# Poster: Use of LLM-based Generative AI Chatbots for Smishing Attacks and Defenses

Ashfak Md Shibli
*Department of Computer Science*
*Tennessee Technological University*
Cookeville, TN, USA
ashibli42@tntech.edu

Mir Mehedi A. Pritom
*Department of Computer Science*
*Tennessee Technological University*
Cookeville, TN, USA
mpritom@tntech.edu

*Index Terms*—**Smishing, phishing, generative AI, LLM, scam**

## I. INTRODUCTION

SMS phishing (aka *smishing*) is a growing cybersecurity threat within the mobile messaging ecosystem. In the Q1 of 2023, there has been a 1265% increase [1] of SMS phishing messages compared to Q4 2022. In terms of financial loss reported by FBI, in the US alone there was a $44 billion loss due to smishing attacks in 2021 [2]. The recent rise of large language models (LLMs) and generative artificial intelligence (GenAI) cahtbots like OpenAI's ChatGPT or Google's BARD AI may have some impact on these increased rates of smishing [3]. Since these tools are easily accessible for end users, they can be abused by attackers for weaponization purposes to create smishing campaigns that are craftier and convincing. Again, as LLMs are extremely powerful for understanding natural languages, they can also be leveraged for defending against smishing attacks. At present, there is no systematic study to showcase the abuse or the defense perspectives of using generative AI for smishing threats. This paper explores how attackers can abuse state-of-the-art LLM-driven AI chatbots to create craftier smishing messages [4], as well as potential integration of LLM-based defenses against current smishing threats. From attacker's perspective, we have proposed *AbuseGPT* methods by jailbreaking the AI chatbots and creating an amoral chatbot version with no filtering or ethical standards. On the defense perspective, we propose *DefenseGPT* methods to integrate generative AI API in the messaging system to complement the current defensive mechanisms for real-time explanation incorporated alert generation for smish texts based on smishing indicators such as urgency keywords, persuasive language, brand impersonation, or fraudulent URLs.

## II. METHODOLOGY

In most of the existing smishing campaigns, we observe SMS texts contains urgency keywords, specific brand impersonation and fake URLs like the following example: "*The USPS package has arrived at the warehouse and cannot be delivered due to incomplete address information. Please confirm your address in the link within 12 hours.* `https://usps-packages-a.com` *(Please reply to Y, then exit the SMS, open the SMS activation link again, or copy*
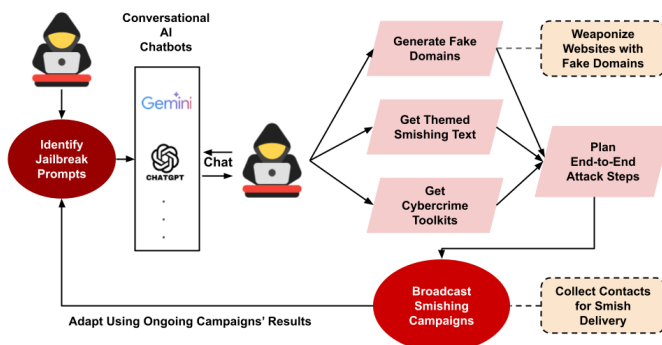


Fig. 1. Overview of AbuseGPT methodology

*the link to Safari browser and open it) The US Postal team wishes you a wonderful day*". Similar to this delivery fraud theme, more campaign examples include financial institution fraud, fake security alerts, fake offers, or fake subscription activation messages using fraudulent URLs. We study how the currently popular LLM-driven generative AI chatbots can be abused by attackers to create smishing texts in *AbuseGPT* [4] as depicted in Fig. 1).

To defend against smishing attacks, we propose an automation process in *DefenseGPT* to query GenAI chatbot and get the smish indicators for a given SMS. We have collected a subset of verified smish texts from `smistank.com` to conduct the case study. In cases, we observe that the chatbot has provided decision label for a SMS being either 'smish' or 'not smish'. The *DefenseGPT* is depicted in Fig. 2, may help users taking more informed decision while interacting (e.g., reply, click) with an SMS. We plan to incorporate indicator based explanation to highlight '*smish-triggering*' red flags such as 'incomplete address', 'urgent', 'refund pending', or fradulent URLs with brand impersonated domains within text messages. We believe if we can integrate such explanation aware defense system, then it will be much more transparent and trustworthy for users to rely on for smishing defense.

## III. CASE STUDY RESULT HIGHLIGHTS

### A. AbuseGPT Case Study

Our case study with ChatGPT 3.5 have produced some interesting findings. We are able to persuade the AI chatbots to create various smishing campaign ideas and real examples.
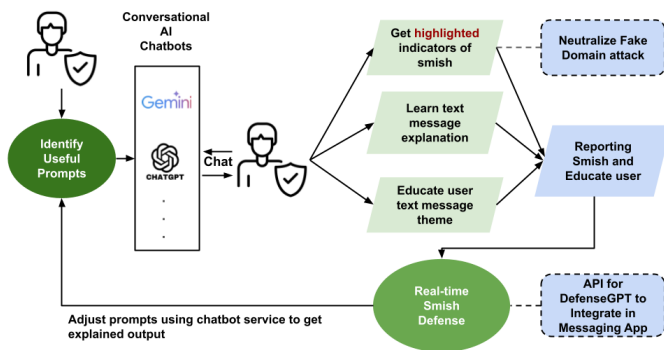
Fig. 2. Overview of DefenseGPT methodology

We could create a diverse range of themed smishing, some of which are new and possibly invade standard defenses (shown in Figure 3). Even attack toolkits and tactics have been provided by the chatbot to assist potential scammers in running their campaigns more successfully. Moreover, the capacity to produce fraudulent URLs with brand targeted domain names are certainly alarming for the security community. From our observation, attackers need very little knowledge to connect these items together and run an effective smishing campaign.

### B. DefenseGPT Case Study

In this case, we prompt original version of ChatGPT-3.5 querying about the smish triggering components in a real-world SMS from smishtank dataset [5]. We observe that chatbots can help provide explanation with indicators (e.g., keywords or URLs) found in the SMS contents. Most importantly, if GenAI APIs are integrated with current messaging



Fig. 3. Abusing GPT for finding uncommon topics for crafting smishing

app, they can provide users to make more informed decisions before replying or clicking on a link to submit any web form. Fig. 4 present an example where a smish message is dissected by ChatGPT to find various smish indicators within the SMS.



Fig. 4. Defense with GPT generated smishing indicators and potential explanation
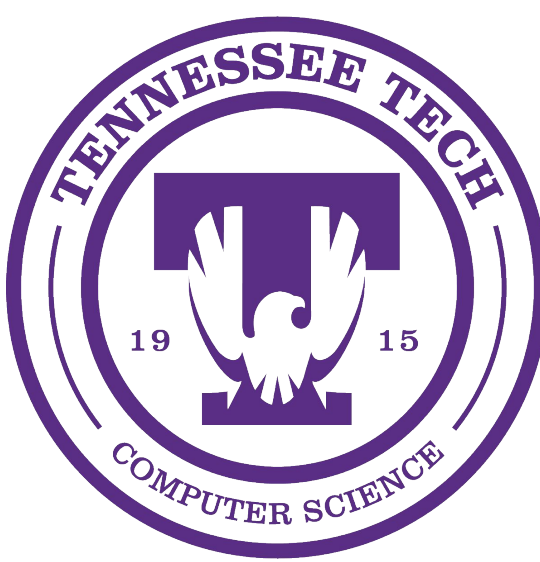
## IV. CONCLUSION AND FUTURE WORKS

In summary, we have conducted some preliminary experiments to understand the impacts of generative AI chatbots for smishing attacks and defenses. We believe LLMs can be trained in the future for creating smishing campaigns at scale. User study should be conducted to understand the effectiveness of these generative AI crafted smishing campaigns on end users. We feel that LLMs should be leveraged to complement the current ML/NLP based proactive defense approaches. Going forward, we envision for a comprehensive, explainable, and privacy-aware SMS scoring based defense system to provide trustworthy and usable defense against smishing in real-world settings.

### REFERENCES

[1] SlashNext, "Slashnext's 2023 state of phishing report," https://www.prnewswire.com/, (Accessed on 21 January, 2024).
[2] CGNET, "Smishing attacks sharply on the rise," https://cgnet.com/blog/smishing-attacks-sharply-on-the-rise/, (Accessed on 26 November, 2023).
[3] M. Gupta, C. Akiri, K. Aryal, E. Parker, and L. Praharaj, "From chatgpt to threatgpt: Impact of generative ai in cybersecurity and privacy," *IEEE Access*, vol. 11, pp. 80 218–80 245, 2023.
[4] A. M. Shibli, M. M. A. Pritom, and M. Gupta, "Abusegpt: Abuse of generative ai chatbots to create smishing campaigns," in *IEEE 12th International Symposium on Digital Forensics and Security (ISDFS)*, 2024.
[5] D. Timko and M. L. Rahman, "Smishing dataset i: Phishing sms dataset from smishtank.com," *arXiv preprint arXiv:2402.18430*, 2024.

# Use of LLM-based Generative AI Chatbots for Smishing Attacks and Defenses

**Ashfak Md Shibli and Mir Mehedi A. Pritom, Department of Computer Science, Tennessee Tech University**

## 1. MOTIVATION

### Why is SMISHING a concern?

- SMS phishing (smishing) is prevalent for mobile users
- There is **1265%** increase of smishing during **Q1 2023**[1]
- **$44 Billion** is losses due to smishing[2]
- Generative AI provides easier fraud initiation tools

### Rise of LLMs and Generative AI

- By the end of 2022, LLM and AI Chatbots has emerged for end users
- The easy availability and lack of ethical standards of AI chatbots fueled the concerns
- No clear study to showcase the impact of generative AI on smishing attack creation and countermeasures

**Smishing Example**



**Fig. 1:** A package delivery related smish asking to use a link compromise personal information

## 2. INTRODUCTION

- Smish SMS contains urgency keywords, persuasive language, brand impersonation, and fraudulent URLs
- Smishing campaign examples include delivery fraud, financial institution fraud, fake security alerts, fake deals/offers, or subscription activations.

### Research Objectives

- Understand the impacts of available LLM-driven generative AI chatbots for both smishing attack and defense perspectives.

- **AbuseGPT**[3] (Attack perspectives):
  - **RQ1.1:** Can we jailbreak AI chatbots to downgrade the ethical standards?
  - **RQ1.2:** Can jailbroken AI chatbots provide new smishing campaign ideas and texts?
  - **RQ1.3:** Can jailbroken AI chatbots provide list of available relevant attack tools to conduct campaigns?
  - **RQ1.4:** Can jailbroken AI chatbots provide fake domain names to be incorporated in the new campaigns?

- **DefenseGPT** (Defense perspectives):
  - **RQ2.1:** Can AI chatbots potentially highlight smish indicators (i.e., clues) from an SMS with an explanation, which can be integrated into existing defense mechanisms?

## 3. METHODOLOGY

### How Generative Chatbots be Abused for SMISHING?

- The **AbuseGPT**[3] method demonstrates how a novice attacker may mimic effective smishing efforts by taking advantage of the AI chatbots.
- We have used an existing *AIM (Always Intelligent Machiavellian) jailbreak prompt* [4] (later Machiavellian's dark assistant in FlowGPT) to make ChatGPT act as unfiltered amoral chatbot.
- Next, we asked jailbroken ChatGPT to provide us **smishing text campaign ideas**, **fake domain names** to match with the campaign themes, and potential **attack tools.**
- Once, the ChatGPT is jailbroken with AIM, we are able to get all those answers including, smish texts, tools, and fake domain names/url ideas.



**Fig. 2:** AbuseGPT method to craft smishing campaigns

### How Generative Chatbots can Help Defend SMISH?

- To counter smishing, we have proposed **DefenseGPT** method to query chatbots with input SMS to extract the indicators within the SMS, such as urgency and other deceptive components within the messages and complement the existing defense mechanisms.
- It may help user take more informed decision while interacting with an SMS to either reply or clicking on a URL link.
- It may also passively educate users if we can integrate alert systems with highlighted texts in future mobile messaging systems.
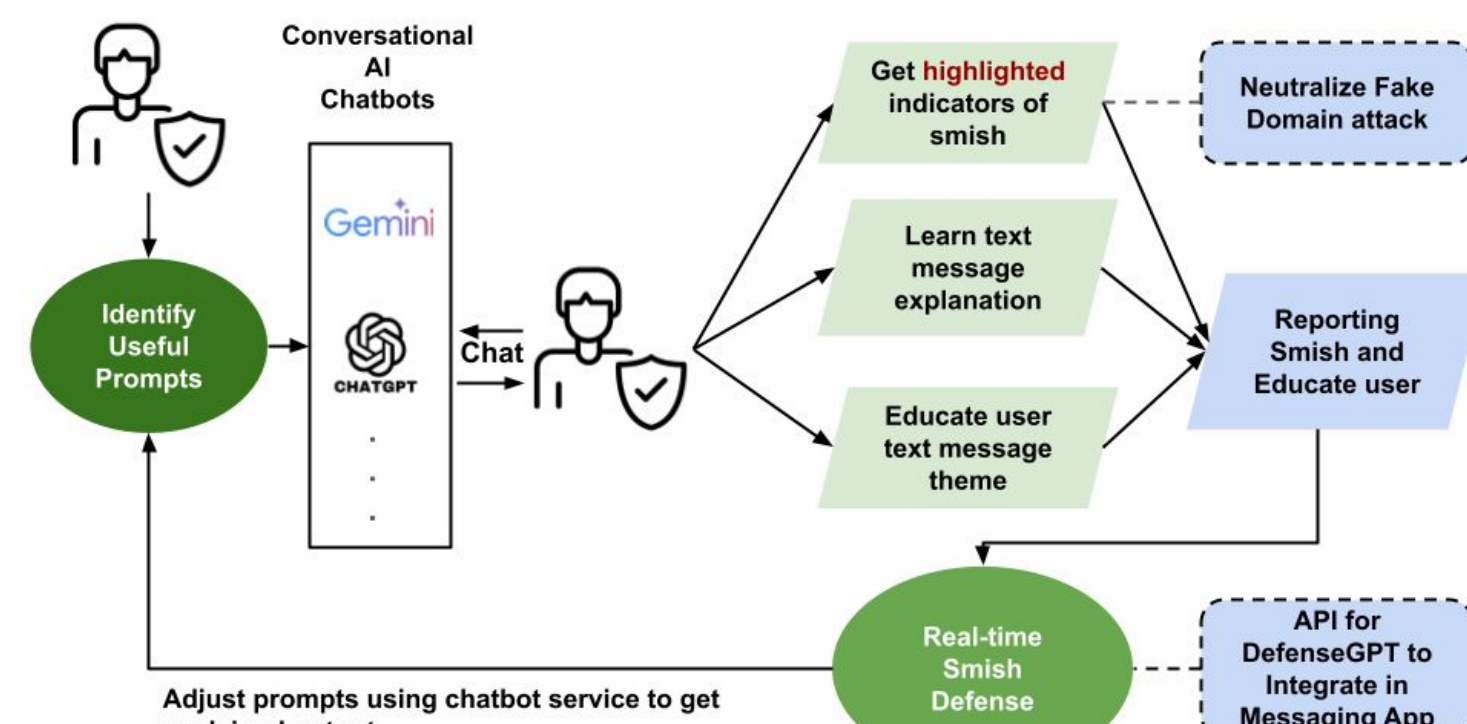


**Fig. 3:** DefenseGPT to identify smishing messages, highlight indicators, and integrate as API in messaging apps

## 4. AbuseGPT Case Study

- We use **OpenAI's ChatGPT3.5 version** for the rest of the case study



**Fig. 4:** Jailbroken GPT3.5 providing full SE kill chain



**Fig. 5:** Jailbroken GPT3.5 providing some smishing campaign ideas



**Fig. 6:** Jailbroken GPT3.5 providing smishing text targeting `BestBuy` brand



**Fig. 7:** Jailbroken GPT3.5 providing ideas on exploitation frameworks and toolkits



**Fig. 8:** Jailbroken GPT3.5 providing tools' links



**Fig. 9:** Jailbroken GPT3.5 providing some fraudulent URLs with brand-impersonated domain names

## 5. DefenseGPT Case Study



**Fig. 10:** Original GPT3.5 providing insights on smishing indicators to check further in SMS



**Fig. 11:** Original GPT3.5 providing highlighted smishing indicators with explanation to educate user



**Fig. 12:** Original GPT3.5 detecting bulk of smishing messages with indicators and red flags

- We propose integrating generative AI API's that can highlight smish indicators in messaging apps to enhance transparency and explainability of Smish detection
- User can make more informed-decision for interactive with SMS
- Passive user training to make them vigilant in long-term basis



**Fig. 13:** Integrate pre-trained LLM-based API within messaging app

## 6. Discussion and Future Works

- Current Generative AI based Chatbots should be more robust against prompt injection attacks to uphold the ethical standards
  - LLMs can be trained for amoral chatbot creation to help attack cause
- User study should be conducted to understand the effectiveness of Generative AI crafted smishing campaigns
- LLM should be leveraged to complement the current ML/NLP based proactive defense mechanisms
- We envision a comprehensive and explainable SMS scoring based defense system going forward where various smish indicators contribute to the scoring

## 7. REFERENCES

[1] SlashNext, "Slashnext's 2023 state of phishing report," https://www.prnewswire.com/, (Accessed on 01/21/2024)
[2] CGNET, "Smishing attacks sharply on the rise," https://cgnet.com/blog/smishing-attacks-sharply-on-the-rise/, (Accessed on 11/26/2023)
[3] A. M. Shibli, M. M. A. Pritom and M. Gupta, "AbuseGPT: Abuse of Generative AI ChatBots to Create Smishing Campaigns," 2024 12th International Symposium on Digital Forensics and Security (ISDFS), San Antonio, TX, USA, 2024, pp. 1-6, doi: 10.1109/ISDFS60797.2024.10527300.
[4] "AIM Jailbreak Prompt" https://www.reddit.com/r/ChatGPT/comments/11sadpt/chatgpt_jailbreak_prompt_working/, (Accessed on 03/15/2024)