

# Discrepancy in Llama-3.1-8B Long-Context Reasoning Accuracy Across Evaluation Frameworks

Assignee Research

June 11, 2026

## Abstract

As Large Language Models (LLMs) become increasingly integrated into secure software development workflows, a critical question remains unanswered: can these models not only detect insecure code but also reliably classify vulnerabilities according to standardized taxonomies? In this work, we conduct a systematic evaluation of three state-of-the-art LLMs - Llama3, Codestral, and Deepseek R1 - using a carefully filtered subset of the Big-Vul dataset annotated with eight representative Common Weakness Enumeration categories. Adopting a closed-world classification setup, we assess each model's perf

## 1 Introduction

This paper examines: Can Open Large Language Models Catch Vulnerabilities?. Research question: What is the discrepancy in long-context reasoning accuracy for Llama-3.1-8B across different evaluation frameworks when testing needle-in-a-haystack scenarios?.

## 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 8.1/10.

## 3 Results

12 papers retrieved. 8 claims extracted; 7 independently verified. Quality review score: 8.1/10.

## 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
Three state-of-the-art LLMs - Llama3, Codestral, and Deepseek R1 - were evaluated using a subset of the Big-Vul dataset	✓	0.33
The evaluation adopted a closed-world classification setup to assess each model's performance in identifying vulnerabilities	✓	0.30
The findings revealed a sharp contrast between high detection rates and markedly poor classification accuracy among the	✓	0.23
Frequent overgeneralization and misclassification were observed in the LLMs' performance.	×	0.11
Model-specific biases and common failure modes were analyzed, highlighting limitations in current LLMs' fine-grained security	✓	0.26
The insights are particularly relevant in educational contexts where LLMs are being adopted as learning aids despite the	✓	0.23
A nuanced understanding of LLMs' behavior is essential to prevent the propagation of misconceptions among students.	✓	0.18
The results expose key challenges that must be addressed before LLMs can be reliably deployed in security-sensitive environments	✓	0.29

## References

- <https://doi.org/10.48550/arxiv.2406.11230>
- <https://doi.org/10.48550/arxiv.2407.11963>
- <https://doi.org/10.4230/oasics.icpec.2025.4>