

Llama3 and Codestral Performance Gaps in Vulnerability Classification Across Programming Languages

Assignee Research

June 4, 2026

Abstract

This report synthesises findings from 2 peer-reviewed papers addressing the following research question: How does the performance gap between Llama3 and Codestral in vulnerability classification (F1-score) vary when evaluated on Big-Vul samples with different programming languages (e.g., C vs. Java). 10 claims were extracted from source literature; 6 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Can Open Large Language Models Catch Vulnerabilities?. Research question: How does the performance gap between Llama3 and Codestral in vulnerability classification (F1-score) vary when evaluated on Big-Vul samples with different programming languages (e.g., C vs. Java) under uniform obfuscation levels?.

2 Methodology

Systematic literature search across multiple databases yielded 2 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.7/10.

3 Results

2 papers retrieved. 10 claims extracted; 6 independently verified. Quality review score: 6.7/10.

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 study evaluates three LLMs: Llama3, Codestral, and Deepseek R1.	×	0.15
The evaluation uses a carefully filtered subset of the Big-Vul dataset.	✓	0.16
The dataset subset is annotated with eight representative Common Weakness Enumeration (CWE) categories.	✓	0.17
The study adopts a closed-world classification setup.	×	0.12
The evaluated models demonstrate high detection rates for the presence of vulnerabilities.	×	0.14
The evaluated models demonstrate markedly poor classification accuracy when mapping vulnerabilities to correct CWE label	✓	0.15
The models exhibit frequent overgeneralization and misclassification of vulnerabilities.	×	0.12
The study analyzes model-specific biases and common failure modes.	✓	0.17
Current LLMs have limitations in performing fine-grained security reasoning.	✓	0.21
LLMs are being adopted as learning aids in educational contexts.	✓	0.18

References

- <https://openalex.org/W7124227854>
- <https://doi.org/10.4230/oasics.icpec.2025.4>