

# Scaling Inference Efficiency of Small Language Models for Code Weakness Detection

Assignee Research

May 30, 2026

## Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: How does the inference efficiency (throughput, latency) of SLMs trained for CWE detection scale with model size when benchmarked on a private codebase, and how does this compare to larger models. Abstract Data scarcity is a major challenge when training deep learning (DL) models. DL demands a large amount of data to achieve exceptional performance. 17 claims were extracted from source literature; 16 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 7.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: A survey on deep learning tools dealing with data scarcity: definitions, challenges, solutions, tips, and applications. Research question: How does the inference efficiency (throughput, latency) of SLMs trained for CWE detection scale with model size when benchmarked on a private codebase, and how does this compare to larger models?.

## 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 7.8/10.

## 3 Results

16 papers retrieved. 17 claims extracted; 16 independently verified. Quality review score: 7.8/10.

## 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
Data scarcity is a major challenge when training deep learning (DL) models.	✓	0.30
Deep learning demands a large amount of data to achieve exceptional performance.	✓	0.20
Many applications have small or inadequate data to train deep learning frameworks.	✓	0.19
Manual labeling typically involves human annotators with a vast background of knowledge.	✓	0.21
The manual annotation process is costly, time-consuming, and error-prone.	✓	0.20
Every deep learning framework is fed by a significant amount of labeled data to automatically learn representations.	✓	0.25
A larger amount of data generates a better deep learning model.	×	0.11
Deep learning model performance is application dependent.	✓	0.18
Data scarcity is the main barrier for many applications dismissing the use of deep learning.	✓	0.22
Having sufficient data is the first step toward any successful and trustworthy deep learning application.	✓	0.21
The paper presents a survey on techniques to overcome challenges including small datasets, imbalanced datasets, and lack	✓	0.20
Transfer Learning (TL) is a state-of-the-art solution to address the issue of lack of training data.	✓	0.20
Self-Supervised Learning (SSL) is a state-of-the-art solution to address the issue of lack of training data.	✓	0.22
Generative Adversarial Networks (GANs) are a state-of-the-art solution to address the issue of lack of training data.	✓	0.21
Model Architecture (MA) is a state-of-the-art solution to address the issue of lack of training data.	✓	0.20
Physics-Informed Neural Network (PINN) is a state-of-the-art solution to address the issue of lack of training data.	✓	0.23
Deep Synthetic Minority Oversampling Technique (DeepSMOTE) is a state-of-the-art solution to address the issue of lack o	✓	0.25

## References

- <https://doi.org/10.3390/computation13020030>
- <https://doi.org/10.1186/s40537-023-00727-2>
- <https://doi.org/10.1038/s43856-024-00717-2>