

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

Assignee Research

June 4, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: What is the effect of instruction fine-tuning on language model mathematical problem-solving accuracy v5. 15 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Improving Large Language Model Fine-tuning for Solving Math Problems. Research question: What is the effect of instruction fine-tuning on language model mathematical problem-solving accuracy v5.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.7/10.

3 Results

13 papers retrieved. 15 claims extracted; 0 independently verified. Quality review score: 3.7/10.

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 experiments use the MATH dataset with 4,500 original test examples for training and validation, and 500 test example | × | 0.04 |
| Two sources of correct step-by-step solutions are used for training: original human-written explanations from the MATH d | × | 0.12 |
| The PRM800K dataset only covers a subset of problems in the original MATH dataset. | × | 0.05 |
| Solution correctness is evaluated using an automatic grading script that checks mathematical equivalence rather than tex | × | 0.03 |
| Pass@1 performance is evaluated using greedy decoding. | × | 0.07 |
| Majority voting performance (Maj1@N) is evaluated using nucleus sampling with a temperature of 0.6 and a top-p value of | × | 0.05 |
| PaLM 2-S* and PaLM 2-L models are fine-tuned using the Maximum Likelihood Estimation (MLE) training objective. | × | 0.06 |
| Three fine-tuning strategies were explored: using only original MATH solutions, using only PRM800K GPT-4 solutions, and | × | 0.10 |
| All fine-tuned models achieved their best performance within two epochs. | × | 0.09 |
| Few-shot baseline results were obtained using a customized 4-shot prompt designed in Lewkowycz et al. (2022). | × | 0.06 |
| Models fine-tuned on PRM800K solutions achieved significantly better performance than those fine-tuned on the original M | × | 0.08 |
| Original solutions in the MATH dataset are characterized as more abstract, while GPT-4 generated solutions are more fine | × | 0.07 |
| In the training process, the step-by-step solution (S) and the final answer (A) are concatenated into a single text sequ | × | 0.08 |
| Two re-ranking strategies are compared: re-ranking all candidate solutions (RR.All) and re-ranking solutions in the top- | × | 0.12 |
| Two loss functions are compared for solution-cluster re-ranking: Lcls-margin and Lcls-xent. | × | 0.09 |

References

- <http://arxiv.org/abs/2310.10047v1>
- <http://arxiv.org/abs/2509.25160v1>
- <http://arxiv.org/abs/2312.10793v3>