

# Open-Weight LLMs Latency-Accuracy Trade-offs in Obfuscated C-C++ Vulnerability Classification

Assignee Research

June 4, 2026

## Abstract

This report synthesises findings from 8 peer-reviewed papers addressing the following research question: What is the correlation between inference latency and vulnerability classification accuracy for open-weight LLMs processing obfuscated C/C++ code. 7 claims were extracted from source literature; 6 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 7.4/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Enriching Location Representation with Detailed Semantic Information. Research question: What is the correlation between inference latency and vulnerability classification accuracy for open-weight LLMs processing obfuscated C/C++ code?.

## 2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.4/10.

## 3 Results

8 papers retrieved. 7 claims extracted; 6 independently verified. Quality review score: 7.4/10.

## 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 open-source Large Language Models (LLMs), including Mistral 7B and Llama3.1:8b-instruct-fp16, for an	✓	0.31
The methodology utilizes retrieval-augmented generation (RAG) techniques incorporating a two-step process where LLMs fir	✓	0.33
The original prompt design yielded strong results for the battery dataset but required modification to improve performan	✓	0.33
An adjusted prompt emphasizing rule inference significantly improved anomaly detection performance for the powertrain da	✓	0.26
Mistral 7B achieved F1-scores up to 0.99 in the experiments.	✓	0.24
Llama3.1:8b-instruct-fp16 reached an F1-score of 1.0 in complex scenarios.	✓	0.21
Gemma 2 reached an F1-score of 1.0 in complex scenarios.	×	0.11

## References

- <https://doi.org/10.1186/s42400-025-00361-w>
- <https://doi.org/10.4230/lipics.giscience.2025.3>
- <https://doi.org/10.48550/arxiv.2403.05530>