue systems via speech (Tudor Car et al., 2020). Text-based systems may also be more likely to violate privacy considerations (Tudor Car et al., 2020). Thus, we recommend that researchers engage in further exploration of multimodal or spoken dialogue systems when applicable and appropriate.

Many of the surveyed systems were also implemented on mobile phones. Although an advantage

of mobile-based systems is that they are readily available using a technology familiar to most users, Lee et al. (2018) found that users significantly reduced their usage over time when engaging long-term with mobile health applications. Tudor Car et al. (2020) suggest that one way to overcome this limitation in mobile-based systems is by directly embedding them in applications or platforms with which users already engage habitually (e.g., Facebook Messenger). This more ambient dissemination approach may facilitate easier and more lasting integration of system use in individuals' daily lives.

Finally, we identified that most systems (84%) target only patients, with research on systems targeted towards clinicians and caregivers remaining limited. We recommend further exploration of systems targeted towards these critical audiences. This may offer broad, high-impact support in understanding, diagnosing, and treating patients' health issues (Valizadeh et al., 2021; Kaelin et al., 2021).

9.2 Replicability Concerns

Data accessibility restrictions reduce the capacity of public health research (Strongman et al., 2019), and these limitations may be partially responsible for the imbalance of pipeline versus end-to-end architectures (§9.1). Only a small percentage of papers surveyed (29%) ventured to discuss the quantity or characteristics of the data used during system development in any way. A lack of data transparency hinders scientific progress and severely impedes replicability. We call upon researchers to publish data when permissible by governing protocol, and descriptive statistics to the extent allowable when circumstances prevent data release. We also view the development of high-quality, publicly available datasets as an important frontier in translational dialogue systems research (§9.1).

Many of the surveyed papers also lack important implementation details, such as evaluation methods (34%). This prevents the research community from replicating developed systems and generalizing study findings more broadly (Walker et al., 2018). Well-established guidelines exist and are being increasingly enforced within the NLP community to prevent reproducibility issues (Dodge et al., 2019). The disregard of reproducibility best practices observed with many healthcare dialogue systems may be partially attributed to the most common target venues for this work, which may place less emphasis on replication. This validates a cen-

539 tral motivator for publishing this survey—without
540 adequate inclusion of target domain *and* techni-
541 cal stakeholders in interdisciplinary, translational
542 research, progress will remain constrained. We
543 strongly urge researchers in this domain to provide
544 implementation details in their publications.

545 9.3 Potential Ethical and Privacy Issues

546 Real-world medical data facilitates the devel-
547 opment of high-quality healthcare applications
548 (Bertino et al., 2005), but protecting the rights
549 and privacy of contributors to the data is critical
550 for ensuring ethical research conduct (Institute of
551 Medicine, 2009), as is proper treatment of copy-
552 right protections. We screened all included papers
553 for coverage of privacy and ethical concerns, and
554 observed that only 27% of the surveyed papers con-
555 sidered participant or patient privacy in the design
556 of their system. Moreover, only 14% of the sur-
557 veyed papers documented any evidence of Institu-
558 tional Review Board (or IRB-equivalent) approval.

559 Research involving healthcare dialogue systems
560 is unquestionably human-centered, and as such the
561 absence of ethical oversight in the design of such
562 systems is a grave concern. Although technical
563 researchers entering this space may be unfamiliar
564 with human subjects research and protocol, we urge
565 all dialogue systems researchers to submit their
566 experimental design and protocol for review by an
567 appropriate external review board. We also ask that
568 researchers consider the potential harms from use
569 or misuse of their systems, following guidelines
570 established by the ACM Code of Ethics.⁸

571 9.4 Room for Increased Language Diversity

572 We observed that most systems (56%) targeted En-
573 glish speakers. Developing multilingual dialogue
574 systems or systems for speakers of low-resource
575 languages brings up various challenges (López-
576 Cózar Delgado and Araki, 2005), but solving this
577 problem could have tremendous benefit for
578 individuals in non-English speaking communities
579 with minimal or unreliable healthcare access. The
580 systems developed by Duggenpudi et al. (2019),
581 Rahman et al. (2019), and Grover et al. (2009) pro-
582 vide case examples for how such systems may be
583 implemented. We also note that while troubling,
584 a 56% share of systems targeted towards English
585 speakers is consistent with linguistic homogeneity
586 in the field in general, and actually slightly low

587 relative to many other NLP tasks (Mielke, 2016;
588 Bender, 2009). Healthcare dialogue systems may
589 on some level offer a case example for how appli-
590 cations originally designed for high-resource (i.e.,
591 English-language) settings can be adapted and re-
592 engineered to provide better coverage of the di-
593 verse, real-world potential user base.

594 9.5 Minimal Investigation of Usability or 595 User Engagement

596 Finally, more than 50% (37/70) of the included
597 papers did not evaluate system usability or gen-
598 eral user experience. Usability testing can improve
599 productivity and safeguard against errors (Rogers
600 et al., 2005), both of which are critical in healthcare
601 tasks. Therefore, we urge the research community
602 to consider and assess usability when designing for
603 this domain. The systems among those surveyed
604 that do this already (e.g., those developed by Wang
605 et al. (2020), Lee et al. (2020b), Wei et al. (2018),
606 or Demasi et al. (2020)) provide case examples for
607 how it might be done.

608 Almost 60% of the surveyed systems were not
609 explicitly designed to engage users, despite this
610 being a common objective in the general domain
611 (Ghazarian et al., 2019). Healthcare dialogue sys-
612 tems may stand to benefit particularly well from
613 such measures, since patient engagement is predic-
614 tive of adoption and adherence to healthcare out-
615 comes (Montenegro et al., 2019). To increase user
616 satisfaction and system performance, we recom-
617 mend that the research community more purpose-
618 fully consider engagement when designing their
619 healthcare-oriented dialogue systems.

620 10 Conclusion

621 In this work, we conducted a systematic techni-
622 cal survey of task-oriented dialogue systems in
623 the healthcare domain, narrowing the translational
624 gap between basic and applied dialogue systems
625 research. We comprehensively searched through
626 4070 papers in computer science, NLP, and AI
627 databases, finding 70 papers that satisfied our inclu-
628 sion criteria. We analyzed these papers based on
629 numerous technical factors, and present evidence-
630 based recommendations stemming from our find-
631 ings. It is our hope that interested researchers find
632 the information provided to be a unique and help-
633 ful resource for developing task-oriented dialogue
634 systems for healthcare applications.

⁸<https://www.acm.org/code-of-ethics>

11 Ethical Considerations

Beyond the concrete changes suggested during the discussion, it is important to consider the broader ethical implications of task-oriented dialogue systems in healthcare settings. Although the goal of such systems may not be to replace human healthcare providers, it is likely that deployed systems would support clinicians, defraying workload for overburdened individuals. In doing so, these systems may have significant impact on healthcare decision-making. Machines are imperfect, and thus a possible harm is that these systems may misinterpret user input or make incorrect predictions—a mistake that in high-stakes healthcare settings could prove detrimental or even dangerous. Researchers and developers should be cognizant of possible harms stemming from the use and misuse of task-oriented dialogue systems for healthcare settings, and should implement both automated (e.g., strict thresholds for diagnostic suggestions) and human (e.g., training to ensure staff awareness of potential system fallibilities) safeguards.

Moreover, a potential benefit of these systems is their potential to meaningfully and beneficially extend healthcare access to underserved populations. As such, it is important to ensure that automated systems do not fall prey to the same biases often observed among human healthcare providers (FitzGerald and Hurst, 2017). Systems trained to perform healthcare tasks using datasets that are not representative of the target population may exhibit poorer performance with users who already experience marginalization or are otherwise vulnerable, impeding or even reversing benefits. We call upon researchers to examine, debias, and curate their training data such that task-oriented dialogue systems for healthcare applications elevate, rather than diminish, outcomes for the historically underserved users which they are best poised to benefit.

References

Parham Aarabi. 2013. [Virtual cardiologist — a conversational system for medical diagnosis](#). In *2013 26th IEEE Canadian Conference on Electrical and Computer Engineering (CCECE)*, pages 1–4.

Yuna Ahn, Yilin Zhang, Yujin Park, and Joonhwan Lee. 2020. [A chatbot solution to chat app problems: Envisioning a chatbot counseling system for teenage victims of online sexual exploitation](#). In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI EA '20, page

1–7, New York, NY, USA. Association for Computing Machinery.

Mohammad Rafayet Ali, Seyedeh Zahra Razavi, Raina Langevin, Abdullah Al Mamun, Benjamin Kane, Reza Rawassizadeh, Lenhart K. Schubert, and Ehsan Hoque. 2020. [A virtual conversational agent for teens with autism spectrum disorder: Experimental results and design lessons](#). In *Proceedings of the 20th ACM International Conference on Intelligent Virtual Agents, IVA '20*, New York, NY, USA. Association for Computing Machinery.

Mohammad Rafayet Ali, Taylan Sen, Benjamin Kane, Shagun Bose, Thomas Carroll, Ronald Epstein, Lenhart K. Schubert, and Ehsan Hoque. 2021. [Novel computational linguistic measures, dialogue system and the development of sophie: Standardized online patient for healthcare interaction education](#). *IEEE Transactions on Affective Computing*, pages 1–1.

Robert J. Amdur and Chuck Biddle. 1997. [Institutional Review Board Approval and Publication of Human Research Results](#). *JAMA*, 277(11):909–914.

Masahiro Araki, Kana Shibahara, and Yuko Mizukami. 2011. [Spoken dialogue system for learning braille](#). In *2011 IEEE 35th Annual Computer Software and Applications Conference*, pages 152–156.

Suket Arora, Kamaljeet Batra, and Sarabjit Singh. 2013. [Dialogue system: A brief review](#). *CoRR*, abs/1306.4134.

Mikel Artetxe, Sebastian Ruder, Dani Yogatama, Gorka Labaka, and Eneko Agirre. 2020. [A call for more rigor in unsupervised cross-lingual learning](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7375–7388, Online. Association for Computational Linguistics.

Lekha Athota, Vinod Kumar Shukla, Nitin Pandey, and Ajay Rana. 2020. [Chatbot for healthcare system using artificial intelligence](#). In *2020 8th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions) (ICRITO)*, pages 619–622.

Saminda Sundeepa Balasuriya, Laurianne Sitbon, Andrew A. Bayor, Maria Hoogstrate, and Margot Brerton. 2018. [Use of voice activated interfaces by people with intellectual disability](#). In *Proceedings of the 30th Australian Conference on Computer-Human Interaction, OzCHI '18*, page 102–112, New York, NY, USA. Association for Computing Machinery.

R. V. Belfin, A. J. Shobana, Megha Manilal, Ashly Ann Mathew, and Blessy Babu. 2019. [A graph based chatbot for cancer patients](#). In *2019 5th International Conference on Advanced Computing Communication Systems (ICACCS)*, pages 717–721.

739	Emily M. Bender. 2009. Linguistically naïve != language independent: Why NLP needs linguistic typology . In <i>Proceedings of the EACL 2009 Workshop on the Interaction between Linguistics and Computational Linguistics: Virtuous, Vicious or Vacuous?</i> , pages 26–32, Athens, Greece. Association for Computational Linguistics.	794
740		795
741		796
742		797
743		798
744		799
745		
746	Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? . <i>Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency</i> .	800
747		801
748		802
749		803
750		804
751	E. Bertino, B.C. Ooi, Y. Yang, and R.H. Deng. 2005. Privacy and ownership preserving of outsourced medical data . In <i>21st International Conference on Data Engineering (ICDE'05)</i> , pages 521–532.	805
752		806
753		807
754		808
755	Timothy Bickmore and Toni Giorgino. 2004. Some novel aspects of health communication from a dialogue systems perspective. <i>AAAI Fall Symposium - Technical Report</i> .	809
756		810
757		811
758		812
759	Vildan Bilici, Emiel Krahmer, Saskia Riele, and Raymond Veldhuis. 2000. Preferred modalities in dialogue systems.	813
760		814
761		815
762	Antoine Bordes, Y-Lan Boureau, and Jason Weston. 2017. Learning end-to-end goal-oriented dialog .	816
763		817
764	Willem-Paul Brinkman, Dwi Hartanto, Ni Kang, Daniel de Vliegheer, Isabel L. Kampmann, Nexhmedin Morina, Paul G.M. Emmelkamp, and Mark Neerinx. 2012a. A virtual reality dialogue system for the treatment of social phobia . In <i>CHI '12 Extended Abstracts on Human Factors in Computing Systems</i> , CHI EA '12, page 1099–1102, New York, NY, USA. Association for Computing Machinery.	818
765		819
766		820
767		821
768		822
769		823
770		824
771		825
772		826
773	Willem-Paul Brinkman, Dwi Hartanto, Ni Kang, Daniel Vliegheer, Isabel Kampmann, Nexhmedin Morina, Paul Emmelkamp, and Mark Neerinx. 2012b. A virtual reality dialogue system for the treatment of social phobia . pages 1099–1102.	827
774		828
775		829
776		830
777		831
778	Jacqueline Brixey, Rens Hoegen, Wei Lan, Joshua Rusow, Karan Singla, Xusen Yin, Ron Artstein, and Anton Leuski. 2017. SHIHbot: A Facebook chatbot for sexual health information on HIV/AIDS . In <i>Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue</i> , pages 370–373, Saarbrücken, Germany. Association for Computational Linguistics.	832
779		833
780		834
781		835
782		836
783		837
784		838
785		839
786	Leonardo Campillos Llanos, Dhouha Bouamor, Éric Bilinski, Anne-Laure Ligozat, Pierre Zweigenbaum, and Sophie Rosset. 2015. Description of the Patient-Genesys dialogue system . In <i>Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 438–440, Prague, Czech Republic. Association for Computational Linguistics.	840
787		841
788		842
789		843
790		844
791		845
792		846
793		847
		848
		849
		850
		851
	Leonardo Campillos-Llanos, Catherine Thomas, Éric Bilinski, Pierre Zweigenbaum, and Sophie Rosset. 2020. Designing a virtual patient dialogue system based on terminology-rich resources: Challenges and evaluation . <i>Natural Language Engineering</i> , 26(2):183–220.	
	Bo-Wei Chen, Po-Yi Shih, Karunanithi Bharanitharan, Po-Chuan Lin, Jhing-Fa Wang, and Chiaming Chen. 2013. Customizable cloud-healthcare dialogue system based on lvcsr with prosodic-contextual post-processing . In <i>2013 1st International Conference on Orange Technologies (ICOT)</i> , pages 246–249.	
	Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. 2017. A survey on dialogue systems: Recent advances and new frontiers . <i>CoRR</i> , abs/1711.01731.	
	Ching-Hua Chuan and Susan Morgan. 2021. Creating and evaluating chatbots as eligibility assistants for clinical trials: An active deep learning approach towards user-centered classification . <i>ACM Trans. Comput. Healthcare</i> , 2(1).	
	Karl Daher, Jacky Casas, Omar Abou Khaled, and Elena Mugellini. 2020. Empathic chatbot response for medical assistance . In <i>Proceedings of the 20th ACM International Conference on Intelligent Virtual Agents, IVA '20</i> , New York, NY, USA. Association for Computing Machinery.	
	Prathyusha Danda, Brij Mohan Lal Srivastava, and Manish Shrivastava. 2016. Vaidya: A spoken dialog system for health domain . In <i>Proceedings of the 13th International Conference on Natural Language Processing</i> , pages 161–166, Varanasi, India. NLP Association of India.	
	Johan Oswin De Nieva, Jose Andres Joaquin, Chaste Bernard Tan, Ruzel Khyvin Marc Te, and Ethel Ong. 2020. Investigating students' use of a mental health chatbot to alleviate academic stress . In <i>6th International ACM In-Cooperation HCI and UX Conference, CHIUXID '20</i> , page 1–10, New York, NY, USA. Association for Computing Machinery.	
	Orianna Demasi, Yu Li, and Zhou Yu. 2020. A multi-persona chatbot for hotline counselor training . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3623–3636, Online. Association for Computational Linguistics.	
	Jan Deriu, Álvaro Rodrigo, Arantxa Otegi, Guillermo Echegoyen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. 2019. Survey on evaluation methods for dialogue systems . <i>CoRR</i> , abs/1905.04071.	
	David DeVault, Kallirroi Georgila, Ron Artstein, Fabrizio Morbini, David Traum, Stefan Scherer, Albert Skip Rizzo, and Louis-Philippe Morency. 2013. Verbal indicators of psychological distress in interactive dialogue with a virtual human . In <i>Proceedings of the SIGDIAL 2013 Conference</i> , pages 193–202, Metz, France. Association for Computational Linguistics.	

852	Alessandro Di Nuovo, Josh Bamforth, Daniela Conti,	analysis of current evaluation protocols . In <i>Proceed-</i>	907
853	Karen Sage, Rachel Ibbotson, Judy Clegg, Anna	<i>ings of the 21th Annual Meeting of the Special Inter-</i>	908
854	Westaway, and Karen Arnold. 2020. An explo-	<i>est Group on Discourse and Dialogue</i> , pages 236–	909
855	rative study on robotics for supporting children with	245, 1st virtual meeting. Association for Computa-	910
856	autism spectrum disorder during clinical procedures.	tional Linguistics.	911
857	In <i>Companion of the 2020 ACM/IEEE International</i>		
858	<i>Conference on Human-Robot Interaction, HRI '20,</i>	Chloë FitzGerald and Samia Hurst. 2017. Implicit	912
859	page 189–191, New York, NY, USA. Association for	bias in healthcare professionals: a systematic review.	913
860	Computing Machinery.	<i>BMC medical ethics</i> , 18(1):1–18.	914
861	Emily Dinan, Gavin Abercrombie, A. Stevie Bergman,	Floyd Garvey and Suresh Sankaranarayanan. 2012. In-	915
862	Shannon Spruit, Dirk Hovy, Y-Lan Boureau, and	telligent agent based flight search and booking sys-	916
863	Verena Rieser. 2021. Anticipating safety issues in	tem . <i>International Journal of Advanced Research in</i>	917
864	e2e conversational ai: Framework and tooling.	<i>Artificial Intelligence</i> , 1(4).	918
865	Francesca Dino, Rohola Zandie, Hojjat Abdollahi,	Sarik Ghazarian, Ralph M. Weischedel, Aram Gal-	919
866	Sarah Schoeder, and Mohammad H. Mahoor. 2019.	styan, and Nanyun Peng. 2019. Predictive en-	920
867	Delivering cognitive behavioral therapy using a con-	gagement: An efficient metric for automatic eval-	921
868	versational social robot . In <i>2019 IEEE/RSJ Interna-</i>	uation of open-domain dialogue systems . <i>CoRR,</i>	922
869	<i>tional Conference on Intelligent Robots and Systems</i>	abs/1911.01456.	923
870	<i>(IROS)</i> , pages 2089–2095.		
871	Jesse Dodge, Suchin Gururangan, Dallas Card, Roy	Nancy Green, William Lawton, and Boyd Davis. 2004.	924
872	Schwartz, and Noah A. Smith. 2019. Show your	An assistive conversation skills training system for	925
873	work: Improved reporting of experimental results.	caregivers of persons with alzheimer’s disease . In	926
874	In <i>Proceedings of the 2019 Conference on Empirical</i>	<i>Proceedings of the AAAI 2004 Fall Symposium on</i>	927
875	<i>Methods in Natural Language Processing and the</i>	<i>Dialogue Systems for Health Communication</i> .	928
876	<i>9th International Joint Conference on Natural Lan-</i>		
877	<i>guage Processing (EMNLP-IJCNLP)</i> , pages 2185–	Aditi Sharma Grover, Madelaine Plauché, Etienne	929
878	2194, Hong Kong, China. Association for Computa-	Barnard, and Christiaan Kuun. 2009. Hiv health	930
879	tional Linguistics.	information access using spoken dialogue systems:	931
880	Suma Reddy Duggenpudi, Kusampudi Siva Subra-	Touchtone vs. speech . In <i>Proceedings of the 3rd In-</i>	932
881	hamanyam Varma, and Radhika Mamidi. 2019.	<i>ternational Conference on Information and Commu-</i>	933
882	Samvaadhana: A Telugu dialogue system in hospi-	<i>nication Technologies and Development, ICTD’09,</i>	934
883	tal domain . In <i>Proceedings of the 2nd Workshop on</i>	page 95–107. IEEE Press.	935
884	<i>Deep Learning Approaches for Low-Resource NLP</i>		
885	<i>(DeepLo 2019)</i> , pages 234–242, Hong Kong, China.	Donghoon Ham, Jeong-Gwan Lee, Youngsoo Jang,	936
886	Association for Computational Linguistics.	and Kee-Eung Kim. 2020. End-to-end neural	937
887	Wilmer Stalin Erazo, Germán Patricio Guerrero, Car-	pipeline for goal-oriented dialogue systems using	938
888	los Carrión Betancourt, and Iván Sánchez Salazar.	GPT-2 . In <i>Proceedings of the 58th Annual Meet-</i>	939
889	2020. Chatbot implementation to collect data on	<i>ing of the Association for Computational Linguis-</i>	940
890	possible covid-19 cases and release the pressure	<i>tics</i> , pages 583–592, Online. Association for Com-	941
891	on the primary health care system . In <i>2020 11th</i>	putational Linguistics.	942
892	<i>IEEE Annual Information Technology, Electronics</i>		
893	<i>and Mobile Communication Conference (IEMCON),</i>	Helen Hastie. 2012. Metrics and Evaluation of Spoken	943
894	pages 0302–0307.	Dialogue Systems , pages 131–150.	944
895	Ahmed Fadhil and Ahmed Ghassan Tawfiq AbuRa’ed.	Samuel Holmes, Anne Moorhead, Raymond Bond,	945
896	2019. Ollobot - towards a text-based arabic health	Huiru Zheng, Vivien Coates, and Michael Mctear.	946
897	conversational agent: Evaluation and results . In	2019. Usability testing of a healthcare chatbot: Can	947
898	<i>RANLP</i> .	we use conventional methods to assess conversa-	948
899	Shahla Farzana, Mina Valizadeh, and Natalie Parde.	tional user interfaces? In <i>Proceedings of the 31st Eu-</i>	949
900	2020. Modeling dialogue in conversational cogni-	<i>ropean Conference on Cognitive Ergonomics, ECCE</i>	950
901	tive health screening interviews . In <i>Proceedings of</i>	2019, page 207–214, New York, NY, USA. Associa-	951
902	<i>the 12th Language Resources and Evaluation Con-</i>	tion for Computing Machinery.	952
903	<i>ference</i> , pages 1167–1177, Marseille, France. Euro-		
904	pean Language Resources Association.	Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu,	953
905	Sarah E. Finch and Jinho D. Choi. 2020. Towards uni-	Semih Yavuz, and Richard Socher. 2020. A simple	954
906	fied dialogue system evaluation: A comprehensive	language model for task-oriented dialogue . <i>CoRR,</i>	955
		abs/2005.00796.	956
		Matthew B. Hoy. 2018. Alexa, siri, cortana, and	957
		more: An introduction to voice assistants . <i>Medical</i>	958
		<i>Reference Services Quarterly</i> , 37(1):81–88. PMID:	959
		29327988.	960

961	Chin-Yuan Huang, Ming-Chin Yang, Chin-Yu Huang,	B. Amir H. Kargar and Mohammad H. Mahoor. 2017.	1016
962	Yu-Jui Chen, Meng-Lin Wu, and Kai-Wen Chen.	A pilot study on the ebear socially assistive robot:	1017
963	2018. A chatbot-supported smart wireless interac-	Implication for interacting with elderly people with	1018
964	tive healthcare system for weight control and health	moderate depression. In <i>2017 IEEE-RAS 17th In-</i>	1019
965	promotion. In <i>2018 IEEE International Conference</i>	<i>ternational Conference on Humanoid Robotics (Hu-</i>	1020
966	<i>on Industrial Engineering and Engineering Manage-</i>	<i>manoids)</i> , pages 756–762.	1021
967	<i>ment (IEEM)</i> , pages 1791–1795.		
968	Tae-Ho Hwang, JuHui Lee, Se-Min Hyun, and KangY-	Fakhri Karray, Milad Alemzadeh, Jamil Saleh, and	1022
969	oon Lee. 2020. Implementation of interactive health-	Mo Nours Arab. 2008. Human-computer interac-	1023
970	care advisor model using chatbot and visualiza-	tion: Overview on state of the art. <i>International</i>	1024
971	tion. In <i>2020 International Conference on Informa-</i>	<i>Journal on Smart Sensing and Intelligent Systems,</i>	1025
972	<i>tion and Communication Technology Convergence</i>	1:137–159.	1026
973	<i>(ICTC)</i> , pages 452–455.		
974	Koji Inoue, Pierrick Milhorat, Divesh Lala, Tianyu	William Kearns, Nai-Ching Chi, Yong Choi, Shih-Yin	1027
975	Zhao, and Tatsuya Kawahara. 2016. Talking with	Lin, Hilaire Thompson, and George Demiris. 2019.	1028
976	ERICA, an autonomous android. In <i>Proceedings</i>	A systematic review of health dialog systems. <i>Method-</i>	1029
977	<i>of the 17th Annual Meeting of the Special Interest</i>	<i>ods of Information in Medicine</i> , 58:179–193.	1030
978	<i>Group on Discourse and Dialogue</i> , pages 212–215,		
979	Los Angeles. Association for Computational Lin-	Liliana Laranjo, Adam Dunn, Huong Ly Tong, A. Baki	1031
980	guistics.	Kocaballi, Jessica Chen, Rabia Bashir, Didi Surian,	1032
981	Institute of Medicine. 2009. <i>Beyond the HIPAA Pri-</i>	Blanca Gallego, Farah Magrabi, Annie Lau, and En-	1033
982	<i>vacancy Rule: Enhancing Privacy, Improving Health</i>	rico Coiera. 2018. Conversational agents in health-	1034
983	<i>Through Research</i> . The National Academies Press,	care: A systematic review. <i>Journal of the American</i>	1035
984	Washington, DC.	<i>Medical Informatics Association</i> , 0.	1036
985	Hifza Javed, Myoungsoon Jeon, Ayanna Howard,	Choong Lee and Hyung-Jin Yoon. 2017. Medical big	1037
986	and Chung Hyuk Park. 2018. Robot-assisted	data: promise and challenges. <i>Kidney Research and</i>	1038
987	socio-emotional intervention framework for chil-	<i>Clinical Practice</i> , 36:3–11.	1039
988	dren with autism spectrum disorder. In <i>Companion</i>		
989	<i>of the 2018 ACM/IEEE International Conference on</i>	Dongkeon Lee, Kyo-Joong Oh, and Ho-Jin Choi. 2017.	1040
990	<i>Human-Robot Interaction, HRI '18</i> , page 131–132,	The chatbot feels you - a counseling service using	1041
991	New York, NY, USA. Association for Computing	emotional response generation. In <i>2017 IEEE Inter-</i>	1042
992	Machinery.	<i>national Conference on Big Data and Smart Com-</i>	1043
993	Daniel Jurafsky and James H. Martin. 2009. <i>Speech</i>	<i>puting (BigComp)</i> , pages 437–440.	1044
994	<i>and Language Processing (2nd Edition)</i> . Prentice-	Ju Yeon Lee, Ju Young Kim, Seung Ju You, You Soo	1045
995	Hall, Inc., USA.	Kim, Hye Yeon Koo, Jeong Hyun Kim, Sohye Kim,	1046
996	Dipesh Kadariya, Revathy Venkataramanan,	Jung Ha Park, Jong Soo Han, Siye Kil, Hyerim	1047
997	Hong Yung Yip, Maninder Kalra, Krishnaprasad	Kim, Ye Seul Yang, and Kyung Min Lee. 2019. De-	1048
998	Thirunarayanan, and Amit Sheth. 2019. kbot:	velopment and usability of a life-logging behavior	1049
999	Knowledge-enabled personalized chatbot for	monitoring application for obese patients. <i>Journal</i>	1050
1000	asthma self-management. In <i>2019 IEEE Inter-</i>	<i>of Obesity and Metabolic Syndrome</i> , 28(3):194–202.	1051
1001	<i>national Conference on Smart Computing</i>	Publisher Copyright: Copyright © 2019 Korean So-	1052
1002	<i>(SMARTCOMP)</i> , pages 138–143.	ciety for the Study of Obesity.	1053
1003	Vera C Kaelin, Mina Valizadeh, Zurisadai Salgado, Na-	Kyunghee Lee, Hyeyon Kwon, Byungtae Lee, Guna	1054
1004	talie Parde, and Mary A Khetani. 2021. Artificial	Lee, Jae Ho Lee, Yu Rang Park, and Soo-Yong Shin.	1055
1005	intelligence in rehabilitation targeting the partici-	2018. Effect of self-monitoring on long-term patient	1056
1006	pation of children and youth with disabilities: Scop-	engagement with mobile health applications. <i>PLOS</i>	1057
1007	ing review. <i>J Med Internet Res</i> , 23(11):e25745.	<i>ONE</i> , 13(7):1–12.	1058
1008	Takeshi Kamita, Atsuko Matsumoto, Boyu Sun, and	Yi-Chieh Lee, Naomi Yamashita, and Yun Huang.	1059
1009	Tomoo Inoue. 2020. Promotion of continuous use	2020a. Designing a chatbot as a mediator for pro-	1060
1010	of a self-guided mental healthcare system by a chat-	moting deep self-disclosure to a real mental health	1061
1011	bot. In <i>Conference Companion Publication of the</i>	professional. <i>Proc. ACM Hum.-Comput. Interact.</i> ,	1062
1012	<i>2020 on Computer Supported Cooperative Work and</i>	4(CSCW1).	1063
1013	<i>Social Computing, CSCW '20 Companion</i> , page	Yi-Chieh Lee, Naomi Yamashita, Yun Huang, and Wai	1064
1014	293–298, New York, NY, USA. Association for	Fu. 2020b. "i hear you, i feel you": Encouraging	1065
1015	Computing Machinery.	deep self-disclosure through a chatbot. In <i>Proceed-</i>	1066
		<i>ings of the 2020 CHI Conference on Human Factors</i>	1067
		<i>in Computing Systems, CHI '20</i> , page 1–12, New	1068
		York, NY, USA. Association for Computing Machin-	1069
		ery.	1070

1071	Peter Ljunglöf, Britt Claesson, Ingrid Mattsson Müller,	<i>Processing Systems</i> , SSPS 2020, page 91–95, New	1128
1072	Stina Ericsson, Cajsa Ottesjö, Alexander Berman,	York, NY, USA. Association for Computing Machin-	1129
1073	and Fredrik Kronlid. 2011. Lekbot: A talking	ery.	1130
1074	and playing robot for children with disabilities .		
1075	In <i>Proceedings of the Second Workshop on Speech</i>	Joao Luis Zeni Montenegro, Cristiano André da Costa,	1131
1076	<i>and Language Processing for Assistive Technolo-</i>	and Rodrigo da Rosa Righi. 2019. Survey of conver-	1132
1077	<i>gies</i> , pages 110–119, Edinburgh, Scotland, UK. As-	sational agents in health . <i>Expert Systems with Appli-</i>	1133
1078	sociation for Computational Linguistics.	<i>cations</i> , 129:56–67.	1134
1079	Peter Ljunglöf, Staffan Larsson, Katarina	Fabrizio Morbini, David DeVault, Kallirroi Georgila,	1135
1080	Heimann Mühlenbock, and Gunilla Thunberg.	Ron Artstein, David Traum, and Louis-Philippe	1136
1081	2009. TRIK: A talking and drawing robot for	Morency. 2014. A demonstration of dialogue pro-	1137
1082	children with communication disabilities .	cessing in SimSensei kiosk . In <i>Proceedings of the</i>	1138
1083	In <i>Proceedings of the 17th Nordic Conference of Com-</i>	<i>15th Annual Meeting of the Special Interest Group</i>	1139
1084	<i>putational Linguistics (NODALIDA 2009)</i> , pages	<i>on Discourse and Dialogue (SIGDIAL)</i> , pages 254–	1140
1085	275–278, Odense, Denmark. Northern European	256, Philadelphia, PA, U.S.A. Association for Com-	1141
1086	Association for Language Technology (NEALT).	putational Linguistics.	1142
1087	A. Loisel, N. Chaignaud, and J-Ph. Kotowicz.	Fabrizio Morbini, Eric Forbell, David DeVault, Kenji	1143
1088	2007. Designing a human-computer dialog sys-	Sagae, David Traum, and Albert Rizzo. 2012. A	1144
1089	tem for medical information search . In <i>2007</i>	mixed-initiative conversational dialogue system for	1145
1090	<i>IEEE/WIC/ACM International Conferences on Web</i>	healthcare . In <i>Proceedings of the 13th Annual Meet-</i>	1146
1091	<i>Intelligence and Intelligent Agent Technology -</i>	<i>ing of the Special Interest Group on Discourse and</i>	1147
1092	<i>Workshops</i> , pages 350–353.	<i>Dialogue</i> , pages 137–139, Seoul, South Korea. As-	1148
1093	Ramón López-Cózar Delgado and Masahiro Araki.	sociation for Computational Linguistics.	1149
1094	2005. <i>Spoken, Multilingual and Multimodal Dia-</i>	Denis Newman-Griffis, Jill Fain Lehman, Carolyn	1150
1095	<i>logue Systems: Development and Assessment</i> . Wi-	Rosé, and Harry Hochheiser. 2021. Translational	1151
1096	ley, Chichester, UK.	NLP: A new paradigm and general principles for nat-	1152
1097	Ryan Lowe, Nissan Pow, Iulian Serban, Laurent Char-	ural language processing research . In <i>Proceedings</i>	1153
1098	lin, Chia-Wei Liu, and Joelle Pineau. 2017. Train-	<i>of the 2021 Conference of the North American Chap-</i>	1154
1099	ing end-to-end dialogue systems with the ubuntu di-	<i>ter of the Association for Computational Linguistics:</i>	1155
1100	alogue corpus . <i>Dialogue and Discourse</i> , 8:31–65.	<i>Human Language Technologies</i> , pages 4125–4138,	1156
1101	Raju Maharjan, Per Bækgaard, and Jakob E. Bardram.	Online. Association for Computational Linguistics.	1157
1102	2019. "hear me out": Smart speaker based conversa-	Kyo-Joong Oh, Dongkun Lee, Byungsoo Ko, and Ho-	1158
1103	tional agent to monitor symptoms in mental health .	Jin Choi. 2017. A chatbot for psychiatric counsel-	1159
1104	In <i>Adjunct Proceedings of the 2019 ACM Interna-</i>	ing in mental healthcare service based on emotional	1160
1105	<i>tional Joint Conference on Pervasive and Ubiqui-</i>	dialogue analysis and sentence generation . In <i>2017</i>	1161
1106	<i>tous Computing and Proceedings of the 2019 ACM</i>	<i>18th IEEE International Conference on Mobile Data</i>	1162
1107	<i>International Symposium on Wearable Computers</i> ,	<i>Management (MDM)</i> , pages 371–375.	1163
1108	UbiComp/ISWC '19 Adjunct, page 929–933, New	Alexandros Papangelis, Robert Gatchel, Vangelis Met-	1164
1109	York, NY, USA. Association for Computing Machin-	sis, and Fillia Makedon. 2013. An adaptive dialogue	1165
1110	ery.	system for assessing post traumatic stress disorder .	1166
1111	Rohit Binu Mathew, Sandra Varghese, Sera Elsa Joy,	In <i>Proceedings of the 6th International Conference</i>	1167
1112	and Swanthana Susan Alex. 2019. Chatbot for dis-	<i>on PErvasive Technologies Related to Assistive En-</i>	1168
1113	ease prediction and treatment recommendation us-	<i>vironments</i> , PETRA '13, New York, NY, USA. As-	1169
1114	ing machine learning . In <i>2019 3rd International</i>	sociation for Computing Machinery.	1170
1115	<i>Conference on Trends in Electronics and Informat-</i>	Falguni Patel, Riya Thakore, Ishita Nandwani, and San-	1171
1116	<i>ics (ICOEI)</i> , pages 851–856.	tosh Kumar Bharti. 2019. Combating depression in	1172
1117	Michael McTear. 2010. Chapter 9 - the role of spo-	students using an intelligent chatbot: A cognitive be-	1173
1118	ken dialogue in user–environment interaction . In	havioral therapy . In <i>2019 IEEE 16th India Council</i>	1174
1119	Hamid Aghajan, Ramón López-Cózar Delgado, and	<i>International Conference (INDICON)</i> , pages 1–4.	1175
1120	Juan Carlos Augusto, editors, <i>Human-Centric Inter-</i>	Frano Petric, Damjan Miklic, and Zdenko Kovacic.	1176
1121	<i>faces for Ambient Intelligence</i> , pages 225–254. Aca-	2017. Robot-assisted autism spectrum disorder di-	1177
1122	ademic Press, Oxford.	agnostics using pomdps . In <i>Proceedings of the</i>	1178
1123	Sabrina J. Mielke. 2016. Language diversity in ACL	<i>Companion of the 2017 ACM/IEEE International</i>	1179
1124	2004 - 2016 .	<i>Conference on Human-Robot Interaction</i> , HRI '17,	1180
1125	Mahdi Naser Moghadasi, Yu Zhuang, and Hashim Gell-	page 369–370, New York, NY, USA. Association for	1181
1126	ban. 2020. Robo: A counselor chatbot for opioid ad-	Computing Machinery.	1182
1127	dicted patients . In <i>2020 2nd Symposium on Signal</i>		

1183	Marco Polignano, Fedelucio Narducci, Andrea Iovine,	Naohiro Shoji, Takayo Namba, and Keiichi Abe. 2020.	1240
1184	Cataldo Musto, Marco De Gemmis, and Giovanni	Proposal of spoken interactive home doctor system	1241
1185	Semeraro. 2020. Healthassistantbot: A personal	for elderly people . In <i>2020 IEEE 9th Global Confer-</i>	1242
1186	health assistant for the italian language . <i>IEEE Ac-</i>	ence on Consumer Electronics (GCCE) , pages 421–	1243
1187	cess , 8:107479–107497.	423.	1244
1188	A. Prange, Margarita Chikobava, P. Poller, Michael	Noppon Siangchin and Taweesak Samanchuen. 2019.	1245
1189	Barz, and D. Sonntag. 2017. A multimodal dialogue	Chatbot implementation for icd-10 recommenda-	1246
1190	system for medical decision support inside virtual	tion system . In <i>2019 International Conference on</i>	1247
1191	reality. In <i>SIGDIAL Conference</i> .	<i>Engineering, Science, and Industrial Applications</i>	1248
1192	Libo Qin, Yijia Liu, Wanxiang Che, Haoyang Wen,	<i>(ICESI)</i> , pages 1–6.	1249
1193	Yangming Li, and Ting Liu. 2019. Entity-consistent	Daneil Sonntag and Manuel Moller. 2010. Prototyping	1250
1194	end-to-end task-oriented dialogue system with kb re-	semantic dialogue systems for radiologists . In <i>2010</i>	1251
1195	triever .	<i>Sixth International Conference on Intelligent Envi-</i>	1252
1196	Juan C. Quiroz, Tristan Bongolan, and Kiran Ijaz.	<i>ronments</i> , pages 84–89.	1253
1197	2020. Alexa depression and anxiety self-tests: A	Daniel Sonntag, Gerhard Sonnenberg, Robert Neßel-	1254
1198	preliminary analysis of user experience and trust .	rath, and Gerd Herzog. 2009. Supporting a rapid	1255
1199	In <i>Adjunct Proceedings of the 2020 ACM Interna-</i>	dialogue engineering process .	1256
1200	<i>tional Joint Conference on Pervasive and Ubiqui-</i>	Prakhar Srivastava and Nishant Singh. 2020. Autom-	1257
1201	<i>tous Computing and Proceedings of the 2020 ACM</i>	atized medical chatbot (medibot) . In <i>2020 Interna-</i>	1258
1202	<i>International Symposium on Wearable Computers</i> ,	<i>ational Conference on Power Electronics IoT Applica-</i>	1259
1203	UbiComp-ISWC '20, page 494–496, New York, NY,	<i>tions in Renewable Energy and its Control (PARC)</i> ,	1260
1204	USA. Association for Computing Machinery.	pages 351–354.	1261
1205	Md. Moshir Rahman, Ruhul Amin, Md Nazmul	H. Strongman, R. Williams, W. Meeraus, T. Murray-	1262
1206	Khan Liton, and Nahid Hossain. 2019. Disha: An	Thomas, J. Campbell, L. Carty, D. Dedman, A. Gal-	1263
1207	implementation of machine learning based bangla	lagher, J. Oyinlola, A. Kousoulis, and J. Valentine.	1264
1208	healthcare chatbot . In <i>2019 22nd International Con-</i>	2019. Limitations for health research with restricted	1265
1209	<i>ference on Computer and Information Technology</i>	data collection from uk primary care . <i>Pharmacoepi-</i>	1266
1210	<i>(ICCIT)</i> , pages 1–6.	demiol Drug Saf .	1267
1211	Michelle L. Rogers, Emily S. Patterson, Roger J. Chap-	Bo-Hao Su, Shih-Pang Tseng, Yu-Shan Lin, and Jhing-	1268
1212	man, and Marta L. Render. 2005. Usability testing	Fa Wang. 2018. Health care spoken dialogue system	1269
1213	and the relation of clinical information systems to	for diagnostic reasoning and medical product recom-	1270
1214	patient safety.	mendation . In <i>2018 International Conference on Or-</i>	1271
1215	Nudtaporn Rosruen and Taweesak Samanchuen. 2018.	<i>ange Technologies (ICOT)</i> , pages 1–4.	1272
1216	Chatbot utilization for medical consultant system .	Konstantinos Tsiakas, Lynette Watts, Cyril Lutterodt,	1273
1217	In <i>2018 3rd Technology Innovation Management</i>	Theodoros Giannakopoulos, Alexandros Papange-	1274
1218	<i>and Engineering Science International Conference</i>	lis, Robert Gatchel, Vangelis Karkaletsis, and Fillia	1275
1219	<i>(TIMES-iCON)</i> , pages 1–5.	Makedon. 2015. A multimodal adaptive dialogue	1276
1220	Sanket Sanjay Sadavarte and Eliane Bodanese. 2019.	manager for depressive and anxiety disorder screen-	1277
1221	Pregnancy companion chatbot using alexa and ama-	ing: A wizard-of-oz experiment . In <i>Proceedings</i>	1278
1222	zon web services . In <i>2019 IEEE Pune Section Inter-</i>	<i>of the 8th ACM International Conference on PErva-</i>	1279
1223	<i>national Conference (PuneCon)</i> , pages 1–5.	<i>sive Technologies Related to Assistive Environments</i> ,	1280
1224	Iulian Vlad Serban, Ryan Lowe, Peter Henderson, Lau-	PETRA '15, New York, NY, USA. Association for	1281
1225	rent Charlin, and Joelle Pineau. 2015. A survey of	Computing Machinery .	1282
1226	available corpora for building data-driven dialogue	Lorainne Tudor Car, Dhakshenya Ardhithy Dhina-	1283
1227	systems . <i>CoRR</i> , abs/1512.05742.	garan, Bhone Myint Kyaw, Tobias Kowatsch, Shafiq	1284
1228	Bhuvan Sharma, Harshita Puri, and Deepika Rawat.	Joty, Yin-Leng Theng, and Rifat Atun. 2020. Con-	1285
1229	2018. Digital psychiatry - curbing depression us-	versational agents in health care: Scoping review	1286
1230	ing therapy chatbot and depression analysis . In <i>2018</i>	and conceptual analysis . <i>J Med Internet Res</i> ,	1287
1231	<i>Second International Conference on Inventive Com-</i>	22(8):e17158.	1288
1232	<i>munication and Computational Technologies (ICI-</i>	A. Vaidyam, Hannah Wisniewski, J. Halamka, M. S.	1289
1233	<i>CCT)</i> , pages 627–631.	Kashavan, and J. Torous. 2019. Chatbots and con-	1290
1234	Tianhao She, Xin Kang, Shun Nishide, and Fuji Ren.	versational agents in mental health: A review of	1291
1235	2018. Improving leo robot conversational abil-	the psychiatric landscape . <i>The Canadian Journal of</i>	1292
1236	ity via deep learning algorithms for children with	Psychiatry , 64:456 – 464.	1293
1237	autism . In <i>2018 5th IEEE International Confer-</i>		
1238	<i>ence on Cloud Computing and Intelligence Systems</i>		
1239	<i>(CCIS)</i> , pages 416–420.		

1294	Mina Valizadeh, Pardis Ranjbar-Noiey, Cornelia Caragea, and Natalie Parde. 2021. Identifying medical self-disclosure in online communities . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4398–4408, Online. Association for Computational Linguistics.	1349
1295		1350
1296		1351
1297		1352
1298		1353
1299		1354
1300		1355
1301		1356
1302	Richard M Walker, Gene A Brewer, M Jin Lee, Nicolai Petrovsky, and Arjen van Witteloostuijn. 2018. Best Practice Recommendations for Replicating Experiments in Public Administration . <i>Journal of Public Administration Research and Theory</i> , 29(4):609–626.	1357
1303		1358
1304		1359
1305		1360
1306		1361
1307		1362
1308	Jinping Wang, Hyun Yang, Ruosi Shao, Saeed Abdullah, and S. Shyam Sundar. 2020. Alexa as coach: Leveraging smart speakers to build social agents that reduce public speaking anxiety . In <i>Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems</i> , CHI '20, page 1–13, New York, NY, USA. Association for Computing Machinery.	1363
1309		1364
1310		1365
1311		1366
1312		1367
1313		1368
1314		1369
1315	J. V. Waterschoot, Iris Hendrickx, Arif Khan, E. Klappers, M. D. Korte, H. Strik, C. Cucchiari, and M. Theune. 2020. Bliss: An agent for collecting spoken dialogue data about health and well-being . In <i>LREC</i> .	1370
1316		1371
1317		1372
1318		1373
1319		1374
1320	Zhongyu Wei, Qianlong Liu, Baolin Peng, Huaixiao Tou, Ting Chen, Xuanjing Huang, Kam-fai Wong, and Xiangying Dai. 2018. Task-oriented dialogue system for automatic diagnosis . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 201–207, Melbourne, Australia. Association for Computational Linguistics.	1370
1321		1371
1322		1372
1323		1373
1324		1374
1325		1375
1326		1376
1327		1377
1328	Charles Welch, Allison Lahkala, Veronica Perez-Rosas, Siqi Shen, Sarah Seraj, Larry An, Kenneth Resnicow, James Pennebaker, and Rada Mihalcea. 2020. Expressive interviewing: A conversational system for coping with COVID-19 . In <i>Proceedings of the 1st Workshop on NLP for COVID-19 (Part 2) at EMNLP 2020</i> , Online. Association for Computational Linguistics.	1378
1329		1379
1330		1380
1331		1381
1332		1382
1333		1383
1334		1384
1335		1385
1336	Tsung-Hsien Wen, Milica Gasic, Nikola Mrksic, Lina Maria Rojas-Barahona, Pei-Hao Su, Stefan Ultes, David Vandyke, and Steve J. Young. 2016. A network-based end-to-end trainable task-oriented dialogue system . <i>CoRR</i> , abs/1604.04562.	1386
1337		1387
1338		1388
1339		1389
1340		1390
1341	Angang Xu, Zhe Liu, Yufan Guo, Vibha Sinha, and Rama Akkiraju. 2017. A new chatbot for customer service on social media . In <i>Proceedings of the 2017 CHI conference on human factors in computing systems</i> , pages 3506–3510.	1391
1342		1392
1343		1393
1344		1394
1345		1395
1346	Jing Xu, Da Ju, Margaret Li, Y-Lan Boureau, Jason Weston, and Emily Dinan. 2021. Recipes for safety in open-domain chatbots .	1396
1347		1397
1348		1398
	L. Xu, Q. Zhou, Ke Gong, Xiaodan Liang, Jianheng Tang, and Liang Lin. 2019. End-to-end knowledge-routed relational dialogue system for automatic diagnosis . <i>ArXiv</i> , abs/1901.10623.	1399
		1400
	Keigo Yabuki and Kaoru Sumi. 2018. Learning support system for effectively conversing with individuals with autism using a humanoid robot . In <i>2018 IEEE International Conference on Systems, Man, and Cybernetics (SMC)</i> , pages 4266–4270.	1401
		1402
	Akihiro Yorita, Simon Egerton, Carina Chan, and Naoyuki Kubota. 2020. Chatbot for peer support realization based on mutual care . In <i>2020 IEEE Symposium Series on Computational Intelligence (SSCI)</i> , pages 1601–1606.	1403
		1404
	Zheng Zhang, Ryuichi Takanobu, Minlie Huang, and Xiaoyan Zhu. 2020. Recent advances and challenges in task-oriented dialog system . <i>CoRR</i> , abs/2003.07490.	1405
		1406
	Yin Jiang Zhao, Yan Ling Li, and Min Lin. 2019. A review of the research on dialogue management of task-oriented systems . 1267:012025.	1407
		1408
	Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum. 2020. The design and implementation of xiaoice, an empathetic social chatbot . <i>Computational Linguistics</i> , 46(1):53–93.	1409
		1410

A Multi-Objective Systems

Multi-Objective System	# Papers
Diagnosis + Assistance	7
Diagnosis + Intervention	2
Diagnosis + Monitoring	1
Diagnosis + Counseling	1
Intervention + Monitoring	2
Intervention + Assistance	1
Assistance + Counseling	2
Intervention + Monitoring + Diagnosis	2
Intervention + Monitoring + Assistance	2
Intervention + Monitoring + Counseling	1
Diagnosis + Monitoring + Counseling	1
Diagnosis + Assistance + Intervention	2
Diagnosis + Intervention + Monitoring + Assistance	1

Table 8: Distribution of varying combinations of multiple system objectives across the surveyed papers.

Conversational agents seek to generate dialogues that have value to their end-users. We categorized included articles as having one or more of the following objectives: diagnosis, monitoring, intervention, counseling, or assistance. We found that 25

out of 70 surveyed systems were designed for more than one target objective, and provide additional details describing these multi-objective systems in Table 8.

B Multi-Device Systems

Multi-Device Category	# Papers
Desktop/Laptop + Mobile-based	8
Desktop/Laptop + VE	5
Desktop/Laptop + Robot	2
Mobile-based + PDA systems	2
Desktop/Laptop + GUI	1
Desktop/Laptop + PDA systems	1
Mobile-based + VE	1

Table 9: Details regarding the distribution of multi-device systems across the surveyed papers (20 total).

User Population	# Papers
Lab Experiments	15
Field Experiments	17
Crowdsourcing	1
Not Specified	4

Table 10: Distribution of user populations across the surveyed papers that conducted a human evaluation.

Human Evaluation Type	# Papers
Interact with the System	8
Rate a Dialogue	1
Both	28

Table 11: Distribution of evaluation types across the surveyed papers that conducted a human evaluation.

Many of the surveyed systems functioned using multiple device types. Table 9 shows the distribution of included devices across all multi-device systems. We found that the most common multi-device pairing was systems operating using computers and mobile devices.

C Additional Evaluation Details

From among the surveyed systems that conducted system and/or human evaluations, we further examined the types of evaluations conducted. Table 10 describes the populations leveraged for human evaluation across the surveyed systems, and Table

Type of System Evaluation	# Papers
Task Completion	4
Task Performance	9
Response Correctness	5
Naturalness	2
Response Time	3
Routing Time	1

Table 12: Type of system evaluation across the surveyed papers.

11 presents broad categories of the types of human evaluations conducted. We found that most human evaluations were conducted in a laboratory or field setting, and often included opportunities for participants to both interact with the system directly, and rate the quality of the dialogue. Table 12 details the various types of system evaluations conducted across the surveyed systems. We found that the most common assessment item in system evaluations was the system’s overall task performance.

D Included Papers

In this systematic review, we investigated 4070 papers involving dialogue systems for healthcare applications, identifying 70 papers that satisfied our defined inclusion criteria. We comprehensively analyzed these papers on the basis of numerous technical factors. We provide aggregated statistics for each of these categories in the main body of the paper. In Table 13 beginning on the following page, we provide a listing of each included paper and its categorization across all included classes. Full references for each included paper can be found in the bibliography.

Paper	DS Arch.	DM Arch.	Mod.	Device	Sys. Obj.	Engagement	Dom. of Research	Target Aud.	Lang.	Eval. Method	Dataset Size
Papangelis et al. (2013)	Pipeline	Intent-based	Multi-Modal	Desk /Lap	Monitoring, Intervention, Diagnosis	Yes	PTSD	Patients	English	Not Specified	Not Specified
Brinkman et al. (2012a)	Pipeline	Rule-based	Speech	Virtual Environment	Monitoring, Diagnosis	No	Social Phobia	Clinicians	English	Human Evaluation	Not Specified
Ali et al. (2020)	Pipeline	Intent-based	Speech	Desk /Lap	Monitoring, Assistance, Intervention	Yes	Autism Spectrum Disorder	Patients	English	Human Evaluation	46 videos
Tsiakas et al. (2015)	Pipeline	Intent-based	Multi-Modal	Desk /Lap, Virtual Environment	Diagnosis, Assistance	Yes	Anxiety Disorders, Depression, PTSD	Patients	English	Human Evaluation	90 speech segments
Wang et al. (2020)	Pipeline	Hybrid	Speech	PDA	Intervention	Yes	Social Phobia	Patients	English	Human Evaluation	Not Specified
Balasuriya et al. (2018)	Pipeline	Hybrid	Speech, GUI	PDA	Monitoring	Yes	Intellectual Disability	Patients	English	Human Evaluation	Not Specified
Chuan and Morgan (2021)	Pipeline	Intent-based	Speech	Desk /Lap	Assistance	No	Clinical Application	Patients	English	Human Evaluation	Not Specified
Grover et al. (2009)	Pipeline	Rule-based	Speech	Telephone	Assistance	No	HIV	Clinicians	Setswana	Human & Automated Evaluation	Not Specified
Petric et al. (2017)	Pipeline	Intent-based	Speech	Robot	Diagnosis	No	Autism Spectrum Disorder	Clinicians	English	Human Evaluation	Not Specified
Javed et al. (2018)	Not Specified	Not Specified	Speech, GUI	Robot	Monitoring	Yes	Autism Spectrum Disorder	Patients	English	Human Evaluation	Not Specified

Di Nuovo et al. (2020)	Not Specified	Not Specified	Speech	Robot	Monitoring	Yes	Autism Spectrum Disorder	Patients, Caregivers	English	Human Evaluation	Not Specified
Quiroz et al. (2020)	Pipeline	Hybrid	Speech	PDA, Mobile	Diagnosis, Intervention	Yes	Depression, Anxiety	Patients	English	Human Evaluation	Not Specified
Maharjan et al. (2019)	Pipeline	Hybrid	Speech	PDA, Mobile	Monitoring	No	Mental Health	Patients	English	Not Specified	Not Specified
Ahn et al. (2020)	Not Specified	Not Specified	Text	Mobile	Intervention, Assistance	Yes	Online sexual exploitation, PTSD	Patients	Korean	Not Specified	Not Specified
Kamita et al. (2020)	Not Specified	Not Specified	Text	Mobile	Intervention	Yes	Cognitive Behavioral Therapy, stress reduction	Patients	Japanese	Human Evaluation	Not Specified
Lee et al. (2020b)	Pipeline	Hybrid	Speech	Mobile	Monitoring	Yes	Health-related Self-disclosure	Patients	English	Human Evaluation	Not Specified
Moghadasi et al. (2020)	Pipeline	Hybrid	Text	Desk/Lap, Mobile	Assistance, Counseling	No	Opioid Addiction	Patients	English	Not Specified	20,494 records
De Nieva et al. (2020)	Pipeline	Hybrid	Text	Mobile	Monitoring, Intervention, Counseling	Yes	Anxiety, Depression	Patients	English	Human & Automated Evaluation	Not Specified
Lee et al. (2020a)	Pipeline	Hybrid	Text	Mobile	Monitoring	Yes	Health-related Self-disclosure	Patients	English	Human Evaluation	Not Specified
Daher et al. (2020)	Pipeline	Rule-based	GUI	Not Specified	Monitoring	No	Empathy for medical Assistance	Patients	English	Human Evaluation	Not Specified
Holmes et al. (2019)	Pipeline	Hybrid	Multi-Modal	Mobile	Assistance	Yes	Weight Loss	Patients	English	Human & Automated Evaluation	Not Specified

Oh et al. (2017)	Pipeline	Intent-based	Multi-Modal	Mobile	Diagnosis, Monitoring, Intervention	Yes	Psychiatric Counseling	Patients	Korean	Not Specified	49,846,477 records
Dino et al. (2019)	Pipeline	Rule-based	Speech	Robot	Intervention	Yes	Depression	Patients	English	Human Evaluation	Not Specified
Patel et al. (2019)	Not Specified	Not Specified	Text	Not Specified	Diagnosis	No	Stress, Depression	Patients	English	Not Specified	7,652 records, ISEAR dataset
Sharma et al. (2018)	Not Specified	Not Specified	Text	Mobile	Diagnosis, Intervention, Assistance	No	Depression	Patients	Not Specified	Not Specified	Not Specified
Belfin et al. (2019)	Pipeline	Intent-based	Multi-Modal	Desk /Lap, Mobile	Assistance	No	Cancer	Patients	English	Not Specified	Not Specified
Yorita et al. (2020)	Pipeline	Rule-based	Multi-Modal	Mobile	Diagnosis, Counseling	No	Stress Management	Clinicians	English	Not Specified	Not Specified
Kargar and Ma-hoor (2017)	Pipeline	Rule-based	Speech	Robot	Intervention	Yes	Depression	Patients	English	Human Evaluation	Not Specified
Hwang et al. (2020)	Pipeline	Rule-based	Text	Not Specified	Diagnosis, Intervention	No	Medical Assistance	Patients	Korean	Not Specified	Not Specified
Srivastava and Singh (2020)	Pipeline	Rule-based	Text	Not Specified	Diagnosis, Assistance	Yes	Disease Diagnosis	Patients	English	Human Evaluation	Not Specified
Mathew et al. (2019)	Pipeline	Rule-based	Text	Mobile	Diagnosis, Assistance	Yes	Disease Diagnosis	Patients	English	Human Evaluation	Not Specified
Athota et al. (2020)	Pipeline	Rule-based	Multi-Modal	Mobile	Diagnosis, Assistance	No	Disease Diagnosis	Patients	English	Not Specified	Not Specified
Sadavarte and Bodanese (2019)	Pipeline	Hybrid	Multi-Modal	PDA	Assistance	No	Pregnancy	Patients	English	Human Evaluation	Not Specified
Lee et al. (2017)	Pipeline	Hybrid	Text	Mobile	Counseling	Yes	Psychiatric Counseling	Patients	Korean	Not Specified	Not Specified

Rahman et al. (2019)	Pipeline	Hybrid	Text	Not Specified	Diagnosis, Monitoring, Counseling	No	Medical Assistance	Patients	Bengali	Automated Evaluation	4,961 records
Yabuki and Sumi (2018)	Not Specified	Not Specified	Speech	Robot	Intervention	No	Autism Spectrum Disorder	Care-givers	English	Not Specified	Not Specified
Su et al. (2018)	Pipeline	Intent-based	Speech	Not Specified	Diagnosis, Assistance	No	Disease Diagnosis	Patients	Chinese	Automated Evaluation	Not Specified
Shoji et al. (2020)	Not Specified	Not Specified	Speech	Desk /Lap, PDA	Diagnosis	No	Pneumonia	Patients	Not Specified	Automated Evaluation	Not Specified
Polignano et al. (2020)	Pipeline	Hybrid	Multi-Modal	Mobile	Diagnosis, Intervention, Assistance, Monitoring	No	Medical Assistance	Patients	Italian	Human & Automated Evaluation	1,865,700 records
Ali et al. (2021)	Pipeline	Hybrid	Speech	Desk /Lap, Virtual Environment	Intervention	No	Cancer	Clinicians	English	Automated Evaluation	382 transcripts of conversations
Aarabi (2013)	Pipeline	Intent-based	Text	Not Specified	Diagnosis	No	Cardiology	Patients	English	Not Specified	Not Specified
Loisel et al. (2007)	Pipeline	Hybrid	Text	Not Specified	Assistance	No	Medical Assistance	Patients	French	Not Specified	Not Specified
Rosruen and Samanchuen (2018)	Pipeline	Hybrid	Multi-Modal	Desk /Lap, Mobile	Assistance	No	Medical Assistance	Patients	Chinese	Automated Evaluation	Not Specified
Sonntag and Moller (2010)	Pipeline	Intent-based	Multi-Modal	Desk /Lap	Assistance	Yes	Radiology	Clinicians	Not Specified	Human & Automated Evaluation	Not Specified
Kadariya et al. (2019)	Pipeline	Hybrid	Multi-Modal	Mobile	Monitoring, Intervention	Yes	Asthma	Patients	English	Human & Automated Evaluation	Not Specified

Siangchin and Samanchuen (2019)	Pipeline	Hybrid	Text	Mobile	Assistance	No	Medical Assistance	Clinicians	Chinese	Human & Automated Evaluation	Not Specified
Erazo et al. (2020)	Pipeline	Rule-based	Text	Desk /Lap, Mobile	Diagnosis, Assistance	No	COVID-19	Patients	Not Specified	Human Evaluation	Not Specified
Huang et al. (2018)	Pipeline	Hybrid	Multi-Modal	Mobile	Monitoring, Intervention	Yes	Weight Loss	Patients	English, Chinese	Not Specified	Not Specified
Chen et al. (2013)	Pipeline	Rule-based	Speech	Desk /Lap, Mobile	Assistance	No	Medical Assistance	Patients, Caregivers	Chinese	Human Evaluation	MAT 400 dataset
Araki et al. (2011)	Pipeline	Intent-based	Multi-Modal	Desk /Lap	Intervention	No	Visually Impaired	Patients	Japanese	Human Evaluation	Not Specified
She et al. (2018)	End-to-End	Not Applicable	Speech	Robot	Intervention	Yes	Autism Spectrum Disorder	Patients	English	Automated Evaluation	Tager-Flusberg, Nadig ASD English, and Rollins Corpus
Yabuki and Sumi (2018)	Not Specified	Not Specified	Speech	Robot	Intervention	Yes	Autism Spectrum Disorder	Caregivers	Japanese	Not Specified	Self-Constructed dataset
Wei et al. (2018)	Pipeline	Intent-based	Text	Not Specified	Diagnosis	No	Medical Assistance	Clinicians	Chinese	Automated Evaluation	Self-Constructed dataset
Fadhil and AbuRa'ed (2019)	Pipeline	Intent-based	Multi-Modal	Mobile	Monitoring, Assistance, Intervention	No	Medical Assistance	Patients	Arabic	Human Evaluation	Not Specified
Demasi et al. (2020)	Pipeline	Intent-based	Text	Not Specified	Counseling	No	Mental Health	Patients	English	Human Evaluation	Self-Constructed dataset
Waterschoot et al. (2020)	Pipeline	Intent-based	Speech	Not Specified	Monitoring	No	Mental Health	Patients	Dutch	Not Specified	Self-Constructed dataset

Danda et al. (2016)	Pipeline	Hybrid	Speech	Desk /Lap, Mobile	Diagnosing, Intervention, Assistance	No	Medical Assistance	Patients	Indian	Human & Automated Evaluation	CMU arctic dataset
Duggenpudi et al. (2019)	Pipeline	Rule-based	Text	Not Specified	Assistance	No	Medical Assistance	Patients	Telugu	Human Evaluation	Self-Constructed dataset
Prange et al. (2017)	Pipeline	Rule-based	Multi-Modal	Mobile	Assistance	No	Medical Assistance	Clinicians	Not Specified	Not Specified	475 records
Campillos Llanos et al. (2015)	Pipeline	Intent-based	Multi-Modal	Not Specified	Intervention	No	Medical Assistance	Clinicians	French	Not Specified	Not Specified
Welch et al. (2020)	Pipeline	Intent-based	Text	Not Specified	Counseling, Assistance	Yes	Mental Health	Patients	Not Specified	Human Evaluation	Not Specified
Ljunglöf et al. (2009)	Pipeline	Intent-based	Speech	Desk /Lap, Robot	Intervention	No	Communication Disorders	Patients	Swedish	Human Evaluation	Not Specified
Ljunglöf et al. (2011)	Pipeline	Intent-based	Speech	Desk /Lap, Robot	Intervention	Yes	Communication Disorders	Patients	Swedish	Human Evaluation	Not Specified
Brixey et al. (2017)	Pipeline	Hybrid	Text	Desk /Lap, Mobile	Assistance	No	HIV	Patients	English	Human Evaluation	Self-Constructed dataset
Morbini et al. (2014)	Pipeline	Rule-based	Speech	Desk /Lap, Virtual Environment	Counseling	Yes	Mental Health	Patients	English	Not Specified	Not Specified
DeVault et al. (2013)	Not Specified	Not Specified	Speech	Desk /Lap, Virtual Environment	Diagnosis	No	Mental Health	Clinicians	English	Not Specified	Not Specified
Inoue et al. (2016)	Pipeline	Rule-based	Multi-Modal	Mobile, Virtual Environment	Counseling	Yes	Mental Health	Patients	Not Specified	Not Specified	Not Specified
Morbini et al. (2012)	Pipeline	Intent-based	Text	Desk /Lap, Mobile	Counseling	Yes	PTSD	Patients	English	Not Specified	Not Specified
Xu et al. (2019)	End-to-End	Not Applicable	Text	Not Specified	Diagnosis	No	Disease Diagnosis	Patients	Chinese	Human & Automated Evaluation	Self-Constructed dataset

Green et al. (2004)	Pipeline	Rule- based	Speech	Desk /Lap	Inter- ven- tion	No	Dementia	Care- givers	English	Human Evalu- ation	Not Spec- ified
---------------------------	----------	----------------	--------	--------------	------------------------	----	----------	-----------------	---------	--------------------------	-----------------------
