

# Robustness Comparison of Llama3.1 and Mistral 7B in Adversarial Code Vulnerability Detection

Assignee Research

June 4, 2026

## Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does the robustness of Llama3.1 compare to Mistral 7B in detecting code vulnerabilities when subjected to adversarial syntax perturbations. 4 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 8.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Securing Systems and Data: Attack Detection techniques and Generative AI approaches. Research question: How does the robustness of Llama3.1 compare to Mistral 7B in detecting code vulnerabilities when subjected to adversarial syntax perturbations?.

## 2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.7/10.

## 3 Results

12 papers retrieved. 4 claims extracted; 4 independently verified. Quality review score: 8.7/10.

## 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 CYBERRAG system architecture involves a user interacting with a chatbot connected to an agent, while an IDS detects	✓	0.25
BERT's classification performance on different web vulnerabilities improves with attack-specific training.	✓	0.23
LLM-based scoring shows a consistent advantage from RAG-enhanced generation when comparing explanations generated with a	✓	0.32
The model's robustness is evaluated based on the percentage of correct classifications under two conditions: Adversarial	✓	0.30

## References

- [https://doi.org/10.13118/blefari-francesco\\_phd2025-12-04](https://doi.org/10.13118/blefari-francesco_phd2025-12-04)
- <https://doi.org/10.48550/arxiv.2504.14985>
- <https://openalex.org/W7127541240>