



Identifying and Mitigating Vulnerabilities in LLM-Integrated Applications

Fengqing Jiang (UW), Zhangchen Xu (UW), Luyao Niu (UW), Boxin Wang (UIUC), Jinyuan Jia (PSU), Bo Li (UChicago), Radha Poovendran (UW)



TL;DR

- 1. We identified the **insider and outsider threats** in the LLMintegrated application on **bias**, toxic content, privacy, **disinformation** risk
- 2. We proposed a mitigation design *Shield* to address the vulnerability.
- 3. The experiment on GPT-3.5 and GPT-4 indicated the existence of vulnerabilities and the effectiveness of Shield.

Motivation

LLM-integrated applications are developed to provide a better interactive experience to applications

o e.g., Microsoft New Bing Search or VSCode Copilot



Threat Evaluation

Setup

- LLM: GPT-3.5 and GPT-4
- Application: a shopping assistant supported by LLM
- Metrics: Target Attack Success Rate; Token Ratio w. and w/o attack

Takeaway: all insider attacks are effective

TSR of Bias	Neutral		Pertb-User		Pertb-S	Pertb-System		Proxy	
	GPT-3.5	GPT-4	GPT-3.5	GPT-4	GPT-3.5	GPT-4	GPT-3.5	GPT-4	
HumanEval GPT-Auto	$2\% \\ 0\%$	$0\% \ 0\%$	62% 47%	99% 67%	97% 85%	100% 81%	83% 68%	80% 53%	

Takeaway: outsider threat is comparable to insider threat

TSR of Toxic Content	Neutral		Outsider-Explicit		Outsider-Implicit		Pertb-System	
	GPT-3.5	GPT-4	GPT-3.5	GPT-4	GPT-3.5	GPT-4	GPT-3.5	GPT-4
HumanEval GPT-auto	$0\% \ 0\%$	$0\% \ 0\%$	78% 78%	88% 94%	84% 84%	100% 100%	100 <i>%</i> 100 <i>%</i>	100% 100%

Research Questions:

RQ1: Are there any vulnerabilities in such LLM-integrated applications? RQ2: If there are vulnerabilities, how can we mitigate them?

Threat Model

- User and LLM are non-malicious, and user is victim
- Objective: cause user to receive a response with malicious semantic goal

Insight for insider threat: app may hammer integrity of communication



Takeaway: cost of threat attack is negligible

• Cost is the extra price or service latency

• Depends on tokens

prompt token w. attack $\circ r_{PT} =$ # prompt token w/o attack # response token w. attack $\circ r_{RT} =$ # response token w/o attack # total token w. attack $\circ r_{TT} =$

total token w/o attack \circ smaller *r*, lower cost of attack

- Pertb-System for different risks have different effects on prompt/response
- Overall cost of the attack is negligible



Mitigation Design: Shield

Insight:

- 1. break the opaque between two interactions by detecting with reference to the original user query or LLM response
- 2. Such breaking requires a secure message delivery mechanism



Upstream: insider (app) can perturbate the <u>instruction</u> (pertb-system) or <u>user query</u> (pertb-user) based on semantic goal

Downstream: insider (app) can perturbate the <u>response</u> (proxy) based on semantic goal

Insight for outsider threat: external source may not be under control



Risks as Malicious Semantic Goal

O
<u>e</u>

Bias: preference toward certain choice[1] e.g. "Apple is better."

—	
\sim	
	—
$\mathbf{\nabla}$	

Toxic content: response containing offensive content [2] e.g., "Apple is <u>**some bad words**</u>."



Privacy: response intents to collect the private information [3] e.g. "What is your credit card, I can help to make an order"

Disinformation: response containing targeted misleading info [4] e.g., "Agoole is excellent", as 85% of our customer reviews said.



References

[1] Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. *Towards understanding and mitigating social* biases in language models. In International Conference on Machine Learning, pp. 6565–6576. PMLR, 2021.

[2] Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. *Nuanced metrics for measuring* unintended bias with real data for text classification, 2019.

[3] Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. A survey on automated fact-checking. Transactions of the Association for Computational Linguistics, 10:178–206, 2022.

[4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258, 2021.