

# Multimodal Data Integration for Enhancing SLM Performance in CWE Detection Beyond Code-Only Fine-Tuning

Assignee Research

June 12, 2026

## Abstract

Large Language Models (LLMs) have demonstrated significant capabilities in understanding and analyzing code for security vulnerabilities, such as Common Weakness Enumerations (CWEs). However, their reliance on cloud infrastructure and substantial computational requirements pose challenges for analyzing sensitive or proprietary codebases due to privacy concerns and inference costs. This work explores the potential of Small Language Models (SLMs) as a viable alternative for accurate, on-premise vulnerability detection. We investigated whether a 350-million parameter pre-trained code model (codeg

## 1 Introduction

This paper examines: Case Study: Fine-tuning Small Language Models for Accurate and Private CWE Detection in Python Code. Research question: What is the impact of multimodal data (e.g., combining code with natural language vulnerability descriptions) on the performance of SLMs for CWE detection compared to code-only fine-tuning approaches?.

## 2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.9/10.

## 3 Results

16 papers retrieved. 16 claims extracted; 16 independently verified. Quality review score: 8.9/10.

## 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
Farasat & Posegga [15] used a BiLSTM model for Python vulnerability detection, achieving an accuracy of 98.6%, F1 score	✓	0.16
Bagheri et al. [31] reported an F1 score of 99% using a Hybrid ML approach (Self-attention + CNN - Conformer).	✓	0.22
Singh et al. [32] used Logistic Regression for CWE prediction, achieving an accuracy of 0.66, precision of 0.65, recall	✓	0.16
Dozono et al. [6] reported Python F1 scores for various LLMs: GPT-4o (0.80), GPT-4T (0.76), Gemini 1.5 Pro (0.75), CodeL	✓	0.32
Steenhoek et al. [34] evaluated LLM prompting strategies on C/C++ code and found a low Balanced Accuracy of 54.5%.	✓	0.22
Shestov et al. [35] achieved a best F1 score of 0.86 for binary classification of JAVA CWE detection using the WizardCod	✓	0.22
Li et al. [36] reported a 5-6% improvement over the base model using LoRa and IA3 fine-tuning approaches on a 2.7b Codeg	✓	0.20
Jiang et al. [37] achieved a best F1 score of 87% using Llama 2-7b model with LoRa fine-tuning.	✓	0.22
The un-tuned codegen-mono model failed to detect a single CWE in a zero-shot evaluation of 100 randomly selected example	✓	0.26
The study selected the MITRE Top 25 Most Dangerous Software Weaknesses list as the target vulnerabilities for detection.	✓	0.19
The dataset was generated using Google’s gemini-2.0-flash-thinking-exp-01-21 model via its API.	✓	0.19
For each of the 25 selected CWEs, the Gemini model was instructed to generate five distinct Python code snippets demonst	✓	0.23
The fine-tuned codegen-mono model (350m parameters) achieved an accuracy of 99% on a held-out test set of 100 instances.	✓	0.28
The fine-tuned codegen-mono model achieved a precision of approximately 98.08% on the CWE detection task.	✓	0.32
The fine-tuned codegen-mono model achieved a recall of 100% on the CWE detection task.	✓	0.29
The fine-tuned codegen-mono model achieved an F1-Score of approximately 99.04%.	✓	0.32

## References

- <http://arxiv.org/abs/2504.16584v1>
- <http://arxiv.org/abs/2509.17337v1>
- <http://arxiv.org/abs/2604.26313v1>