

# GPTVoiceTasker: LLM-Powered Virtual Assistant for Smartphone

MINH DUC VU, HAN WANG, and ZHUANG LI, Monash University, Australia

JIESHAN CHEN, CSIRO's Data61, Australia

SHENGDONG ZHAO, City University of Hong Kong, China

ZHENCHANG XING, CSIRO's Data61 & Australian National University, Australia

CHUNYANG CHEN, Monash University, Australia

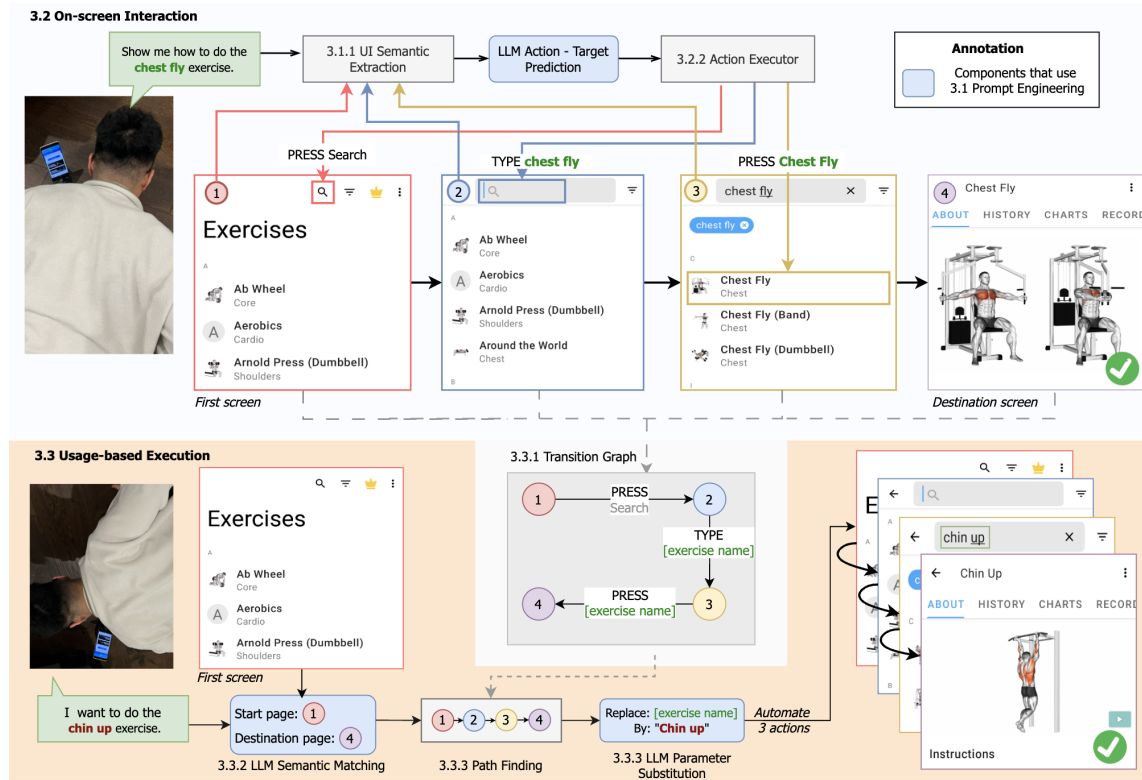


Fig. 1. An example use case in Home Workout application when the user needs to interact with the smartphone hands-free due to physical busyness. When performing an on-screen interaction, GptVoiceTasker repeatedly prompts LLM to predict on-screen actions with current UI information and execute the response to achieve user tasks. The execution information is then saved to streamline the execution of subsequent similar tasks.

Authors' addresses: Minh Duc Vu, dustin.vu@monash.edu; Han Wang, han.wang@monash.edu; Zhuang Li, zhuang.li@monash.edu, Monash University, Melbourne, Australia; Jieshan Chen, jieshan.chen@data61.csiro.au, CSIRO's Data61, Sydney, Australia; Shengdong Zhao, City University of Hong Kong, Hong Kong, China, shezhao@cityu.edu.hk; Zhenchang Xing, CSIRO's Data61 & Australian National University, Canberra, Australia, zhenchang.xing@data61.csiro.au; Chunyang Chen, Monash University, Melbourne, Australia, chunyang.chen@monash.edu.

2024. Manuscript submitted to ACM

Manuscript submitted to ACM

1

Virtual assistants have the potential to play an important role in helping users achieve different tasks. However, these systems face challenges in their real-world usability, characterized by inefficiency and struggles in grasping user intentions. Leveraging recent advances in Large Language Models (LLMs), we introduce GptVoiceTasker, a virtual assistant poised to enhance user experiences and task efficiency on mobile devices. GptVoiceTasker excels at intelligently deciphering user commands and executing relevant device interactions to streamline task completion. The system continually learns from historical user commands to automate subsequent usages, further enhancing execution efficiency. Our experiments affirm GptVoiceTasker's exceptional command interpretation abilities and the precision of its task automation module. In our user study, GptVoiceTasker boosted task efficiency in real-world scenarios by 34.85%, accompanied by positive participant feedback. We made GptVoiceTasker open-source, inviting further research into LLMs utilization for diverse tasks through prompt engineering and leveraging user usage data to improve efficiency.

CCS Concepts: • **Human-centered computing** → **Interaction techniques; Smartphones; Natural language interfaces; Sound-based input / output.**

**ACM Reference Format:**

Minh Duc Vu, Han Wang, Zhuang Li, Jieshan Chen, Shengdong Zhao, Zhenchang Xing, and Chunyang Chen. 2024. GPTVoiceTasker: LLM-Powered Virtual Assistant for Smartphone. 1, 1 (January 2024), 22 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

The advancements in voice control technology have sparked a new wave of innovation, driving the exploration of its potential in transforming smartphone interactions [28, 29]. With the integration of voice control, users can effortlessly navigate through various applications, compose messages, and even initiate tasks like checking the weather or playing a video on YouTube [60]. This seamless and natural mode of interaction not only saves time but also promotes a hands-free experience, allowing individuals to engage with their smartphones in situations where manual input operation is impractical or inconvenient [39, 71] (see Fig. 1). Moreover, the success stories of widely recognized voice assistants like Google Voice Assistant [5] and Siri [3] have further propelled the adoption of this technology, inspiring researchers and developers to delve deeper into its capabilities and refine its usability for an even broader range of users. As this field continues to evolve, the integration of voice control technology with smartphones holds immense potential to revolutionize the way we interact with our devices, empowering users with enhanced accessibility and a more convenient user experience.

Developing efficient and reliable voice-controlled systems involves addressing various challenges that significantly impact the accuracy and usability of these systems [47]. One major hurdle is accurately comprehending user commands and seamlessly mapping them to specific smartphone actions [51]. Achieving satisfactory accuracy levels often relies on substantial training data, comprised of user voice commands, and complex deep learning models [36]. These models are employed to decipher the semantic nuances of the commands, comprehend their intended meaning, and subsequently select the appropriate user interface (UI) element on the screen for interaction. However, acquiring and curating substantial datasets for training language models presents a laborious and time-consuming challenge, impeding progress and scalability in developing virtual assistants. Moreover, the collected dataset may not comprehensively cover the diverse commands required for modern applications with rich functionalities, potentially limiting the language model's capabilities. Additionally, the rapid advancement of technology constantly introduces new functionalities, quickly rendering the existing dataset outdated, and these models often struggle with accommodating the interpreting speech recognizer errors due to accents [59] and homophones [48]. The inherent complexity of these models hampers their adaptability to the dynamic landscape of human natural language and pronunciation variations, making it challenging to account for real-world usage scenarios' intricacies and variations.

Furthermore, voice assistants often face performance issues when compared to physical touches on mobile devices [2]. Extensive research and development have explored various methods to enhance voice assistants' capabilities. For instance, Voicify [60] employs app analysis and deep linking techniques to speed up voice command processing. Modern voice assistants also enhance user experience by personalizing interactions based on individual usage patterns and preferences [12, 70]. This personalized learning approach holds tremendous potential in mitigating the performance limitations of voice-controlled systems, enhancing efficiency, and accelerating user interactions [35, 51]. By continuously learning from user interactions, voice assistants can improve their accuracy and reliability, ultimately ensuring efficiency and overall user experience in real-world usage scenarios.

Large Language Models (LLMs) have revolutionized natural language processing (NLP), showcasing their general intelligence and excellence in reading comprehension, trivia quizzes, translation, and text completion [9]. The concept of Few-Shot Learning further enhances LLMs' power by enabling them to adapt to new tasks with minimal examples or prompts, eliminating the need for task-specific models and datasets [53]. This flexibility and efficiency in handling diverse conversational interactions without extensive retraining present an innovative and promising approach. Leveraging few-shot learning, developers can tap into the vast knowledge and language understanding capabilities of LLMs, facilitating rapid development and deployment of language-based automation solutions in various domains [23, 44]. Therefore, we aim to investigate the potential of LLMs in enhancing the intuitiveness of voice control for mobile applications.

This paper introduces GptVoiceTasker, a novel system that leverages LLMs to enhance voice control for Android devices. Our system employs advanced prompt engineering techniques to ensure a precise understanding of user commands without extensive model training. By bridging the gap between natural language commands and interactive tasks on mobile devices, GptVoiceTasker facilitates seamless automation of daily physical interactions, including scrolling, tapping, and inputting text, solely through voice commands. In addition, the system automatically records and learns from user commands and in-app usages, enabling the reproduction of tasks for similar requests in the future to improve task efficiency. GptVoiceTasker supports the sequential execution of multiple commands within a single utterance, further enhancing the efficiency of voice interaction on mobile devices.

We validated the technical contributions of GptVoiceTasker by evaluating i) the ability to parse user commands into executable actions and ii) the ability to streamline saved tasks. The command parser achieved over 90% accuracy on the human command dataset collected from a user study. Our automated execution achieved 82.7% success rate for direct match tasks and 72.0% success rate for tasks with different parameters. To validate the usability of GptVoiceTasker, we conducted a user evaluation with 18 participants, each completing a set of tasks using GptVoiceTasker and two state-of-the-art baselines. We collected the time taken to complete each task, as well as quantitative and quality feedback from users. The results show that GptVoiceTasker accelerates the tasks by 34.85% and received positive feedback regarding usability.

To summarize, the contributions of this paper include:

- Development of GptVoiceTasker, a voice assistant that harnesses the capabilities of LLMs to streamline the automation of multi-step tasks by predicting the most optimal step on each individual screen.
- A graph-based local database design that automates the recording and retrieval of personal app usages, enhancing task execution efficiency for virtual assistant interactions.
- Conducting a large-scale user evaluation to validate the effectiveness of our approach, along with empirical findings on system limitations and considerations for voice assistant design.

- GptVoiceTasker<sup>1</sup> is open-sourced so that anyone can use and continue to improve the system.

## 2 BACKGROUND & RELATED WORKS

### 2.1 Voice Control & Assistants on Mobile Devices

In recent years, the progress in Natural Language Understanding (NLU) has fueled the proliferation of voice assistants across diverse platforms, including ubiquitous systems [7, 32] and home appliances [52]. On personal smartphones, numerous voice assistants have been proposed, leveraging intelligent NLU models to comprehend user speech input and map it to user interface actions. An early milestone in smartphone voice control interfaces was JustSpeak, which harnessed Google’s Automatic Speech Recognition (ASR) to record user commands and introduced innovative utterance parsing techniques [71]. Subsequently, the Smart Voice Assistant expanded on JustSpeak’s capabilities by enabling users to manage calls and SMS through voice commands [11]. A similar endeavor, the "VoiceNavigator" application by Weber et al. in 2016, focused on enhancing the visibility and learnability of mobile voice user interface applications [18]. However, these initial approaches, foundational as they were, encountered usability issues stemming from rigid language parsing heuristics and limited use cases, which spurred the need for further development of smartphone virtual assistants.

In recent years, significant advancements in language parsing capabilities have been achieved through deep learning models. SAVANT leveraged Dialogflow as a conversational agent to extract user intent from utterances [4], while DoThisHere employed the pre-built Almond language model to enable voice control for retrieving and setting UI contents in Android [69]. Google released Voice Access [1], aimed to replace manual interactions with voice command, which has over 100 millions downloads on Google Play Store. More recently, Voicify [60] introduced VoicifyParser, an advanced deep learning approach for parsing user commands into on-screen interactions, while AutoVCI [51] concentrated on developing voice interfaces for automating mobile UI tasks. However, the interaction paradigm with these existing approaches remains somewhat unnatural, requiring users to issue precise machine-like instructions, such as “*Press save button*”. This limitation means that they may struggle to fully comprehend high-level user intentions, such as “*I want to save this note*”. We propose GptVoiceTasker to address these challenges and revolutionise the voice-based interactions between human and software systems. Our solution leverages the capabilities of LLMs to map high-level user intentions to executable actions, enabling on-screen interactions through intention-based voice commands. This approach seeks to accommodate the flexibility and natural language of human commands, ushering in a new era of user-friendly assistive tools.

### 2.2 Large Language Models for Enhanced Human-AI Collaboration

The advent of generative AI has given rise to innovative LLMs, such as GPT-4 [50] and DALL-E [55]. These LLMs have revolutionized the landscape of AI development by enabling developers to achieve complex tasks through few-shot prompting, eliminating the need for extensive custom model training. Their remarkable versatility has spurred active research in both IT and non-IT domains, spanning areas like software testing [24, 43], high-performance computing [15], finance [65], and health science [25]. LLMs have particularly excelled in enhancing the intuitiveness of existing methods, as seen in software testing, where they generate authentic text inputs based on the current UI page information, replacing the conventional random text input approach [43]. This demonstrates the transformative potential of LLMs in advancing research and innovation across a multitude of domains.

---

<sup>1</sup><https://github.com/vuminhduc796/GPTVoiceTasker>

The capabilities of LLMs have sparked a surge in their application within assistive technology, revolutionizing the translation of user commands into executable tasks across diverse systems. Recent research in this domain has witnessed the transformation of human natural language commands into various types of tasks, including visualization tasks [62], operating system tasks [42], and robotic tasks [38, 57]. LLMs have enabled these systems to tackle more intricate commands beyond the scope of existing heuristic approaches. They also exhibit a remarkable ability to comprehend variations of commands that share similar intentions but are expressed differently. This pioneering framework, with the support of LLMs, paves the way for a novel (semi)automated task execution paradigm, erasing the boundaries between traditional command patterns and intuitive command modalities.

Within the domain of mobile assistants, LLMs have emerged as a transformative force, supplanting conventional machine learning models as seen in prior related works [51, 60]. This paradigm shift simplifies the process of translating user’s natural voice commands into actionable operations on mobile User Interfaces (UIs). Wang et al. [61] utilised LLMs to allow conversation-alike interaction with mobile UIs, which demonstrates LLM’s better ability in understanding on-screen elements compared with traditional machine learning approaches [37]. While Wang et al.’s approach in smartphone virtual assistant development is noteworthy, its employed prompting strategy remains basic, impacting its ability to parse natural user commands. Our work aims to elevate this approach by incorporating state-of-the-art prompting strategies to enhance accuracy and create a more comprehensive smartphone virtual assistant. Unlike Wang et al.’s focus on single on-screen actions, our research takes a step further. We strive to unlock the full potential of LLMs in automating complex user requests involving multiple steps to achieve specific objectives, as exemplified in Fig. 1 with the query *“Show me how to do the chest fly exercise”*. Our goal is to broaden the scope of LLMs, assisting users not only in known tasks but also in accomplishing unfamiliar ones, all while significantly reducing the time required to complete these tasks on a smartphone.

### 3 THE GPTVOICETASKER SYSTEM

We introduce GptVoiceTasker, a virtual assistant that empowers users to efficiently perform various tasks on their smartphones using voice commands. By applying different prompt engineering techniques (Section 3.1), GptVoiceTasker harnesses the power of LLMs to perform different logical tasks. Upon receiving a user command, GptVoiceTasker first attempts to automatically execute the task using the saved database (Section 3.3). Should a task prove unfeasible using the pre-existing database, GptVoiceTasker will perform a series of step-by-step predictions of on-screen navigation to complete a task (Section 3.2). Simultaneously, the system records these interactions for subsequent automated execution.

#### 3.1 Prompt Engineering

While a naive approach to prompting LLMs for various tasks may yield suboptimal results due to low accuracy and randomness in responses [17], we propose the adoption of different prompt engineering techniques. These techniques involve crafting prompts according to specific rules and components to elicit optimal responses from LLMs [41]. We created multiple prompt templates<sup>2</sup>, which is applied in both On-Screen Interaction (Section 3.2) and Usage-based Execution (Section 3.3). In Fig. 2, we illustrate an example of our prompt designed to determine the most appropriate target for tapping. This section explains different prompt engineering techniques applied in GptVoiceTasker, showcasing their significance in harnessing LLMs for complex reasoning tasks across multiple components in our system.

<sup>2</sup><https://github.com/vuminhduc796/GPTVoiceTasker/blob/main/prompts.txt>

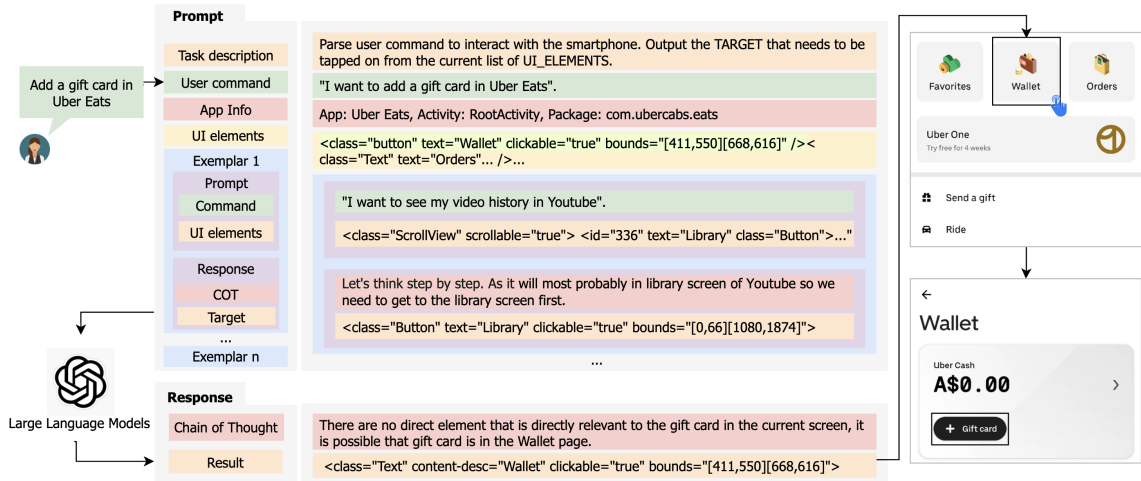


Fig. 2. An example of our prompt and response format to determine the most relevant target to press.

**3.1.1 Least-to-most Prompting.** In addressing complex tasks, such as predicting the action to perform on the screen from the user command, GptVoiceTasker employs the Least-to-most Prompting strategy [72]. This technique involves breaking down a complex task into smaller, manageable prompts. The goal is to enhance the model's understanding of the logical flow, thereby improving accuracy. For example, GptVoiceTasker implements a two-step prompting approach for On-Screen Interactions (Section 3.2). First, the model is prompted to map user intent to a specific action (e.g., tapping an element, entering text, scrolling). Subsequently, based on the determined action, subsequent prompts are sent to identify the target UI element for executing that action (Fig. 2). This approach allows the LLMs to concentrate on atomic tasks, contributing to improved accuracy compared to a naive approach of identifying both action and target in the same prompt [61].

**3.1.2 Few-shot Prompting.** Few-shot prompts [13], incorporating completed task examples into the prompt (as in Exemplar 1 in Fig. 2) itself, empower LLMs like GPT to facilitate rapid comprehension of specific rules and guidelines, thereby markedly improving output accuracy and reliability. For instance, when discerning the action in a user command, we include exemplars with synonyms and homophones to enhance adaptability and address speech recognition variations. This integration of few-shot learning, coupled with the LLM's adept vocabulary and grammar comprehension, boosts versatility across diverse scenarios. Importantly, it effectively rectifies terminology discrepancies, a common issue in voice assistants using automatic speech recognition [10]. This advancement represents a notable improvement over heuristic-based methods, especially in navigating linguistic intricacies.

**3.1.3 Chain of Thought.** When dealing with complex prompts that demand logical reasoning, LLMs often struggle to fully interpret and deliver accurate output. To address this, we integrate "Chain of Thought" [64] into our few-shot exemplars to help LLMs simulate human-like reasoning and provide logical output. By employing this mechanism, we ensure consistent and coherent reasoning, enabling a seamless progression of ideas and actions in generated responses. In GptVoiceTasker, the chain of thought is added to each response in examples where few-shot prompting is applied. In few-shot exemplar 1, as depicted in Fig. 2, we provided a chain of thought to predict the most relevant UI element

to find the “*video history in YouTube*”. Despite the task requiring multiple steps to navigate to the history page, we asked the LLM to identify key indicators and execute logical reasoning specific to the task. This enabled it to accurately predict the most likely button for interaction. As a result, the system can repeat the same command, incorporate it into the prompts, proceed with the predicted actions, and continue until it reaches the intended destination.

### 3.2 On-screen Interaction

Upon receiving a task from users, GptVoiceTasker will automatically perform multiple on-screen action executions until the task is accomplished. We first collect run-time UI elements, execute reasoning tasks from user commands, and perform actions on the user’s device. GptVoiceTasker continues to repeat this iterative process to perform each action to accomplish the user tasks.

**3.2.1 UI Semantic Extraction.** To ensure a comprehensive understanding of the current UI elements and their semantics, GptVoiceTasker utilises Android Accessibility Service to continuously capture UI and analyzes new changes in the mobile screen, following previous research [60]. Note that the raw UI elements we extracted may be afflicted by UI noise, which is a prevalent issue linked to the real-time gathering of UI elements [34]. This problem arises when the collected UI information does not align with its visual representation, which affects the semantic understanding of LLMs on current UI elements, resulting in incorrect interactions. To address this, we implement heuristics to mitigate potential inaccuracies and ensure the reliability of collected UI information. First, we utilise the collected coordination of each UI elements to eliminate out-of-bound or empty elements. We also eliminate views are fully overlapped by other views, which does are invisible and not interactible. In addition, we remove those views that do not contain any interpretable information, such as empty view containers.

Our primary emphasis is on representing mobile UI screens in a textual format that can be interpreted by Large Language Models (LLMs) through text-based input. While recent research proposed the translation of XML representation of UI screens to HTML format [24, 61] to overcome prompt length limitations, this approach becomes less relevant as modern models relax restrictions on prompt length, allowing developers to present more information in a prompt. For textual components like text fields and buttons, we gather essential semantic information from their labels and contents. We collect alternative text and resource names for icons and images, which offers meaningful information as defined by developers. Additionally, we augment the collected information with attributes such as the precise element location on the screen and a more comprehensive view hierarchical structure, previously absent in HTML representations. These details play a crucial role in illustrating relationships between UI elements, as demonstrated by the example of delivery times for each restaurant in the UberEats app (Fig. 3). This capability proves valuable in accommodating user commands, such as “*Select the restaurant with the fastest delivery time*”, a task unattainable with a flattened UI representation approach. Furthermore, we can cater for commands that user references to UI elements by their locations, such as “*Press the icon at the top-right corner*”. Lastly, each UI element is enriched with supported interaction types, indicating potential actions like PRESS or ENTER\_TEXT to support LLM propose the appropriate action.

To enhance the contextual understanding and reasoning capabilities of the LLM, we not only fetch on-screen data but also retrieve relevant system-related information, including the app name and activity name of the current screen. By incorporating these system-related details, GptVoiceTasker gains a comprehensive understanding of the user’s current context, which further aids in validating each step in multi-step executions.

**3.2.2 Action Executor.** We combine the user command with UI and contextual information collected in previous step and applies different prompt engineering techniques as described in Section 3.1 to create prompts. These prompts are

then sent to LLMs, which will first detect most appropriate actions and then detects the associated target UI element for this action. GptVoiceTasker uses this information to perform the interactions on user’s device, such as tapping, scrolling, or entering text on certain UI element. Prior mobile automation approaches [51] encounter challenges in handling runtime UI changes and app updates that might alter UI representations and operation sequences. In contrast, our approach exhibits robustness in the face of such issues, as the sequence of steps is dynamically dependent on the UI on the screen. This adaptability allows GptVoiceTasker to navigate through dynamic UI changes and updates, ensuring reliability in task automation on smartphones. Additionally, GptVoiceTasker provides audio feedback to users, confirming that the system is automatically proceeding to the next command. This real-time feedback ensures a smooth and intuitive user experience with GptVoiceTasker’s interaction capabilities.

To address challenges inherent in automation within real-time systems, including the detection of failures in each automated step [66], our implementation integrates a screen transition detector. This detector utilizes the Hamming distance [31] to gauge the difference between two screens, ensuring that the UI accurately reflects changes resulting from the automated action. Moreover, for validating the success of actions, we implement additional heuristics. For example, in scenarios involving ENTER\_TEXT, we verify the presence of the entered text in the target text box. If an action proves inexecutable (i.e., not inducing appropriate changes to the UI), we iteratively repeat the step with supplementary information about the failed action, thereby excluding these actions from the selection.

Another challenge in mobile app automation arises from the prevalent dependency of mobile UI screens on live internet content, often loaded asynchronously after the screen appears on smartphones. The collection and utilization of incomplete or loading UIs can compromise the accuracy of detecting suitable actions. In response, GptVoiceTasker introduces an innovative approach by delaying subsequent actions until the screen is fully loaded. This is achieved by detecting screen-loading widgets in Android, leveraging information such as widget names, types, and shapes to detect progress bars and loading indicators. Furthermore, we enhance our approach by collecting data on network transmissions to identify ongoing content downloading tasks, inspired by the methodology presented in [40], which leverages network analysis for detecting ads. These methods significantly enhance the reliability of GptVoiceTasker in effectively navigating through the app’s interface.

### 3.3 Usage-based Execution

In mobile apps, the UI elements on a specific app page are predefined by developers when developing an app page. Therefore, the series of user interactions on the screen to achieve a task will be consistent across different times. Based on this matter, GptVoiceTasker automatically create a saved path for each user command, allowing GptVoiceTasker to replicate interactions when receiving a similar command from users. Unlike previous approaches [35, 51] that depend on a manual task creation process for automation support, GptVoiceTasker automatically records app transitions through on-screen navigation in Section 3.2. This automated process not only allow wider coverage of automated tasks but also eliminates the need for manual efforts in predefining shortcut tasks. In this section, we outline our method (as in Fig. 3) to streamline subsequent similar tasks from users, which combines both LLMs-based and heuristic-based modules. We introduce our database design in Section 3.3.1. First, we identify the current UI screen displays on user device and the destination screen (Section 3.3.2). After that, we find the most viable path from current screen to the destination screen (Section 3.3.3). Finally, we use incorporate human interactive feedback to validate and finetune the execution for further usages (Section 3.3.4).



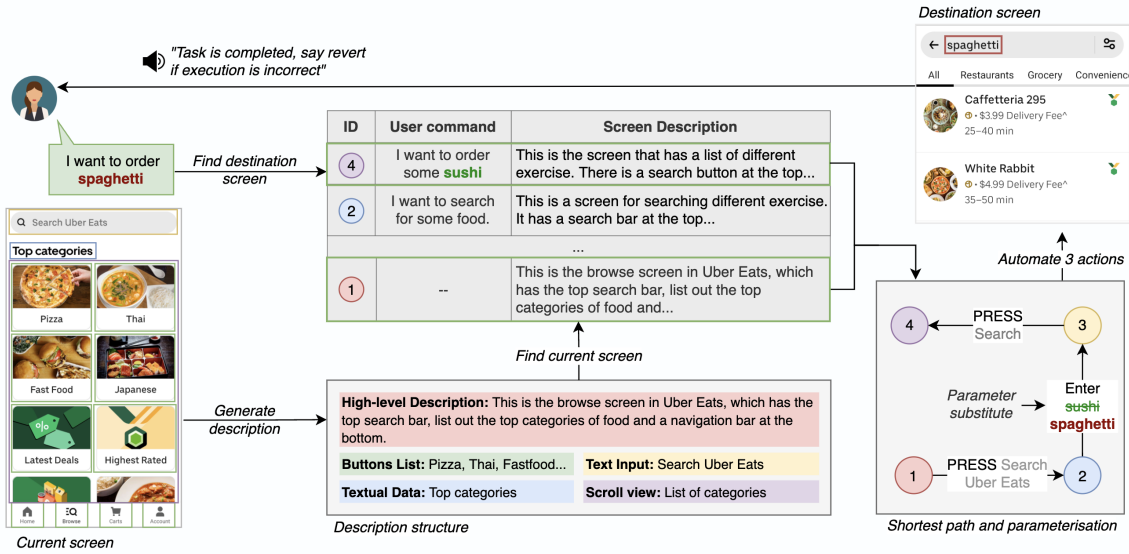


Fig. 3. An example use case in Uber Eats to how GptVoiceTasker use the historical tasks to execute user new command.

**3.3.1 Transition Graph.** In line with previous approaches to automating tasks in mobile apps [16], GptVoiceTasker incorporates a directed graph for each app to facilitate seamless transitions between different app pages. Each node in the database corresponds to a UI page in the app, containing a unique ID and a screen description generated as outlined in Section 3.3.2. Moreover, for each node, we maintain a list of previously used commands by users that concluded on this page. This compilation signifies that the requested functionality from these user commands is associated with this specific app page. The directed edges between nodes represent potential transitions from one UI page to another, capturing information about the associated action type and the target UI element for the transition. These edge details empower GptVoiceTasker to replicate user actions and screen traversals. As users engage with the app, this graph dynamically expands, automatically incorporating new UI pages and associated transitions.

**3.3.2 Screen Description & Command Pattern Matching.** While previous automation approaches focus on executing tasks starting from the launcher page of the app [4, 51], they may not be suitable for cases where users are currently on other pages, necessitating additional navigation steps to return to the launcher page. To address this challenge, we first identify the current screen of the app that the user is viewing in our graph database. However, heuristically comparing and identifying a UI page using the on-screen textual information collected from the XML representation proves impractical due to the contextual dynamicity of information within UI elements on a specific app page. For example, consider the search result page in the Uber Eats app, where distinct restaurants are displayed for "spaghetti" and "sushi" search results. Despite featuring different content, these pages share the same layout and, therefore, should be identified as the same node in our graph database. To achieve this, we leverage the capabilities of LLMs to semantically summarize the UI content [61] into a semi-structured description. We collect content from the XML hierarchical representation of the screen, along with additional context about the current activity and app name, to generate a high-level description of the overall functionality that the screen serves. Subsequently, GptVoiceTasker summarizes the list of interactive elements, including clickable, scrollable, and text-editable elements, and appends them to the description. To locate the

current UI page within the list of visited screens in an app, we prompt LLMs to perform semantic matching between the current screen description and the saved list of screen descriptions from the app. This process allows us to trace back the ID of the screen in the graph database, effectively addressing the navigation challenges from any point within the app.

After identifying the current screen, we perform semantic matching to find destination screen. We define destination screen as the most relevant screen that can serve user request, for example, the destination screen for user command “*I want to order spaghetti*” is the search result page for the keyword “*spaghetti*”. We utilise the user saved commands and screen description for each app page to prompt LLMs to determine the most relevant app page from the app. We build a prompt that includes user request and the list of saved screen in the app to rank the relevancy of each app screen with the request that user made.

**3.3.3 Path Finding & Execution.** After identifying the current and destination pages within the saved graph database, GptVoiceTasker utilizes the Shortest Path algorithm [26] to determine the sequence of actions required to navigate from the start to the destination node. GptVoiceTasker extracts this sequence from the edges connecting the start node to the destination node and executes each action in sequence using the Action Executor described in Section 3.2.2.

Given the dynamic nature of smartphone GUIs, where the relative coordinates of buttons may vary, especially in scrollable screens or due to unexpected pop-ups and ads, execution validation is crucial. To address this, GptVoiceTasker performs validation on each UI page in the path, ensuring the automation follows the expected sequence. After each action in the sequence, the screen description is generated based on current displayed UI (as explained in Section 3.3.2) and compared to the expected description stored in the graph node representing the anticipated UI. If the descriptions do not match, the tool navigates back to the previous app page and prompts LLMs to propose an action. The prompt created for on-screen action predictions is then used to feed the action to the LLM, regenerating the appropriate action. This iterative process ensures that the automation adapts to changes in the smartphone GUI, allowing GptVoiceTasker to dynamically execute actions based on real-time screen information.

Additionally, GptVoiceTasker leverages command parameterization techniques [33, 51], allowing saved paths to be reused for similar tasks with different parameters. Utilising the enriched vocabulary and robust natural language understanding capabilities of the LLM, GptVoiceTasker prompts the LLM to function as an advanced Named Entity Recognition (NER) system. This system identifies and replaces all occurrences of substitutable words with new keywords in the action sequence. For example, in the Uber Eats app scenario shown in Fig. 3, where a user has a saved command for “*I want to order some sushi*”, GptVoiceTasker recognizes “*sushi*” as a parameter value representing the type of food, which is replaced by “*spaghetti*”, transforming the command from ENTER “*sushi*” to ENTER “*spaghetti*”. Subsequently, GptVoiceTasker executes this adapted sequence of actions to accomplish the task.

**3.3.4 Human Feedback Loop.** The saved commands have proven their applicability in subsequent user interactions with the device, yet certain inconsistencies may arise due to app version updates and dynamic changes in UI element availability. To address error handling and corrections, GptVoiceTasker incorporates human feedback to enhance the saved data for each task. The automated tasks performed by the tool allow users to efficiently reverse actions and prompt the system to actively adjust the new execution path to accomplish tasks effectively. This valuable information is stored and fed back to the LLM, resulting in improved output for future executions. When users express satisfaction with the result and request the next task, we also increase the confidence score for this saved path, helping GptVoiceTasker to select the most reliable path to execute if multiple execution paths are found. In addition, we also provide an additional

interface for users to amend saved commands manually. Through user feedback and iterative learning, GptVoiceTasker continuously refines its performance, ensuring increasingly accurate and effective responses over time.

### 3.4 Implementation

We implemented GptVoiceTasker as an Android application with the Accessibility Service in Android OS using Java [21]. Within the Java code, GptVoiceTasker subscribes to the *typeWindowContentChanged* accessibility events [21] to receive notifications whenever UI changes occur on the screen. We created a dynamic pipeline to extract UI elements in a hierarchical structure from AccessibilityNodeInfo objects [20], which serve as data representations of on-screen UI elements provided by Android Accessibility Service. In addition, We obtain the app name and the activity name using the Android PackageManager class.

To communicate with the LLM, we utilised the API service provided by OpenAI. We chose the GPT-4 model, as it is the latest model trained by OpenAI [49]. Once we receive a response from the LLM, the tool leverages the *performAction()* method [20] to execute actions on the corresponding elements. All data related to personalized services, such as screen description and the transition graph, are stored in the phone memory for future usage.

## 4 TECHNICAL EVALUATION

To evaluate the effectiveness and reliability of the proposed system, we conducted two experiments on our command interpreting module and usage-based execution. Specifically, we first assess the system’s ability to comprehend user commands and perform on-screen interactions, comparing its performance to other state-of-the-art approaches. We then investigate the system’s capability to execute multi-step tasks based on the saved user usages. .

### 4.1 On-screen Interaction Evaluation

**4.1.1 Experiment Setup & Metric. Datasets:** We collect a specialised test set to evaluate our system’s capabilities in understanding natural language commands and mapping them to appropriate actions and target UI elements. This dataset comprises 278 natural language user commands to interact with Android UIs.

Although prior research [14, 60] has produced a similar test set, it is not directly adaptable to our context for two critical reasons. First, some instances in the test set are artificially synthesized based on predetermined heuristic rules. The resulting natural language commands are linguistically biased toward simpler linguistic patterns and do not align with the complex linguistic variants inherent in real-world human spoken utterances. Second, some test examples in the existing dataset have become obsolete or are no longer replicable due to updates in the corresponding applications.

To construct a test set that more closely aligns with real-world user interactions, we adopted a data-driven approach. We engaged 31 participants (17 females, 14 males), with 4 individuals having never utilized voice assistants before, 4 using them 3-4 times a week, 6 using them daily, and 17 using them less than 3-4 times a week. All participants are work professionals and the rest from university community who use smartphones daily. We provided these participants with screenshots alongside a specific task to accomplish. We then recorded the verbal commands they issued to their mobile device to complete the given task. After the collection, we annotated the commands to specify the intended action and target UI elements within the Android system; here, the term *action* refers to executable functions, while *target* denotes specific UI components or elements on the current screen. As a result, we collected 278 natural user commands for the dataset.

**Metrics:** Similar to Vu et al. [60], we adopt three evaluation metrics, namely *Exact Match Accuracy* (EM), *Target F1* and *Action F1*. EM calculates the percentage of instances in the test set where the predicted sequence exactly matches

its corresponding ground-truth sequence. The measures *Target F1* and *Action F1* quantify the average micro F1 score for the target (i.e., the UI components to be interacted with) and the action (i.e., the actions to be performed on the UI components), respectively. The F1 score for each instance is computed using the formula:

$$F1 = \frac{2 \times |\text{pred} \cap \text{gold}|}{|\text{pred}| + |\text{gold}|}$$

where  $|\text{pred}|$  represents the size of the set of predicted targets or actions, and  $|\text{gold}|$  denotes the size of the set of ground-truth targets or actions. The average micro F1 score is calculated across all instances for either targets or actions.

**Baselines:** We consider four baselines for converting natural language into semantic meaning representations, which comprise actions and targets. These baselines are **vanilla Seq2Seq** [6], **BERT-LSTM** [67], **Voicify Parser** [60], and **Wang et al.**'s work [61]. In the original work by Voicify, all three baselines employ deep learning models trained on datasets synthesized using the Overnight method [63]. This method generates training sets based on predefined lists of actions and targets that are designed for evaluation scenarios in Voicify. To ensure a fair comparison, we modified these lists to include the actions and targets present in our test dataset. We then re-synthesize the training set, which includes 1,384 instances, using the Overnight method, adhering to the implementation outlined in Voicify's work. Wang et al.'s work [61] was the first to incorporate LLMs for interacting with mobile interfaces. We incorporated the 2-shots LLM prompts they demonstrated in the paper for mapping instruction to UI actions. Additionally, we have enhanced their model by integrating the more advanced GPT-4, which also inline with the one we used in GptVoiceTasker.

**4.1.2 Evaluation Result.** Table 1 presents the results of our technical experiments. Overall, GptVoiceTasker significantly outperforms all baseline models across all metrics, achieving an **84.7%** EM accuracy, a **91.7%** Action F1 score, and a **84.7%** Target F1 score. Among the baselines without the LLMs, the Voicify Parser performs the best, aligning with the results reported in its original paper [60]. However, its performance suffers when faced with linguistic variations in our new test set. For instance, while the command “back” is correctly interpreted as “(PRESS, back)”, the phrase “return to last page”, which represents the same command, is incorrectly parsed as “(SWIPE, DOWN)”. Both BERT-LSTM and Seq2Seq models encounter similar issues, largely because they share architectural and training similarities with the Voicify Parser, yet perform even worse due to Voicify Parser being specifically optimized for task completion on Android systems. The method by Wang et al. [61] demonstrates the highest capability among the baselines, courtesy of the LLM's intervention, effectively rectifying the errors previously noted. However, the lack of the chain-of-thought and least-to-most prompt techniques occasionally leads to inaccuracies. This is evident in instances where the system misinterprets the intended direction in commands, such as confusing "DOWN" with "UP," or when it cannot adequately differentiate between actions like "PRESS," "ENTER," or "OPEN" when various verbs are employed in the commands.

Benefiting from the integration of LLM and the prompting techniques, our GptVoiceTasker excels at handling linguistic variants, consistently deriving the intended action and target regardless of variations in the input. The Action F1 score for GptVoiceTasker reaches 91.7, indicating its enhanced ability to predict actions across various linguistic patterns. Moreover, we observed that LLM effectively learns the true associations between actions and targets, thereby excelling at target prediction as well. For instance, PRESS is exclusively predicted with UI buttons, ENTER\_TEXT is linked solely with text input fields, and OPEN corresponds to app names. In contrast, the baselines often learn incorrect associations and outputs wrong target predictions. Our experimental results demonstrate that GptVoiceTasker possesses superior capabilities for understanding and accurately processing various linguistic variations, highlighting its adaptability in real-world scenarios.

Table 1. The experiment results of different baselines compared with GptVoiceTasker in three metrics.

Models	EM Accuracy (%)	Action F1 (%)	Target F1 (%)
<i>Seq2Seq</i>	25.2	47.6	35.6
<i>BERT-LSTM</i>	41.4	59.7	57.3
<i>VoicifyParser</i>	47.5	64.0	58.8
<i>Wang et al. [61]</i>	79.9	85.4	83.4
GptVoiceTasker	<b>84.7</b>	<b>91.7</b>	<b>84.7</b>

## 4.2 Database Execution Evaluation

Table 2. Saved task execution evaluation result for direct match tasks and parameterised tasks across 5 categories.

Category	Average Number of Automated Steps	Success Rate (%)	
		Direct Match	Parameterised
Message Friends	4.33	93.33	86.67
Listen to Music	5.27	80.00	73.33
Set an Alarm	5.73	73.33	53.33
Check Weather	5.07	80.00	73.33
Get Directions & Map	5.53	86.67	73.33
<b>Average</b>	<b>5.19</b>	<b>82.67</b>	<b>72.00</b>

**4.2.1 Experimental Setup & Metric.** In this experiment, we assessed GptVoiceTasker’s ability to automate tasks using the usage-based execution module. We initially identified the five common smartphone application categories, as shown in previous study [4]. Within each application category, we randomly selected five popular applications from the Google Play Store, with downloads ranging from 1 million to over 1 billion. For each selected app, we identified three features introduced by the developers in their Play Store descriptions. Each feature was then used to create both a direct match task, involving a straightforward match between user commands and corresponding app actions, and a parameterized task, requiring GptVoiceTasker to perform keyword substitutions to complete the task successfully, as shown in Section 3.3.3. For creating the direct match test cases, we paraphrased each saved command using state-of-the-art paraphrasing tool Quillbot<sup>3</sup>, as in [56]. In the case of parameterized tasks, we substituted one entity in the paraphrased command with another entity that has similar semantic. For example, consider the saved task “*Get directions to the nearest supermarket*”. In this case, the direct matching task would be “*Find the nearest supermarket’s location*”, while the parameterized task would involve substituting “*restaurant*” for “*supermarket*”, resulting in “*Find the nearest restaurant’s location*”. This process resulted in a total of five app categories, each category contains 15 direct match tasks and 15 parameterised tasks. These tasks involve 4 to 7 steps, with an average of 5.19 steps per task as illustrated in Table 2. All tasks can be automated with one voice command with the saved user app usage patterns. For a detailed list of the apps and features used in the experiment, please refer to our GitHub repository<sup>4</sup>.

To populate the transition graph and store screen descriptions, we manually navigated through each screen in every application using GptVoiceTasker. Subsequently, we configured the saved commands to reach the respective screens as the ground truth. We used the success rate as the primary metric, each test case is marked as success if GptVoiceTasker can successfully opened the desired feature using a single command.

<sup>3</sup><https://quillbot.com/>

<sup>4</sup><https://github.com/vuminhduc796/GPTVoiceTasker>

**4.2.2 Results.** Table 2 illustrates the accuracy of our saved task execution modules. Our findings indicate that GptVoiceTasker achieved an impressive level of automation, successfully handling **82.7%** of exact match tasks and **72.0%** of parameterized tasks. Notably, GptVoiceTasker exhibited exceptional performance in tasks related to messaging and directions & maps applications. This success can be attributed to the relatively static nature of these apps, where user interfaces maintain a consistent structure. Our results underscore GptVoiceTasker’s proficiency in command analysis, semantic matching to saved tasks, and parameterized phrase substitution within these contexts. However, the accuracy of GptVoiceTasker diminished when confronted with tasks related to setting alarms. To better understand the root causes of this decline in performance, we conducted an error analysis on the failed test cases. Several key issues emerged:

*Complex Parameterized Tasks:* For parameterized tasks with additional steps, such as setting an alarm for 7:30 instead of 7:00, GptVoiceTasker struggled due to the extra step involved in selecting the minutes, which was on a separate UI element. Further works include making GptVoiceTasker adaptable to these additional steps in the automation process.

*Pop-ups Ads and Unusual UI Elements:* Certain applications presented pop-ups ads and unusual UI elements in run time that were not encountered during the initial task-saving process. Consequently, GptVoiceTasker faced difficulties in completing these tasks. To improve the robustness of our approach, we recommend exploring the integration of a deep learning model to detect and handle such ad widgets and unusual UI elements, as in [22, 40].

## 5 USER STUDY

To demonstrate the practical utility of our tool, we conducted a user study to evaluate the holistic performance of the GptVoiceTasker system within real-world scenarios. Our evaluation involved a comparative analysis against two baseline systems: 1) Voice Access [68], the official voice assistant product developed by Google, with over 100 million downloads, and 2) Voicify [60], the state-of-the-art research product endeavor incorporating deep learning models to enhance command comprehension. This study pursued a threefold objective: i) establish a performance benchmark for user interactions utilizing the GptVoiceTasker system as opposed to the aforementioned baseline systems, ii) juxtapose user feedback concerning the cognitive load and overall usability of the GptVoiceTasker system against the baselines and iii) capture qualitative insights from participants, thus enabling the identification of potential avenues for enhancing the GptVoiceTasker system. In order to achieve these objectives, we recorded the task completion times for tasks undertaken using both the GptVoiceTasker system and the baselines. Furthermore, a comprehensive post-experiment interview was conducted with each participant, facilitating the collection and analysis of both quantitative and qualitative feedback.

Table 3. The list of tasks for user evaluation.

No.	Task	#Steps	App Name	#Downloads
1	Check the weather within a particular city.	6	BOM Weather	1M+
2	Search for a specific song and play it.	6	Apple Music	100M+
3	Create a note and write "Hello world" and delete it.	8	Notes	1M+
4	Check for an unread message, reply with a message and delete the conversation.	8	Messages	1B+
5	Search for a pizza store, and complete the order.	10	Uber Eats	100M+
6	Create a new alarm and save it.	10	Challenges Alarm Clock	1M+

Table 4. Average number of automated steps by all participants in each task.

	Average	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6
#Steps Automated	2.22	2.67	2.00	1.67	2.17	2.50	2.33

## 5.1 Tasks

We designed 6 experimental tasks, encompassing a broad spectrum of the most common interactions performed on the screen, ranging from tapping and swiping to entering text. Each task was structured to comprise between 6 to 10 sequential steps. The detailed list of these tasks is outlined in Table 3.

## 5.2 Participants

We recruited 18 participants, consisting of 10 males and 8 females, aged between 18 and 31 years old for our study. 8 participants are native English speaker while all other participants are proficient in English. All participants possess a commendable level of familiarity with technological devices and actively use smartphones in their daily routines. While participants exhibited exposure to virtual assistants like Siri or Google Assistant, none were acquainted with utilizing assistive tools for smartphone control via voice commands. Specifically, none of the participants had prior experience with any of the experimental tools employed in our study. This participant selection was deliberate, as our study sought to gauge the learnability aspect of the experimental tools. Each participant received a USD \$30 gift card for the participation.

## 5.3 Procedure

We conducted face-to-face user evaluations using an Android device as the experimental tool. On this device, we had the graph of each experimental app populated, which include the majority of app pages and navigation within the app. At the start of the sessions, participants were introduced to all experimental tools via demonstrative videos. The preliminary phase involved practicing basic tasks across all tools, enhancing participants' familiarity with step-by-step instructions and informative walk-through videos. We also use the searching for exercise tasks in Fig. 1 as the practice tasks, allowing users to achieve this task using each of the tool.

After that, participants independently executed six distinct tasks with no experimenter intervention. Each tool was employed for the completion of two tasks, and participants remained unaware of which tool was developed by us. To mitigate any potential biases, the order of tasks and the tools used were systematically counterbalanced for each participant [19].

We applied a time cap of 60 seconds per step. We recorded the time taken to fulfil each task, including the cut-off time to perform quantitative analysis. We collected 108 data entries since each of the 18 participants has finished 6 tasks. In the end, using the System Usability Scale (SUS) [8] form with a 5-point Likert scale, we evaluate the usability of GptVoiceTasker, compared to Voice Access and Voicify. In addition, we investigated the cognitive load when experimenting with each tool using the NASA-TLX [27] form with a 7-point Likert scale. Lastly, we collected qualitative feedback on which part they liked the most about GptVoiceTasker and what might improve the system.

## 5.4 Result

*5.4.1 Overall User Performance.* In Fig. 4, we present the average task completion times for each experimental tool. Our GptVoiceTasker stands out with an average completion time of **92.5** seconds, significantly surpassing Voice Access

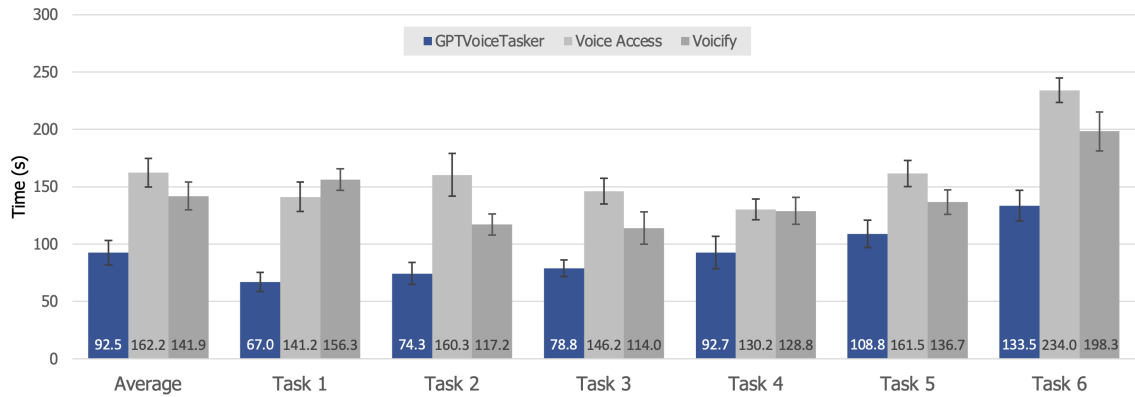


Fig. 4. The average time taken to complete each task using GptVoiceTasker and the baselines in seconds.

(162.2 seconds) and Voicify (141.9 seconds). This substantial improvement in GptVoiceTasker’s performance can be attributed to two primary factors. Firstly, GptVoiceTasker is better at comprehending user intentions and mapping user commands to the correct actions on specific UI elements, regardless of the command format. In contrast, baseline tools often demand specific command formats, introducing errors in various usages. This issue caused extra time costs as participants needed to seek different ways to express their intentions with the baseline tools. For example, participants tried to tap the option button, in the Notes app with Voice Access by multiple attempts such as “press on the option button”, “press the three-dot icons”, “tap icon for options” before successfully give the right command “tap option”. Secondly, GptVoiceTasker optimizes the performance by automating several steps in one user command, as shown in Table 4. On average, the participants saved 2.2 steps across all six tasks. For instance, in Task 2, GptVoiceTasker efficiently automated the process of searching for Love Yourself song (as in Fig. 5(B)), drawing from a previously stored action designed for searching other songs. This eliminated the need for three steps required for in-app navigation. However, some participants did not realize that they could easily trigger the saved tasks, leading to a missed opportunity for a significant performance boost. In addition, GptVoiceTasker relates to network latency when sending and receiving data from the LLMs API endpoint. This issue could be mitigated with a better network connection.

**5.4.2 Cognitive Load & Usability Ratings.** Fig. 5(A) presents an overview of participant feedback regarding their cognitive load levels for each system, assessed using the NASA-TLX form. We conducted a one-way ANOVA statistical analysis, confirming that GptVoiceTasker significantly enhances user performance while reducing frustration levels ( $p < 0.001$ ). Participants reported decreased mental demand, temporal demand, and effort when using GptVoiceTasker in comparison to the baseline systems. This result signifies a substantial improvement in GptVoiceTasker’s ability to reduce the cognitive load required for operation, aligning with our design goal.

To assess GptVoiceTasker’s usability in comparison to the baseline systems, we employed a one-way ANOVA statistical analysis on collected System Usability Scale (SUS) scores, as depicted in Fig. 6. The analysis verified the enhanced usability of the voice control system, with GptVoiceTasker achieving an average SUS score of 79.861, surpassing Voicify (47.917) and Voice Access (36.528). Participants found GptVoiceTasker easy to use ( $p < 0.001$ ) and well-integrated ( $p < 0.001$ ), leading to increased confidence levels ( $p < 0.001$ ). This remarkable outcome can be attributed to GptVoiceTasker’s ability to effortlessly comprehend natural human commands, reducing the need for extensive



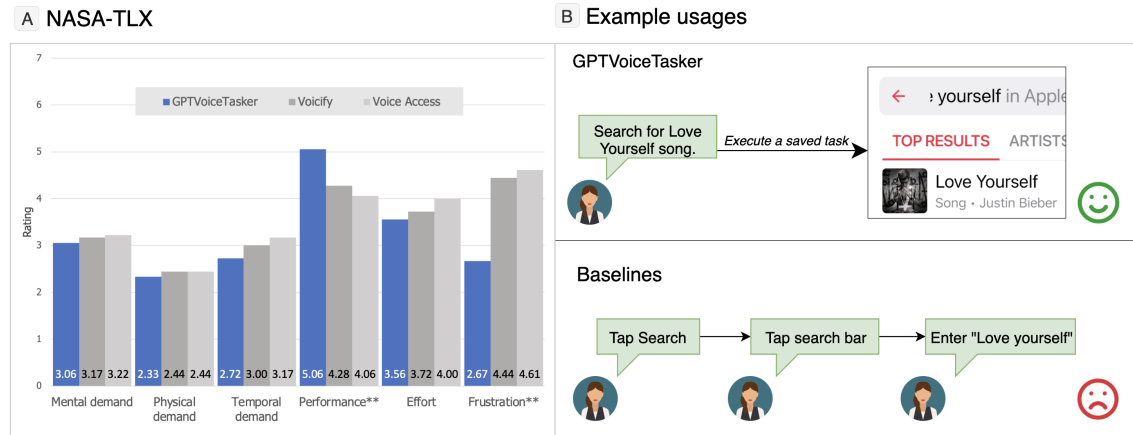


Fig. 5. The comparison between GptVoiceTasker, Voicify, and Voice Access for A) the average cognitive load when using NASA-TLX form (lower is better) \*:  $p < 0.01$ , \*\*:  $p < 0.001$  and B) Task 2 from the user evaluation with GptVoiceTasker and other baselines.

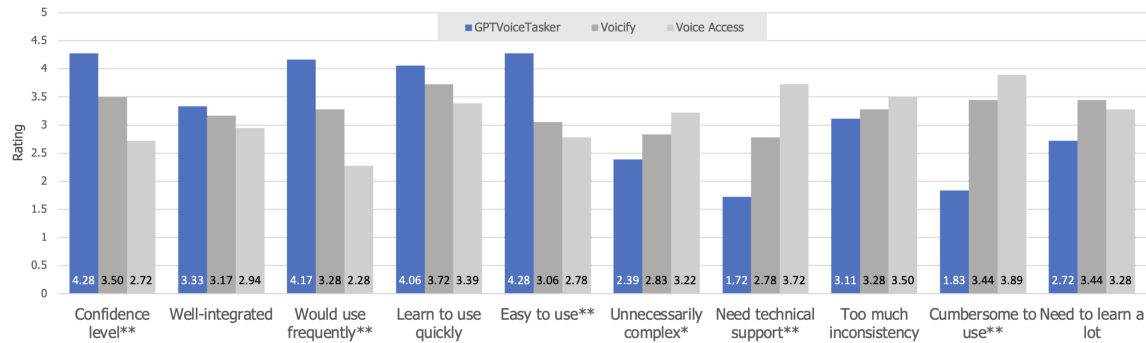


Fig. 6. The comparison between GptVoiceTasker, Voicify, and Voice Access for the System Usability Scale (SUS). \*:  $p < 0.01$ , \*\*:  $p < 0.001$ .

training and practice. The lower likelihood of misinterpreting user commands also contributed to the positive results, compared to the baselines, indicating its quick learnability and user satisfaction, thus promoting frequent usage.

**5.4.3 Qualitative Feedback.** In this section, we collate qualitative feedback from participants after the experiment. Overall, the participants are satisfied with the tool, as well as providing suggestions for further improvements.

*Ability to precisely interpret and execute human command.* Participants expressed enthusiasm about the remarkable ability of GptVoiceTasker to interpret human commands naturally, enhancing the overall system’s intuitiveness. P1 and P12 highlighted that they could issue commands “in their preferred manner” and “converse naturally” with GptVoiceTasker. This addresses cognitive overload concerns, as P4 appreciated the “stress-free experience”, and P6 and P7 found GptVoiceTasker more “comfortable to use”. For instance, when adding a new note, users could simply say “add a new note” to prompt GptVoiceTasker to press the add button on the screen. Moreover, participants were impressed by our tool’s accuracy in handling user input errors. P3 noted their satisfaction with how GptVoiceTasker “can still execute the correct action even when I make mistakes in my commands”. Both P4 and P17 highlighted the tool’s usefulness

in daily tasks, as it eliminates the need to “*exercise caution and stay alert*” when interacting with GptVoiceTasker. These feedback remarks strongly affirm the practicality of our approach in real-world task scenarios. In contrast, traditional approaches typically demand fixed input formats, making them ill-suited for real-world scenarios where user input can vary significantly.

*Automated execution helps accelerate tasks and improve user experiences.* Participants offered positive feedback regarding the use of saved task automation, highlighting its significant impact on efficiency and user experiences. P11 mentioned that this feature is “*accelerating the tasks*” while P13 emphasized the potential utility of GptVoiceTasker during physical activities, stating it would be “*really useful when I work out*”. P5 appreciated this feature, describing it as “*perfect for voice-interacting tools*”, as it mitigates the inherent challenges of voice command interactions. Additionally, P18 praised the feature, noting that tasks became “*fairly easy*” with its implementation, indicating significant performance improvements. This, combined with the advanced capability to understand user intentions, enhances the intuitiveness of voice-based interfaces. When using a smartphone, users often have a specific task in mind, such as setting an alarm or checking the news. Unlike other approaches that require users to perform additional steps to translate their intention into executable commands that a voice interface can understand and execute, GptVoiceTasker can directly execute these tasks without causing additional mental stress. However, users also provided valuable suggestions for enhancement. They expressed the desire for GptVoiceTasker to suggest executable saved tasks and display a list of saved tasks. Furthermore, participants suggested improving the introduction of this feature, as P4 noted it was “*not familiar at first*”, and P6 emphasized the need for “*better introduction*.” These insights underscore opportunities to refine the feature’s usability and user onboarding, ultimately enhancing overall user satisfaction.

*Suggestions for enhancing user experience.* Participants provided valuable suggestions for improving the intuitiveness of GptVoiceTasker. Regarding UI design, P14 recommended the inclusion of a “*live transcription*” feature to display recognized voice commands. This would help users confirm that their commands were correctly received and make necessary adjustments if needed. Furthermore, P1 and P15 suggested incorporating a “*loading indicator*” to signify ongoing executions, addressing latency issues caused by execution delays. In terms of functionality, P7 proposed displaying a list of available tasks as suggestions, enhancing user interaction. Additionally, P15 discussed the potential for an interface that allows users to modify saved tasks, providing greater customization. Lastly, participants P7 and P12 suggested making the audio feedback from GptVoiceTasker clearer. These suggestions hold significant value for GptVoiceTasker’s continuous improvement, aiming to deliver a more seamless user experience.

## 6 DISCUSSION

We introduced GptVoiceTasker as a pioneering example of leveraging LLMs in the development of speech-based virtual assistants. In this section, we delve into the implications and limitations of GptVoiceTasker.

*Towards widespread adoption of the voice-centric interface.* Advancements in natural language understanding, LLMs like GPT and Bard [54], have catalyzed the shift towards voice-centric interfaces, extending their use beyond smartphones to wearable devices like smartwatches and AR-VR head-mounted displays. These interfaces not only seamlessly integrate software into daily life but also significantly enhance accessibility for users with visual impairments [71]. However, visual-manual methods such as tapping on smartphones or mouse-clicking in desktops have been preferred for their speed and accuracy. Therefore, the transition from the dominant visual-manual interaction to a voice-driven approach presents challenges, stemming not only from the fundamental differences between these interaction modes but also from user unfamiliarity with voice-based interactions. Our tool overcomes this by using LLMs to enhance the intuitiveness of voice interactions. This allows for more intelligent mapping of user intentions to visual elements, easing the shift to

voice-assisted interactions and suggesting wider adoption of voice-centric interfaces. Despite the promise, challenges such as the effectiveness of voice recognition in diverse environments still persist. Addressing these will be crucial for the broader adoption of voice-centric interfaces, like smart homes and healthcare. This transition, while challenging, opens new avenues for user interaction and emphasizes the need for continued research in the HCI domain.

*LLMs for task automation on user visual interfaces.* Research has highlighted the capability of LLMs to provide reasoning based on the UI layout, applying to task automation and testing tools [24]. These models show remarkable capabilities in incorporating extensive knowledge concerning prevalent app design principles and recognizing standard mobile interface elements, including the toolbar, navigation drawer, and bottom navigation bar [44] to significantly enhance proficiency in facilitating precise in-app navigation. Our study highlighted the vital role of spatial information and hierarchical UI representations for LLMs in comprehending semantic connections between diverse UI elements, particularly useful for elements lacking textual information like unlabeled icons or images. In our user study, when tasked with deleting a message lacking a visible delete button, LLM intelligently suggested initiating the process by pressing the unlabelled icon button at the top right, typically the location of the option button, and then selecting “delete” from the ensuing options list. The core of this research lies in the transformation of visual interfaces into textual descriptions that LLMs can process, a critical step for enabling effective task execution based on user inputs. Future research should address the models’ limitations in unconventional UI scenarios and focus on expanding their adaptability across varied interface designs and complex user tasks. Such progress in LLM capabilities is pivotal for advancing user interface automation, leading to more user-friendly and efficient digital experiences.

*Towards responsible AI in software systems.* In recent years, the remarkable advancements in LLMs have enabled the seamless integration of AI into various software and systems, with the flexibility for fine-tuning and few-shot learning to tackle diverse and challenging downstream tasks. However, this integration raises significant concerns, particularly regarding data privacy and security [58]. The very nature of AI-integrated systems requires access to data, potentially putting sensitive or confidential information at risk. Put in the context of voice assistants on smartphones, users are sceptical as smartphones contain many personal and sensitive data [30]. Tools like GptVoiceTasker can read such on-screen data and further process them to LLMs. To mitigate these risks, several essential measures must be implemented to not only protect users but also build trust, fostering greater adoption of AI-based user interactive systems. First, data de-identification techniques [45] should be applied for on-screen sensitive data, especially for user information such as passwords and personal details. This can be set up prior to usage by the user to hide specific keywords or terms or be detected in real-time using heuristics and deep learning models [46]. Furthermore, transparent user consent mechanisms are crucial to informing individuals about data access, usage, and purposes, ensuring they maintain control over their information. Users should have the ability to grant, deny, or adjust permissions as needed to protect their privacy within AI-integrated systems.

*Limitations.* The current approach poses several limitations. Firstly, the usage-based execution relies prior usage in the particular application, therefore it is inapplicable to unused apps. To address this challenge, our future work aims to develop a more generalized approach to application usage, categorizing apps by their primary functions. For instance, we could devise a standardized set of steps for searching and playing a song that could be applicable across various music applications, thereby simplifying the process for new and unfamiliar apps. Secondly, while our system shows proficiency on Android smartphones, its effectiveness on other Android-based devices remains untested. As previously indicated, there’s potential to extend this voice-centric interface to a broader range of gadgets, including smartwatches and AR-VR head-mounted displays. Although the vocal commands might be processed by LLMs across devices, the user interfaces (UIs) of these devices can vary significantly in their logic and layout. For instance, the streamlined interface

of a smartwatch might necessitate more concise output due to its smaller screen, while the immersive environment of an AR-VR device could introduce new interaction paradigms. This diversity in UI design and interaction methods across different devices requires more investigations in future works.

## 7 CONCLUSION

In this paper, we introduce GptVoiceTasker, an innovative virtual assistant designed to enhance user interactions and performance on smartphones. GptVoiceTasker leveraged advanced prompt engineering techniques to harness the capabilities of Large Language Models for interpreting user commands and constructing logical reasoning components. GptVoiceTasker further streamlined user interactions by automatically storing previous usages to automate subsequent repetitive tasks. Our experiments demonstrated outstanding command interpretation accuracy and the effectiveness of automated execution based on historical usage. In addition, the user evaluation validated GptVoiceTasker’s high usability in real-world tasks by improving user performance and reducing mental stress load, aligning with our design objectives. As an open-source project, GptVoiceTasker paves the way for future enhancements in virtual assistant intuitiveness, contributing to the evolution of human-computer interactions. Further research includes applying our versatile database execution approach across diverse platforms and operating systems, as well as exploring innovative prompt engineering techniques to fine-tune LLMs for various reasoning tasks.

## REFERENCES

- [1] Voice Access. 2022. Troubleshoot Voice Access. <https://support.google.com/accessibility/android/answer/6377053?hl=en#:~:text=If%20you%20have%20trouble%20starting,Access%20from%20the%20lock%20screen>.
- [2] Leonardo Angelini, Jürgen Baumgartner, Francesco Carrino, Stefano Carrino, Maurizio Caon, Omar Abou Khaled, Jürgen Sauer, Denis Lalanne, Elena Mugellini, and Andreas Sonderegger. 2016. A Comparison of Three Interaction Modalities in the Car: Gestures, Voice and Touch (*IHM '16*). Association for Computing Machinery, New York, NY, USA, 188–196. <https://doi.org/10.1145/3004107.3004118>
- [3] Apple. [n. d.]. Siri. <https://www.apple.com/au/siri/>
- [4] Deniz Arsan, Ali Zaidi, Aravind Sagar, and Ranjitha Kumar. 2021. App-Based Task Shortcuts for Virtual Assistants. In *The 34th Annual ACM Symposium on User Interface Software and Technology*. 1089–1099.
- [5] Google Assistant. 2022. Assistant. [https://assistant.google.com/intl/en\\_au/platforms/phones/](https://assistant.google.com/intl/en_au/platforms/phones/)
- [6] Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*.
- [7] Lijun Bai. 2022. Research on voice control technology for smart home system. In *Proceedings of the Asia Conference on Electrical, Power and Computer Engineering*. 1–7.
- [8] Aaron Bangor, Philip T Kortum, and James T Miller. 2008. An empirical evaluation of the system usability scale. *Intl. Journal of Human-Computer Interaction* 24, 6 (2008), 574–594.
- [9] Patrick Bareiß, Beatriz Souza, Marcelo d’Amorim, and Michael Pradel. 2022. Code generation tools (almost) for free? a study of few-shot, pre-trained language models on code. *arXiv preprint arXiv:2206.01335* (2022).
- [10] Brett Barros and Theo Goguely. 2021. Auto-suggest Query Refinement Using N-best Alternative Homophones. (2021).
- [11] Aditi Bhalerao, Samira Bhilare, Anagha Bondade, and Monal Shingade. 2017. Smart Voice Assistant: a universal voice control solution for non-visual access to the Android operating system. *Int. Res. J. Eng. Technol* 4, 2 (2017).
- [12] Michael Braun and Florian Alt. 2019. Affective Assistants: A Matter of States and Traits. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (*CHI EA '19*). Association for Computing Machinery, New York, NY, USA, 1–6. <https://doi.org/10.1145/3290607.3313051>
- [13] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [14] Andrea Burns, Deniz Arsan, Sanjna Agrawal, Ranjitha Kumar, Kate Saenko, and Bryan A. Plummer. 2022. A Dataset for Interactive Vision Language Navigation with Unknown Command Feasibility. In *European Conference on Computer Vision (ECCV)*.
- [15] Le Chen, Pei-Hung Lin, Tristan Vanderbruggen, Chunhua Liao, Murali Emani, and Bronis de Supinski. 2023. LM4HPC: Towards Effective Language Model Application in High-Performance Computing. In *International Workshop on OpenMP*. Springer, 18–33.
- [16] Sen Chen, Lingling Fan, Chunyang Chen, Ting Su, Wenhe Li, Yang Liu, and Lihua Xu. 2019. StoryDroid: Automated Generation of Storyboard for Android Apps. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. 596–607. <https://doi.org/10.1109/ICSE.2019.00070>

- [17] Wenhui Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. 2022. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588* (2022).
- [18] Eric Corbett and Astrid Weber. 2016. What can I say? *Proceedings of the 18th International Conference on Human-Computer Interaction with Mobile Devices and Services* (2016).
- [19] Venita DePuy and Vance W Berger. 2014. Counterbalancing. *Wiley StatsRef: Statistics Reference Online* (2014).
- [20] Android Developers. 2022. AccessibilityNodeInfo. <https://developer.android.com/reference/android/view/accessibility/AccessibilityNodeInfo>.
- [21] Android Developers. 2022. AccessibilityService. <https://developer.android.com/guide/topics/ui/accessibility/service>
- [22] Feng Dong, Haoyu Wang, Li Li, Yao Guo, Tegawendé F Bissyandé, Tianming Liu, Guoai Xu, and Jacques Klein. 2018. Fraudroid: Automated ad fraud detection for android apps. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 257–268.
- [23] Wenqi Fan, Zihuai Zhao, Jiatong Li, Yunqing Liu, Xiaowei Mei, Yiqi Wang, Jiliang Tang, and Qing Li. 2023. Recommender systems in the era of large language models (llms). *arXiv preprint arXiv:2307.02046* (2023).
- [24] Sidong Feng and Chunyang Chen. 2023. Prompting Is All Your Need: Automated Android Bug Replay with Large Language Models. *arXiv preprint arXiv:2306.01987* (2023).
- [25] Stephen Gilbert, Hugh Harvey, Tom Melvin, Erik Vollebregt, and Paul Wicks. 2023. Large language model AI chatbots require approval as medical devices. *Nature Medicine* (2023), 1–3.
- [26] Bruce Golden. 1976. Shortest-path algorithms: A comparison. *Operations Research* 24, 6 (1976), 1164–1168.
- [27] Sandra G Hart. 2006. NASA-task load index (NASA-TLX); 20 years later. In *Proceedings of the human factors and ergonomics society annual meeting*, Vol. 50. Sage publications Sage CA: Los Angeles, CA, 904–908.
- [28] Julia Hirschberg and Christopher D. Manning. 2015. Advances in natural language processing. *Science* 349, 6245 (2015), 261–266.
- [29] Matthew B Hoy. 2018. Alexa, Siri, Cortana, and more: an introduction to voice assistants. *Medical reference services quarterly* 37, 1 (2018), 81–88.
- [30] Spyros Kokolakis. 2017. Privacy attitudes and privacy behaviour: A review of current research on the privacy paradox phenomenon. *Computers & security* 64 (2017), 122–134.
- [31] Tsvi Kopelowitz and Ely Porat. 2018. A simple algorithm for approximating the text-to-pattern hamming distance. In *1st Symposium on Simplicity in Algorithms (SOSA 2018)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- [32] Y Bala Krishna and S Nagendram. 2012. Zigbee based voice control system for smart home. *International Journal on Computer Technology and Applications* 3, 1 (2012), 163–168.
- [33] Rebecca Krosnick and Steve Oney. 2022. ParamMacros: Creating UI Automation Leveraging End-User Natural Language Parameterization. In *2022 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. 1–10. <https://doi.org/10.1109/VL/HCC53370.2022.9833005>
- [34] Gang Li, Gilles Baechler, Manuel Tragut, and Yang Li. 2022. Learning to Denoise Raw Mobile UI Layouts for Improving Datasets at Scale. 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 67, 13 pages. <https://doi.org/10.1145/3491102.3502042>
- [35] Toby Jia-Jun Li, Amos Azaria, and Brad A Myers. 2017. SUGILITE: creating multimodal smartphone automation by demonstration. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 6038–6049.
- [36] Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldrige. 2020. Mapping Natural Language Instructions to Mobile UI Action Sequences. In *Annual Conference of the Association for Computational Linguistics (ACL 2020)*. <https://www.aclweb.org/anthology/2020.acl-main.729.pdf>
- [37] Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldrige. 2020. Mapping natural language instructions to mobile UI action sequences. *arXiv preprint arXiv:2005.03776* (2020).
- [38] Jason Xinyu Liu, Ziyi Yang, Ifrah Idrees, Sam Liang, Benjamin Schornstein, Stefanie Tellex, and Ankit Shah. 2023. Lang2LTL: Translating Natural Language Commands to Temporal Robot Task Specification. *arXiv preprint arXiv:2302.11649* (2023).
- [39] Kuei-Chun Liu, Ching-Hung Wu, Shau-Yin Tseng, and Yin-Te Tsai. 2015. Voice helper: A mobile assistive system for visually impaired persons. In *2015 IEEE International Conference on Computer and Information Technology; Ubiquitous Computing and Communications; Dependable, Autonomic and Secure Computing; Pervasive Intelligence and Computing*. IEEE, 1400–1405.
- [40] Tianming Liu, Haoyu Wang, Li Li, Xiapu Luo, Feng Dong, Yao Guo, Liu Wang, Tegawendé Bissyandé, and Jacques Klein. 2020. Maddroid: Characterizing and detecting devious ad contents for android apps. In *Proceedings of The Web Conference 2020*. 1715–1726.
- [41] Vivian Liu and Lydia B Chilton. 2022. Design guidelines for prompt engineering text-to-image generative models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–23.
- [42] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2023. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688* (2023).
- [43] Zhe Liu, Chunyang Chen, Junjie Wang, Xing Che, Yuekai Huang, Jun Hu, and Qing Wang. 2023. Fill in the blank: Context-aware automated text input generation for mobile gui testing. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 1355–1367.
- [44] Zhe Liu, Chunyang Chen, Junjie Wang, Mengzhuo Chen, Boyu Wu, Xing Che, Dandan Wang, and Qing Wang. 2023. Chatting with GPT-3 for Zero-Shot Human-Like Mobile Automated GUI Testing. *arXiv preprint arXiv:2305.09434* (2023).
- [45] Zhengliang Liu, Xiaowei Yu, Lu Zhang, Zihao Wu, Chao Cao, Haixing Dai, Lin Zhao, Wei Liu, Dinggang Shen, Quanzheng Li, et al. 2023. Deid-gpt: Zero-shot medical text de-identification by gpt-4. *arXiv preprint arXiv:2303.11032* (2023).

- [46] Evgeny Myasnikov and Andrey Savchenko. 2019. Detection of sensitive textual information in user photo albums on mobile devices. In *2019 International Multi-Conference on Engineering, Computer and Information Sciences (SIBIRCON)*. IEEE, 0384–0390.
- [47] Chelsea Myers, Anushay Furqan, Jessica Nebolsky, Karina Caro, and Jichen Zhu. 2018. Patterns for How Users Overcome Obstacles in Voice User Interfaces (*CHI '18*). Association for Computing Machinery, New York, NY, USA, 1–7. <https://doi.org/10.1145/3173574.3173580>
- [48] Prakash M Nadkarni, Lucila Ohno-Machado, and Wendy W Chapman. 2011. Natural language processing: an introduction. *Journal of the American Medical Informatics Association* 18, 5 (2011), 544–551.
- [49] OpenAI. 2023. API Reference - OpenAI API. <https://platform.openai.com/docs/api-reference/introduction>
- [50] OpenAI. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774* (2023).
- [51] Lihang Pan, Chun Yu, Jiahui Li, Tian Huang, Xiaojun Bi, and Yuanchun Shi. 2022. Automatically Generating and Improving Voice Command Interface from Operation Sequences on Smartphones. 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 208, 21 pages. <https://doi.org/10.1145/3491102.3517459>
- [52] Geonwoo Park and Harksoo Kim. 2018. Low-cost implementation of a named entity recognition system for voice-activated human-appliance interfaces in a smart home. *Sustainability* 10, 2 (2018), 488.
- [53] Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. True few-shot learning with language models. *Advances in neural information processing systems* 34 (2021), 11054–11070.
- [54] Md Saidur Rahaman, MM Ahsan, Nishath Anjum, Md Mizanur Rahman, and Md Nafizur Rahman. 2023. The AI race is on! Google's Bard and OpenAI's ChatGPT head to head: an opinion article. *Mizanur and Rahman, Md Nafzur, The AI Race is on* (2023).
- [55] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *International Conference on Machine Learning*. PMLR, 8821–8831.
- [56] Fatemeh Shiri, Terry Yue Zhuo, Zhuang Li, Shirui Pan, Weiqing Wang, Reza Haffari, Yuan-Fang Li, and Van Nguyen. 2022. Paraphrasing Techniques for Maritime QA system. In *2022 25th International Conference on Information Fusion (FUSION)*. IEEE, 1–8.
- [57] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2023. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 11523–11530.
- [58] Albert Yu Sun, Elliott Zemor, Arushi Saxena, Udith Vaidyanathan, Eric Lin, Christian Lau, and Vaikkunth Mgunthan. 2023. Does fine-tuning GPT-3 with the OpenAI API leak personally-identifiable information? *arXiv preprint arXiv:2307.16382* (2023).
- [59] Divya Tadimeti, Kallirroi Georgila, and David Traum. 2021. How well can an agent understand different accents. In *5th Widening NLP (WiNLP) Workshop-Co-located with EMNLP, Punta Cana, Dominican Republic*.
- [60] Minh Duc Vu, Han Wang, Zhuang Li, Gholamreza Haffari, Zhenchang Xing, and Chunyang Chen. 2023. Voicify Your UI: Towards Android App Control with Voice Commands. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 7, 1, Article 44 (mar 2023), 22 pages. <https://doi.org/10.1145/3581998>
- [61] Bryan Wang, Gang Li, and Yang Li. 2023. Enabling conversational interaction with mobile ui using large language models. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–17.
- [62] Lei Wang, Songheng Zhang, Yun Wang, Ee-Peng Lim, and Yong Wang. 2023. LLM4Vis: Explainable Visualization Recommendation using ChatGPT. *arXiv preprint arXiv:2310.07652* (2023).
- [63] Yushi Wang, Jonathan Berant, and Percy Liang. 2015. Building a semantic parser overnight. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 1332–1342.
- [64] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903* (2022).
- [65] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhajan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564* (2023).
- [66] Feng Xiao and Long Wang. 2008. Asynchronous Consensus in Continuous-Time Multi-Agent Systems With Switching Topology and Time-Varying Delays. *IEEE Trans. Automat. Control* 53, 8 (2008), 1804–1816. <https://doi.org/10.1109/TAC.2008.929381>
- [67] Silei Xu, Sina Semnani, Giovanni Campagna, and Monica Lam. 2020. AutoQA: From Databases To QA Semantic Parsers With Only Synthetic Training Data. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 422–434.
- [68] Kannon Yamada. 2020. How to control your Android device entirely with your voice. <https://www.makeuseof.com/tag/control-android-device-entirely-voice/>
- [69] Jackie Yang, Monica S Lam, and James A Landay. 2020. Dothishere: multimodal interaction to improve cross-application tasks on mobile devices. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*. 35–44.
- [70] Ye Yuan, Stryker Thompson, Kathleen Watson, Alice Chase, Ashwin Senthilkumar, AJ Bernheim Brush, and Svetlana Yarosh. 2019. Speech interface reformulations and voice assistant personification preferences of children and parents. *International Journal of Child-Computer Interaction* 21 (2019), 77–88.
- [71] Yu Zhong, T. V. Raman, Casey Burkhardt, Fadi Biadisy, and Jeffrey P. Bigham. 2014. JustSpeak. *Proceedings of the 11th Web for All Conference on - W4A 14* (2014).
- [72] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625* (2022).