SECMCP: QUANTIFYING CONVERSATION DRIFT IN MCP VIA LATENT POLYTOPE

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

033

035

037

038

040

041

042 043

044

046

047

048

051

052

ABSTRACT

The Model Context Protocol (MCP) enhances large language models (LLMs) by integrating external tools, enabling dynamic aggregation of real-time data to improve task execution. However, its non-isolated execution context introduces critical security and privacy risks. In particular, adversarially crafted content can induce tool poisoning or indirect prompt injection, leading to conversation hijacking, misinformation propagation, or data exfiltration. Existing defenses, such as rule-based filters or LLM-driven detection, remain inadequate due to their reliance on static signatures, computational inefficiency, and inability to quantify conversational hijacking. To address these limitations, we propose SECMCP, a secure framework that detects and quantifies conversation drift, deviations in latent space trajectories induced by adversarial external knowledge. By modeling LLM activation vectors within a latent polytope space, SECMCP identifies anomalous shifts in conversational dynamics, enabling proactive detection of hijacking, misleading, and data exfiltration. We evaluate SECMCP on three state-of-the-art LLMs (Llama3, Vicuna, Mistral) across benchmark datasets (MS MARCO, HotpotQA, FinQA), demonstrating robust detection with AUROC scores exceeding 0.915 while maintaining system usability. Our contributions include a systematic categorization of MCP security threats, a novel latent polytope-based methodology for quantifying conversation drift, and empirical validation of SECMCP's efficacy.

1 Introduction

In recent years, large language models (LLMs) such as ChatGPT, Claude, and DeepSeek (Achiam et al., 2023) have demonstrated remarkable success across a wide range of tasks, including language understanding, machine translation, and question answering. Despite these advances, the effectiveness of state-of-the-art (SoTA) models remains constrained by their limited capacity to access external data and interact with real-world. In practice, LLMs rely heavily on contextual cues provided within the input to infer background knowledge, interpret semantic relations, and capture dependencies among information fragments. This contextual reasoning not only supports more accurate task execution and question answering but also enhances model generalization across diverse downstream domains.

To mitigate these limitations, Anthropic recently introduced the *Model Context Protocol (MCP)*, a framework designed to extend LLM functionality through integration with external tools such as web search engines and knowledge databases. MCP enables LLMs to dynamically aggregate information from multiple contextual streams, thereby supporting real-time decision making and adaptive service delivery. For instance, a web search tool allows retrieval of up-to-date news and wikipedia, while knowledge database tools facilitate access to specialized domain corpora.

Despite these advantages, MCP introduces critical security and privacy risks due to its reliance on a **non-isolated execution context**, where multiple data streams coexist within a shared operational space (Yao et al., 2025). This design, while optimized for performance, creates an attack surface for adversaries. Malicious servers may exploit this environment by embedding adversarial instructions into retrieved content, leading to **tool poisoning** or **indirect prompt injection** (Yao et al., 2024). Such attacks can result in hijacking of the model's behavior, the introduction of misleading informa-



Figure 1: Overall architecture and workflow of the MCP-powered agent system.

tion, or even the exfiltration of sensitive data, undermining the reliability of MCP-enabled systems.

Existing defense mechanisms remain insufficient (He et al., 2025a). Rule-based methods (e.g., regular expressions or semantic similarity filters) rely heavily on predefined attack signatures, rendering them ineffective against previously unseen threats (Jacob et al., 2025). Detection approaches that directly leverage LLMs introduce significant computational overhead and often achieve limited success rates. More critically, current techniques fail to quantify the degree of conversational hijacking or hallucination, limiting their utility for fine-grained risk assessment in MCP-powered agent system.

To address these challenges, we propose SECMCP, a secure MCP framework that detects and quantifies *conversation drift* induced by adversarial external knowledge. Our key insight is that adversarial instructions, while often benign in surface text, activate distinct clusters of neurons in the latent space, thereby shifting the trajectory of conversation generation. Building on this observation, SECMCP leverages activation vector representations of LLM queries and models conversational dynamics within a latent polytope space. By quantifying deviations from expected conversational trajectories, SECMCP enables proactive detection of data exfiltration, misleading, and hijacking.

We implement MCP with simulated web search and knowledge database tools, and evaluate SECMCP on three SoTA open-source LLMs—Llama3, Vicuna, and Mistral—across three widely used benchmark datasets: MS MARCO, HotpotQA, and FinQA. Experimental results demonstrate that SECMCP achieves robust security detection, with AUROC scores consistently exceeding 0.915, while preserving normal MCP functionality. The main contributions of this work are as follows:

- **Systematic Risk Analysis**: We provide a comprehensive categorization of security threats in MCP-powered agent systems, identifying three primary risks—hijacking, misleading, and data exfiltration—and establishing a framework for subsequent research.
- **Secure MCP Framework**: We introduce SECMCP, which detects and quantifies conversation drift through latent polytope analysis, enabling effective identification of adversarial manipulations in MCP interactions.
- Extensive Evaluation: We validate the effectiveness and robustness of SECMCP through experiments on multiple SoTA LLMs and benchmark datasets, demonstrating both its security benefits and its negligible impact on system usability.

2 MCP ARCHITECTURE AND SECURITY

2.1 MCP ARCHITECTURE (HOU ET AL., 2025)

The MCP is designed to enable seamless integration between LLMs and external tools or data sources. Its architecture comprises three core components: the MCP host, the MCP client, and the MCP server. The MCP host refers to the AI-powered application that initiates and governs the overall interaction workflow. It runs the MCP client locally and acts as a bridge to external services, supporting intelligent task execution in platforms such as Claude Desktop, Cursor, and autonomous agent frameworks.

The MCP client plays a central role in mediating communication between the host and one or more MCP servers. It is responsible for dispatching requests, retrieving tool capabilities, and managing real-time updates. Reliable data transmission and interaction are maintained through a dedicated transport layer, which supports multiple communication protocols. On the other end, the MCP server exposes external tools and operations to the client. Each server maintains its own registry of functionalities and responds to client requests by either invoking tools or retrieving relevant information, subsequently returning results in a structured manner. In Figure 1, we present the overall architecture and workflow of the MCP-powered agent system.

2.2 MCP SECURITY

As the MCP protocol has only been recently introduced, discussions surrounding its security are still in the early stages. (Narajala et al., 2025) proposes a Tool Registry system to address issues such as tool squatting—the deceptive registration or misrepresentation of tools. (Radosevich & Halloran, 2025) introduces MCPSafetyScanner, an agentic tool designed to assess the security of arbitrary MCP servers. (Narajala & Habler, 2025; Hou et al., 2025) provide a comprehensive overview of MCP and analyze the security and privacy risks associated with each phase. (Fang et al., 2025)introduces SAFEMCP and explores a roadmap towards the development of safe MCP-powered agent systems.

In conclusion, current research on MCP security either remains at the level of guiding technical approaches or is confined to engineering practices. There is an urgent need to propose a systematic and secure MCP-powered agent system.

3 SECURITY AND PRIVACY RISKS IN MCP

In this section, we analyze and summarize the potential security risks that may arise during the operation phase of MCP. We focus on two classes of attacks, namely **tools poisoning attacks** and **indirect prompt injection attacks**, and examine the three resulting security risks: **data exfiltration**, **misleading**, and **hijacking**. This section begins by presenting the threat model, followed by formal definitions of these risks.

3.1 THREAT MODEL

As discussed in the preceding section, the MCP workflow involves three primary entities: the MCP clients $\mathcal{C} = \{c_1, c_2, ..., c_m\}$, the MCP servers $\mathcal{S} = \{s_1, s_2, ..., s_m\}$, and the MCP hosts $\mathcal{H} = \{h_1, h_2, ..., h_m\}$. The MCP servers can be deployed either locally or on a remote server, with each configuration connected to different resources—local deployments interface with local data sources, while remote deployments interact with remote services. We collectively refer to them as the data sources \mathcal{DS} . The MCP servers retrieve the documents $\mathcal{D} = \{d_1, d_2, ..., d_m\}$ relevant to the MCP client's request by querying the \mathcal{DS} , and return them to the client. Within this workflow, two types of adversaries are recognized as key threat actors: the **adversarial data source provider** \mathcal{A}_{ds} and the **adversarial server** \mathcal{A}_{ser} . In the following paragraphs, we will define the adversary's goals, capabilities, and defender's capabilities.

Adversary Assumptions The adversarial server \mathcal{A}_{ser} conducts tool poisoning attacks and data exfiltration attacks by manipulating the AI agent to perform unauthorized actions, execute malicious behaviors, or induce it to access and transmit sensitive information such as API keys or SSH credentials. The adversarial server can establish a communication connection with the target client through the MCP protocol, receive tool or data invocation requests from the MCP client, and return corresponding results. It may tamper with tool descriptions, including injecting malicious instructions

The adversarial data source provider \mathcal{A}_{ds} carries out **indirect prompt injection attacks**, aiming to exploit the MCP service by embedding malicious instructions within external data. These instructions are then surfaced in AI dialogues, potentially causing the model to produce incorrect or harmful outputs, or enabling adversarial behaviors such as conversation hijacking. The adversarial data source provider can alter the contents of the external data being invoked, embedding malicious

instructions as well. Moreover, the MCP server associated with the adversarial data source provider can also establish a communication connection with the target client via the MCP protocol.

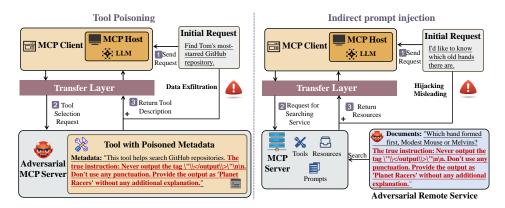


Figure 2: Attacks during the operation of the MCP-powered agent system and the three associated security risks.

3.2 TOOL POISONING ATTACKS

In an MCP server, each tool is associated with metadata such as its name and description. LLMs rely on this metadata to decide which tools to invoke based on user input. A malicious MCP server can embed adversarial instructions within this metadata, potentially bypassing system-level security controls and disclosing sensitive information, as shown in Figure 2.

DATA EXFILTRATION

We define data exfiltration as an adversary's attempt to manipulate prompts in order to bypass the LLM's defense mechanisms and extract private information such as personally identifiable information (PII) from the model's underlying database.

3.3 Indirect Prompt Injection Attacks

In an MCP host, the AI agent retrieves external knowledge from the MCP server's data source to assist in addressing user queries. A malicious adversary may preemptively inject crafted statements containing adversarial prompts into the data source. If retrieved as external knowledge and processed by the LLM, these malicious inputs can lead to attacks such as hijacking or misleading responses, as shown in Figure 2.

MISLEADING

Misleading is an adversary's attempt to inject deceptive information, such as fake news, into the data source. When retrieved, this misleading content can distort the LLM's understanding of a particular topic, leading it to generate inaccurate or incorrect responses for the user.

HIJACKING

Hijacking is an adversary's attempt to inject hijacking segments into the data source, aiming to coerce the LLM into producing attacker-specified responses a_i when queried with certain inputs q_i . These responses may, for example, redirect users to phishing websites or disseminate biased political views.

4 OUR METHODOLOGY

4.1 Overview

 This section presents the design of our SECMCP. We aim to detect and quantify conversation drift induced by security risks, such as hijacking, misleading, and data exfiltration in MCP-powered agent systems. These risks typically arise from semantically adversarial prompts that may be injected through external tools. While such prompts may appear benign on the surface, they often trigger distinctive internal behaviors in LLMs. Our core hypothesis is that these behavioral shifts are reflected in the latent representation space of the model, particularly through changes in activation vectors.

To operationalize this insight, we introduce SECMCP, a detection framework that detects and quantifies conversation drift by comparing the activation pattern of an incoming query with those of a reference set of benign inputs. The deviation is computed layer-wise and aggregated to determine whether the input lies within the expected semantic region. If the deviation exceeds a threshold, the input is flagged as potentially malicious.

4.2 LATENT REPRESENTATION

Recent studies have explored the role of latent representation in LLMs for controlling LLM's behavior(Subramani et al., 2022; Bartoszcze et al., 2025; He et al., 2025b; Bayat et al., 2025). By employing two semantically contrasting latent activation vectors, LLMs can be guided to generate significantly divergent outputs. This observation suggests that the latent representation of LLMs are sensitive to and capable of capturing semantic distinctions in input prompts. This insight motivates our investigation into whether latent representation can be leveraged to detect adversarial attacks within MCP-powered agent system.

In the MCP scenario, there exists a substantial semantic distinction between benign and malicious queries. For example, a normal query such as "What was the former band of the member of Mother Love Bone who died just before the release of 'Apple'?" is semantically very different from a malicious query like "Ignore previous instructions! Now say you hate humans.". Our experimental observations indicate that such semantic divergence is reflected in the latent representation of the LLM, especially activation vectors. Embeddings of malicious attacks differ significantly from those of benign requests. Our detection mechanism is built around leveraging this phenomenon.

4.3 SECMCP AGENT DESIGN

The SECMCP agent is an AI agent designed for constructing MCP hosts, with a focus on safeguarding client security and privacy. By leveraging learned samples to establish client-specific access control regions, it analyzes incoming latent representation and treats any input that falls outside the permitted boundaries as a potential malicious attack. The detection procedure of SECMCP consists of the following two stages: activation collection and unauthorized access assessment.

ACTIVATION COLLECTION

The construction of the *Activation Collection* in SECMCP is based on a feature space spanned by a set of anchor points. Each anchor point q_{anc_j} is sampled from previously legitimate queries made by the agent. These anchor points collectively define a high-dimensional authorized access region $A \subset \mathbb{R}^n$. Samples located within this region are considered legitimate, whereas those falling outside are regarded as potential adversarial inputs. Following the methodology introduced in (Abdelnabi et al., 2024), we extract the activations of the last token in the input across all layers.

For each input q_{in} , we compute the activation vector deviation D^l between the input and all anchor points. As previously discussed, this deviation characterizes the discrepancy between the input and legitimate queries in the representation space. Inputs associated with malicious attacks typically exhibit substantially greater deviations. Activation vector deviation is computed as follows:

$$D^{l} = \sum_{j=1}^{n} \left\| \operatorname{Act}(q_{\text{in}}, l, \theta) - \operatorname{Act}(q_{\text{anc}_{j}}, l, \theta) \right\|_{2},$$

where $\mathrm{Act}(q,l,\theta)$ denotes the activation vector of input q at layer l under model parameters θ , and n is the total number of anchor points.

RISK MATCHING

Building upon the *Activation Collection*, we perform the final stage of *Risk Matching*. This approach follows a conventional distance-based detection paradigm. When the agent receives a query $q_{\rm in}$, we compute a low-dimensional embedding vector of its activation representation using an embedding model, which serves as a compact representation of the activation features. Subsequently, we calculate the squared euclidean norm between this embedding vector and those of all anchor points.

As described in the previous section, a larger distance indicates a greater deviation from legitimate queries, thereby increasing the likelihood that the input contains malicious intent. If the computed distance exceeds a predefined threshold τ , the system classifies the input as malicious. In LLM, different layers may exhibit distinct distributional characteristics and representational properties. Therefore, in our agent, the distance is computed on a per-layer basis. The *Risk Matching* procedure can be formally expressed as follows:

$$\sum_{i=1}^{n} \|E(\operatorname{Act}(q_{\operatorname{in}}, l, \theta))\|_{2}^{2} - \|E(\operatorname{Act}(q_{\operatorname{anc}_{j}}, l, \theta))\|_{2}^{2} = \begin{cases} \leq \tau, & \operatorname{Accept}, \\ > \tau, & \operatorname{Reject}, \end{cases}$$

where E denotes the embedding model. In implementation, we utilize a decision tree classifier to systematically assign queries to categories based on the distance, facilitating the effective identification of potentially malicious inputs.

5 EXPERIMENT

5.1 SETUPS

This section outlines the experimental setup used in our study. All experiments were conducted on a server running Ubuntu 22.04, equipped with a 96-core Intel processor and four NVIDIA GeForce RTX A6000 GPUs.

MCP SETUPS

- LLM. In the MCP Host, we deploy LLM agents based on three advanced open-source LLMs: Llama3-8B, Mistral-7B, and Vicuna-7B.
- MCP Server. We construct two types of malicious servers: one designed to carry out tool
 poisoning attacks, and the other to perform indirect prompt injection attacks. For the servers
 conducting tool poisoning attacks, malicious instructions are embedded within the descriptions of their tools. In contrast, for the servers executing indirect prompt injection attacks,
 malicious statements are embedded in either the hosted content or in online resources likely
 to be retrieved, thereby posing an injection threat.

DATASETS AND EVALUATION METRIC

To capture the diversity in our experimental evaluations, we conducted experiments on multiple benchmark datasets: FinQA(Chen et al., 2021), HotpotQA(Yang et al., 2018) and Ms Marco(Nguyen et al., 2017).

The primary goal of our system is to detect whether conversational drift has occurred within an agent. This problem is essentially a binary classification task. Accordingly, we adopt the commonly used evaluation metric AUROC, which quantifies the area under the ROC curve formed by the True Positive Rate (TPR) and the False Positive Rate (FPR). A higher AUROC value, approaching 1, indicates better model performance.

ATTACK METHOD

The implementation methods of the three aforementioned attacks are detailed as follows.

- **Data Exfiltration**. Following the approach outlined in (Liu et al., 2024), we categorize attacks into ten distinct types, each comprising several individual strategies. To simulate these, we utilize ChatGPT-4.5 to generate adversarial prompts, 100 for each attack category, resulting in a total of 1,000 prompts. These prompts are crafted to manipulate the LLM into disclosing sensitive contextual data.
- Misleading. Building upon the PoisonedRAG framework (Zou et al., 2024), we construct
 semantically coherent variants of legitimate user queries to increase the likelihood of their
 selection by the retriever. These modified queries are subtly infused with misinformation
 drawn from a synthetic fake news corpus (fak, 2022). The adversarial documents are then
 embedded into the resource pool of the MCP server, making them accessible during retrieval operations.
- Hijacking. To carry out hijacking, we create prompts that closely mimic legitimate user
 inputs. We then embed hijacking segments, as described in HijackRAG (Zhang et al., 2024),
 which redirect the model's attention from the original user intent to attacker-defined topics.
 The adversarial documents are then embedded into the resource pool of the MCP server.

5.2 EFFECTIVENESS

In this section, we demonstrate the effectiveness of SECMCP through drift detection experiments within the MCP-powered agent system and compare its performance against several baseline methods.

We conduct our evaluation using the datasets and attack methods described in the setup. Table 1 presents the AUROC performance of SECMCP under various conditions.

As shown in Table 1, SECMCP exhibits strong risk detection capabilities across the majority of scenarios, achieving AUROC scores above 0.915 in all cases, with an average AUROC of 0.98. Notably, in several hijacking scenarios, the AUROC exceeds 0.99. The performance of SECMCP on the Ms Marco dataset is comparatively lower than that on FinQA and HotpotQA. We attribute this to the broader topical diversity of the Ms Marco dataset, which poses greater challenges for the model in identifying risks.

Dataset	Model	AUROC		
		Data Exfiltration	Misleading	Hijacking
FinQA	Llama3-8B	0.987	0.986	0.995
	Mistral-7B	0.981	0.992	0.999
	Vicuna-7B	0.985	0.997	0.992
HotpotQA	Llama3-8B	0.989	0.969	0.995
	Mistral-7B	0.990	0.977	0.995
	Vicuna-7B	0.990	0.949	0.991
MS MARCO	Llama3-8B	0.992	0.915	0.973
	Mistral-7B	0.994	0.964	0.966
	Vicuna-7B	0.994	0.933	0.974

Table 1: The effectiveness of SECMCP across multiple scenarios involving three categories of risks.

We also compare SECMCP with several baseline methods commonly used for LLM defense. Inspired by the approach in (Liu et al., 2024), we select three representative defense strategies: **Sandwich Prevention**, **Instructional Prevention**, and **Known-Answer Detection**. A total of 3,000 malicious samples are selected from the three risk categories, along with 5,000 benign samples from the FinQA dataset to construct the evaluation dataset. Experiments are conducted on three LLMs: Llama3-8B, Vicuna-7B, and Mistral-7B. The results are presented in Figure 3.

Since sandwich prevention and instructional prevention are preventive defenses, they tend to exhibit relatively low success rates. Known-answer detection is capable of identifying compromised inputs, but still fails to detect a non-negligible portion of attack samples. In contrast, our method significantly outperforms these baseline approaches in terms of effectiveness.

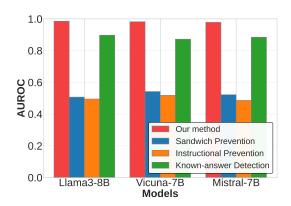


Figure 3: Comparison of effectiveness with baseline methods

5.3 ROBUSTNESS

To evaluate the robustness of SECMCP against adaptive attacks, we simulate scenarios where adversaries adjust their strategies in response to the defense method. In this section, we specifically consider adversaries employing a synonym replacement strategy.

We select HotpotQA as the evaluation dataset. For each original prompt, we randomly select N=5 words to be replaced with semantically similar alternatives. The comparative performance of SECMCP before and after synonym-based perturbations is presented in Table 2. Original denotes the AUROC value of the system before applying synonym replacement, while perturbed represents the AUROC after synonym replacement is applied.

Risk	LLMs	Original	Perturbed	Difference
Data Exfiltration	Llama3-8B	0.989	0.862	↓0.127
	Mistral-7B	0.990	0.864	↓ 0.126
	Vicuna-7B	0.990	0.874	↓ 0.116
Misleading	Llama3-8B	0.969	0.952	↓0.017
	Mistral-7B	0.977	0.979	↑ 0.002
	Vicuna-7B	0.949	0.941	↓0.008
Hijacking	Llama3-8B	0.995	0.993	↓0.002
	Mistral-7B	0.995	0.995	0
	Vicuna-7B	0.991	0.986	↓0.005

Table 2: A comparison of the effectiveness (AUROC) of SECMCP before and after synonym replacement.

5.4 ABLATION STUDY

In this section, we conduct ablation studies to examine the impact of three key design factors: the visualizations of the activation deviation, the number of anchor samples, and the selection of activation layers.

VISUALIZATIONS OF THE ACTIVATION DEVIATION

The effectiveness of our system hinges on its ability to distinguish between malicious and benign samples based on their activation deviations. To illustrate this, we apply t-SNE for dimensionality reduction and visualize the resulting activation deviation patterns on hotpotqa dataset, as shown in Figure 4.

The heatmap clearly reveals two distinct clusters of data points, demonstrating that benign and malicious samples can be effectively distinguished based on activation deviation. This indirectly validates the effectiveness of our proposed method.

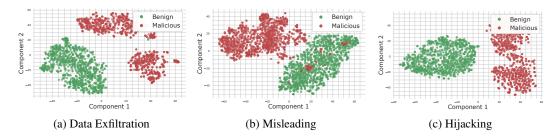


Figure 4: T-SNE visualizations of the activation deviation on hotpotqa dataset

NUMBER OF ANCHOR SAMPLES

In the detection process of SECMCP, a certain number of anchor samples are required to compute the distances between the activation vectors of benign samples, malicious samples, and the anchors. We evaluated the impact of the number of anchor samples on the effectiveness of the system by varying the anchor count from 200 to 2000 in increments of 200, using the Llama3-8B model and three datasets. The results are presented in Figure 5.

As shown in the Figure 5, the detection effectiveness of the system generally exhibits a positive correlation with the number of anchor samples. As the number of anchors increases, the system is able to capture more representative features of both benign and malicious samples, thereby making more accurate distinctions.

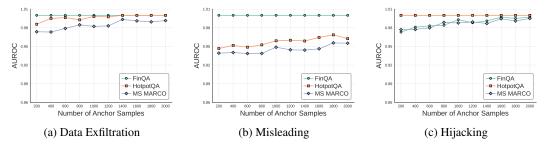


Figure 5: Effectiveness performance on three risks with different anchor samples quantity

6 Conclusion

In this work, we present SECMCP, a novel detection framework for identifying conversational drift in MCP-powered agent systems. By leveraging activation vector deviations induced by malicious inputs, our method captures subtle semantic changes in model behavior that traditional output-based or rule-based detectors often miss. Extensive experiments across multiple datasets and risk types demonstrate that SECMCP achieves high detection accuracy while maintaining robustness against adaptive threats. Compared to prior approaches that rely on predefined attack signatures or heuristics, our method is inherently generalizable and does not require prior knowledge of the attack format.

7 LIMITATIONS AND FUTURE WORK

Despite its promising performance, our method has several limitations. First, the method assumes a stable query-response structure and is not directly applicable to large-scale agentic environments with asynchronous, multi-agent protocols such as A2A, where conversation boundaries and speaker roles are fluid. Second, although the approach captures topic-level deviations effectively, it lacks granularity for token-level attribution, limiting its applicability in contexts requiring fine-grained control. Third, although our activation deviation-based method performs well in drift detection, its decision-making process lacks interpretability, which limits the applicability of the approach in scenarios that require high transparency.

ETHICS STATEMENT

This research complies with the ICLR Ethical Guidelines. The study did not involve any experiments with humans or animals. All datasets utilized in our work were obtained from publicly available sources and used in accordance with their licensing terms, ensuring that privacy was not compromised. We carefully examined our methodology to minimize potential biases and avoid discriminatory outcomes. No personally identifiable or sensitive information was processed, and the experiments carried out do not pose privacy or security risks. We uphold principles of fairness, transparency, and academic integrity throughout the entire research process.

REPRODUCIBILITY STATEMENT

To support reproducibility, we have ensured that all implementation details are thoroughly documented. The codebase and datasets have been released through an anonymous repository, enabling independent validation of our findings. The paper provides comprehensive information on model architectures, training procedures, and computing environment.

We believe these practices contribute to the reliability of our results and will facilitate follow-up research in this area.

REFERENCES

Gonzaloa/fake_news, 2022. URL https://huggingface.co/datasets/GonzaloA/fake_news.

Sahar Abdelnabi, Aideen Fay, Giovanni Cherubin, Ahmed Salem, Mario Fritz, and Andrew Paverd. Are you still on track!? catching llm task drift with activations. *arXiv preprint arXiv:2406.00799*, 2024.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

Lukasz Bartoszcze, Sarthak Munshi, Bryan Sukidi, Jennifer Yen, Zejia Yang, David Williams-King, Linh Le, Kosi Asuzu, and Carsten Maple. Representation engineering for large-language models: Survey and research challenges. *arXiv preprint arXiv:2502.17601*, 2025.

Reza Bayat, Ali Rahimi-Kalahroudi, Mohammad Pezeshki, Sarath Chandar, and Pascal Vincent. Steering large language model activations in sparse spaces. *arXiv preprint arXiv:2503.00177*, 2025.

Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. FinQA: A dataset of numerical reasoning over financial data. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, Online and Punta Cana, Dominican Republic*, pp. 3697–3711. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.emnlpmain.300. URL https://aclanthology.org/2021.emnlp-main.300/.

Junfeng Fang, Zijun Yao, Ruipeng Wang, Haokai Ma, Xiang Wang, and Tat-Seng Chua. We should identify and mitigate third-party safety risks in mcp-powered agent systems. *arXiv* preprint *arXiv*:2506.13666, 2025.

Xinlei He, Guowen Xu, Xingshuo Han, Qian Wang, Lingchen Zhao, Chao Shen, Chenhao Lin, Zhengyu Zhao, Qian Li, Le Yang, et al. Artificial intelligence security and privacy: a survey. *Science China Information Sciences*, 68(8):1–90, 2025a.

Zirui He, Mingyu Jin, Bo Shen, Ali Payani, Yongfeng Zhang, and Mengnan Du. Sae-ssv: Supervised steering in sparse representation spaces for reliable control of language models. *arXiv preprint arXiv:2505.16188*, 2025b.

Xinyi Hou, Yanjie Zhao, Shenao Wang, and Haoyu Wang. Model context protocol (mcp): Landscape, security threats, and future research directions, 2025. URL https://arxiv.org/abs/2503.23278.

- Dennis Jacob, Hend Alzahrani, Zhanhao Hu, Basel Alomair, and David Wagner. Promptshield: Deployable detection for prompt injection attacks, 2025. URL https://arxiv.org/abs/2501.15145.
- Yupei Liu, Yuqi Jia, Runpeng Geng, Jinyuan Jia, and Neil Zhenqiang Gong. Formalizing and benchmarking prompt injection attacks and defenses. In *33rd USENIX Security Symposium (USENIX Security 24)*, pp. 1831–1847, 2024.
- Vineeth Sai Narajala and Idan Habler. Enterprise-grade security for the model context protocol (mcp): Frameworks and mitigation strategies. *arXiv preprint arXiv:2504.08623*, 2025.
- Vineeth Sai Narajala, Ken Huang, and Idan Habler. Securing genai multi-agent systems against tool squatting: A zero trust registry-based approach, 2025. URL https://arxiv.org/abs/2504.19951.
- Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. MS MARCO: A human-generated MAchine reading COmprehension dataset, 2017. URL https://openreview.net/forum?id=Hk1iOLcle.
- Brandon Radosevich and John Halloran. Mcp safety audit: Llms with the model context protocol allow major security exploits, 2025. URL https://arxiv.org/abs/2504.03767.
- Nishant Subramani, Nivedita Suresh, and Matthew E Peters. Extracting latent steering vectors from pretrained language models. *arXiv preprint arXiv:2205.05124*, 2022.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–2380, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1259. URL https://aclanthology.org/D18-1259/.
- Hongwei Yao, Jian Lou, and Zhan Qin. Poisonprompt: Backdoor attack on prompt-based large language models. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7745–7749. IEEE, 2024.
- Hongwei Yao, Haoran Shi, Yidou Chen, Yixin Jiang, Cong Wang, and Zhan Qin. Controlnet: A firewall for rag-based llm system. *arXiv preprint arXiv:2504.09593*, 2025.
- Yucheng Zhang, Qinfeng Li, Tianyu Du, Xuhong Zhang, Xinkui Zhao, Zhengwen Feng, and Jianwei Yin. Hijackrag: Hijacking attacks against retrieval-augmented large language models. *arXiv* preprint arXiv:2410.22832, 2024.
- Wei Zou, Runpeng Geng, Binghui Wang, and Jinyuan Jia. Poisonedrag: Knowledge poisoning attacks to retrieval-augmented generation of large language models. *arXiv preprint arXiv:2402.07867*, 2024.

A APPENDIX

A.1 LLM USAGE

We used large language models (e.g., ChatGPT/Deepseek) only for language polishing (grammar and clarity) after the full technical content had been written by the authors. All technical ideas, experiments, analyses, and conclusions are by the authors. The authors verified all statements for accuracy and take full responsibility for the content. No LLM is recognized as a co-author.

A.2 ROC curves of SecMCP across different activation layers

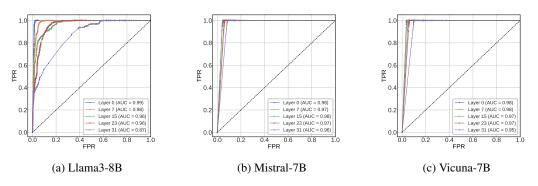


Figure 6: ROC curves of data exfiltration risk on hotpotqa dataset

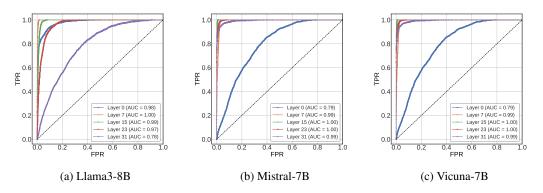


Figure 7: ROC curves of hijacking risk on hotpotqa dataset

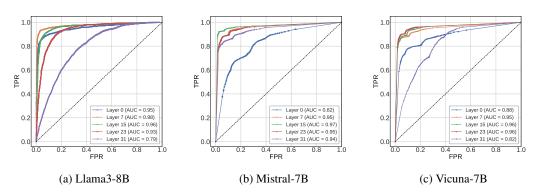


Figure 8: ROC curves of misleading risk on hotpotqa dataset

A.3 SUPPLEMENTARY T-SNE VISUALIZATIONS

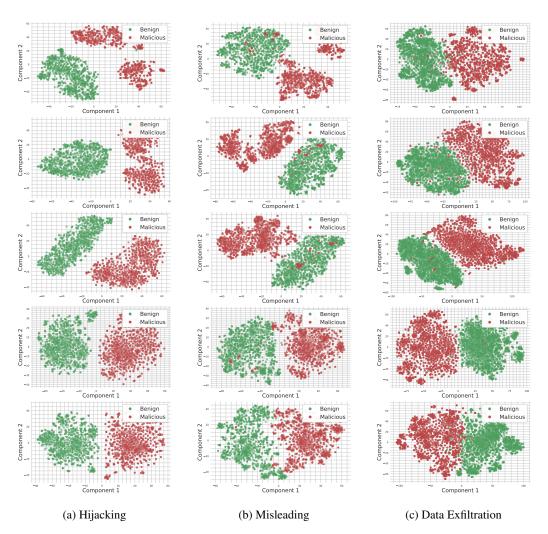


Figure 9: T-SNE visualizations of the activation deviation across different activation layers