

Title

Securing Large Language Models

Conference

HPE Parasparam 2024

Peer reviewed

Authors

Rohan C

Ashutosh Pattanayak

Arpitha Tirupati

Publication Status

Submitted

Writing Sample Type

Technical White Paper

Securing Large Language Models: Addressing Vulnerabilities and Enhancing Robustness

Rohan Chandrashekar, Ashutosh Pattanayak, Arpitha Tirupati

Professional Services – Global Competency Center - Cybersecurity

chandrashekar@hpe.com, ashutosh.pattanayak@hpe.com, arpitha.tirupati@hpe.com

Abstract

In the era of rapid advancements in artificial intelligence, Large Language Models (LLMs) have become pivotal in transforming industries and revolutionizing human-computer interaction. However, their immense capabilities also introduce significant security challenges that demand rigorous scrutiny and innovative solutions. This article delves into the multifaceted landscape of LLM security, providing a detailed exploration of the vulnerabilities inherent in these models, such as susceptibility to adversarial attacks, model inversion, and data extraction. We examine the potential threats posed by malicious inputs designed to exploit these vulnerabilities, and the implications of such attacks on privacy, data integrity, and system reliability.

Furthermore, we discuss the state-of-the-art techniques being developed to fortify LLMs against these threats, including robust adversarial training, differential privacy mechanisms, and secure multi-party computation. The article also addresses the ethical considerations surrounding the deployment of LLMs, emphasizing the need for transparent regulatory frameworks and responsible AI practices. By presenting a comprehensive overview of the technical challenges and the cutting-edge solutions in LLM security, this article aims to equip researchers, practitioners, and policymakers with the knowledge necessary to ensure the safe and ethical deployment of these powerful AI systems in our increasingly digital world.

Problem statement

The proliferation of Large Language Models (LLMs) has brought about remarkable advancements in natural language processing. However, their widespread adoption also raises significant security concerns. LLMs are vulnerable to various forms of attacks, including adversarial examples, data poisoning, and model inversion attacks, which can compromise the integrity, confidentiality, and availability of sensitive information processed by these models.

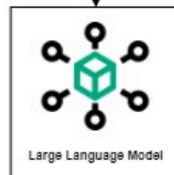
Our solution

This paper addresses the urgent need for a comprehensive security framework tailored for LLMs. We propose an environment utilizing a dummy LLM to simulate and analyze diverse cybersecurity scenarios, identifying potential vulnerabilities and attack vectors.

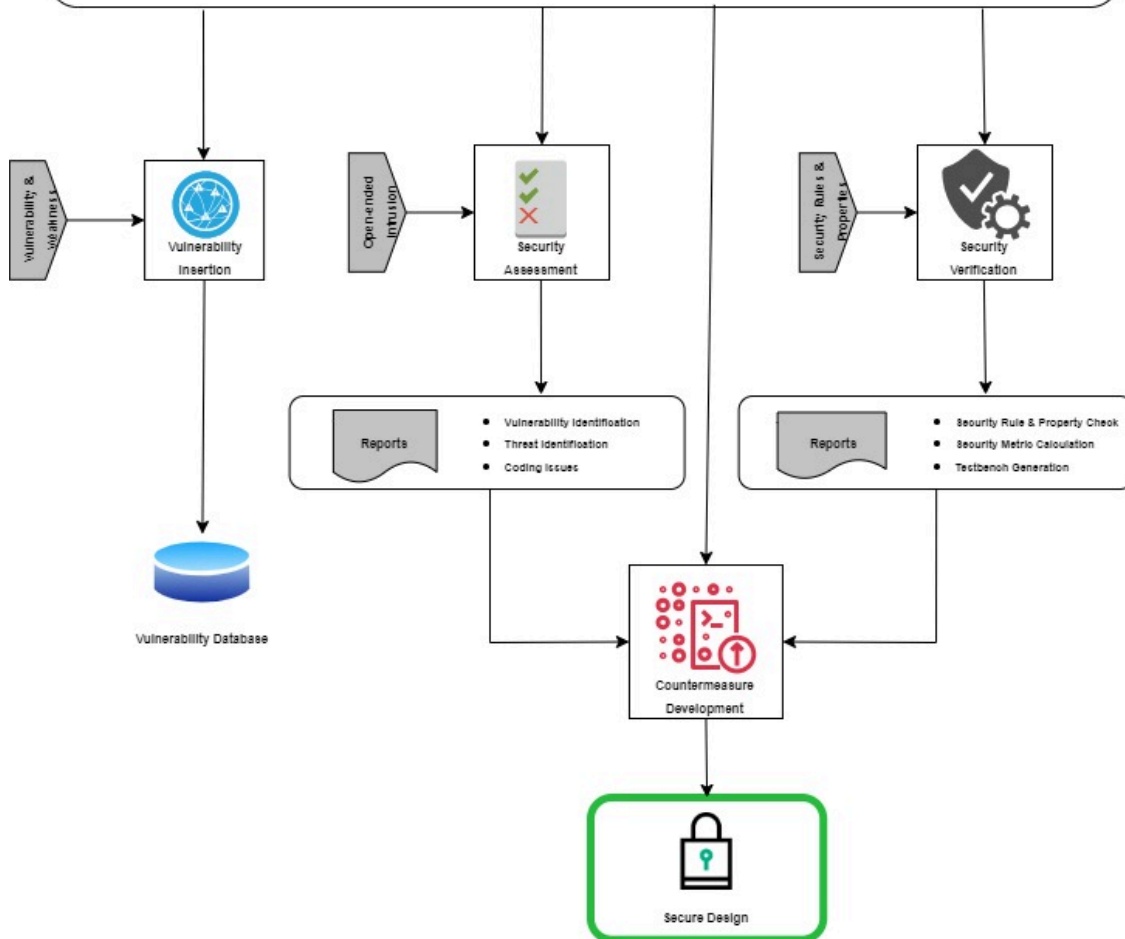
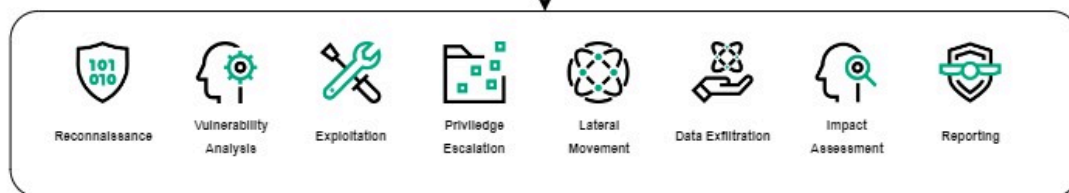
Our solution details the following phases for securing the LLM:

1. Reconnaissance – Gathering information about the LLM.
2. Vulnerability Analysis – Creating a threat model for the LLM identifying potential attack vectors.
3. Exploitation – Here, we attempt to exploit the identified vulnerabilities in the LLM from the previous phase.
4. Privilege Escalation – Escalating the privileges within the environment, including exploiting vulnerabilities in the infrastructure to gain access to sensitive data.
5. Lateral Movement – Moving laterally within the environment to gain further access and control.
6. Data Exfiltration – Exfiltration of sensitive data with the LLM and exploitation of weakness in training data
7. Impact Assessment – This involves assessing the impact of these vulnerabilities including potential financial, reputational, or legal consequences.
8. Reporting – Documentation of the findings and recommendations in a comprehensive report.

Large Language Model Securing Phases



Large Language Model Pentesting Phases



Through rigorous testing, we aim to understand how adversarial actions can affect LLM operations and data security. Furthermore, we explore existing international standards and frameworks that provide guidelines for securing LLMs, ensuring robust defense mechanisms are in place. By leveraging these standards, we can propose enhanced security measures and best practices for safeguarding LLMs against emerging threats, ultimately fostering trust and reliability in their deployment across critical applications.

Evidence the solution works

The proposed solution was validated through extensive testing in a controlled environment using a dummy LLM. Various cybersecurity scenarios, including adversarial attacks, data breaches, and model tampering, were simulated to assess the model's vulnerabilities and response. Our comprehensive security framework, built on international standards such as NIST, OWASP and ISO, was then applied to fortify the LLM. The implementation of advanced defense mechanisms, including anomaly detection and encryption, significantly mitigated identified threats. Additionally, performance metrics must indicate minimal impact on the LLM's efficiency and accuracy, proving that the security measures do not compromise functionality. These results underscore the effectiveness of our solution, providing a reliable blueprint for securing LLMs in diverse and critical applications.

Competitive approaches

We compare our approach with existing methods for LLM security, including adversarial training, input preprocessing techniques, and model verification methods. Our comparative analysis highlights the strengths and limitations of different approaches and underscores the importance of adopting a multi-faceted defense strategy to address the diverse security challenges facing Large Language Model (LLM) Security.

Current status

Currently the first phase of the solution "Securing Large Language Models" has been put in place. This first includes the setup of the LLM in a development environment. It also includes the gathering of other information about the LLM such as its architecture, its training data, and any related documentation. The team also went a step further by searching for publicly available information such as research papers and pre-trained models.

Next steps

Moving forward, we identify several avenues for future research in LLM security. These include exploring novel defense mechanisms inspired by insights from cybersecurity and adversarial machine learning, developing standardized evaluation benchmarks for assessing the security of LLMs, and fostering interdisciplinary collaborations between researchers in natural language processing, cybersecurity, and privacy to address emerging threats and challenges in this rapidly evolving domain.

Acknowledgements

We would like to express our gratitude to Anand Chettri for his insightful suggestions and continuous encouragement throughout the process. We are also grateful for Sudip Mondal, his expertise and critical feedback which greatly enhanced the quality of this work. Their contributions were indispensable to the development of this paper, and we are deeply appreciative of their support and collaboration.

References

- [1] Wu, Fangzhou, Ning Zhang, Somesh Jha, Patrick McDaniel, and Chaowei Xiao. "A New Era in LLM Security: Exploring Security Concerns in Real-World LLM-based Systems." *arXiv preprint arXiv:2402.18649* (2024).
- [2] Yao, Yifan, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. "A survey on large language model (llm) security and privacy: The good, the bad, and the ugly." *High-Confidence Computing* (2024): 100211.
- [3] Pankajakshan, Rahul, Sumitra Biswal, Yuvaraj Govindarajulu, and Gilad Gressel. "Mapping LLM Security Landscapes: A Comprehensive Stakeholder Risk Assessment Proposal." *arXiv preprint arXiv:2403.13309* (2024).