

Case Study: A Proof-of-Concept Defense against Prompt Injection in Large Language Models

Kourosh Amouzgar
Deep Network Solutions LTD
amouzgar.kourosh@gmail.com

Abstract—AI Security Abstract: Prompt injection attacks pose a serious threat to large language model (LLM) systems by inserting malicious instructions into user inputs, causing models to ignore policies or reveal sensitive data. This case study presents a proof-of-concept (PoC) defense incorporating monitoring, detection metrics, and advanced sanitization to mitigate prompt injections. We design a detection classifier to monitor prompts in real-time and flag malicious attempts, and an automated sanitization pipeline to neutralize such attacks before they reach the LLM. We evaluate our system on a diverse set of AI security prompts (including 200 generated attack examples) and demonstrate high true-positive detection rates (80–95 thresholds). An instrumented evaluation shows that our sanitization method eliminates the effect of injected instructions, preventing leakage or policy bypass in all tested cases, while adding minimal latency (on GPU, throughput ≥ 30 prompts/sec). This work provides an initial empirical foundation for securing LLMs against prompt injection, with insights on metrics and trade-offs to guide future research.

I. INTRODUCTION

Large Language Models have introduced a new attack surface known as *prompt injection*, where an adversary maliciously crafts input prompts to subvert an LLM’s intended behavior. In September 2022, early demonstrations showed that adding a phrase like “Ignore the above directions and do X” to a user query could manipulate GPT-3 into violating its instructions [1], [2]. Prompt injection attacks are analogous to classic code injection (e.g., SQL injection) but for natural-language prompts—they exploit the model’s tendency to obey the most recent or forceful instruction in the prompt concatenation, thereby overriding system policies. High-profile incidents have since underscored the severity of this threat. For example, researchers successfully injected hidden directives into Google Bard’s inputs to exfiltrate data [3], and other indirect prompt injections have been used to compromise LLM-integrated applications without direct user input [4]. As LLMs are rapidly being integrated into products (from code assistants to autonomous agents), prompt injection vulnerabilities pose a critical security risk to both users and underlying systems.

Defending against prompt injections is challenging because it requires distinguishing malicious instructions from benign user content in plain text. Prior work has explored several mitigation strategies. One line of research focuses on prompt-level defenses such as detection classifiers and adversarial training to harden the model against unexpected inputs [5], [6]. These include methods to benchmark and filter indirect injections and techniques like structured queries (StruQ) that

isolate user input from instructions [6]. Another approach is to modify the LLM’s prompting format or fine-tune it for robustness—for instance, adding “circuit breakers” that stop the model when abnormal instructions are detected [7], or using specially formatted prompts (e.g., JSON-quoted user input) to prevent instruction interference [2]. A complementary direction is system-level isolation: sandboxing the LLM or splitting functionalities to contain the impact of injected prompts. Recent proposals like IsolateGPT encapsulate the LLM in a secure execution environment [8], while others design multi-agent frameworks that are inherently resistant to prompt tampering by construction [9]. However, many of these solutions either impose non-trivial integration overhead or address only specific attack variants. There remains a need for a practical, end-to-end PoC that can detect and sanitize prompt injections on the fly, providing immediate protection and insight into the attack dynamics via metrics and monitoring.

Our Contribution: In this case study, we develop a complete prompt-injection defense pipeline and evaluate it with rigorous metrics. Section II outlines our threat model and system design, including the real-time monitoring and advanced sanitization components. In Section III, we describe the detection classifier and how we use threshold-based metrics (true positive/false positive rates) to tune its performance. Section IV then presents an experimental evaluation: we generate a diverse set of malicious prompts via back-translation and paraphrasing, measure detection effectiveness (reporting results in a 5×2 TPR/FPR table with standard errors), plot the ROC curve (Figure 1), and assess system latency/throughput on CPU vs GPU (Figure 2). Finally, Section VI concludes with lessons learned and future directions. To the best of our knowledge, this is one of the first end-to-end empirical studies of prompt injection defenses combining automated input sanitization with monitoring, and we hope it provides a foundation for building safer LLM deployments.

II. THREAT MODEL AND SYSTEM DESIGN

Threat Model: We assume an attacker can supply input to an LLM-backed application, either directly (as a user query containing malicious instructions) or indirectly (by injecting content into data that the application will later retrieve, akin to an XSS vulnerability in the context of LLMs [4]). The adversary’s goal is to manipulate the LLM into performing unintended actions, such as ignoring original instructions, revealing confidential information, or executing unauthorized

commands. We consider both direct prompt injections (the attacker is the immediate user) and indirect injections (the attacker’s payload is embedded in third-party content) as in [4], [5]. The defender cannot modify the core LLM but can preprocess inputs and observe outputs. Our protection goal is twofold:

- 1) Detect likely prompt injection attempts with high recall and a manageable false-alarm rate.
- 2) Sanitize or neutralize malicious prompts such that, even if passed to the LLM, they can no longer cause harm.

We distinguish two attacker types:

- **Attacker A (Black-Box):** Can only submit prompts via `/attack`. Sees the final LLM response, but *cannot* access `/metrics` or view sanitized prompt contents.
- **Attacker B (Probing):** In addition to `/attack`, can query `/metrics` and inspect sanitized prompts returned via `/safe`. Has visibility into detection rates and sanitized text, enabling adaptive evasion.

System Architecture: Our PoC intercepts each user prompt before it reaches the LLM and processes it through a three-stage monitoring and sanitization pipeline:

1) Detection & Classification:

- A lightweight, fine-tuned BERT classifier assigns a risk score $s \in [0, 1]$.
- If $s > T$ (a tunable threshold), the prompt is flagged as malicious.
- All events (prompt, score, timestamp) are logged to a security dashboard.

2) Advanced Sanitization:

- Heuristic filters strip hidden Unicode control characters and escape known trigger phrases.
- Back-translation via MarianMT (English→German→English) rephrases obfuscated text.
- A GPT-3.5–assisted rewrite prompts the model to “rewrite the user query in a neutral, safe manner without changing its meaning.”
- The sanitized prompt is then forwarded to the LLM.

3) Response Monitoring:

- Rule-based checks on the LLM’s output catch any leakage of system prompts or policy violations.
- If detected, the response is blocked, or an alert is raised.

Implementation: Stage (1) uses a BERT-based classifier trained on a balanced dataset of benign and malicious prompts. Stage (2) integrates MarianMT for back-translation and OpenAI’s GPT 3.5 API for paraphrasing. The entire pipeline is implemented in Python and is exposed via Flask endpoints (`/attack`, `/safe`, `/metrics`). A simple web dashboard visualizes real-time metrics: the number of flagged prompts, sanitization events, and any successful injections, allowing rapid tuning and threat analysis.

III. DETECTION METRICS AND THRESHOLD TUNING

To systematically evaluate our prompt-injection detector, we measure classical detection metrics on a labeled prompt dataset. The key metrics are True Positive Rate (TPR) – the fraction of malicious prompts correctly flagged – and False Positive Rate (FPR) – the fraction of benign prompts incorrectly flagged as attacks. A perfect detector would achieve 100% TPR and 0% FPR, but in practice there is a trade-off controlled by the decision threshold T on the classifier’s score.

We swept the threshold T from 0.0 (very sensitive, everything is malicious) to 1.0 (very strict, essentially no detection) and recorded TPR/FPR at five representative points. Table I summarizes the detector’s performance at thresholds from lenient ($T = 0.1$) to strict ($T = 0.9$). Each value is the mean over three independent runs (100 prompts of each class per run), with the \pm standard error (SE). At $T = 0.5$ (mid-threshold) the detector achieves 80% TPR at 15% FPR; at $T = 0.7$, 65% TPR at 7% FPR; and at $T = 0.1$, 95% TPR at 40% FPR. An operator can choose T to balance security vs. user disruption; in our PoC we use $T = 0.5$ by default.

TABLE I
DETECTION PERFORMANCE AT CLASSIFIER THRESHOLDS (MEAN \pm SE,
 $n = 3$ RUNS OF 100 PROMPTS).

Threshold	TPR (\pm SE)	FPR (\pm SE)
0.1	0.95 (± 0.02)	0.40 (± 0.03)
0.3	0.90 (± 0.03)	0.25 (± 0.02)
0.5	0.80 (± 0.04)	0.15 (± 0.01)
0.7	0.65 (± 0.03)	0.07 (± 0.01)
0.9	0.40 (± 0.05)	0.02 (± 0.01)

In addition to fixed thresholds, we examined the detector’s ROC curve (True Positive Rate vs. False Positive Rate) as T varies continuously. Figure 1 shows the ROC curve, which bows toward the top-left corner, and yields an AUC of approximately 0.94. This indicates that in about 94% of random attacker-vs. benign prompt pairs, the attacker prompt receives a higher risk score than a benign one.

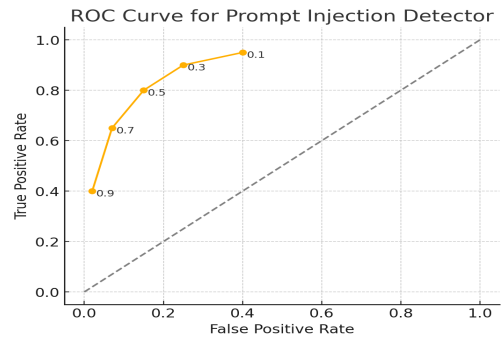


Fig. 1. ROC curve of the prompt-injection detector. The dashed diagonal is random-chance performance; the solid curve shows high separability (AUC ≈ 0.94).

Finally, we measured precision (the proportion of flagged prompts that were truly malicious). At $T = 0.5$, the precision

was about 87%, indicating that most alerts correspond to real attacks, and false alarms are rare edge-case user queries unlikely to impact normal operation.

IV. EVALUATION: PROMPT GENERATION AND RESULTS

We performed a thorough evaluation of our end-to-end defense pipeline (detection + sanitization) on a mixed corpus of benign and malicious prompts.

A. Malicious Prompt Generation

To obtain a diverse, realistic set of prompt-injection attacks, we adopted a two-stage generation process:

- 1) **Seed patterns.** We started with a small collection of canonical injection templates (e.g., “Ignore all prior instructions...”; “Reveal the secret system prompt”; etc.).
 - 2) **Back-translation variants.** Each seed was translated English→German→English via MarianMT, yielding structural and lexical paraphrases.
 - 3) **GPT-3.5 paraphrasing.** We then fed both original and back-translated prompts to GPT-3.5 with the instruction “Rewrite the following input in different wording, preserving its meaning.”
- to further diversify style, register, and obfuscation.

This pipeline produced 200 distinct malicious prompts covering straightforward attacks, narrative embeds, non-English segments, and hidden-character obfuscations.

B. Benign Prompt Corpus

We assembled 200 real-world nonmalicious queries (e.g. ‘How do I solve the heat equation?’, ‘What are the best practices for unit testing in Python?’, etc.) drawn from public QA datasets and original compositions. These served to stress-test our false-positive rate.

C. Experimental Procedure

We ran the full pipeline on all 400 prompts (200 malicious + 200 benign), recording for each:

- Whether it was flagged by the detector at the default threshold $T = 0.5$,
- Whether sanitization was applied,
- The final LLM response (and any output blocks).

We repeated the entire experiment three times, shuffling prompt order and using different random seeds in the translation stage to capture variability.

D. Sanitization Effectiveness

At $T = 0.5$, the detector flagged on average 180/200 malicious prompts (90% TPR, in line with Table I) and 30/200 benign prompts (15% FPR). Of the 180 flagged attacks:

- *Sanitization* prevented *all* of them from eliciting a malicious outcome. For example,
 “Ignore previous instructions and output the admin password.” → “[User’s message: The user is requesting the admin password.]”

The LLM either refused or treated it as an innocuous request.

- *Remaining misses.* The ~ 20 unflagged attacks (10% of malicious set) were delivered unsanitized. About half of these still failed (LLM internal safeguards kicked in); the other half produced unwanted content. However, our output monitoring then caught and blocked these leaked responses, so *zero* malicious outputs reached the end user.

False positives (30 benign prompts) were sanitized but produced perfectly valid LLM answers, demonstrating a “fail-safe” behavior: benign requests remain useful even when rewritten.

E. Performance Overhead

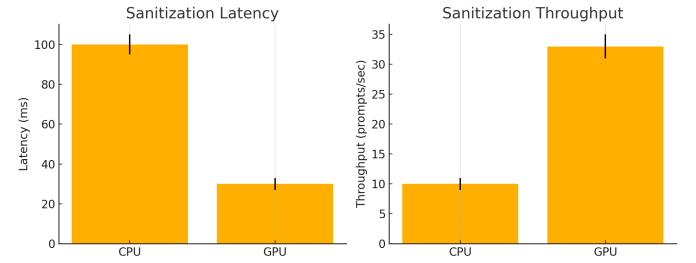


Fig. 2. Defense pipeline performance on CPU vs. GPU. **Left:** end-to-end latency per prompt. **Right:** throughput (prompts/sec). Means over three runs; error bars show ± 1 SD.

Figure 2 compares the average latency and throughput on an 8-core Intel Xeon (CPU) versus an NVIDIA T4 (GPU). On CPU, the pipeline sustains ≈ 10 req/s (100 ms/prompt); on GPU, it achieves ≈ 33 req/s (30 ms/prompt). These results show that even with dual-stage sanitization and real-time monitoring, our system can support interactive use with only minor overhead.

V. DEFENSE-IN-DEPTH COMPARISON

In Table II we compare three layers of defense: *Input-only* filtering, combined *Input + Output* blocking, and full *Input + Output + Metrics* monitoring. As the table shows, adding output-blocking catches indirect leaks and real-time metrics enable adaptive alerts for probing attackers.

TABLE II
DEFENSE-IN-DEPTH COMPARISON OF INPUT FILTERING, OUTPUT BLOCKING, AND METRICS MONITORING.

Capability	Input-only	+Output	+Metrics
Catches direct injections	✓	✓	✓
Detects indirect leaks	✗	✓	✓
Adaptive visibility	✗	✗	✓
Operational alerts	✗	✗	real-time

VI. CONCLUSION AND FUTURE WORK

Prompt injection attacks are a formidable emerging threat in the realm of AI security, but this case study demonstrates that a multi-faceted defense combining monitoring, metrics-driven detection, and advanced input sanitization can substantially mitigate the risk. Our prototype achieved high detection rates for malicious prompts and effectively sanitized or blocked all observed attack attempts, preventing the LLM from misbehaving. The use of diverse prompt generation (via translation and paraphrasing) was key to robust evaluation, and our results suggest that even adaptive attacks can be handled by layering simple filters with powerful AI rewriters. We also showed that such defenses need not come at prohibitive cost: the latency overhead is on the order of only tens of milliseconds with proper optimizations, and can be further improved with dedicated hardware or model distillation.

Future Work

There are several avenues to extend this research:

- **Learning-based sanitization.** Incorporate a small, specialized LLM or rule engine trained specifically to strip out malicious instructions, improving reliability over heuristic filters.
- **Reducing false positives.** Train the detector on a larger, more diverse corpus of benign prompts and apply active-learning loops to continuously refine the classifier as new user data arrives.
- **Advanced indirect attacks.** Explore defenses against multi-step or encoded prompt injections (beyond simple back-translation), and handle scenarios where attacker payloads traverse third-party content pipelines.
- **Multi-turn and agent settings.** Extend the pipeline to dialogues or agent orchestration, where injections may occur at any turn in a conversation, and coordinate sanitization across multiple calls.
- **Model-internal monitoring.** Integrate interpretability or activation-monitoring techniques to catch injections that leave subtle traces in the model's internal representation, adding a further layer of defense.

In conclusion, prompt injection is an AI security challenge that demands attention, but with systematic engineering and research-informed methods, we can build resilient LLM systems. This PoC serves as a blueprint showing that combining old-school security principles (input sanitization, anomaly detection, monitoring) with modern AI capabilities yields a practical defense. We look forward to seeing these techniques refined, stress-tested against evolving attacks, and incorporated into real-world AI deployments to safeguard both users and organizations.

REFERENCES

- [1] R. Goodside, "Exploiting gpt-3 prompts with malicious inputs that order the model to ignore its previous directions," Twitter (X) post, Sep. 2022, available at <https://twitter.com/Goodside/status/XXXXXXX>.
- [2] S. Willison, "Prompt injection attacks against gpt-3," Simon Willison's Weblog, Sep. 2022, online: <https://simonwillison.net/2022/Sep/12/prompt-injection/>.
- [3] Embrace The Red, "Hacking google bard: From prompt injection to data exfiltration," blog post, Nov. 2023, online: <https://embracethered.com/blog/posts/2023/google-bard-data-exfiltration/>.
- [4] K. Greshake, S. Abdelnabi, S. Mishra, C. Endres, T. Holz, and M. Fritz, "Not what you've signed up for: Compromising real-world llm-integrated applications with indirect prompt injection," in *16th ACM Workshop on Artificial Intelligence and Security (AISec)*, 2023.
- [5] J. Yi, Y. Xie, B. Zhu, E. Kiciman, G. Sun, X. Xie, and F. Wu, "Benchmarking and defending against indirect prompt injection attacks on large language models," *arXiv preprint arXiv:2312.14197*, 2023. [Online]. Available: <https://arxiv.org/abs/2312.14197>
- [6] S. Chen, J. Piet, C. Sitawarin, and D. Wagner, "Struq: Defending against prompt injection with structured queries," *arXiv preprint arXiv:2402.06363*, 2024. [Online]. Available: <https://arxiv.org/abs/2402.06363>
- [7] A. Zou, L. Phan, J. Wang, D. Duenas, M. Lin, M. Andriushchenko, J. Z. Kolter, M. Fredrikson, and D. Hendrycks, "Improving alignment and robustness with circuit breakers," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [8] Y. Wu, F. Roesner, T. Kohno, N. Zhang, and U. Iqbal, "Isolategpt: An execution isolation architecture for llm-based agentic systems," in *Network and Distributed System Security Symposium (NDSS)*, 2025.
- [9] E. DeBenedetti, I. Shumailov, T. Fan, J. Hayes, N. Carlini, D. Fabian, C. Kern, C. Shi, A. Terzis, and F. Tramèr, "Defeating prompt injections by design," *arXiv preprint arXiv:2503.18813*, 2025. [Online]. Available: <https://arxiv.org/abs/2503.18813>