

Instruction Fine-Tuning Boosts Language Model Mathematical Problem-Solving Accuracy

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: What is the effect of instruction fine-tuning on language model mathematical problem-solving accuracy v10. 9 claims were extracted from source literature; 4 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: When Reasoning Beats Scale: A 1.5B Reasoning Model Outranks 13B LLMs as Discriminator. Research question: What is the effect of instruction fine-tuning on language model mathematical problem-solving accuracy v10.

2 Methodology

Systematic literature search across multiple databases yielded 12 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

12 papers retrieved. 9 claims extracted; 4 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
A 1.5B distilled DeepSeek-R1 model achieves 87% higher F1 score than CodeLlama-7B in SQL query discrimination.	✓	0.21
A 1.5B distilled DeepSeek-R1 model achieves 3.7% better discrimination accuracy than CodeLlama-7B.	✓	0.24
A 1.5B distilled DeepSeek-R1 model achieves 3.7% higher execution accuracy than CodeLlama-13B.	✓	0.25
Using logit-based soft scoring versus binary true/false discrimination yields performance differences of less than 1.5%.	×	0.03
Increasing the compute budget beyond 1024 tokens yields less than 0.4% performance gain for reasoning models.	×	0.10
Using extremely low compute budgets results in less than 2% accuracy and greater than 94% failure rate for reasoning mod	×	0.08
DeepSeek-R1 underperforms as a generator compared to smaller non-reasoning LLMs.	✓	0.26
The study evaluates LLMs' ability to discriminate correct and incorrect SQL queries by re-labeling oracle-generated outp	×	0.04
Discriminator performance was tested in a naive setting and an enhanced setting that filters by executability via enviro	×	0.02

References

- <http://arxiv.org/abs/2505.03786v1>
- <http://arxiv.org/abs/2310.04793v2>

- <http://arxiv.org/abs/2312.10793v3>