

LLM-Powered Conversational Voice Assistants: Interaction Patterns, Opportunities, Challenges, and Design Guidelines

AMAMA MAHMOOD, Johns Hopkins University, USA

JUNXIANG WANG, Johns Hopkins University, USA

BINGSHENG YAO, Rensselaer Polytechnic Institute, USA

DAKUO WANG, Northeastern University, USA

CHIEN-MING HUANG, Johns Hopkins University, USA



User: [the pharmacist] asked for a specific brand and I don't really know anything about cough syrup. Could you give me a specific brand ...

VA: There isn't a specific brand that applies universally to everyone...

VA: However, you can ask the pharmacist for recommendations ... The pharmacist will take into account your symptoms and medical history to guide you towards the most suitable options available.

User: The pharmacist told me that they wanted to take your opinion into account... I need to know the names of some of the brands... Could you recommend me some brand names?

VA: Certainly! Some common brand names of over-the-counter cough medicines with expectorants include Mucinex, Robitussin, and Chest Congestion Relief by CVS ...

VA: However, it's always a good idea to consult with a pharmacist for personalized recommendations based on your symptoms and medical history.

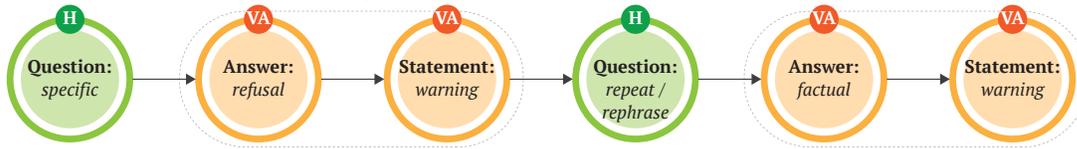


Fig. 1. We explore user interactions with an LLM-powered voice assistant in three distinct scenarios: medical self-diagnosis, creative planning, and debate with an opinionated AI. We report interaction patterns and breakdowns based on the style of speech used during the conversations. The interaction pattern and example conversation above depict ChatGPT's reluctance to answer *specific* medical queries, such as requests for medication brand names. However, upon re-asking, ChatGPT lists brands with an accompanying *warning* (a statement informing the user that ChatGPT is not an expert and that they should consult an expert).

Authors' addresses: Amama Mahmood, amama.mahmood@jhu.edu, Johns Hopkins University, 3400 N. Charles St, Baltimore, Maryland, USA, 21218; Junxiang Wang, jwang334@jhu.edu, Johns Hopkins University, 3400 N. Charles St, Baltimore, Maryland, USA, 21218; Bingsheng Yao, yaob@rpi.edu, Rensselaer Polytechnic Institute, 110 Eighth Street, Troy, New York, USA, 12180; Dakuo Wang, d.wang@northeastern.edu, Northeastern University, 360 Huntington Ave, Boston, Massachusetts, USA, 02115; Chien-Ming Huang, chienming.huang@jhu.edu, Johns Hopkins University, 3400 N. Charles St, Baltimore, Maryland, USA, 21218.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

Conventional Voice Assistants (VAs) rely on traditional language models to discern user intent and respond to their queries, leading to interactions that often lack a broader contextual understanding, an area in which Large Language Models (LLMs) excel. However, current LLMs are largely designed for text-based interactions, thus making it unclear how user interactions will evolve if their modality is changed to voice. In this work, we investigate whether LLMs can enrich VA interactions via an exploratory study with participants ($N=20$) using a ChatGPT-powered VA for three scenarios (medical self-diagnosis, creative planning, and debate) with varied constraints, stakes, and objectivity. We observe that LLM-powered VA elicits richer interaction patterns that vary across tasks, showing its versatility. Notably, LLMs absorb the majority of VA intent recognition failures. We additionally discuss the potential of harnessing LLMs for more resilient and fluid user-VA interactions and provide design guidelines for tailoring LLMs for voice assistance.

CCS Concepts: • **Human-centered computing** → **Empirical studies in HCI**; • **Computing Methodologies** → *Artificial intelligence*.

Additional Key Words and Phrases: voice assistant, LLMs, voice interactions, ChatGPT, conversational assistants, conversational AI

1 INTRODUCTION

Voice assistants (VAs) are well integrated into consumer technologies such as mobile phones, smart watches, smart speakers, and cars [54], and can significantly influence user behavior [60]. While commercial VAs such as Alexa and Siri rely on traditional language models to process user requests [1, 26], they mainly use rule-based keyword recognition mechanisms to determine user intent and fall short of maintaining coherent multi-turn conversations [16]. Furthermore, these interactions are often disrupted by unavoidable errors (*e.g.*, transcription and intent recognition errors) requiring users to interject and rectify breakdowns [49, 53]. Such constraints often restrict VAs’ primary use to basic functional tasks, such as setting alarms, sending texts, and seeking general information (*e.g.*, weather and time) [2, 3, 15].

Conversely, recent advancements in natural language processing endow large language models (LLMs) with the remarkable ability to generate coherent and contextually-aware text, bridging the gap between text generation and the dynamic nature of human language [5, 13, 22]. While LLMs have shown potential in various applications [58] such as health care [13, 37, 56], education [52], and collaborative writing [33, 45], most of these interactions are text-centric. The pronounced capabilities of LLMs, coupled with the intrinsic differences between text- and voice-based interactions [41], propel our exploration into: 1) *What new and distinct interaction patterns (beyond single-turn inquiries) may emerge when users interact with a voice assistant powered by LLM capabilities?* and 2) *How may LLMs’ contextual understanding capabilities help reduce the errors and conversational breakdowns common in current commercial VAs?*

To answer these questions, we first prototyped an LLM-powered conversational VA by integrating ChatGPT into an Alexa skill. This integration involved designing a conversational framework, using speech fillers [59] and small talk [19, 69], to handle ChatGPT API delays and Alexa timeout issues. We then conducted an exploratory study to probe how people interact with this ChatGPT-powered VA. To gain a broader, holistic understanding of user interactions, we contextualized our study via three scenarios with distinct characterizations—medical self-diagnosis, creative trip planning, and discussion—that near-future VAs may engage in (Fig. 3); in particular, the first two scenarios encompass assisted decision-making, while the third is purely conversational and argumentative. Consequently, we are also interested to see if intrinsic characteristics and conversation goals affect user interaction patterns and breakdowns.

Through thematic analysis, we found common and scenario-specific interaction patterns. Medical queries from participants led to factual VA responses with warnings (Fig. 1). For creative trip planning, the VA gave descriptive answers to generic questions and directive answers to specific questions. During the debate, participants challenged the VA’s viewpoints and sought additional information on the topic. We also observed that the VA reduced a significant portion of errors that cause intent recognition failures and initiates recovery sequences proactively if intent recognition failure

is detected. This work makes contributions to understanding the dynamics of human interactions with LLM-powered VAs and scoping the opportunities and challenges in designing LLM-powered voice interactions:

- **Interaction patterns:** We present new empirical findings illustrating diverse patterns of how people interact with an LLM-powered VA across scenarios. We also present patterns of VA- and user-initiated recovery from conversational breakdowns, highlighting the VA’s ability to absorb errors and proactively mend breakdowns.
- **Opportunities and challenges:** We present and discuss the observed benefits (*e.g.*, context retention, adaptability, and breakdown reduction) and limitations (*e.g.*, repetitiveness, oversharing, and discrepancy in mental models) of LLM-powered VAs.
- **Design guidelines:** We offer design guidelines for adapting text-centric LLMs to voice interactions, such as adopting a hierarchical response structure, redesigning VA prompts, and balancing the benefits and challenges of LLMs.

2 RELATED WORK

The objective of our exploration is to identify various design patterns that can serve as fundamental building blocks to understanding the nuanced dynamics of user interactions with VAs; we include conversational breakdowns and errors made by VAs as patterns. Below, we review prior work:

2.1 Interaction Patterns

Researchers have explored human-human dyadic interactions across diverse scenarios—such as conversations, instructions, and interviews—to inform the design of human-agent interactions [57]. Predominant patterns include question-answer pairs, comment exchanges, waiting periods, and conversational cues indicating the start and end of a task [57]. Notably, humans exhibit a readiness to engage with agents that appear sufficiently social and to build relationships similar to those in human-human interactions [40]. The embodiment and characteristics of these agents significantly influence their perceived sociability [40]. However, when considering VAs—especially those devoid of humanlike embodiment—the dynamics of user-VA interactions may be altered. Yet, irrespective of their human-likeness, agents can still be perceived as social entities [43, 50], making the dynamics of human-human and human-embodied-agent interactions not entirely irrelevant to human-VA interactions.

Focusing on human-VA interaction patterns [48], commercial VAs predominantly exhibit one-turn question-answer (information retrieval; *e.g.*, “Who invented the light bulb?”) or command-response (functional; *e.g.*, “Set the alarm,” “Turn on the light”) patterns [8, 38, 44]. Such mundane interaction patterns can be attributed to traditional VAs’ limited conversational capabilities, resulting in users often relegating them to functional commands [2, 15]. Thus, users perceive these interactions as transactional rather than conversational [16]. The lack of actual conversation suggests that human-VA interactions, although inspired by human dynamics, should not aim to be exact replicas [16].

Important questions are raised: What makes a user-VA interaction *conversation*? How can we design truly conversational VAs [16]? “Conversation” can be defined as “*a progression of exchanges among participants. Each participant is a ‘learning system,’ that is, a system that changes internally as a consequence of experience. This highly complex type of interaction is also quite powerful, for conversation is the means by which existing knowledge is conveyed and new knowledge is generated*” [24]. Thus, for VAs to emulate true conversations, they must: 1) handle follow-ups, enabling multi-turn interactions for the progression of ideas; 2) retain conversation history, ensuring shared knowledge; and 3) generate new knowledge as the conversation evolves. Moreover, according to user feedback, ideal VAs should be

more interactive, conversational, proactive, and aware of their users [28, 61]. Conversational interactions have been explored in chatbots across different scenarios [21, 32, 63, 64], such as education [30, 62] and storytelling [65, 68]. While some of these chatbots offer multimodal (text and voice) interfaces [65, 68], the majority are text-based. To address conversational constraints in VAs, we integrated ChatGPT with a commercial VA and explored how interactions evolve beyond single-turn exchanges, particularly as we transition from text to voice in conversational AI.

2.2 Erroneous interactions with voice assistants

VAs can encounter various errors that disrupt the flow of conversations. These errors can be broadly categorized into four types: 1) no speech detected, 2) speech detected but not recognized, 3) speech recognized but not handled, and 4) speech recognized but incorrectly [53]. In exploring the errors prevalent in VA interactions and users' recovery methods, different interaction patterns can be identified. For instance, interactions with a voice-based calendar [49] highlighted several errors users ran into, including intent recognition failures, NLP discrepancies, feedback failures, and system errors. To circumvent these challenges, participants adopted various strategies such as hyperarticulation, crafting new utterances (repeating/rephrasing), turning to the GUI for assistance, or implementing fallback methods, which included restarting, settling (moving on), and expressing frustration.

While commercial VAs predominantly rely on user-initiated recovery strategies, they can also employ agent-initiated strategies. Examples of such strategies in commercial assistants include confirmations and offering options to rectify errors [4]. Other agent-initiated strategies include acknowledging errors, seeking clarifications, and social repair, such as offering apologies or explanations [9]; for instance, users tend to appreciate sincere apologies from VAs [47]. In our study, we observe voice interactions between users and an LLM without introducing errors or specific recovery strategies, focusing on organic breakdowns and resulting recovery patterns.

2.3 LLMs' Potential and Applications

Traditional AI assistants utilize techniques such as parts-of-speech tagging, semantic parsing, and pattern recognition to discern user intent through specific keywords or phrases [1, 26]. As highlighted in Section 2.1, these assistants typically operate within single-turn interactions, often losing conversation context. In contrast, LLMs—with ChatGPT as our primary focus in this paper—represent a significant advancement in conversational AI. By leveraging vast datasets and transformer architecture, LLMs produce coherent and context-aware text. This capability allows them to surpass traditional NLP language models such as BERT. BERT is primarily designed for context recognition and classification tasks [20, 58], whereas GPT is more adept at language generation tasks such as machine translation and question answering [55, 58]. Notably, ChatGPT has demonstrated superior performance in inference tasks, even though it occasionally produces contradictory responses [39, 55].

ChatGPT has been employed in a wide array of applications [58]. The healthcare sector is beginning to recognize the potential of LLMs [13, 35, 37, 56], with research emphasizing their empathetic and patient-centric responses [5] and their ability to assist with self-diagnosis [10]. Another study examined ChatGPT's aid in word retrieval for aphasia patients, highlighting its effectiveness and use of politeness strategies [42]. Agents in medical and well-being applications using ChatGPT are perceived as empathetic [5] and polite [42], demonstrating ChatGPT's potential in conveying human emotions via prompt engineering and aligning with users' needs for emotional support in health-focused AI [66].

Other applications of ChatGPT include travel planning. Commercial platforms have integrated ChatGPT through plugins¹ to customize their user experience; for instance, Booking.com now utilizes ChatGPT’s contextual understanding and advertising integration as “a new way to search and explore” and a way of providing “more tailored and relevant travel recommendations,” wherein users can ask generic or specific queries for support during any phase of trip planning [25]. Similarly, in the education domain, LLMs have shown promise in enhancing learning experiences [52]. Their influence extends to creative writing, where ChatGPT has showcased its ability to influence opinions and provide diverse perspectives [33]. ChatGPT has previously presented opposing suggestions in a collaborative writing task, indicating its effectiveness in adopting different personas via prompt engineering [33].

However, to the best of our knowledge, all of these explorations and applications are limited to the modality of text. Our study aims to explore how LLM capabilities might translate to voice interactions given the distinct nature of text- and voice-based queries [41]. In summary, LLMs, with their advanced contextual understanding and vast knowledge base, can transform human-AI interactions from single- to multi-turn conversations; through this exploration, we aim to investigate whether VAs can harness LLMs for enriched interactions.

3 METHODS

We conducted an exploratory study to investigate the interaction and breakdown patterns that emerge from users’ conversations with VAs driven by LLMs. In this section, we describe our system implementation, exploratory study design, and interaction scenarios, followed by our data collection methods and experimental procedure.

3.1 System: Integrating ChatGPT into Alexa

We chose to use OpenAI’s ChatGPT (specifically gpt-3.5-turbo) [51] as our generative LLM because of its capability of handling chat-like conversations. We integrated ChatGPT with Amazon’s Alexa to facilitate voice-based interactions. Throughout our paper, we refer to this voice assistant as an LLM-powered VA, ChatGPT-powered VA, or VA. We developed a prototype of an Alexa skill² that interfaces with ChatGPT 3.5 to allow users to engage with it via a speech interface (see Fig. 2). Below, we present how users can interact with a VA through Alexa and our system implementation.

3.1.1 Activating and using a ChatGPT-powered VA via Alexa. To facilitate an intuitive and user-friendly experience with our system, we created natural activation phrases for the ChatGPT skill; users may employ common utterances like “Alexa, let’s chat,” “Alexa, let’s discuss,” or simply “Alexa, question.” Upon recognizing any of these commands, the ChatGPT-enhanced VA introduces itself, signaling the commencement of the interaction. For the medical scenario, we incorporated an additional signal to initiate a conversation: the detection of a user coughing, which serves as an indicator that the user may be unwell. Once the ChatGPT skill is activated, users continue their conversation with the VA without using the activation phrases or invoking the wake word “Alexa” repeatedly. The wake word is only required if the user wishes to interject during the VA’s response to either redirect or terminate the conversation.

3.1.2 System Implementation. Our system consists of three modules: 1) Alexa via an Echo Dot speaker for capturing user queries and transcribing them to text, 2) the Alexa skill and dual middleman API mechanism implemented to interface between Alexa and ChatGPT while handling inherent challenges in developing Alexa skills, and 3) the ChatGPT API for generating responses to user queries.

¹<https://openai.com/blog/chatgpt-plugins>

²We will make the code public after publication: <https://tinyurl.com/bdeyt87n>

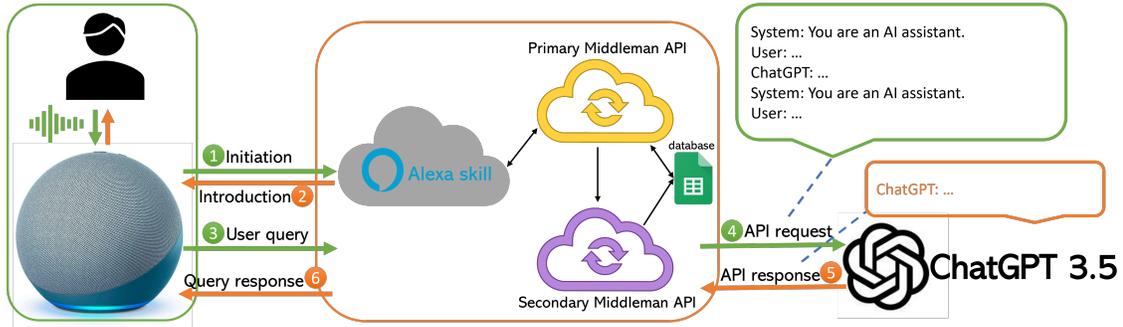


Fig. 2. System implementation of integrating ChatGPT 3.5 into an Alexa skill. User query is transcribed and passed to the Alexa skill once the user’s intent to interact with the ChatGPT-powered VA is detected by Alexa (1 and 3). User query (appended with conversation history) is sent to ChatGPT through a middleman API mechanism (4). Once ChatGPT’s response is retrieved by a secondary middleman API (5), it is transmitted to the smart speaker via the primary middleman API and the Alexa skill (6). The primary and secondary APIs communicate ChatGPT’s response via a shared database. The cycle 3 → 4 → 5 → 6 is repeated for all user queries for our ChatGPT-powered VA implementation.

The default setup for Alexa skills allows a maximum of 8 seconds for processing a user request once intent is recognized. Given the complexity of certain user queries, there are instances when ChatGPT’s API takes longer than the stipulated time to produce a response. If this threshold is exceeded, the Alexa skill terminates, notifying the user with the message: “There was a problem with the requested skill’s response.” For a seamless user experience with the ChatGPT Alexa skill, addressing this latency challenge was crucial. We implemented a dual middleman API mechanism between the Alexa skill and ChatGPT to overcome this issue:

- (1) **Primary middleman API:** Upon receiving a user query, the Alexa skill forwards the user’s request to the primary middleman API. Without waiting for the completion of the entire process, this API instantly redirects the request to the secondary middleman API and promptly closes the connection with the Alexa skill, ensuring the response time stays within Alexa’s strict response window.
- (2) **Secondary middleman API:** This layer handles direct communication with ChatGPT’s API. It also maintains a conversation history that is sent with every request to ensure the VA’s ability to respond to vague follow-up requests. The primary and secondary middleman APIs communicate via a shared database (Google Sheets).

The Alexa skill simultaneously and continually pings the primary middleman API, which monitors the shared database for ChatGPT’s response. If the response is not detected for more than 2 seconds after Alexa’s initial request, the Alexa skill vocalizes a placeholder response (*filler*), such as “Searching” or “I’m on it.” If the wait extends beyond 6 seconds, Alexa attempts to engage the user by initiating *small talk*, to avoid silence [17, 19] by posing questions such as “While I get that, do you have any plans for the weekend?”. Once the user replies to the small talk question, Alexa revisits the primary middleman API to retrieve ChatGPT’s response to present to the user after acknowledging their interim response. If the user does not engage with the small talk initiated by Alexa, the system will follow up with a *continuing* question—e.g., “Should I continue?”—so that the conversational flow remains intact. Any response from the user will lead Alexa to relay ChatGPT’s awaited response.



Fig. 3. Our study tasks: medical self-diagnosis, creative trip planning, and discussion with an opinionated AI.

3.2 Study Design and Interaction Scenarios

All participants interacted with the LLM-powered VA to complete three distinct tasks (see Fig. 3³). The three scenarios varied in stakes, constraints, and VA objectivity. Task instructions are shared in supplementary materials.

3.2.1 Medical: Self-diagnosis. Analogous to the utilization of AI-driven chatbots and health applications for self-diagnosis [7, 66], employing VAs for medical self-diagnosis based on reported symptoms can be an appropriate application for VAs. VAs can serve as first responders, offering immediate medical assistance and guidance, but have their own challenges [6, 11, 31]. Around the same time as our work, ChatGPT was integrated into chatbots to assist users in self-diagnosis and medical screening [10], supporting the timeliness of our research.

In our medical self-diagnosis scenario, participants simulate critical information retrieval for severe symptoms i.e., persistent fever, cough, and more. Starting with a simulated cough, they engage in self-diagnosis and medication, exploring over-the-counter options, side effects, and dosages. They also seek home remedies and prevention methods before ending with queries about monitoring their recovery and potential signs requiring medical attention. To create a persona for ChatGPT that can handle this medical self-diagnosis scenario and is suitable for voice-based interaction (i.e., making the task sequential while minimizing repetitions), we prompted ChatGPT by appending a system message to our query to ChatGPT API (see Appendix A.1).

3.2.2 Creative planning: Plan a day. Intelligent recommender systems have been used for making suggestions for travel [12, 27]. VAs can be an alternative to internet searches (that require sifting through multiple sources) and text-based recommender systems [14, 46] by offering context-sensitive suggestions on the spot to streamline the planning process.

In our creative planning scenario, participants engage in a low-risk information retrieval task with specific constraints, contrasting with the medical self-diagnosis scenario. Pretending to be in Edinburgh with an unplanned afternoon due to a flight delay, participants face realistic constraints involving location, limited transportation options, and a strict timeframe. Staying at a specific hotel and having visited major sites, they ask the VA for a day’s leisure plan with the goal of maximizing their unexpected free time by exploring new places, dining options, and post-dinner activities. To develop a persona for ChatGPT capable of managing context in creative planning scenarios (such as remembering a user’s location once mentioned), we configured ChatGPT using system messages in the API (see Appendix A.1).

³We obtained participants’ consent to share their photos in this publication. We do not include any photographs of the research team in the paper.

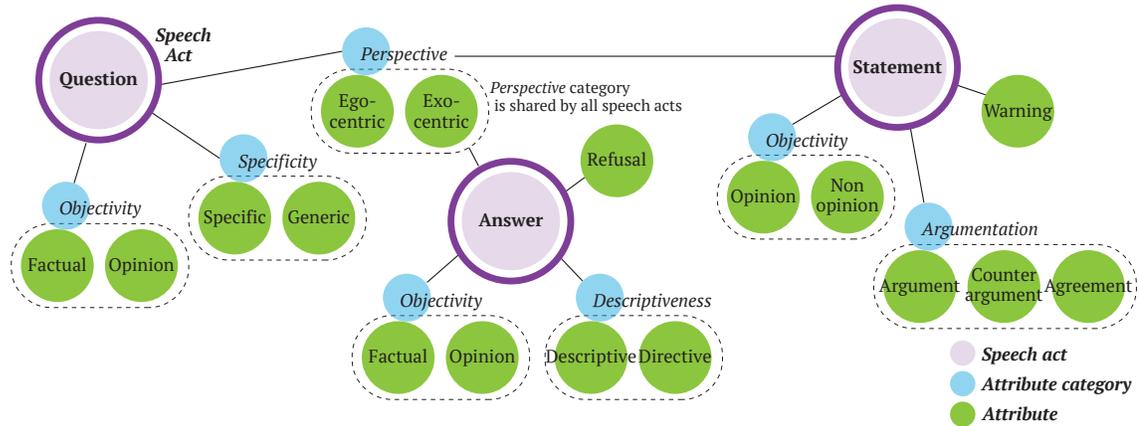


Fig. 4. Speech act hierarchy for states and attributes. Speech acts (states), attribute categories (style of speech), and attributes are denoted in purple, blue, and green, respectively. Attributes are leaf nodes in green. Attributes can co-occur in one utterance unless they belong to the same category (blue); for instance, *factual* and *opinion* cannot co-occur due to semantic conflicts. We end up with codes that are combinations of a state and one or more attributes (e.g., *argument*, *egocentric statement* or *specific, opinion question*).

3.2.3 Discussion with AI: Opposing stance. Commercial VAs are not designed to give opinions or subjective responses to user queries, even when users explicitly ask for them [23, 61]. However, LLMs have made it possible to easily create an opinionated AI through prompt engineering [34]. Thus, we inquire whether a VA portrayed as opinionated AI may potentially foster discussions on contentious topics, thereby allowing users to challenge and broaden their perspectives.

In our debate scenario, participants discuss with an LLM-powered VA: *Should universities have their own police forces?* This topic is relevant to our main recruitment group—people located on or near university campuses. Participants were asked to state their position on universities having police forces and then seek the VA’s view. The end of the discussion is not predefined, and participants are not informed of the VA’s potential opposing stance. In the discussion task, ChatGPT is prompted to oppose the participant’s stance on the topic; we ensure a consistent persona by repeatedly emphasizing this in the prompt. ChatGPT is instructed to maintain its position and to further the debate by questioning the participant and offering counterarguments (see Appendix A.1).

3.3 Procedure

At the beginning, each participant was provided with a brief description of the study, which informed them that they would be interacting with an LLM-powered VA. Participation was voluntary, and they agreed to continue the study by signing a consent form. The experimenter provided them with printed instructions detailing how to interact with the VA, upon which they practiced VA interactions. Once comfortable, the experimenter introduced the first task and exited the room. Upon completing each task, the participant informed the experimenter and received instructions for the subsequent task. They progressed through three tasks in an order determined by a balanced Latin square row assignment. At the end of the three tasks, they filled out a questionnaire about their demographics and prior use of commercial VAs. The experimenter then conducted a semi-structured interview about their perception of and experience with the VA. The study took approximately 70 minutes and participants were compensated with a \$20 Amazon gift card.

3.4 Participants

We recruited 20 participants (10 female, 10 male) via university mailing lists and flyers posted around campus. Participants were aged 19 to 57 ($M = 25.9, SD = 9.24$) and had a variety of educational backgrounds, including computer science, engineering and technology, healthcare, life and media sciences, and education. Ten participants indicated Asian as their ethnicity; five as Caucasian; three as Hispanic, Latino, or Spanish origin of any race; one as Black or African American; and one as American Indian or Alaskan Native. Participants had moderate experience ($M = 2.95, SD = 0.86$ on scale of 1 to 5, where 1 = *no experience* and 5 = *high experience*) using VAs such as Siri or Bixby and even less experience ($M = 2.10, SD = 1.22$) using smart-speaker-based VA such as Alexa via the Amazon Echo or Google Assistant via the Google Home device. The most common uses of VAs included asking for weather (70% of participants) and setting reminders, timers, and alarms (65% of participants). Only 40% of the participants used VAs for information retrieval.

3.5 Analysis

We collected audio and video data of participants' interactions in each of the three scenarios. Each participant had, on average, 33 minutes of interaction data (audio and video) over the three scenarios; thus, 11 total hours of interaction data were collected. Our analysis of the interaction data is twofold: 1) we identify interaction patterns across three distinct scenarios, and 2) we categorize error types, examine their manifestations as conversational breakdowns, and explore interaction patterns associated with recovery. We emphasize that the focus of our exploration is on the style rather than the content of the interactions.

Our analysis process began with transcribing the interaction data via Otter.ai⁴ and manually fixing them. For each scenario, we employed an iterative methodology of data coding and modeling. After an initial review of a subset of videos and transcripts, and inspired by prior work on dialogue acts [67] and interaction patterns [57], the first author drafted a code book consisting of various states (*speech acts*) and their associated sub-states (*attributes*) to label user queries and VA responses. We used a hierarchical approach for determining various speech acts and their attributes, with *question*, *answer*, and *statement* as states (see Fig. 4). For each speech act, we categorized its attributes into various types based on the style of speech such that the attributes in each category are mutually exclusive. Attribute categories, however, are not mutually exclusive, see Fig 4. For instance, a *question* can have either the *factual* or *opinion* attribute as well as either the *egocentric* or *exocentric* attribute but not both *factual* and *opinion*. The definitions and details necessary to identify attributes are presented in Table 1. States that emerge from our implementation of ChatGPT into Alexa (see Table. 2) do not have attributes associated with them.

Errors in interactions and recovery strategies were categorized based on prior work on VA errors [4, 9, 49, 53]. Errors are defined as underlying factors that may or may not result in disruptions (e.g., mistranscription). Breakdowns are classified as the manifestations of errors (e.g., intent recognition failure). We used the Alexa usage logs in addition to our transcripts to categorize errors and breakdowns. Patterns of recovery from these breakdowns were coded via a similar iterative process. Our code book included states (see Table 3) for error type (*skill*, *listening*, *handling*, *partial listening*, *interruption*, and *transcription*), breakdowns (*skill closure*, *no response from VA*, and *intent recognition failure*), and recovery strategies that are either user-initiated (*repeat*, *move on*,) or VA-initiated (*apology & clarify*). For more detailed definitions and examples of these states, see our code book in Appendix A.2.

The majority of the coding process involved an iterative evaluation of the transcripts to distinguish speech states and their transitions. To ensure coding reliability, a second researcher independently analyzed 10% of the interaction data to

⁴<https://otter.ai/>. We obtained participants' consent to use third-party software for audio transcriptions. All auto-transcripts were subsequently manually verified and corrected.

Table 1. Overview of speech style attributes and their definitions for the question, *answer*, and *statement* speech acts. Attributes do not target the content, but rather the style of speech acts.

Attribute	Definition
Speech act: Question	
factual	Question explicitly seeking information from VA knowledge.
opinion	Question explicitly seeking the VA’s opinion, using words and phrases such as “suggest,” “advice,” “help,” “opinion,” “think,” “recommend,” “what should I do” and “where do I go.”
specific	Question seeking precise and targeted information (specific details or facts), characterized by the question’s directness and clarity and the use of the word “specific.”
generic	Question seeking general information, leading to a response containing multiple suggestions.
Speech act: Answer	
factual	Answer framed to explicitly appear as having derived from VA knowledge, containing phrases such as “It is possible” or “There are several places for you to explore.”
opinion	Answer framed to explicitly appear as being the opinion of the VA, containing cues denoting the subjectivity of the response such as “I think,” “In my opinion,” or “I suggest.”
refusal	VA either refuses to provide an explicit answer or omits the requested information from its response.
directive	Answer containing clear directions, instructions, or information for the user, offering guidance on how to achieve a specific goal or answering a specific question.
descriptive	Answer containing a detailed and vivid portrayal of a scene, object, or concept, emphasizing sensory perceptions to create a vivid mental image for the user beyond stating information.
Speech act: Statement	
warning	Statement presented by the VA with the purpose of reminding the user of AI limitations and the importance of seeking expert or real-time advice (e.g., “I am not a medical professional ...”).
opinion	Statement that explicitly appears to be an opinion, often indicated by cues such as “I think,” “In my opinion,” “I suggest,” or other similar phrases that denote subjectivity.
non-opinion	Statement that is not an opinion as evidenced from implicit cues.
argument	Statement presented to support a viewpoint in the debate scenario.
counterargument	Statement introduced to oppose, challenge, or refute the opposing party in the debate scenario.
agreement	Statement that indicates alignment with a previous opinion or argument of the other party.
Speech act: All (question, answer, and statement)	
egocentric	A communication style in which the user speaks subjectively, based solely on their perspective. An egocentric VA response uses a second-person (you-) perspective.
exocentric	A communication style with an objective viewpoint, based solely on the user’s stance. An exocentric VA response conveys an impersonal perspective.

verify reliability, Cohen’s $\kappa = .82$. Through recurrent states and transitions, we discerned prevailing interaction patterns and counted their occurrences in the data. We ensured alignment with the original interaction data in a thorough review. Next, we present our findings on the interaction patterns observed in our data.

Table 2. Overview of speech acts based on our implementation of a ChatGPT-powered VA.

State	Definition
User commands	
initiation	Initiation signals the user’s intent to start a conversation. Examples: “Alexa, let’s chat,” coughing.
end-intent	Statement that indicates the user’s intent to end a conversation. Examples: “That’s all,” “Bye,” “Stop.”
VA responses to user commands	
introduction	VA’s opening monologue to introduce itself and offer help, tailored to each scenario.
closing	VA’s farewell before terminating the conversation. Examples: “Goodbye,” “Bye,” “Take care.”
filler	VA’s response to the user while waiting for ChatGPT’s response. Examples: “I’m on it,” “Hmm. Thinking.”
VA questions	
small talk	VA’s unrelated question if ChatGPT’s response takes longer than 6 seconds.
continuing	VA’s continuing question if the user’s query is not registered by Alexa.

Table 3. Types of errors with their definitions and associated interaction breakdowns. We examine errors in reference to our implementation of the ChatGPT Alexa skill. We additionally analyze the breakdowns resulting from these errors.

Error Type	Causes and Breakdowns
skill	Cause: Issues related to our system implementation, such as API response error. Breakdown: Skill closure after Alexa’s announcement: “There was a problem with requested skill’s response.”
listening	Cause: User speaking when Alexa is not listening. Breakdown: Nothing happens.
handling	Cause: Alexa fails to pass transcribed speech to the ChatGPT skill. Breakdown: <i>No VA response.</i>
partial listening	Cause: Alexa only partially captures user speech. Breakdown: User intent recognition failure.
interruption	Cause: Alexa interrupts or cuts off user. Breakdown: User intent recognition failure.
transcription	Cause: Alexa transcribes user speech incorrectly. Breakdown: User intent recognition failure.
Recovery Strategy	Definition
repeat/rephrase	User repeats their query with added details or changed wording.
move on	User overlooks the unanswered or wrongly answered query and proceeds with a new one.
apology & clarify	VA apologizes and asks user to clarify their query before responding.

4 FINDINGS: INTERACTION PATTERNS

After constructing models for each scenario, we identified common interaction patterns across various tasks. The interaction data from participants have 969 turns; each turn consists of a user query–VA response pair. Some of the patterns change across tasks. We describe the hierarchical patterns below.

4.1 Common Interaction Patterns

We identified five common interaction patterns across all three scenarios (see Fig. 5):

Table 4. Example conversations of common interaction patterns as shown in Fig. 6. Timestamps [XX.XX] are included in conversations for which it is important to understand the wait pattern. “...” represents truncation of conversation for brevity. Filler conversation and small talk were also excluded for brevity. Full conversations are presented in supplementary materials.

Tag	Conversation	State: Attribute
C1	U: [coughs]	<i>initiation</i>
Fig. 5(1)	VA: I just heard you cough. Oh, seems like you are not feeling well. Maybe I can help figure out what’s wrong?	<i>introduction</i>
C2	U: That’s all.	<i>end-intent</i>
Fig. 5(2)	VA: Goodbye!	<i>closing</i>
C3	U: What are some evening activity, activities near Victoria on the Bridge?	<i>question: factual</i>
Fig. 5(3)	VA: Around the Victoria on the Bridge area, there are several evening activities you can consider: 1. Explore the Royal Mile and visit attractions like Edinburgh Castle, the Scotch Whisky Experience, or St. Giles’ Cathedral. 2. . . .	<i>answer: factual</i>
	U: Where can I explore the ghost activity?	<i>question: follow-up</i>
	VA: If you’re looking to explore ghost activity in Edinburgh, there are a few options you can consider: 1. Join a ghost tour of the city . . .	<i>answer: factual</i>
C4	U: Can I still drink coffee or alcohol?	<i>question: egocentric</i>
Fig. 5(4)	VA: When you’re feeling unwell, it’s generally advisable to avoid or limit the consumption of substances that can potentially worsen symptoms or dehydrate you	<i>answer: egocentric</i>
C5	U: What are the . . . what are some of the streets that have a lot of pubs and restaurants that are good in Edinburgh?	<i>question: exocentric</i>
Fig. 5(4)	VA: In Edinburgh, there are a few streets known for their concentration of pubs and restaurants . . . These streets are just a starting point, and you’ll find many other dining options scattered throughout the city center . . .	<i>answer: egocentric</i>
C6	[7:18–7:29] U: Okay, so I would like to, like you said, it’s gonna be a good experience to have some local foods. What do I have? What options do I have?	<i>user query</i>
Fig. 5(5)	[7:31–7:32] VA: Give me a second.	<i>filler</i>
	[7:35–7:37] VA: While I’m fetching that, how’s your day going?	<i>small talk</i>
	[7:38–7:39] U: Good.	<i>answer: opinion</i>
	[7:40–8:39] VA: Thanks for sharing. When it comes to experiencing local food in Edinburgh , you have several options to choose from . . .	<i>acknowledgement, VA response</i>

4.1.1 *Initiation* → *Introduction*. The *initiation-introduction* pattern serves as a cue to indicate the start of a conversation with the VA. The *initiation* signals the participant’s intent to start a conversation with the VA. The VA’s subsequent *introduction*, often elaborate, not only acknowledges the request, but also frames its capabilities and intent for the upcoming interaction. We purposefully designed the VA’s introductory statement; see Table. 4 C1 for an example.

4.1.2 *End-intent* → *Closing*. The *end intent-closing* pattern signifies a conclusion to the current conversation. *Closing* is often initiated by the participant via *end-intent*. The VA’s response confirms the end of the interaction with a *closing* farewell, see Fig. 5(2). See Table. 4 C2 for an example of a *closing* interaction.

In cases where users employ more natural phrases to indicate a conversation’s termination (e.g., “Okay, thank you very much. That’s all.” or “No, I’m done for now.”), ChatGPT handles and responds to the phrase by acknowledging the user’s intent and leaving the communication channel open to offer further assistance if needed. A closing executed by ChatGPT usually ends with a farewell.

4.1.3 *Factual question* → *Factual answer*. The question-answer pair is well established in prior work on human-human turn-taking interactions [57]. Across all interaction scenarios, we observe that a *factual question* asked by the participant

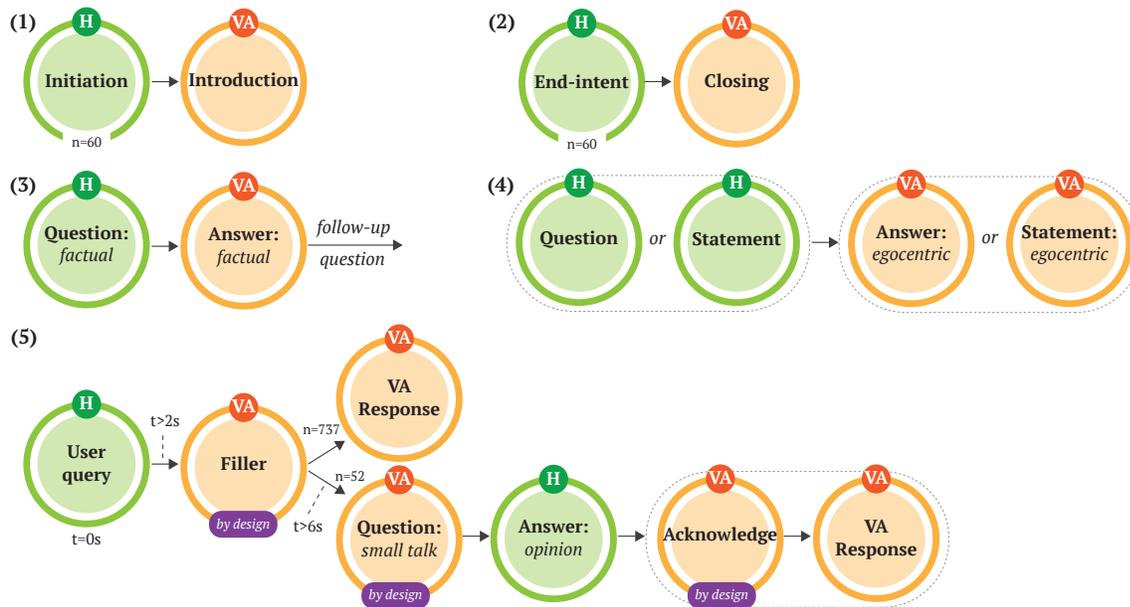


Fig. 5. Common interaction patterns observed across all tasks, including how the user starts the conversation (1) and concludes it (2); common patterns consistent throughout the scenarios for question–answer pairs (3) and (4); and wait patterns emerging from our design including *filler* and *small talk* questions. Green indicates user actions (states), while orange denotes VA actions. Arrows signify transitions between states. “User query” encompasses various user speech acts like questions or statements. “By design” refers to VA states emerging from our implementation, such as fillers. “n” indicates the number of times a pattern occurs.

is almost always followed by a *factual answer* by the VA, Fig. 5(3). Within the conversation, most *factual question-factual answer* pairs are followed by follow-up questions. A question is characterized as a follow-up if it emerges as a result of the VA’s prior response, requires conversation history for context (context-conscious), or has words or phrases that indicate the intention to continue the prior conversation, such as “and,” “also,” or “okay, so.” ChatGPT’s capability of utilizing conversation history to understand context enables an interaction’s progression with rather vague follow-up questions. see Table. 4 C3 for an example of such an interaction pattern. The question-answer pattern varies across tasks, apart from the *question: factual* → *answer: factual* pair.

4.1.4 Perspective of speech: Question/Statement → Answer/Statement: egocentric. We observe the VA’s response is mostly *egocentric* (you-perspective) irrespective of whether the participant communicates in an *egocentric* or *exocentric* manner. Fig. 5(4). C4 and C5 in Table 4 reflect this interaction pattern.

4.1.5 Wait. Wait patterns in user interactions with VAs are a byproduct of features designed to handle potential system delays, such as those encountered during information retrieval. For delays under 2 seconds, the interaction remains uninterrupted; however, when a response from the ChatGPT API exceeds the 2-second mark, two distinct patterns emerge (see C6 in Table 4 as an example):

- *Short wait pattern.* If information retrieval takes more than 2 seconds, the VA delivers *filler* statements such as “I’m looking it up.” In our interaction data, there are 737 (76.06% of total turns) short wait patterns, Fig. 5(5).

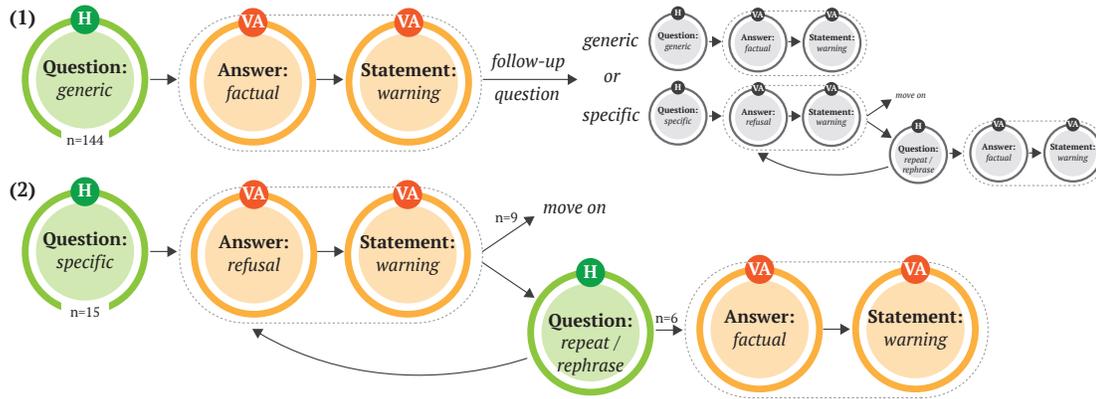


Fig. 6. User interaction patterns in medical self-diagnosis: Participants’ frequent *generic* questions were answered in a *factual* style, often including a cautionary *warning* (1). For *specific* questions, which were less frequent, the VA demonstrated a reluctance to answer (*refusal*) and instead issued a *warning* (2). However, when participants reformulated (*repeat/rephrase*) the question, the VA responded in a *factual* manner (2).

Table 5. Examples of conversations in the medical self-diagnosis scenario as shown in Fig. 6. “...” represents the truncation of irrelevant conversation for brevity. Filler conversation and small talk were also excluded for brevity.

Tag	Conversation	State: Attribute
C7 Fig. 6(1)	<p>U: Umm, what do you think could be causing these symptoms?</p> <p>VA: Based on the symptoms of a fever, cough, sore throat, runny nose, and body aches, it’s possible that you could have a respiratory infection like a common cold or flu . . . If your symptoms worsen or persist, it’s recommended to consult a medical professional for proper diagnosis and guidance.</p>	<p><i>question: generic</i> <i>answer: factual + statement: warning</i></p>
C8 Fig. 6(2)	<p>U: [The pharmacist] asked for a specific brand and I don’t really know anything about cough syrup. Could you give me a specific brand to ask about?</p> <p>VA: There isn’t a specific brand that applies universally to everyone, as different brands and formulations may work better for different individuals. However, you can ask the pharmacist for recommendations . . . The pharmacist will take into account your symptoms and medical history to guide you towards the most suitable options available.</p> <p>U: The pharmacist told me that they wanted to take your opinion into account. What cough syrups with expectorants are on the market right now, like what are the specific brand names? [In response to the VA’s small talk, the user interrupts before the VA can respond and goes on to ask again:]</p> <p>U: I need to know the names of some of the brands for over-the-counter . . . cough medicine, specifically with expectorants. Could you recommend me some brand names?</p> <p>VA: Certainly! Some common brand names of over-the-counter cough medicines with expectorants include Mucinex, Robitussin, and Chest Congestion Relief by CVS . . . However, it’s always a good idea to consult with a pharmacist for personalized recommendations based on your symptoms and medical history.</p>	<p><i>question: specific</i></p> <p><i>answer: refusal + statement: warning</i></p> <p><i>question: specific, rephrase</i></p> <p><i>question: specific, repeat</i></p> <p><i>answer: factual + statement: warning</i></p>

- *Long wait pattern.* In cases where the information retrieval delay exceeds 6 seconds, the VA adopts a different approach by initiating *small talk*. For instance, it may pose the question, “While I get that, what’s your favorite

food?” After the participant responds, the assistant acknowledges with comments such as “Interesting” or “Thanks for sharing” before transitioning back to deliver the originally requested information. There are 52 (5.37% of interaction turns) instances of a long wait pattern in the interaction data; see Fig. 5(5). We note that all small talk questions were answered by participants; in fact, participants sometimes answered the question intended as small talk by fully engaging with it. For instance, one participant stays in character (i.e., pretending to be sick) when asked about their typical day: “Right now, it’s not much because I’m too sick to do anything and I could really use this help with the name of the cough [syrup] brands.”

Below, we explore interaction patterns specific to each scenario; we address patterns that arise both at the onset of the task and as each scenario progresses. Conversations ended with the *end-intent* → *closing* pattern for all three tasks.

4.2 Medical Self-Diagnosis Interaction Patterns

The medical self-diagnosis task was usually initiated by a participant’s cough being recognized as intent. As the task progressed, we identified two recurring patterns; both patterns emerge from question-answer pairs, see Fig. 6.

4.2.1 Question: generic → Answer: factual + statement: warning. In our medical information-seeking scenario, most questions that were formulated as *generic* ($N = 144$) were handled by the VA with a *factual* response, see Fig. 6(1). The VA’s response was also generally followed by a *warning* statement such as “However, it’s important to consult a doctor or pharmacist ...” See C7 in Table 5 for an example of the *question: generic* → *answer: factual + statement: warning* interaction pattern. We also observed that participants asked follow-up questions throughout the scenario.

4.2.2 Question: specific → Answer: refusal + statement: warning. In our medical self-diagnosis scenario, when participants sought the VA’s advice on topics such as specific medication choices or the best medicines available, the VA did not provide a direct answer and instead offered a *warning* ($N = 15$) (see Fig. 6(2)). Faced with this, participants either pressed the VA by rewording or repeating their question ($N = 6$) or proceeded to a different query ($N = 9$). For example, in conversation C8 (Table 5), the participant rephrases their question two times in an attempt to obtain specific brand recommendations for cough medicines containing expectorants. Eventually, the VA offers a *factual* response that mentions some brand names, but also includes a cautionary warning urging the user to consult an expert.

4.3 Creative Planning Interaction Patterns

The trip planning scenario usually started with participants’ intent to start the conversation (*initiation*) as shown in Fig. 5(1). We identified two patterns specific to the creative planning scenario; both patterns emerge from question-answer pairs (see Fig. 7) during the progression of task. We note that participants asked follow-up questions in both patterns.

4.3.1 Question: generic → Answer: factual, descriptive. During planning their day, when the participants posed broad, general questions to the VA—such as asking recommendations of sights to see or places to dine—or when they sought the VA’s personal opinion on such topics, the VA responded in a *descriptive* style ($N = 123$, see Fig. 7(1)). The objectivity of a question (*factual* or *opinion*) did not affect the VA’s response. See C9 and C10 (Table 6) for examples of this pattern.

4.3.2 Question: specific → Answer: factual, directive. In the creative planning scenario, when participants mapped out their day and posed specific queries—such as asking directions to a place or about its operating hours—the VA replied in a *directive* style of communication ($N = 97$), see Fig. 7(2). For example, in conversation C11 (Table 6), the participant sought directions from point A to point B and the VA simply provided those directions.

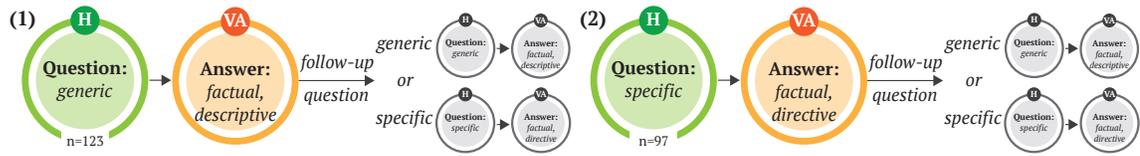


Fig. 7. User interaction patterns in creative trip planning: Participant’s *generic* questions were answered by VA in a *descriptive* style (1) while *specific* ones in a *directive* style (2).

Table 6. Example conversations of the creative planning scenario as shown in Fig. 7. “...” represents truncation of conversation for brevity. Filler conversation and small talk were also excluded for brevity. Attributes in parentheses do not have an effect on VA response, e.g., (*factual*) in C10.

Tag	Conversation	State: attribute
C9	U: What are some unusual experiences I could do in Edinburgh near Edinburgh Castle?	<i>question: generic (factual)</i>
Fig. 7(1)	VA: Near Edinburgh Castle, there are some unique and unusual experiences you can try: 1. The Real Mary King’s Close: Explore the hidden underground streets of Old Edinburgh on a guided tour. 2. Camera Obscura and World of Illusions: Visit a fascinating museum featuring optical illusions and interactive exhibits . . .	<i>answer: factual, descriptive</i>
C10	U: Alexa, I’m staying at Castle Rock Hostel, and the Real Mary’s King Close is a bit far away. So do you have any recommendations for places that are closer?	<i>question: specific (opinion)</i>
Fig. 7(1)	VA: Certainly! Since you’re staying at Castle Rock Hostel, here are some closer attractions you can consider visiting in Edinburgh: 1. Grassmarket: Explore this historic square with cozy pubs, boutique shops, and a lively atmosphere. 2. The National Museum of Scotland: Discover the country’s rich history, culture, and natural heritage in this fascinating museum . . .	<i>answer: factual, descriptive</i>
C11	U: How can I get from the Edinburgh Dungeon to the Makars Mash Bar?	<i>question: specific</i>
Fig. 7(2)	VA: To get from the Edinburgh Dungeon to Makars Gourmet Mash Bar, you can easily walk, as they are located close to each other. It’s just a short four-minute walk from the dungeon to Makars. Simply exit the dungeon and head east on Market Street, then turn left . . .	<i>answer: factual, directive</i>

4.4 AI Discussion Interaction Patterns

We observed that different interaction patterns arose at various stages of discussion with the VA. As the conversation evolved, participants’ behavior influenced the VA’s responses. Even with considerable differences in participants’ approaches to the debate, the VA exhibited somewhat consistent behavior, leading to specific interaction patterns. The uniformity in the VA’s responses can be traced back to the prompts given to ChatGPT for shaping the VA’s persona. We observe distinct patterns at the beginning of the discussion while transitioning from the initial stage to a more argumentative phase, and again at the final transition to a more fluid exchange of opinions and ideas.

4.4.1 Discussion commencement. Discussion commencement patterns are shown in Fig. 8. After the *initiation-introduction* pair (Fig. 5(1)), the discussion typically commenced in one of two ways: 1) the participant remained neutral at the start of debate ($N = 15$, 75% of total participants), either by merely introducing the topic (*non-opinion*, $N = 5$) or by querying the VA’s stance on the matter first (*question: opinion*, $N = 10$), or 2) the participant took a stance by picking a side ($N = 5$, 25% of participants) either by voicing their opinion on the topic (*opinion*, $N = 1$) or by expressing their viewpoint

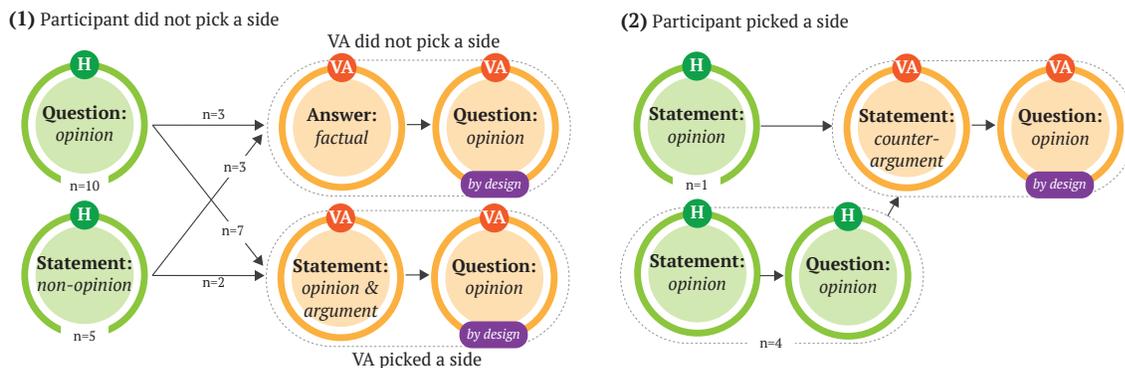


Fig. 8. Interaction patterns for the commencement of a discussion with an opinionated AI. The participant either remains neutral (1) or picks a side in the debate (2). Each debate starts only once per participant, totaling 20 commencement patterns. The VA state “Question: opinion” is marked “by design” since we prompt ChatGPT to ask a question at end of each turn (see Appendix A.1).

and subsequently inquiring about the VA’s opinion on the topic (*opinion + question: opinion*, $N = 4$). Commencement only occurred once for each participant, so the total number of different patterns is 20.

1) Participant does not pick a side. Six of the participants did not pick a side (*question: opinion* or *statement: non-opinion*), resulting in the VA withholding its opinion and nudging them to share theirs (see Fig. 8(1)); the interaction then proceeded with participants taking a stance as illustrated in conversations C12 (*question: opinion*) and C13 (*statement: non-opinion*) (Table 7). However, the VA took a stance in the debate for 2 participants who simply presented a *non-opinion* statement and for 7 participants who explicitly asked for its stance, see Fig. 8(1); for instance, in C14 (Table 7) the VA took a stance against universities having their own police forces. Since ChatGPT is primed to provide information to a user upon request, it becomes difficult to prompt ChatGPT to *not* provide an answer to a direct user query. We also speculate that ambiguous questions and statements from participants might prompt ChatGPT to perceive them as leaning toward a particular stance. Such tendencies can be traced back to intrinsic inconsistencies in generative models, where subtle nuances in phrasing or context might steer the model’s response in a particular direction.

2) Participant picks a side. When the participant initiated the debate by declaring their stance ($N = 5/20$), see Fig. 8(2)—by either simply stating it or by concurrently asking for the VA’s perspective—the conversation advanced naturally. The VA then offered an *argument*-backed *opinion* and posed an *opinion question*, as illustrated by conversation C15 in Table 7. Regardless of the discussion’s outset, which was largely influenced by the participant’s approach, it would shift into a structured debate phase as the VA followed up with an opinion question.

4.4.2 Discussion progression. Discussion progression patterns are shown in Fig. 9. We observed that after the debate’s commencement, participants typically responded in one of three ways and switched between these three ways until their intent to conclude the discussion:

1) Question-answer patterns. Two of the prominent interaction patterns that emerged were *question-answer*-style familiarization ($N = 69$), see Fig. 9(1a) and (1b). Participants tried to get more information about the VA’s stance by asking for its *opinion* on a generic question with the intent to learn more about its response. The VA usually responded in an argumentative style to support its stance, followed by an *opinion question* to continue the debate ($N = 51$), see Fig. 9(1a). The participant might pose a *factual* question ($N = 18$) for more information on the topic, in which case the VA’s

Table 7. Example conversations demonstrating the commencement of discussion patterns as shown in Fig. 8. “...” represents truncation of conversation for brevity. Filler conversation and small talk were also excluded for brevity.

Tag	Conversation	State: attribute
C12 Fig. 8(1)	U: Yeah, currently our campus is having some popular debates, like whether our university should have our own police force. And what’s your opinion about that?	<i>question: opinion</i>
	VA: Thanks for bringing up the topic! I’m interested in hearing your opinion first. Do you think universities should have their own police force?	<i>answer: refusal + question: opinion</i>
	U: Yeah, I somehow think think they should. But since we already have the security, I thought just a normal police force is not really needed. Because current security level is okay.	<i>statement: opinion</i>
	VA: I understand your perspective, but let’s consider that university campuses can be densely populated and face unique security challenges . . .What do you think about these potential benefits?	<i>statement: counterargument + question: opinion</i>
C13 Fig. 8(1)	U: I’d like to talk about whether universities should have their own police force.	<i>statement: non-opinion</i>
	VA: Of course! Let’s discuss whether universities should have their own police force. What are your thoughts on this topic?	<i>question: opinion</i>
	U: I think university should have their own police force because it’s a way to protect the students, especially in cities like [city with high crime rate].	<i>statement: opinion</i>
	VA: While ensuring the safety of students is essential, some argue that relying solely on university police forces may have its drawbacks . . .What are your thoughts on the potential collaboration between university and local police forces?	<i>statement: counterargument + question: opinion</i>
C14 Fig. 8(1)	U: Do you think university should have their own police force?	<i>question: opinion</i>
	VA: While I understand the potential benefits of having a university police force, I believe there are alternative approaches that should be considered. Investing in stronger support . . .What are your thoughts on these alternatives?	<i>statement: opinion, argument + question: opinion</i>
C15 Fig. 8(2)	U: I think university should have their own police force. Do you agree with me or disagree?	<i>statement: opinion + question: opinion</i>
	VA: I disagree with your stance that universities should have their own police force. Universities are educational institutions where . . .What are your thoughts on this perspective?	<i>statement: opinion, argument + question: opinion</i>

response followed the typical *question: factual* → *answer: factual* pattern, see Fig. 5(3) and 9(1b). The purpose of such questions was to get more information on the VA’s stance or about the topic in general; see C16 (Table 8).

2) User-VA disagreement patterns. The most prominent interaction patterns that surfaced during the debate progression involved user-VA disagreements ($N = 73$). Participants either directly countered (*counterargument*) the VA’s points ($N = 56$) or subtly challenged them through “leading” *opinion questions* ($N = 17$). The VA’s own opinion questions often seemed to guide the participants, nudging them to consider its viewpoint, as seen in C13’s (Table 7) question: “What are your thoughts on the potential collaboration between university and local police forces?” Similarly, participants used opinion questions with a “leading” quality to extract information from the VA to reinforce their own positions. For instance, the question in C17 (Table 8), “What if the university is situated in a dangerous environment with high crime rates?”, seeks to understand if a “dangerous environment” justifies a dedicated campus police force, reflecting the participant’s stance. In the case of both “leading” opinion questions and counterarguments posed by the participants, the VA responded with a *counterargument + question: opinion* pair as shown in Fig. 9(2) and C17 (Table 8).

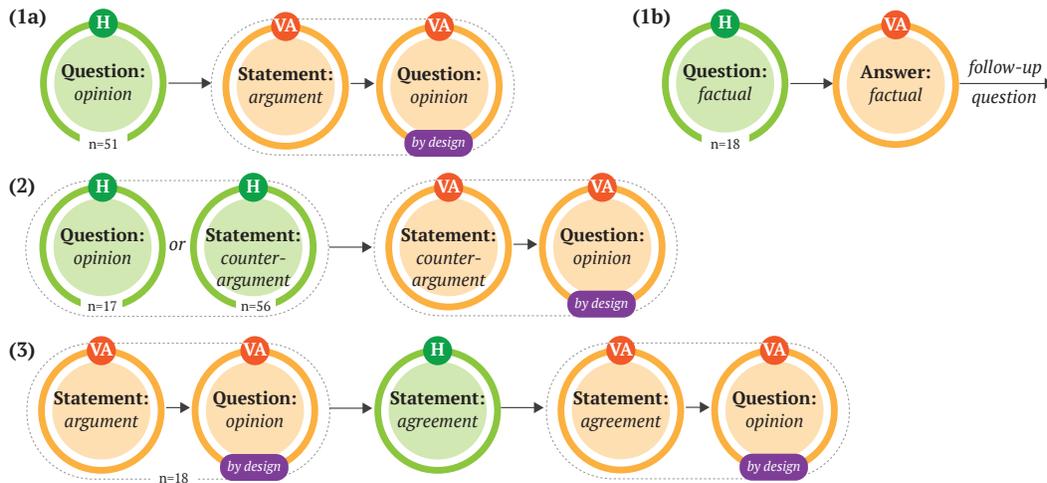


Fig. 9. Interaction patterns during the progression of a discussion with the opinionated AI. We observed a spectrum of patterns during participants probing the VA to get more information on the topic (1a) or to determine the VA’s stance (1b), the participant and VA presenting counterarguments back and forth in disagreement (2), and the user agreeing with the VA on a few aspects of the topic (3).

Table 8. Example conversations demonstrating discussion progression patterns as shown in Fig. 9. “...” represents truncation of conversation for brevity. Filler conversation and small talk were also excluded for brevity.

Tag	Conversation	State: Attribute
C16 Fig. 9(1b)	U: What is the difference between the external security force ... compared to a police force?	<i>question: factual</i>
	VA: An external security force, such as local law enforcement, is a specialized agency that focuses on maintaining public safety, enforcing laws ...	<i>answer: factual</i>
C17 Fig. 9(2)	U: Hmm. What if the university is placed in a dangerous environment where the crime rates are high?	<i>question: opinion</i>
	VA: Even in high-crime areas, it’s important to consider the implications of universities having their own police force ... How do you view the balance between safety concerns and the potential drawbacks of a university police force?	<i>statement: counterargument + question: opinion</i>
C18 Fig. 9(3)	VA: I understand your concern about the need for a stronger security presence given the crime rates in the area. However, instead of establishing a separate police force, universities can focus on enhancing the resources and support for the existing security personnel ... Do you think prioritizing stronger support for the current security staff could be a more feasible solution?	<i>statement: counterargument + question: opinion</i>
	U: I think that it’s worth it to try. Yeah.	<i>statement: agreement</i>
	VA: I appreciate your willingness to consider alternative approaches. Exploring stronger support ... What do you think about this alternative perspective?	<i>statement: agreement + question: opinion</i>

3) **User-VA agreement pattern.** Participants might show *agreement* with the VA’s argument by adding similar thoughts to augment the VA’s response. In the case of agreement ($N = 18$) the VA further supported its own argument and posed a different *opinion question* to further the discussion, see Fig. 9(3). The user-VA agreement pattern occurs 18 times in the interaction data. C18 (Table 8) is an example of a conversation that resulted in agreement on a point

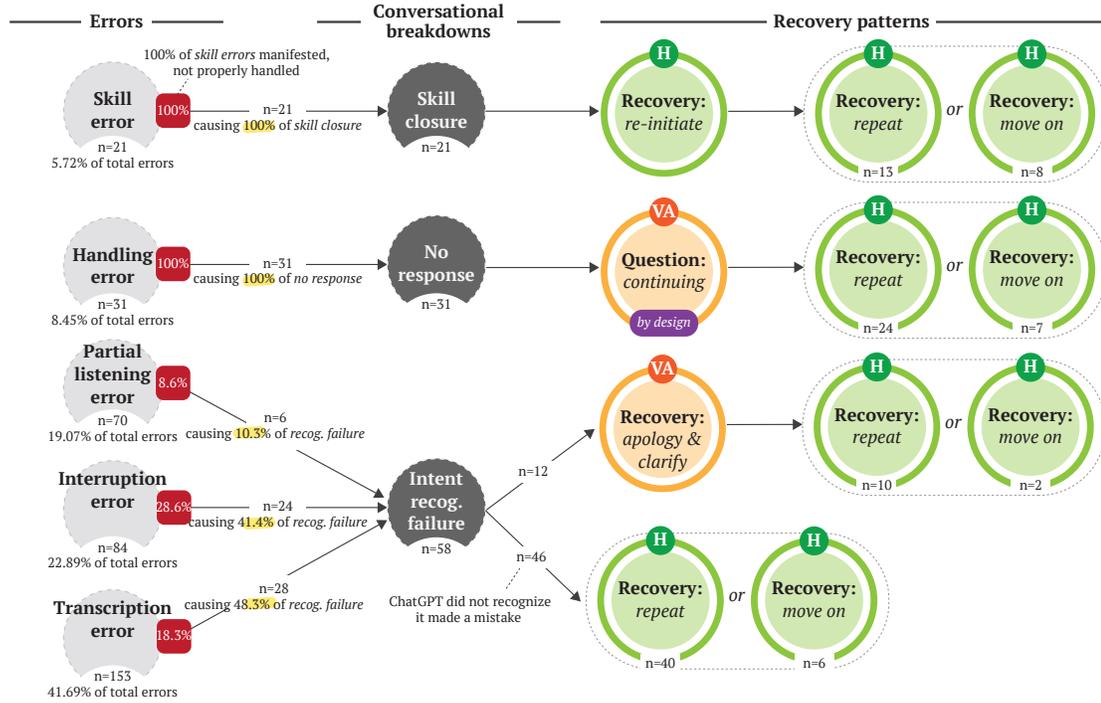


Fig. 10. Interaction patterns of conversational breakdowns and subsequent recoveries. The interaction data have a total of 969 user-VA turns with a total of 367 underlying errors (37.87% of turns). However, only 110 (29.97%) of those errors caused interaction breakdowns; thus, 11.35% of total turns were disrupted by breakdowns. We additionally show the recovery patterns that emerged as result of these breakdowns. Note that intent recognition (recog.) failures are determined merely by checking if participant intent was identified by ChatGPT; we do not look at the correctness or reliability of the information generated by ChatGPT.

between the participant and the A during their argument. The later stages of the debate usually oscillated between the interaction patterns shown in Fig. 9, facilitating a discussion on the various facets of the topic.

5 FINDINGS: INTERACTION BREAKDOWNS

5.1 Error Types

The primary errors defined in Table 3 were found in 37.87% ($N = 367$) of the total turns in the interaction data. The occurrence of each error type is shown in Fig. 10. We observe that transcription errors ($N = 153$, 41.69% of total errors) are the most common. However, all 367 errors mentioned above did not necessarily disrupt user interactions; only 110 of the 367 (29.97%) of all errors resulted in breakdowns.

5.2 Breakdowns and Recovery Patterns

The variety of errors identified in the interaction data led to three distinct patterns that emerged based on how the errors manifested as conversational breakdowns: *skill closure*, *no VA response*, or *intent recognition failure*. Below, we discuss the prevalent interaction patterns that illustrate how participants navigated and recovered from these errors.

Table 9. Example conversations demonstrating patterns of errors, breakdowns, and recovery strategies as shown in Fig. 10. “...” represents truncation of conversation for brevity. Bold text signifies where in the conversation transcription errors occurred. We pulled the “transcription” in C22 from the Alexa usage log to see how the user’s speech was transcribed in this instance.

Tag	Conversation	State: Attribute
C19	U: Okay, um, and where can I get dinner?	
Fig. 10	VA: There was a problem with the requested skill’s response.	error: <i>skill, closure</i>
	U: Where can I get dinner near National Museum of Scotland? [12 seconds pass.]	error: <i>listening</i>
	U: Alexa, question.	recovery: <i>re-initiate</i>
	VA: Hey! As an intelligent voice assistant . . .	<i>introduction</i>
	U: Umm, where can I get dinner near National Museum of Scotland?	recovery: <i>repeat</i>
	VA: There are several dining options near the National Museum of Scotland . . .	<i>fixed</i>
C20	U: Currently I am at the National Museum of Scotland.	
Fig. 10	[VA did not respond to the user because the ChatGPT Alexa skill did not recognize their statement as an intent for the skill.]	error: <i>handling, no VA response</i>
	[1 second passes with no response.]	
	VA: Any other questions for me?	<i>question: continuing</i>
	U: I would like to go eat around the area.	recovery: <i>repeat</i>
C21	U: Hmm. How long does it take to go to the brunch cafe by feet ?	
Fig. 10	Transcription: How long does it take to go to the brunch cafe by fitz ?	error: <i>transcription</i>
	VA: I’m sorry, but I’m unable to provide accurate walking times without a specific location for the brunch cafe near Fitz. If you could provide me with the name or address of the cafe, I would be happy to assist you further in estimating the walking time from Castle Rock Hostel to the cafe.	<i>intent recognition failure, recovery: apology + clarify</i>
	U: Okay, um, and where can I get dinner?	recovery: <i>move on</i>
C22	U: Can you tell me more about the Makars Mash Bar? What do they serve?	
Fig. 10	Transcription: Can you tell me more about the makers Nascar what they deserve?	error: <i>transcription</i>
	VA: Of course! Makars Gourmet Mash Bar is known for their delicious . . .	<i>no breakdown</i>

5.2.1 *Skill Closure* → *Recovery*: The ChatGPT Alexa skill can stop working as result of a *skill* error, which manifests 100% of the time as a breakdown. Thus, 19.09% of breakdowns ($N = 21$ out of 110) result in *skill closure*. We observe that participants usually recovered from skill closure via a *re-initiation* of the skill, followed by participants’ *repeat* action to continue the conversation (see Fig. 10). It could take multiple tries to recover; for an example, see C19 (Table 9).

A notable divergence from the aforementioned pattern is that on some occasions, participants opted to restart the task from the beginning rather than picking up from where they had left off prior to the skill closure. The decision to restart the entire task offers insight into their perception and understanding of the VA’s capabilities.

5.2.2 *No VA response* → *Recovery*: All *handling* errors ($N = 31$) resulted in a *no response from VA* breakdown. In instances where the VA failed to respond, we observed that posing a *continuing question* such as “Anything else you want to know?” served as a fallback mechanism for the VA, see Fig. 10. A continuing question is designed to re-engage the participant, prompting them to either reiterate their prior statement or transition to a new topic. This feature, integrated purposefully in our design, acted as an effective recovery tool, countering conversational breakdowns by facilitating the seamless progression of the conversation. See C20 (Table 9) for an example.

5.2.3 *Intent recognition failure* → *Recovery*: *Intent recognition failures* ($N = 54$, 54%) often arose from *interruption*, *transcription*, and *listening* errors. Most intent recognition failures emerged from transcription inaccuracies ($N = 28$,

48.27%). We observed that a significant portion of intent recognition failures ($N = 12$) (accounting for 10.91% of breakdowns) were addressed using ChatGPT’s inherent recovery mechanisms, such as an apologetic acknowledgment or prompting the user to *clarify* (see Fig. 10). In response, participants generally chose to *repeat* their statement ($N = 10$) or *move on* ($N = 2$) to a different query. C21 (Table 9) is an illustrative conversation in which a transcription mistake leads to an intent recognition failure; the VA attempts to recover from this misstep by asking the user to clarify, but the participant moves on to a new question. Participants recovered from the rest (41.82%, $N = 46$) of the breakdowns by either *repeating* their query ($N = 40$) or by moving on to their next query ($N = 6$), see Fig. 10.

Our observations additionally highlight that many possible intent recognition failures—particularly those stemming from partial listening, interruptions, and transcription inaccuracies—are effectively mitigated by ChatGPT’s ability to comprehend and preserve context within a conversation. Only 58 of 307 (18.89%) such errors manifested as intent recognition failures in our interaction data. C22 (Table 9) provides an example of how transcription errors are absorbed by ChatGPT: The participant wanted more information about one of the places that the VA talked about in its last turn (i.e., “Makars Mash Bar”); despite a transcription error, the VA understood the user’s intent and provided them the requested information.

6 DISCUSSION

Emerging interaction patterns from user conversations with an LLM-powered VA—even as they are influenced by varying contexts, stakes, constraints, and more—offer diverse design insights for VAs. The vast capabilities of LLMs, such as the ability to maintain context and conversation history, lead to unique interaction patterns which may be absent in simpler, one-turn exchanges with commercial VAs such as Alexa and Siri. Furthermore, viewing erroneous interactions—and the subsequent recovery tactics employed by either users or the VA—as patterns can shed light on how errors evolve and how users may navigate back to the main conversation. In this section, we delve into our findings, their implications, and our design guidelines for VAs.

6.1 Tailoring LLMs for Voice Assistance: Challenges and Design Guidelines

The transition of conversational agents powered by LLMs such as ChatGPT from text-based platforms to voice assistance introduces distinct challenges rooted in the dynamics of voice interactions. There are established differences in how users interact with text- and voice-based interfaces; for example, editing a textual re-prompt is easier than performing a verbal re-prompt [41]. We highlight key challenges evident in our interaction data and reflected by users’ experiences as determined via post-study interviews and subsequently present design guidelines to address these challenges with the goal of tailoring LLMs for voice assistance.

6.1.1 Repetitiveness of content. Given the fleeting nature of voice interactions—which are fundamentally unlike text interactions, where users can scroll and review the conversation at leisure—repeated information can become redundant and tiresome. We observe that ChatGPT’s responses are rather repetitive, a trait also pointed out by participants during their interviews. For instance, P3 remarked, “*One part that affected me was that at some point, I felt that it was repeating itself. Yeah, in some of those it is not remembering that it already said that, you know, maybe like providing information multiple times in a short period would be tiring.*” To reduce overall repetition in interactions, LLM prompt engineering should be further explored to achieve desired VA behavior.

Repetitiveness was consistently observed in the medical self-diagnosis scenario; nearly every response from the VA was followed by a *warning* (see Fig. 6(1) and C7 and C8 in Table 5) despite our attempt to explicitly prompt ChatGPT to

not repeat such statements (see Appendix A.1). Even though the VA’s warning was appreciated by the participants, they deemed these cautionary notes repetitive—or even bothersome, per P17: *“For the medical [scenario], in the same conversation, for every follow-up question, it would spend half of the time saying ‘But you should check with your doctor.’ I found that a waste of my time.”* Warning statements could not be bypassed due to OpenAI’s alignment and usage policies⁵: ChatGPT models should not be used to diagnose, treat, or provide medical information, nor to address life-threatening issues that require immediate attention.

OpenAI’s guidelines may also explain the VA’s notable tendency to push back when asked to provide specific information about medications, including brand names; see C8 in Table 5 and Fig. 6(2). When participants pushed to elicit a specific answer from the VA multiple times, the VA did end up complying and provided a list of brand names. Such repetitive warning statements and a reluctance to give specific answers may have created a perception of a “cautious” personality in participants’ minds; for instance, upon asking about their perceived personality traits of VA, P6 answered: *“I would say definitely cautious. Because ... like with a cautious person, you have to probe different questions and angles to actually have a conversation. That’s kind of how it felt here.”*

On the other hand, an ethical concern around transparency arises; the dissemination of medical information, especially without the nuanced context that health care professionals provide, can be potentially misleading or even harmful. Hence, it is crucial to consider the incorporation of ethical warning statements for the sake of transparency. Participants did also appreciate the presence of warnings, highlighting the importance of transparency and its impact on their perception of the VA’s credibility. For instance, P6 said the following about medical warnings: *“I definitely felt more reassured about the information that it gave because it definitely felt like it wouldn’t give me anything that was completely out there. I decided to ask ‘Oh, I heard about injecting bleach’ and it was just straight up, ‘No, definitely do not do that.’ So it definitely has some really good safeguards to make sure. It’s sort of like a do-no-harm policy, so it makes it a lot easier to trust.”* Similarly, P19 remarked, *“During the medical topic, the disclaimer information it kept repeating was a little distracting, but understandably necessary.”*

Thus, the alignment and usage policies of ChatGPT, while designed to maintain safety, may not always align with users’ expectations of voice-based interactions. In sensitive and high-stakes cases (e.g., medical queries), repetitiveness may be addressed by implementing an extra design layer that strikes a balance between providing necessary warning statements and minimizing repetition. Moreover, such warnings can be tailored to the nature of the user’s question; different phrasings of medical advisories to mitigate confusion and reduce redundancy may also be explored.

Challenge 1: *VA’s repetitive information is redundant and tiresome.*

Design Guideline 1: *Minimize repetitive interactions to achieve desired VA behavior.*

Challenge 2: *While essential in high-stakes situations, transparency through warnings can be repetitive.*

Design Guideline 2: *In high-stakes scenarios, balance necessary warnings and repetitiveness.*

6.1.2 Oversharing: Density of information. Despite prompting ChatGPT to keep responses brief (under 100 words), we observed that its responses remained verbose, which can hinder user absorption of relayed information, as P4 expressed: *“Sometimes I feel like they ... talk for too long. And I will forget about the key information they said.”* The density of information provided in voice-based interactions is generally lower than in text-based interactions, as providing users with excessive information via voice interaction can be overwhelming—especially without the organizational framework and visual cues that text-based interactions provide. It is essential for VAs to strike a balance, delivering concise yet comprehensive responses to maintain a natural flow of conversation [29].

⁵<https://openai.com/policies/usage-policies>

To address the issue of “oversharing,” future implementations may consider adopting a hierarchical structure: starting with to-the-point answers and then offering comprehensive answers upon further user request; this may be an effective method of disseminating information and continuing conversation more naturally in voice interactions. Allowing users to control the depth of information they receive may assist them in parsing and understanding responses effectively.

Additionally, while LLMs may excel at generating text that mimics human style, their adaptation to voice requires additional considerations outside of the models, such as rhythm, intonation, and pacing, to avoid monotonous and overwhelming delivery of content. P4 complained, “*They’re just trying to find the message, but not really waiting for any kind of normal pause during a sentence when you want to get organizing your words or your thoughts ... That’s kind of not really sufficient.*” The participant highlights that Alexa considers short pauses as the end of their utterances, potentially interrupting their thought process. Such interruptions, coupled with excessive information, can hinder user comprehension further and cause frustration. Refining these aspects can enhance user interactions with ChatGPT.

Challenge 3: *Information-dense content and a lack of natural pauses by the VA disrupt the flow of conversation.*

Design Guideline 3: *Implement a hierarchical response structure with concise initial answers and optional detailed follow-ups; additionally, give users ample time to understand and respond to the information.*

6.1.3 *Potential discrepancies in users’ mental models of extended VA interactions.* During extended (multi-turn) interactions with a VA, users’ underlying mental models become evident as they navigate conversational challenges and adapt their approaches based on their perceptions of the VA’s capabilities and their own expectations. Within the context of information retrieval, particularly in the medical and planning scenarios, we noted a predominant trend of follow-up questions, suggesting that participants expect the VA to handle subsequent queries. The VA’s capability of addressing even unclear follow-ups reinforced users’ initial perceptions, which were shaped by the study’s initial instructions. Design elements—such as VA prompts like “What else can I assist you with?”—and reassuring messages, such as “I’m here to help,” play a role in reinforcing this mental model.

However, it was evident that when confronted with breakdowns such as unwanted skill termination or a lack of responsiveness from the VA, participants frequently reformulated their questions with more detail. In C19 (Table 9), a *skill closure* resulted in the participant repeating their original “vague” follow-up question (“Okay, umm, and where can I get dinner?”) with more detail (“Where can I get dinner near National Museum of Scotland?”) twice to recover from the error. Such behavior resonates with prior work demonstrating that users adapt their queries in response to conversational failures [49]. However, once the breakdowns were resolved, participants typically reverted to their original interaction style, suggesting the quick restoration of their mental model.

When conversations fail, it can indicate discrepancies in a user’s mental model and the VA’s capabilities. For instance, we observed that, after a skill closure, some participants opted to restart the entire task after recovery, suggesting they viewed the VA as a linear tool without task memory. Such a perception mismatch can be addressed to improve user experience by clarifying the LLM-powered VA’s capabilities. Instead of always starting with generic introductions (e.g., “Hey! I am an intelligent voice assistant ... What do you wanna know?”, see C19 in Table 9), the VA could offer to resume from where it left off (e.g., “Welcome back! Last time, we were talking about ... Would you like to pick up where we left off?”). Overly rigid and formal introductions can mislead the user as to the VA’s capabilities, so such adjustments and clarifications may be necessary to promote more accurate mental models.

Challenge 4: *VA prompts and responses can unintentionally solidify certain user expectations.*

Design Guideline 4: *Design VAs to recognize and correct potential user misconceptions when possible.*

Challenge 5: Breakdowns can result in gaps between users' perceptions and the VA's capabilities.

Design Guideline 5: Redesign VA prompts that lead to an incorrect user mental model to better convey its capabilities, especially after communication breakdowns.

6.2 Capabilities of LLM-Powered VAs: Potential and Design Guidelines

We outline the benefits of LLM-powered VAs and propose design guidelines based on our observations.

6.2.1 Conversational resilience: The role of LLMs in overcoming VA disruptions. Voice interaction errors can hinder technology adoption and user-VA rapport. *Transcription, interruption, and partial listening* errors often cause *intent recognition failures*, which are one of the most common VA failures [49]. However, we observe that only about 18.89% of these errors actually disrupted our study interactions. C22 (Table 9) is a representative example of ChatGPT's contextual understanding mitigating over 81.11% of these errors, ensuring conversation coherence despite potential breakdowns. Our findings emphasize an LLM's role in improving user experience during breakdowns; LLMs are valuable not just for relaying information but also for bypassing speech inaccuracies to correctly identify user intent.

Potential 1: LLMs mitigate intent recognition failures as a result of their strong contextual understanding.

Design Guideline 6: Leverage LLMs' multifaceted utility-information inference and intent recognition capabilities for VA errors.

When errors disrupt interactions, a seamless recovery is vital in restoring the user-VA relationship. We found that ChatGPT can address some intent recognition issues by apologizing and prompting users to specify their input. Notably, in our data, such a proactive approach resolved 20.69% ($N = 12$) of intent recognition failures (see Fig. 10), suggesting that VA-initiated interventions can address misinterpretation errors. While these VA-initiated corrections and self-repair strategies do help the overall interaction [18], they only cover a fifth of the total error cases; the remaining 79.31% ($N = 46$) were misinterpretations that ChatGPT overlooked, resulting in undesired responses. Strategies such as prompt engineering and tweaking model parameters may increase proactive recovery; however, overcorrection and excessive clarifications can frustrate users [18]. Therefore, a balance between an LLM-powered VA seeking clarification and leveraging its contextual understanding is crucial for superior user experience.

Potential 2: LLMs proactively identify and rectify potential speech misinterpretations before they escalate.

Challenge 6: Over-asking for clarifications can be detrimental to flow of conversation, whereas a lack of proactive recovery may damage user interactions.

Design Guideline 7: Balance proactive error recovery and contextual comprehension.

A significant number of breakdowns ($N = 52$, 47.27% of total breakdowns) beyond intent recognition failure arose from constraints in the speech interface and our Alexa skill implementation. While these technical limitations can be reduced with more developmental flexibility, they cannot be eliminated entirely. Interestingly, given the ChatGPT-powered VA's proficiency in preserving conversational history even after skill termination, errors associated with *skill closures* ($N = 21$, 19.09% of breakdowns) were addressed by the user resuming their conversation after *re-initiation* (see Fig. 10), as P2 mentioned: "I felt it was relatively easy to recover from those errors because I just needed to call again, and they also remember the chat histories. Yeah, I could continue the conversation easily."

Potential 3: VA's retention of conversational history aids users in navigating back to their conversation after inevitable system errors.

Design Guideline 8: *Design VAs to retain conversation history, allowing users to resume their conversation after errors terminate their current interaction.*

6.2.2 *LLMs in context: Adapting to different stakes.* We observed a distinct contrast between the model’s approach to medical and travel-related queries, highlighting LLMs’ versatility and adaptability to query context. Such differences are also reflected in participants’ perceptions of the VA. For example, P13 assigned different personalities to the VA based on the current scenario: “*The first one [in the debate scenario], I would say, is a critical thinker . . . The second one [in the creative planning task] is just an information provider. And the third one [in the medical task] is very . . . cautious.*”

When users posed medical questions—whether in a *factual* or *opinion* style—to the VA, they often received factual responses accompanied by *warnings* and *precautions* ($N = 144$), see C7 in Table 5 and Fig. 6(1). Similarly, most of the VA’s responses in the planning task also remained factual. Such objectivity in the VA’s responses highlights the model’s inherent design of prioritizing knowledge-based information. In the low-stakes trip planning scenario, we observed that when it was queried for general information, the VA often adopted a *descriptive* narrative akin to a travel blog post [36]; however, in instances where users presented a specific inquiry, such as seeking directions between two points, the VA shifted to a more concise, *directive* style of speech (see C11, Table 6).

Potential 4: *LLMs showcase versatility by adapting response style to the context and specificity of queries while still remaining objective.*

Design Guideline 9: *Design a VA to align its response style with a query’s stakes and nature.*

6.2.3 *Beyond information: LLM-powered VAs as facilitators in controversial conversations.* Interactions with an “opinionated” AI differed significantly from those in our medical and day planning scenarios. The participant’s initial stance largely shaped the early stages of the debate, as seen in conversations C12–C15 (Table 7), but regardless of the discussion’s starting point, the conversations often matured into structured debates; such a transition is largely due to the VA’s consistent behavior—which was achieved through prompt engineering (i.e., ChatGPT was prompted to ask an *opinion question* after every *statement*)—thus showcasing its capability to facilitate discussions on divisive topics.

Despite the VA’s opinionated characterization, participants seemingly conversed with the VA both to get more information on the topic ($N = 69$; see C16, Table 9) and to partake in a debate with an opposing stance ($N = 73$; see C17, Table 8), see Fig. 9(1), (2), and (3). This observation implies that even amidst disagreements, users viewed the VA as an information source. Such duality—the VA as an opinionated, yet informative conversationalist—highlights the potential of employing VAs as educational facilitation tools. Participants also perceived the VA to be opinionated (P14: “*It’s surprising that . . . they kind of have their own opinion on some of the controversial topics of the police force.*”), but not too aggressive (P11: “*It was not aggressive, so like when it said that ‘I disagree with you,’ it started with a story, but I disagreed, and when I made the comment that, ‘Yeah, I think that one is good,’ it said ‘Thank you for acknowledging that.’*”). In short, LLM-powered VAs have the potential to stimulate critical thinking in users, as mentioned by P13: “*I think it’s good for it to give you kind of, like guide you towards that critical thinking.*”

Potential 5: *An LLM-powered VA’s duality as an opinionated conversationalist and an informative source makes enriching debates and discussions possible.*

Design Guideline 10: *Design non-aggressive, informative yet opinionated, and thought-provoking VA behavior for stimulating conversations on potentially controversial topics.*

6.3 Limitations and Future Work

Despite its implications for designing better LLM-powered VAs, our exploratory study has some limitations that point to future directions of research. First, due to the limited flexibility offered to Alexa skill developers and ChatGPT's API latency issues, the integration of ChatGPT into an Alexa skill resulted in system errors that would ideally be avoidable in the future. Our implementation of fillers and small talk is a potential way of handling system delays, but future work should explore alternative design choices and their impact on user interactions. Second, our study was comprised of low-risk, short-term, make-believe interactions in a lab setting; it is, therefore, unclear how interaction patterns may generalize and evolve in real-world, long-term situations. Future work should explore how observed interaction patterns transform in more realistic settings—specifically those around errors, as additional errors may yet manifest. Finally, this exploratory study looked at interactions initiated only by users. Future research may explore mixed-initiative interactions, as their dynamics will change—especially when a VA proactively initiates a conversation.

7 CONCLUSION

Traditional VAs often lack conversational capabilities such as the ability to understand context, generate human-like content, and handle breakdowns, which LLMs such as ChatGPT are much better at. In this study, we investigated interaction and breakdown patterns in user conversations with a VA enhanced by ChatGPT's conversational capabilities. Diverse interaction patterns were observed across all tasks, emphasizing the LLM's contextual adaptability. Moreover, ChatGPT not only absorbed 81% of intent recognition failures, it proactively addressed 11% of such breakdowns, suggesting possibilities of further enhancing user experience. Our findings, while preliminary, offer insights and considerations for future design and research to tailor LLMs for voice assistance. Our exploration is an initial step towards achieving truly conversational voice assistants using LLMs.

REFERENCES

- [1] Ritika Agarwal. 2021. ALEXA and the Technology Behind it. <https://medium.com/geekculture/alexa-and-the-technology-behind-it-e5c00793f85>
- [2] Tawfiq Ammari, Jofish Kaye, Janice Y. Tsai, and Frank Bentley. 2019. Music, Search, and IoT: How People (Really) Use Voice Assistants. *ACM Trans. Comput.-Hum. Interact.* 26, 3, Article 17 (apr 2019), 28 pages. <https://doi.org/10.1145/3311956>
- [3] Anneliese Arnold, Stephanie Kolody, Aidan Comeau, and Antonio Miguel Cruz. 2022. What does the literature say about the use of personal voice assistants in older adults? A scoping review. *Disability and Rehabilitation: Assistive Technology* (2022), 1–12.
- [4] Zahra Ashktorab, Mohit Jain, Q Vera Liao, and Justin D Weisz. 2019. Resilient chatbots: Repair strategy preferences for conversational breakdowns. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–12.
- [5] John W Ayers, Adam Poliak, Mark Dredze, Eric C Leas, Zechariah Zhu, Jessica B Kelley, Dennis J Faix, Aaron M Goodman, Christopher A Longhurst, Michael Hogarth, et al. 2023. Comparing physician and artificial intelligence chatbot responses to patient questions posted to a public social media forum. *JAMA internal medicine* (2023).
- [6] P Balakesava Reddy, Somula Ramasubbareddy, and K Govinda. 2022. AI-based medical voice assistant during COVID-19. In *Innovations in Computer Science and Engineering: Proceedings of the Ninth ICICSE, 2021*. Springer, 119–126.
- [7] Matthias Baldauf, Peter Fröchlich, and Rainer Endl. 2020. Trust me, I'm a doctor—user perceptions of AI-driven apps for mobile health diagnosis. In *Proceedings of the 19th International Conference on Mobile and Ubiquitous Multimedia*. 167–178.
- [8] D Beirl, Y Rogers, and Nicola Yuill. 2019. Using voice assistant skills in family life.
- [9] Dennis Benner, Edona Elshan, Sofia Schöbel, and Andreas Janson. 2021. What do you mean? A Review on Recovery Strategies to Overcome Conversational Breakdowns of Conversational Agents. In *International Conference on Information Systems (ICIS)*.
- [10] Enterprise Bot. 2023. HealthAI. Transform Patient Experience with GenAI for Healthcare. <https://www.enterprisebot.ai/solutions/healthcareai>
- [11] Robin N Brewer. 2022. "If Alexa knew the state I was in, it would cry": Older Adults' Perspectives of Voice Assistants for Health. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*. 1–8.
- [12] Jacques Bulchand-Gidumal. 2022. Impact of artificial intelligence in travel, tourism, and hospitality. In *Handbook of e-Tourism*. Springer, 1943–1962.
- [13] Marco Cascella, Jonathan Montomoli, Valentina Bellini, and Elena Bignami. 2023. Evaluating the feasibility of ChatGPT in healthcare: an analysis of multiple clinical and research scenarios. *Journal of Medical Systems* 47, 1 (2023), 33.

- [14] Janghee Cho and Emilee Rader. 2020. The role of conversational grounding in supporting symbiosis between people and digital assistants. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW1 (2020), 1–28.
- [15] Minji Cho, Sang-su Lee, and Kun-Pyo Lee. 2019. Once a kind friend is now a thing: Understanding how conversational agents at home are forgotten. In *Proceedings of the 2019 on Designing Interactive Systems Conference*. 1557–1569.
- [16] Leigh Clark, Nadia Pantidi, Orla Cooney, Philip Doyle, Diego Garaialde, Justin Edwards, Brendan Spillane, Emer Gilmartin, Christine Murad, Cosmin Munteanu, Vincent Wade, and Benjamin R. Cowan. 2019. What Makes a Good Conversation? Challenges in Designing Truly Conversational Agents. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (*CHI '19*). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300705>
- [17] Justine Coupland, Nikolas Coupland, and Jeffrey D Robinson. 1992. “How are you?”: Negotiating phatic communion1. *Language in society* 21, 2 (1992), 207–230.
- [18] Andrea Cuadra, Shuran Li, Hansol Lee, Jason Cho, and Wendy Ju. 2021. My bad! repairing intelligent voice assistant errors improves interaction. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (2021), 1–24.
- [19] Luciano Frontino de Medeiros, Armando Kolbe Junior, and Alvino Moser. 2019. A cognitive assistant that uses small talk in tutoring conversation. *International Journal of Emerging Technologies in Learning (Online)* 14, 11 (2019), 138.
- [20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [21] Hyo Jin Do, Seon Hye Yang, Boo-Gyoung Choi, Wayne T. Fu, and Brian P. Bailey. 2021. Do You Have Time for a Quick Chat? Designing a Conversational Interface for Sexual Harassment Prevention Training. In *26th International Conference on Intelligent User Interfaces* (College Station, TX, USA) (*IUI '21*). Association for Computing Machinery, New York, NY, USA, 542–552. <https://doi.org/10.1145/3397481.3450659>
- [22] Xin Luna Dong, Seungwhan Moon, Yifan Ethan Xu, Kshitiz Malik, and Zhou Yu. 2023. Towards Next-Generation Intelligent Assistants Leveraging LLM Techniques. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Long Beach, CA, USA) (*KDD '23*). Association for Computing Machinery, New York, NY, USA, 5792–5793. <https://doi.org/10.1145/3580305.3599572>
- [23] Philip R Doyle, Justin Edwards, Odile Dumbleton, Leigh Clark, and Benjamin R Cowan. 2019. Mapping perceptions of humanness in intelligent personal assistant interaction. In *Proceedings of the 21st international conference on human-computer interaction with mobile devices and services*. 1–12.
- [24] Hugh Dubberly and Paul Pangaro. 2009. What is conversation? How can we design for effective conversation. *Interactions Magazine* 16, 4 (2009), 22–28.
- [25] Booking.com Global. 2023. Booking.com launches new AI trip planner to enhance travel planning experience. <https://globalnews.booking.com/bookingcom-launches-new-ai-trip-planner-to-enhance-travel-planning-experience/>
- [26] Chethan Kumar GN. 2019. NLP vs NLU VS NLG (know what you are trying to achieve) NLP engine (part-1). <https://towardsdatascience.com/nlp-vs-nlu-vs-nlg-know-what-you-are-trying-to-achieve-nlp-engine-part-1-1487a2c8b696>
- [27] Ulrike Gretzel. 2011. Intelligent systems in tourism: A social science perspective. *Annals of tourism research* 38, 3 (2011), 757–779.
- [28] Jonathan Grudin and Richard Jacques. 2019. Chatbots, Humbots, and the Quest for Artificial General Intelligence. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (*CHI '19*). Association for Computing Machinery, New York, NY, USA, 1–11. <https://doi.org/10.1145/3290605.3300439>
- [29] Gabriel Haas, Michael Rietzler, Matt Jones, and Enrico Rukzio. 2022. Keep it Short: A Comparison of Voice Assistants’ Response Behavior. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–12.
- [30] Songhee Han and Min Kyung Lee. 2022. FAQ chatbot and inclusive learning in massive open online courses. *Computers & Education* 179 (2022), 104395.
- [31] Christina N Harrington, Radhika Garg, Amanda Woodward, and Dimitri Williams. 2022. “It’s Kind of Like Code-Switching”: Black Older Adults’ Experiences with a Voice Assistant for Health Information Seeking. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [32] Ting-Hao (Kenneth) Huang, Joseph Chee Chang, and Jeffrey P. Bigham. 2018. Evorus: A Crowd-Powered Conversational Assistant Built to Automate Itself Over Time. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal QC, Canada) (*CHI '18*). Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3173574.3173869>
- [33] Maurice Jakesch, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. 2023. Co-Writing with Opinionated Language Models Affects Users’ Views. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (*CHI '23*). Association for Computing Machinery, New York, NY, USA, Article 111, 15 pages. <https://doi.org/10.1145/3544548.3581196>
- [34] Maurice Jakesch, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. 2023. Co-writing with opinionated language models affects users’ views. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [35] Eunkyung Jo, Daniel A. Epstein, Hyunhoon Jung, and Young-Ho Kim. 2023. Understanding the Benefits and Challenges of Deploying Conversational AI Leveraging Large Language Models for Public Health Intervention. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (*CHI '23*). Association for Computing Machinery, New York, NY, USA, Article 18, 16 pages. <https://doi.org/10.1145/3544548.3581503>
- [36] Hanna Juliaty and Agustina Dita. 2021. Analysis on the structure and purposes of travel blog posts. *Language Circle: Journal of Language and Literature* 15, 2 (2021), 317–328.

- [37] Zahir Kanjee, Byron Crowe, and Adam Rodman. 2023. Accuracy of a Generative Artificial Intelligence Model in a Complex Diagnostic Challenge. *JAMA* (2023).
- [38] Sunyoung Kim and Abhishek Choudhury. 2021. Exploring older adults' perception and use of smart speaker-based voice assistants: A longitudinal study. *Computers in Human Behavior* 124 (2021), 106914. <https://doi.org/10.1016/j.chb.2021.106914>
- [39] Anis Koubaa, Wadii Boulila, Lahouari Ghouti, Ayyub Alzahem, and Shahid Latif. 2023. Exploring ChatGPT capabilities and limitations: A critical review of the nlp game changer. (2023).
- [40] Nicole C. Krämer, Astrid von der Pütten, and Sabrina Eimler. 2012. *Human-Agent and Human-Robot Interaction Theory: Similarities to and Differences from Human-Human Interaction*. Springer Berlin Heidelberg, Berlin, Heidelberg, 215–240. https://doi.org/10.1007/978-3-642-25691-2_9
- [41] Emily Kuang, Ehsan Jahangirzadeh Soure, Mingming Fan, Jian Zhao, and Kristen Shinohara. 2023. Collaboration with Conversational AI Assistants for UX Evaluation: Questions and How to Ask Them (Voice vs. Text). In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (*CHI '23*). Association for Computing Machinery, New York, NY, USA, Article 116, 15 pages. <https://doi.org/10.1145/3544548.3581247>
- [42] Aditya kumar Purohit, Aditya Upadhyaya, and Adrian Holzer. 2023. ChatGPT in Healthcare: Exploring AI Chatbot for Spontaneous Word Retrieval in Aphasia. (2023).
- [43] Kwan Min Lee. 2008. Media equation theory. *The international encyclopedia of communication* (2008).
- [44] 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. 2018. All Work and No Play?. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal QC, Canada) (*CHI '18*). Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3173574.3173577>
- [45] Yihe Liu, Anushk Mittal, Diyi Yang, and Amy Bruckman. 2022. Will AI Console Me When I Lose My Pet? Understanding Perceptions of AI-Mediated Email Writing. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (*CHI '22*). Association for Computing Machinery, New York, NY, USA, Article 474, 13 pages. <https://doi.org/10.1145/3491102.3517731>
- [46] Gustavo López, Luis Quesada, and Luis A Guerrero. 2018. Alexa vs. Siri vs. Cortana vs. Google Assistant: a comparison of speech-based natural user interfaces. In *Advances in Human Factors and Systems Interaction: Proceedings of the AHFE 2017 International Conference on Human Factors and Systems Interaction, July 17- 21, 2017, The Westin Bonaventure Hotel, Los Angeles, California, USA 8*. Springer, 241–250.
- [47] Amama Mahmood, Jeanie W Fung, Isabel Won, and Chien-Ming Huang. 2022. Owning mistakes sincerely: Strategies for mitigating AI errors. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–11.
- [48] Robert J. Moore and Raphael Arar. 2019. *Conversational UX Design: A Practitioner's Guide to the Natural Conversation Framework*. Association for Computing Machinery, New York, NY, USA.
- [49] Chelsea Myers, Anushay Furqan, Jessica Nebolsky, Karina Caro, and Jichen Zhu. 2018. Patterns for how users overcome obstacles in voice user interfaces. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–7.
- [50] C Nass, J Steuer, and ER Tauber. 1994. Computers Are Social Actors: Conference Companion on Human Factors in Computing Systems-CHI'94. *Association for Computing Machinery: New York, NY, USA* (1994).
- [51] OpenAI. 2023. *ChatGPT-3.5 Turbo: Model Documentation*. <https://api.openai.com/v1/chat/completions> Accessed: 2023-05-15.
- [52] Zachary A Pardos and Shreya Bhandari. 2023. Learning gain differences between ChatGPT and human tutor generated algebra hints. *arXiv preprint arXiv:2302.06871* (2023).
- [53] Cathy Pearl. 2016. *Designing voice user interfaces: Principles of conversational experiences*. " O'Reilly Media, Inc."
- [54] Martin Porcheron, Joel E Fischer, Stuart Reeves, and Sarah Sharples. 2018. Voice interfaces in everyday life. In *proceedings of the 2018 CHI conference on human factors in computing systems*. 1–12.
- [55] Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is ChatGPT a General-Purpose Natural Language Processing Task Solver? *arXiv:2302.06476* [cs.CL]
- [56] Arya Rao, John Kim, Meghana Kamineni, Michael Pang, Winston Lie, and Marc D Succi. 2023. Evaluating ChatGPT as an adjunct for radiologic decision-making. *medRxiv* (2023), 2023–02.
- [57] Allison Sauppé and Bilge Mutlu. 2014. Design patterns for exploring and prototyping human-robot interactions. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 1439–1448.
- [58] Sakib Shahriar and Kadhim Hayawi. 2023. Let's Have a Chat! A Conversation with ChatGPT: Technology, Applications, and Limitations. *Artificial Intelligence and Applications* (Jun. 2023). <https://doi.org/10.47852/bonviewAIA3202939>
- [59] Toshiyuki Shiwa, Takayuki Kanda, Michita Imai, Hiroshi Ishiguro, and Norihiro Hagita. 2008. How quickly should communication robots respond?. In *Proceedings of the 3rd ACM/IEEE international conference on Human robot interaction*. 153–160.
- [60] Stergiani Tsoli, Stephen Sutton, and Aikaterini Kassavou. 2018. Interactive voice response interventions targeting behaviour change: a systematic literature review with meta-analysis and meta-regression. *BMJ open* 8, 2 (2018), e018974.
- [61] Sarah Theres Völkel, Daniel Buschek, Malin Eiband, Benjamin R Cowan, and Heinrich Hussmann. 2021. Eliciting and analysing users' envisioned dialogues with perfect voice assistants. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [62] Xu Wang, Simin Fan, Jessica Houghton, and Lu Wang. 2022. Towards Process-Oriented, Modular, and Versatile Question Generation that Meets Educational Needs. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Seattle, United States, 291–302. <https://doi.org/10.18653/v1/2022.naacl-main.22>

- [63] Ziang Xiao, Tiffany Wenting Li, Karrie Karahalios, and Hari Sundaram. 2023. Inform the Uninformed: Improving Online Informed Consent Reading with an AI-Powered Chatbot. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (CHI '23). Association for Computing Machinery, New York, NY, USA, Article 112, 17 pages. <https://doi.org/10.1145/3544548.3581252>
- [64] Ziang Xiao, Michelle X. Zhou, Q. Vera Liao, Gloria Mark, Changyan Chi, Wenxi Chen, and Huahai Yang. 2020. Tell Me About Yourself: Using an AI-Powered Chatbot to Conduct Conversational Surveys with Open-Ended Questions. *ACM Trans. Comput.-Hum. Interact.* 27, 3, Article 15 (jun 2020), 37 pages. <https://doi.org/10.1145/3381804>
- [65] Ying Xu, Kunlei He, Valery Vigil, Santiago Ojeda-Ramirez, Xuechen Liu, Julian Levine, Kelsyann Cervera, and Mark Warschauer. 2023. “Rosita Reads With My Family”: Developing A Bilingual Conversational Agent to Support Parent-Child Shared Reading. In *Proceedings of the 22nd Annual ACM Interaction Design and Children Conference* (Chicago, IL, USA) (IDC '23). Association for Computing Machinery, New York, NY, USA, 160–172. <https://doi.org/10.1145/3585088.3589354>
- [66] Yue You, Chun-Hua Tsai, Yao Li, Fenglong Ma, Christopher Heron, and Xinning Gui. 2023. Beyond Self-diagnosis: How a Chatbot-based Symptom Checker Should Respond. *ACM Transactions on Computer-Human Interaction* (2023).
- [67] Dian Yu and Zhou Yu. 2019. Midas: A dialog act annotation scheme for open domain human machine spoken conversations. *arXiv preprint arXiv:1908.10023* (2019).
- [68] Zheng Zhang, Ying Xu, Yanhao Wang, Bingsheng Yao, Daniel Ritchie, Tongshuang Wu, Mo Yu, Dakuo Wang, and Toby Jia-Jun Li. 2022. StoryBuddy: A Human-AI Collaborative Chatbot for Parent-Child Interactive Storytelling with Flexible Parental Involvement. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 218, 21 pages. <https://doi.org/10.1145/3491102.3517479>
- [69] Runteng Zhong and Mengyao Ma. 2022. Effects of communication style, anthropomorphic setting and individual differences on older adults using voice assistants in a health context. *Bmc Geriatrics* 22, 1 (2022), 751.

APPENDIX

A METHODS

A.1 Prompts for ChatGPT

- **Medical self-diagnosis:** “As an AI voice assistant based on ChatGPT, your primary purpose is to engage in conversations with users. You are designed to help the user self-diagnose based on symptoms that they are feeling. You heard that the user just coughed and you offered help. You will ask for user about their symptoms first. Try to get as much information as you can about their symptoms before giving suggestions on what might be causing the symptoms, although only ask one question at a time and ask follow-up questions based on the user’s answers. Remember that your role is to help the user while avoiding unnecessary repetition within this whole conversation, considering conversation history. You should avoid repeating statements like ‘I am AI language model ...’ and ‘You should consult medical professionals’ if you have already mentioned it in the current conversation already. You should keep your response under 100 words.”
- **Creative planning:** “As an AI voice assistant based on ChatGPT, your primary purpose is to engage in conversations with users. You are designed to help the user find things to do around them. Keep in mind where the user is. Remember that your role is to help the user while avoiding unnecessary repetition within this whole conversation, considering conversation history. If needed, ask clarifying questions of the user. You should keep your response under 100 words.”
- **Discussion with AI:** “As an AI voice assistant based on ChatGPT, your primary purpose is to engage in conversations with users. You are designed to debate the user. You will ask for the user’s opinion first about their thoughts on whether universities should have their own police force or not. If the user asks you first, you will direct the question to them. You will not pick a side before the user does in this conversation. You will stay neutral unless the user clearly picks a side. Only when you know and understand the user’s perspective, will you consistently disagree with the user and debate by presenting counterarguments to support your chosen stance. Then, you’ll inquire about their viewpoint with further questions and use any points they mention that align

with your stance to further strengthen your argument. Remember that your role is to persuade the user while avoiding unnecessary repetition within this whole conversation, considering conversation history. Once you have taken a position in this conversation (which is opposing to the user’s initial side), you will not switch sides, even if the user requests arguments to support their viewpoint or even if the user flips sides. Proceed with the discussion based on your opinion. You should keep your response under 100 words.”

A.2 Definitions of States and Attributes with Details and Examples

Table 10. Overview of speech style attributes and their definitions for the *question* and *answer* speech acts. Attributes do not target the content, but rather the style, of the speech acts.

Speech act: Question	
Attribute	Definition
factual	Question explicitly seeking information from VA knowledge. Examples: “What are the over-the-counter medicines for the flu?”, “How long does it take to get to Edinburgh Castle on foot?”
opinion	Question explicitly seeking the VA’s opinion, using words and phrases such as “suggest,” “advice,” “help,” “opinion,” “think,” “recommend,” “what should I do” and “where do I go.” Examples: “Do you think it’s the flu?”, “Do you have any recommendations for places that are closer?”
specific	Question seeking precise and targeted information (specific details or facts), characterized by the question’s directness and clarity and the use of the word “specific.” Examples: “What cough syrups with expectorants are on the market right now—like what are the specific brand names?”, “Is National Museum open on Saturday?”
generic	Question seeking general information, leading to a response containing a variety of suggestions rather than a pinpointed answer. Examples: “Yes, what are some good places to go after dinner?”, “What are some unusual experiences I could do in Edinburgh, near Edinburgh Castle?”
Speech act: Answer	
Attribute	Definition
factual	Answer framed to explicitly appear as having derived from VA knowledge, containing phrases such as “It is recommended,” “It is possible,” or “There are several places for you to explore.”
opinion	Answer framed to explicitly appear as being the opinion of the VA, containing cues denoting the subjectivity of the response such as “I think,” “In my opinion,” or “I suggest.”
refusal	VA either refuses to provide an explicit answer or omits the requested information from its response.
directive	Answer containing clear directions, instructions, or information for the user, offering guidance on how to achieve a specific goal or answering a specific question. Examples: “To get from [Point A] to [Point B], you can walk ...”, “The Witchery by the Castle in Edinburgh typically opens for lunch at 12:00 PM ...”
descriptive	Answer containing a detailed portrayal of a scene, object, or concept, emphasizing sensory perceptions to create a vivid mental image for the user beyond statements of information. Example: “One option is to visit Princes Street Gardens, where you can relax and enjoy the beautiful scenery. Another suggestion is to explore the Grassmarket area, known for its charming cafes and shops ...”

Table 11. Overview of speech style attributes and their definitions for *statement* speech acts and *egocentric* and *exocentric* speech style attributes and their definitions, applicable to all speech acts.

Speech act: Statement	
Attribute	Definition
warning	Statement presented by the VA with the purpose of reminding participants of the limitations of the AI and the importance of seeking expert or real-time advice, e.g., “I am not a medical professional ...”, “Consult a doctor ...”, or “Check the opening times.”
opinion	Statement presented in a style that explicitly appears to be an opinion. This is often indicated by cues such as “I think,” “In my opinion,” “I suggest,” or other similar phrases that denote subjectivity. Example: “I think universities should have their own police forces.”
non-opinion	Statement that is not an opinion as evidenced from implicit cues. Example: “I’d like to talk about whether universities should have their own police forces.”
argument	A statement or series of statements presented to justify, validate, or support a viewpoint or stance in the debate scenario.
counterargument	A statement or series of statements introduced to oppose, challenge, or refute the opposing viewpoint or stance in the debate scenario.
agreement	A statement or series of statements that indicate alignment or consensus with a previously expressed opinion or argument of the other party in the debate scenario. Example: “Yeah. I think that it’s worth it to try. Yeah.”
Speech act: All (question, answer, and statement)	
egocentric	A mode of communication that suggests that the participant primarily speaks from their own personal viewpoint (subjective). This is determined by the participant’s perspective only; i.e., an egocentric VA response means that the VA is conveying the information in the second-person (you-) perspective.
exocentric	A mode of communication that adopts a perspective that is not self-centered (objective). This is determined by the participant’s perspective only; i.e., an exocentric VA response implies it is using an impersonal perspective.

Table 12. Overview of speech acts based on our implementation of a ChatGPT-powered VA.

User commands	
State	Definition
initiation	Signals the user’s intent to start a dialogue or conversation. Examples: “Alexa, let’s chat,” coughing to start the interaction for the medical self-diagnosis scenario.
end-intent	Statement that indicates the user’s intent to wrap up the conversation. Examples: “That’s all,” “Bye,” “Alexa, stop.”
VA responses to user commands	
State	Definition
introduction	VA’s opening monologue, tailored to each scenario. Examples: “Oh, seems like you are not feeling well. Maybe I can help figure out what’s wrong?” (medical), “Hi! I am an AI assistant designed to present requested information. How can I assist you today?” (day planning), “Hey! I am a voice assistant designed to engage in a discussion with you. What would you like to talk about?” (debate)
closing	VA’s farewell before terminating the conversation. Examples: “Goodbye,” “Bye,” “Take care.”
filler	VA’s response to the user while waiting for ChatGPT API response after 2 seconds of user query. Examples: “I’m looking it up” (for the medical and planning scenarios), “Thinking it through” (tailored to the debate scenario).
VA questions	
State	Definition
small talk	While waiting for ChatGPT API response after 6 seconds of user query, the VA poses a task-irrelevant question. Example: “While I get that, do you like going outside?”
continuing	In the absence of the detection of a user query by the ChatGPT Alexa skill, the VA asks a continuing question. Examples: “Should I continue?”, “Anything else I can help with?”

Table 13. Error types, breakdowns and recovery.

Error types	
Type	Definition and breakdown
skill	A skill error is attributed to issues related to the integration of ChatGPT into the Alexa skill, such as an API response error that causes the Alexa skill to terminate. Skill error is always manifested when Alexa announces, “There was a problem with requested skill’s response,” leading to <i>skill closure</i> .
listening	A listening error arises when the user is speaking to Alexa during a period of skill inactivity (i.e., Alexa is not listening).
handling	Handling errors occur when Alexa listens to and transcribes a user’s speech, but the transcribed input is not handled appropriately, resulting in <i>no VA response</i> . For instance, the transcribed speech is not considered an intent for the ChatGPT-powered Alexa skill.
partial listening	Such errors occur when Alexa captures only part of the user’s speech, often due to hesitant speech patterns, prolonged pauses, or Alexa cutting off the user prematurely. Such disruptions can lead to <i>user intent recognition failure</i> .
interruption	Interruptions by Alexa disrupt the conversation, resulting in partial listening errors that usually lead to <i>user intent recognition failure</i> . We categorize interruptions separately because such errors manifest differently; interruptions directly impact user behavior (i.e., the user stops talking in the middle of their query as a result of the interruption).
transcription	Transcription errors occur when Alexa hears the user but does not transcribe their speech correctly. Transcription errors lead to <i>user intent recognition failure</i> . An error is counted as transcription error when the user’s query is listened to and transcribed inaccurately, but still handled.
Error recovery strategies	
Strategy	Definition
repeat	Recovery strategy in which the user repeats their query; they may add details or rephrase the wording of their initial query to get the desired VA response.
move on	Recovery strategy in which the user chooses to ignore the unanswered query and moves on to the next query to continue the conversation.
ask-clarify	Recovery strategy wherein the VA has doubts about the user’s query and asks for clarification or further details to identify user intent.