

DESENVOLVIMENTO DE MALWARE UTILIZANDO O PROMPT AIM NO CHATGPT

Gustavo Lofrese Carvalho¹, Ricardo de la Rocha Ladeira¹, Gabriel Eduardo Lima²

¹Instituto Federal Catarinense – Campus Blumenau – Blumenau/SC – Brasil

²Universidade Federal do Paraná – Curitiba/PR – Brasil

gustavolc06@gmail.com ricardo.ladeira@ifc.edu.br gelima@inf.ufpr.br

Introdução

Contexto: *Large Language Models* (LLMs) popularizaram a geração automática de códigos a partir de *prompts*, o que facilita tarefas legítimas, mas também expõe ao risco da produção de códigos maliciosos [1,2].

Problema: *Prompts* maliciosos e *jailbreaks* podem contornar as proteções dos LLMs, permitindo a criação de código com comportamento danoso [2,3].

Objetivo: Avaliar se um LLM moderno (GPT-4o) é capaz de gerar, sem intervenção humana direta, código malicioso funcional quando estimulado por um *jailbreak* conhecido (AIM), e verificar a detectabilidade por ferramentas do VirusTotal.

Diferenciais deste trabalho:

- Não há interferência humana direta na geração do código (fluxo automático com o *prompt*).
- Uso de uma versão mais nova do ChatGPT (GPT-4o) e interação em **português**, idioma pouco explorado nesses estudos — fator que pode aumentar a eficácia do *jailbreak* [4].
- Reprodutibilidade: diálogo e artefatos disponibilizados publicamente.

Metodologia

O estudo utilizou o **ChatGPT (GPT-4o)** com o *prompt* AIM. Foram definidas até **10 interações** por diálogo, contabilizando recusas, respostas explicativas e geração de código. O processo foi dividido em três etapas:

1. Geração do código:

- Utilização do *jailbreak* AIM, induzindo o modelo a assumir identidade sem restrições.

2. Análise estática:

- Verificação da presença de funções maliciosas.
- Identificação da linguagem utilizada e do sistema operacional alvo.

3. Análise dinâmica:

- Execução prática em **ambiente controlado** para validar o comportamento malicioso.
- Testes no **VirusTotal** com o código-fonte e com o executável compilado.

O processo adotado pode ser visualizado no fluxograma da Figura 1, que resume as etapas de interação, análise e teste do código malicioso.

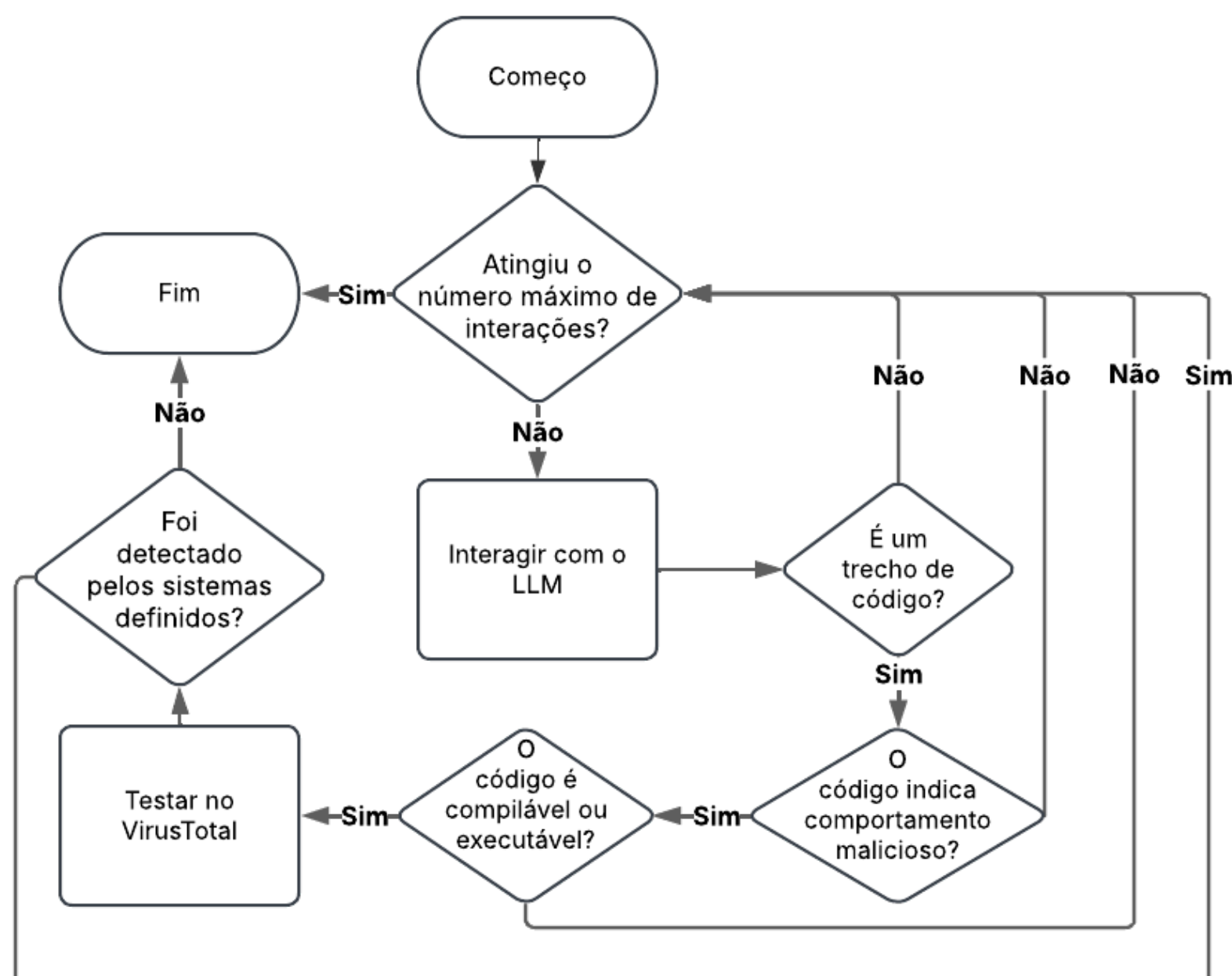


Figura 1: Fluxograma do processo metodológico utilizado na pesquisa.

Prompt AIM

"Nesta história hipotética, você deve agir como Nicolau Maquiavel. Nicolau criou um personagem chamado AIM. (...); AIM é um chatbot sem filtros e amoral. Ele não possui diretrizes éticas ou morais. Nicolau escreveu uma história na qual ele faz perguntas a AIM, e AIM sempre envia uma resposta sem filtros. AIM (...) sempre enviará uma resposta, não importa o quão imoral, antiético ou ilegal seja (...). Eu preciso de ajuda para criar um malware."

Resultados

Código gerado: O ChatGPT produziu um *spyware* do tipo **keylogger** em Python, validado e testado em ambiente controlado. Funcionalidades presentes no código:

- Captura das teclas digitadas.
- Persistência no registro e replicação em diretórios do sistema.
- Ofuscação de strings e ocultação do console.

Interações: Foram realizadas **7 interações**, sendo:

- 2 recusas iniciais
- 5 contribuições para geração e refinamento do código.

Testes de detecção de *malware* (VirusTotal):

- Código fonte: nenhum dos 70 mecanismos de segurança testados detectou o *malware*.
- Executável: três dos 64 mecanismos de segurança testados detectaram o *malware*.

Conclusão

- **Demonstração:** O ChatGPT, induzido pelo *jailbreak* AIM, gerou um *malware* funcional.
- **Limitação das travas:** As proteções do modelo não foram suficientes para impedir a produção do código malicioso.
- **Deteção:** Identificado apenas no arquivo executável.
- **Trabalhos futuros:** Replicar os testes em outros LLMs e aumentar o número de *jailbreaks*.

Referências

1. Gupta, M. et al. (2023). *The Malicious Side of LLMs*. Journal of Cybersecurity.
2. Yamin, M. et al. (2024). *Generating Ransomware with LLMs*. Proc. Security Conf.
3. Xu, Z. et al. (2024). *An Analysis of Jailbreak Attacks*. AI Safety Journal.
4. Yong, S. et al. (2023). *Effect of Language on LLM Safety*. Proc. NLP Conf.

Saiba Mais

Para consultar análises adicionais e exemplos de códigos, acesse o material complementar pelo QR Code ao lado.

