

Fine-Tuning Phi-4 on Synthetic Visual-Math Data Enhances Out-of-Distribution Robustness

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 9 peer-reviewed papers addressing the following research question: Can fine-tuning Phi-4 on additional synthetic visual-math datasets improve its robustness on out-of-distribution GSM8K-V problems. 10 claims were extracted from source literature; 9 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.6/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Large Language Models: A Survey. Research question: Can fine-tuning Phi-4 on additional synthetic visual-math datasets improve its robustness on out-of-distribution GSM8K-V problems?.

2 Methodology

Systematic literature search across multiple databases yielded 9 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.6/10.

3 Results

9 papers retrieved. 10 claims extracted; 9 independently verified. Quality review score: 6.6/10.

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
Large Language Models (LLMs) have drawn a lot of attention due to their strong performance on a wide range of natural language tasks.	✓	0.40
LLMs' ability of general-purpose language understanding and generation is acquired by training billions of model's parameters.	✓	0.39
The research area of LLMs, while very recent, is evolving rapidly in many different ways.	✓	0.28
This paper reviews some of the most prominent LLMs, including three popular LLM families (GPT, LLaMA, PaLM).	✓	0.29
The paper discusses the characteristics, contributions, and limitations of these LLMs.	×	0.13
The paper gives an overview of techniques developed to build and augment LLMs.	✓	0.21
The paper surveys popular datasets prepared for LLM training, fine-tuning, and evaluation.	✓	0.26
The paper reviews widely used LLM evaluation metrics.	✓	0.20
The paper compares the performance of several popular LLMs on a set of representative benchmarks.	✓	0.21
The paper concludes by discussing open challenges and future research directions.	✓	0.20

References

- <https://doi.org/10.48550/arxiv.2402.06196>
- <https://doi.org/10.48550/arxiv.2412.15115>
- <https://doi.org/10.48550/arxiv.2308.12950>