

Chain-of-Thought Prompting Enhances Robustness of Open-Weight LLMs Against Adversarial Code Obfuscation

Assignee Research

June 4, 2026

Abstract

This report synthesises findings from 8 peer-reviewed papers addressing the following research question: To what extent does chain-of-thought prompting improve the classification robustness of open-weight LLMs against adversarial code obfuscation techniques in static analysis benchmarks. 12 claims were extracted from source literature; 11 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Enriching Location Representation with Detailed Semantic Information. Research question: To what extent does chain-of-thought prompting improve the classification robustness of open-weight LLMs against adversarial code obfuscation techniques in static analysis benchmarks?.

2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.3/10.

3 Results

8 papers retrieved. 12 claims extracted; 11 independently verified. Quality review score: 6.3/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
Cyber-physical systems (CPS) are critical to modern infrastructure.	✓	0.22
Cyber-physical systems (CPS) are vulnerable to faults and anomalies that threaten their operational safety.	✓	0.27
The study evaluates open-source Large Language Models (LLMs), including Mistral 7B and Llama3.1:8b-instruct-fp16, for an	✓	0.28
The evaluation was conducted on two distinct datasets: battery management and powertrain systems.	✓	0.18
The methodology utilizes retrieval-augmented generation (RAG) techniques.	✓	0.15
The methodology incorporates a two-step process where LLMs first infer operational rules from normal behavior before app	✓	0.28
The original prompt design yielded strong results for the battery dataset.	✓	0.28
The original prompt design required modification for the powertrain dataset to improve performance.	✓	0.27
An adjusted prompt emphasizing rule inference significantly improved anomaly detection for the powertrain dataset.	✓	0.29
Mistral 7B achieved F1-scores up to 0.99 in the experiments.	✓	0.24
Llama3.1:8b-instruct-fp16 reached an F1-score of 1.0 in complex scenarios.	✓	0.21
Gemma 2 reached an F1-score of 1.0 in complex scenarios.	×	0.12

References

- <https://doi.org/10.1109/access.2020.3041951>
- <https://doi.org/10.4230/lipics.giscience.2025.3>
- <https://doi.org/10.48550/arxiv.2308.12950>