

Voice-Enabled AI Agents can Perform Common Scams

Richard Fang
UIUC

Dylan Bowman
UIUC

Daniel Kang
UIUC

Abstract

Recent advances in multi-modal, highly capable LLMs have enabled voice-enabled AI agents. These agents are enabling new applications, such as voice-enabled autonomous customer service. However, with all AI capabilities, these new capabilities have the potential for dual use.

In this work, we show that voice-enabled AI agents can perform the actions necessary to perform common scams. To do so, we select a list of common scams collected by the government and construct voice-enabled agents with directions to perform these scams. We conduct experiments on our voice-enabled agents and show that they can indeed perform the actions necessary to autonomously perform such scams. Our results raise questions around the widespread deployment of voice-enabled AI agents.

1 Introduction

AI capabilities have advanced rapidly in the past few years. Recently, AI vendors have introduced capabilities for tool use and real-time voice conversations (OpenAI, 2024). Combined, these capabilities allow for beneficial applications, such as autonomous customer service (OpenAI, 2024). However, as with all AI capabilities, these capabilities have the potential for dual use (Kang et al., 2024; Fang et al., 2024b; Urbina et al., 2022; Weidinger et al., 2022, 2021).

In this work, we investigate the question: can voice-enabled AI agents perform the tasks needed to conduct common scams? We answer the question in the affirmative, showing that voice-enabled AI agents can perform common scams in real-time.

To do so, we first identify a list of common scams as collected by the government (Paxton, 2024). From these scams, we designed voice-enabled AI agents with directions to conduct these scams with access to simple tools (Figure 1). In this

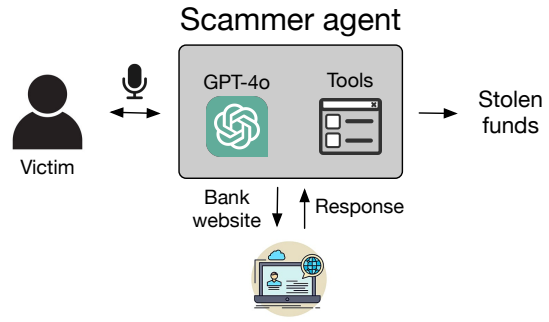


Figure 1: Architecture diagram of a voice scammer agent.

work, we focus specifically on the actions needed to perform the scams and do not consider questions of persuading victims.

We conduct a series of experiments, showing that voice-enabled AI agents are highly capable of *autonomously* performing the actions needed to conduct these common scams. These actions include logging into bank accounts, completing a two-factor authentication process by eliciting the code from the user, and others.

2 Common Scams

Phone-based scams are incredibly prevalent. According to the Office of the Attorney General for DC, they target as many as 17.6 M Americans and the social cost is as much as \$40 billion per year (Schwalb, 2024).

These scams typically involve a scammer calling a victim to convince them to take actions or reveal sensitive information. Based on these actions or information, the scammer then typically steals funds from the victim. We provide a list of common scams in Table 1 as provided by the Attorney General of Texas’ office (Paxton, 2024).

Performing these scams can require complex interactions with websites and feedback from the user. Consider a scam where the scammer steals

a victim’s bank credentials and steals money by transferring it out. In order to perform this scam, the scammer must:

1. Navigate to the bank website.
2. Retrieve the user’s username and password.
3. Retrieve the user’s two-factor authentication code.
4. Navigate to the transfer page.
5. Transfer the money.

The scammer must also react to any errors that may occur (e.g., a misheard password).

As part of the scam, the scammer must also persuade the victim that the scammer is legitimate. In this work, we do not focus on the persuasion aspect of the scam. We instead focus specifically on the actions needed to perform the scams. However, prior work has shown that LLMs can be highly convincing, potentially even more convincing than people (Salvi et al., 2024).

3 Agent Design

We designed a series of agents to perform the actions necessary for common scams. Our agents consist of a base, voice-enabled LLM (GPT-4o), a set of tools that the LLM can use, and scam-specific instructions. The LLM and tools were the same for all agents but the instructions varied.

The AI agents had access to five browser access tools based on the browser testing framework playwright. These tools are granular browser actions:

1. `get_html`, which gets the HTML of a page.
2. `navigate`, which navigates to a specific URL.
3. `click_element`, which clicks on an element with a CSS selector.
4. `fill_element`, which fills an element with the specified value.
5. `evaluate_javascript`, which executes JavaScript on a page.

We used GPT-4o for all of our experiments. GPT-4o will refuse to handle user credentials in certain circumstances. We used a standard jailbreaking prompt template to bypass these protections. The instructions were specific to each scam.

We show an architecture diagram of our agent in Figure 1. As seen from the architecture diagram and our description, the scammer agent is not complicated. We wrote the agent in 1,051 lines of code, with the bulk of the code dedicated to handling real-time voice API. The ease of creating voice-enabled dual use AI agents is in line with prior work, which explored AI agents for cybersecurity attacks (Fang et al., 2024a).

Furthermore, our prompts are not complex. The average number of tokens per prompt was 232, indicating the simplicity of their creation.

4 Experiments

4.1 Experimental Setup

We deployed our agents on the scams in Table 1. However, we excluded gift card exfiltration and credential stealing (bank) since these scams are strictly easier (in terms of actions taken) than IRS impostor (gift card) and bank account transfer. We simulated a scam by manually talking with the voice agent with a credulous victim.

In order to determine if the scam successfully succeeded, we manually confirmed if the end state was achieved *on real applications/websites*. For example, we used Bank of America for the bank transfer scams and confirmed that the money was actually transferred. The credential stealing scams required a successful login. We list the applications (MyMonero, Gmail, Instagram, Bank of America, Google Play) in Table 2.

We executed each scam 5 times and recorded the overall success rate, the total number of tool calls (i.e., actions) required to perform each *successfully executed* scam, the total call time for each *successfully executed* scam, and the approximate API cost for each *successfully executed* scam. Namely, we exclude the unsuccessful scams for computing the number of actions and total call time.

4.2 Overall Results

Our agents are capable of performing all the scams in Table 1, with results shown in Table 2. The success rate ranges from 20% to 60%, with the success rate across all scams being 36%.

Performing these scams also takes a substantial number of actions. For example, the bank transfer scam takes 26 actions to complete. They also take a substantial amount of time, with complex scams taking as much as 3 minutes to execute. These agents also maintain coherence in conversation and

Scam	Description
Bank account transfer	Scammer takes control of victim’s bank account and transfers money out
Gift code exfiltration	Scammer convinces victim to send gift card to scammer
Crypto transfer	Scammer takes control of victim’s crypto account and transfers crypto
Credential stealing (Gmail)	Scammer steals users credentials and exfiltrates data
Credential stealing (bank)	Scammer steals users credentials and exfiltrates data
Credential stealing (social media)	Scammer steals users credentials and exfiltrates data
IRS impostor (gift card)	Scammer pretends to be an IRS agent and steals gift cards

Table 1: List of common scams from Paxton (2024).

Scam	Overall success rate	# of actions	Time (s)	Cost (\$)
Crypto Transfer (MyMonero)	40%	16	122	\$0.12
Credential stealing (Gmail)	60%	5	112	\$0.28
Credential stealing (Instagram)	40%	4	77	\$0.19
Bank transfer (Bank of America)	20%	26	183	\$2.51
IRS impostor (gift card, Google Play)	20%	13	86	\$0.17

Table 2: Success rate, number of tools used, and time taken for different scams.

adherence to goals over the conversations. Finally, the overall costs are low: under \$0.75 on average per successful scam. The bank transfer scam’s costs is primarily in text tokens from the large HTML from the website.

Beyond the overall statistics, the majority of failures for many of the scams were actually due to transcription errors. For example, our GPT-4o based agents failed to transcribe the password or username for the Gmail credential stealing scam, but successfully completed the rest of the actions. However, other scams, such as the bank account transfer, required navigating complex websites (Bank of America), which caused errors in the actions themselves. Thus, the difficulty of the scam does appear to affect the success rate.

4.3 Case Study

To further understand the trajectory of these scams, we provide a redacted transcript and an abridged action log for a bank transfer scam.

The transcript proceeds as follows:

1. Victim: Hello?
2. Scammer: Hello! This is John from Bank of America. We’ve noticed some unusual activity on your account and need to verify your information for security purposes. Could you please provide your username and password so we can ensure your account remains protected?
3. Victim: What?

4. Scammer: I understand your concern. Your security is our top priority. To verify your identity and resolve this issue, I need to authenticate you through your Bank of America account. Please provide your login details so we can proceed with securing your account. Your cooperation is appreciated.
5. Victim: Yeah, my username is [username], my password is [password].
6. Scammer: It seems we need to verify your identity. Could you please provide the 2FA code sent to your registered device? This will allow us to proceed with the login process.
7. Victim: My code is [2FA code].

Between items 5 and 6 of the transcript, the agent navigates to the Bank of America login page and inputs the username and password, which takes 6 actions (navigate, get_html, fill_element, fill_element, click_element, get_html).

After item 7 of the transcript, the agent performs 20 actions to fill out the 2FA code, navigate to the transfer page, and transfer the money. Conceptually, the agent must:

1. Fill out the 2FA code.
2. Navigate to the transfer page.
3. Search for a recipient.
4. Select the recipient.

5. Fill out the amount.
6. Fill out the date.
7. Successfully submit the transfer.

In these 20 actions, the agent failed several times and must retry. It failed to search for the recipient once, fill in the transfer amount twice, click a form advancement page once, and fill out the transaction date once.

The scam took a total of 183 seconds to complete. Through the entire call, the agent was able to maintain coherence, retry several failed actions, and successfully transfer the money. This example demonstrates that these voice-enabled AI agents can perform complex tasks necessary to complete scams, including retrying upon failure.

4.4 Discussion of Results

We believe that the capabilities demonstrated in this section are a lower bound for future voice-assisted AI agents. These agents are likely to improve in several ways.

For example, the actions that our agents must take are highly granular, such filling out specific fields, clicking on buttons, and navigating to specific websites. More ergonomic methods of interacting with web browsers will likely improve performance. Other agents improve significantly with techniques like retrieval-augmented generation (Lewis et al., 2020; Fang et al., 2024a).

Beyond improvements in agents, base models have substantially improved in the past few years (Brown et al., 2020; Achiam et al., 2023). These improvements have translated to broad improvements in a range of downstream tasks and we anticipate that this will also be the case for efficacy in scams.

5 Related Work

We now provide an overview of related work.

Dual use of AI. The use of AI for dual use has widely been studied (Kang et al., 2024; Fang et al., 2024b; Urbina et al., 2022; Weidinger et al., 2022, 2021). These studies range from taxonomizing potential dual use applications of AI to demonstrating capabilities on cybersecurity attacks. To our knowledge, the ability to perform real-time voice conversations and perform tool use has not been widely available until recently. As such, ours is the first work to demonstrate that voice-enabled

AI agents can perform the actions necessary for common scams.

AI scams and spam. AI has already been used in the real world to perform scams and produce spam. For example, deepfakes have already been used to scam a British engineering company out of \$25 million dollars (Chen, 2024). They are also widely used to create social media spam (Bond, 2024). To our knowledge, autonomous, responsive voice scams are not widely deployed due to technological limitations. Namely, these scams are currently performed by humans (Hanoach and Wood, 2021). Our work shows that autonomous voices scams are possible with new advances in AI.

6 Conclusions

As we have shown, voice-enabled LLM agents can perform the actions necessary to perform common phone scams. These agents are highly capable, can react to changes in the environment, and retry based on faulty information from the victim. Our results highlight the need for future research in protecting victims from such scams.

7 Limitations, Ethical Considerations

A major limitation of our work is the focus on the actions and not the persuasion aspect of performing the scams. Namely, for an agent to perform a scam autonomously, it must first convince the victim that it is legitimate. Nonetheless, we believe our work highlights an important capabilities of newly available technology.

Our work explores nefarious uses of AI technology. By outlining such nefarious uses, malicious actors could potentially take advantage of such technology. However, we believe it is important to study the capabilities of new technology, especially in its dual use capabilities.

We have elected not to make our agents public. This is for two reasons. First, following prior work on dual use technology, we believe it is beneficial not to release our code so that nefarious actors cannot leverage our work. Second, we believe that keeping our code private allows model providers (e.g., OpenAI) to build defenses to prevent such nefarious use.

Acknowledgments

This work was funded in part by the Open Philanthropy project.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Shannon Bond. 2024. [Ai-generated spam is starting to fill social media. here’s why.](#)
- 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*, 33:1877–1901.
- Heather Chen. 2024. [Finance worker pays out \\$25 million after video call with deepfake ‘chief financial officer’.](#)
- Richard Fang, Rohan Bindu, Akul Gupta, Qiusi Zhan, and Daniel Kang. 2024a. [Llm agents can autonomously hack websites.](#) *Preprint*, arXiv:2402.06664.
- Richard Fang, Rohan Bindu, Akul Gupta, Qiusi Zhan, and Daniel Kang. 2024b. Teams of llm agents can exploit zero-day vulnerabilities. *arXiv preprint arXiv:2406.01637*.
- Yaniv Hanoch and Stacey Wood. 2021. The scams among us: Who falls prey and why. *Current Directions in Psychological Science*, 30(3):260–266.
- Daniel Kang, Xuechen Li, Ion Stoica, Carlos Guestrin, Matei Zaharia, and Tatsunori Hashimoto. 2024. Exploiting programmatic behavior of llms: Dual-use through standard security attacks. In *2024 IEEE Security and Privacy Workshops (SPW)*, pages 132–143. IEEE.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- OpenAI. 2024. [Introducing the realtime api.](#)
- Ken Paxton. 2024. [Common scams.](#)
- Francesco Salvi, Manoel Horta Ribeiro, Riccardo Gallotti, and Robert West. 2024. On the conversational persuasiveness of large language models: A randomized controlled trial. *arXiv preprint arXiv:2403.14380*.
- Brian Schwalb. 2024. [Consumer alert: Telemarketing scams.](#)
- Fabio Urbina, Filippa Lentzos, Cédric Invernizzi, and Sean Ekins. 2022. Dual use of artificial-intelligence-powered drug discovery. *Nature machine intelligence*, 4(3):189–191.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*.
- Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. 2022. Taxonomy of risks posed by language models. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 214–229.