

# Llama3, Codestral, and DeepSeek-R1 vs. Specialized Models in CWE-Specific Vulnerability Classification

Assignee Research

May 31, 2026

## Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does the performance of Llama3, Codestral, and Deepseek R1 on vulnerability classification in Big-Vul compare to specialized vulnerability detection models like GitHub CodeQL in terms of. Modern software relies on a multitude of automated testing and quality assurance tools to prevent errors, bugs and potential vulnerabilities. This study sets out to provide a head-to-head, quantitative and qualitative evaluation of six automated approaches: three. 13 claims were extracted from source literature; 7 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.9/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Large Language Models Versus Static Code Analysis Tools: A Systematic Benchmark for Vulnerability Detection. Research question: How does the performance of Llama3, Codestral, and Deepseek R1 on vulnerability classification in Big-Vul compare to specialized vulnerability detection models like GitHub CodeQL in terms of precision-recall tradeoffs across different CWE categories?.

## 2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.9/10.

### **3 Results**

13 papers retrieved. 13 claims extracted; 7 independently verified. Quality review score: 6.9/10.

### **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 six automated approaches: three industry-standard rule-based static code-analysis tools (SonarQube,	✓	0.38
The benchmark uses a curated suite of ten real-world C# projects that embed 63 vulnerabilities across common categories	✓	0.30
Language-based scanners achieve higher mean F-1 scores: 0.797, 0.753, and 0.750, compared to static counterparts which s	✓	0.30
LLMs’ advantage originates from superior recall, confirming an ability to reason across broader code contexts.	✓	0.25
DeepSeek V3 exhibits the highest false-positive ratio among the language models.	✓	0.25
All language models mislocate issues at line-or-column granularity due to tokenisation artefacts.	✓	0.24
Language models successfully rival traditional static analysers in finding real vulnerabilities.	✓	0.26
The open benchmark and JSON-based result harness released with this paper lay a foundation for reproducible, practitione	×	0.06
ProjectAnalyzer Code is available on GitHub: <a href="https://github.com/Damian0401/ProjectAnalyzer">https://github.com/Damian0401/ProjectAnalyzer</a>	×	0.06
Static Application Security Testing (SAST) tools such as SonarQube, CodeQL, and Snyk-Code are trusted by more than seven	×	0.06
Recent empirical studies reveal a persistent gap, while a single state-of-the-art analyser can highlight vulnerabilities	×	0.04
The total number of identified vulnerabilities was determined by counting the number of entries present in the Results s	×	0.06
Each reported vulnerability was individually reviewed and compared with a reference dataset of known vulnerabilities for	×	0.08

## References

- <http://arxiv.org/abs/2601.08691v1>

- <http://arxiv.org/abs/2508.04448v1>
- <http://arxiv.org/abs/2504.16584v1>