

Scaling Model Parameters from 7B to 32B and Zero-Shot Performance on CLUE Benchmark

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: How does scaling model parameters from 7B to 32B affect zero-shot performance on the CLUE benchmark compared to few-shot settings. 13 claims were extracted from source literature; 6 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.1/10. 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 scaling model parameters from 7B to 32B affect zero-shot performance on the CLUE benchmark compared to few-shot settings?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.1/10.

3 Results

15 papers retrieved. 13 claims extracted; 6 independently verified. Quality review score: 6.1/10.

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 |
|---|----------|------------|
| The study evaluates three LLMs: Llama3, Codestral, and Deepseek R1. | × | 0.14 |
| The evaluation uses a carefully filtered subset of the Big-Vul dataset. | ✓ | 0.17 |
| The dataset subset is annotated with eight representative Common Weakness Enumeration (CWE) categories. | ✓ | 0.18 |
| The study adopts a closed-world classification setup. | × | 0.13 |
| The models were assessed on their ability to identify the presence of vulnerabilities. | × | 0.07 |
| The models were assessed on their ability to map vulnerabilities to the correct CWE label. | × | 0.10 |
| The evaluated models demonstrated high detection rates for vulnerabilities. | × | 0.11 |
| The evaluated models demonstrated markedly poor classification accuracy for CWE labels. | × | 0.13 |
| The models exhibited frequent overgeneralization and misclassification of vulnerabilities. | × | 0.12 |
| The study analyzes model-specific biases and common failure modes. | ✓ | 0.18 |
| Current LLMs have limitations in performing fine-grained security reasoning. | ✓ | 0.20 |
| LLMs are being adopted as learning aids in educational contexts. | ✓ | 0.17 |
| Key challenges must be addressed before LLMs can be reliably deployed in security-sensitive environments. | ✓ | 0.27 |

References

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

- <https://doi.org/10.1038/s41467-024-45563-x>
- <https://doi.org/10.48550/arxiv.2404.00287>