

Adversarial Robustness of CodeLLMs on RoundTripCodeEval Benchmark

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 3 peer-reviewed papers addressing the following research question: What is the accuracy drop of `codellambda-7b-hf-float16` on RoundTripCodeEval when subjected to adversarial code perturbations, and how does this robustness compare to Llama-2-13b and WizardCoder-13b. 8 claims were extracted from source literature; 8 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 8.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: LLMs Caught in the Crossfire: Malware Requests and Jailbreak Challenges. Research question: What is the accuracy drop of `codellambda-7b-hf-float16` on RoundTripCodeEval when subjected to adversarial code perturbations, and how does this robustness compare to Llama-2-13b and WizardCoder-13b?.

2 Methodology

Systematic literature search across multiple databases yielded 3 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.3/10.

3 Results

3 papers retrieved. 8 claims extracted; 8 independently verified. Quality review score: 8.3/10.

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
The widespread adoption of Large Language Models (LLMs) has heightened concerns about their security, particularly their	✓	0.43
Prior research has been conducted on general security capabilities of LLMs, but their specific susceptibility to jailbre	✓	0.42
MalwareBench is a benchmark dataset containing 3,520 jailbreaking prompts for malicious code-generation, designed to eva	✓	0.38
MalwareBench is based on 320 manually crafted malicious code generation requirements, covering 11 jailbreak methods and	✓	0.41
Experiments show that mainstream LLMs exhibit limited ability to reject malicious code-generation requirements.	✓	0.35
The combination of multiple jailbreak methods further reduces the model’s security capabilities.	✓	0.30
The average rejection rate for malicious content is 60.93%, dropping to 39.92% when combined with jailbreak attack algor	✓	0.35
The code security capabilities of LLMs still pose significant challenges.	✓	0.30

References

- <https://doi.org/10.13140/rg.2.2.27046.59201>
- <https://doi.org/10.18653/v1/2025.acl-long.1350>
- <https://doi.org/10.48550/arxiv.2510.12399>