

Impact of Concept Graph Depth in MathScale on SMBI Benchmark Performance

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 11 peer-reviewed papers addressing the following research question: To what extent does the depth of the concept graph in MathScale influence its performance on the SMBI benchmark compared to shallow concept graphs or no structured knowledge at all. 17 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.6/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: MathScale: Scaling Instruction Tuning for Mathematical Reasoning. Research question: To what extent does the depth of the concept graph in MathScale influence its performance on the SMBI benchmark compared to shallow concept graphs or no structured knowledge at all?.

2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.6/10.

3 Results

11 papers retrieved. 17 claims extracted; 0 independently verified. Quality review score: 3.6/10.

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
The MWPBENCH training set comprises around 20K questions.	×	0.03
GPT-3.5-Turbo-0613 was used for concept extraction.	×	0.04
The concept extraction process obtained 2,018 topics and 8,892 knowledge points.	×	0.05
The edge weight in the concept graph is smoothed using Equation (1) with $\epsilon = 1e-5$.	×	0.04
The concept composition process was repeated for approximately 1K epochs, resulting in 2 million unique concept composit	×	0.03
GPT-3.5-Turbo-0613 was instructed to create 2 million question-answer pairs with these compositions.	×	0.10
The generated datasets were decontaminated by excluding all math questions in the test set of MWPBENCH.	×	0.08
The MathScaleQA dataset was created by combining the generated data with the training set of MWPBENCH.	×	0.05
The validation step (Section 3.4) was excluded from the final pipeline because it did not improve results.	×	0.01
MathScale-7B achieves a 35.0% (micro) and 37.5% (macro) accuracy across MWPBENCH.	×	0.12
MathScale-7B surpasses its best counterparts of equivalent size by 42.9% and 43.7%, respectively.	×	0.13
MathScale-Mistral demonstrates performance parity in both micro and macro averages relative to GPT-3.5-Turbo.	×	0.06
When scaling the size of the MathScaleQA dataset, a nearly logarithmic growth in the performance of the MathScale-7b mod	×	0.10
GSM8K and MATH each contain around 7.5K training examples.	×	0.06
WizardMath introduces an array of operations for GPT-3.5 to generate math questions with increased complexity.	×	0.08
MetaMath bootstraps questions in GSM8K and MATH through answer augmentation, question rephrasing, self-verification, and	×	0.07
The newly generated examples by WizardMath and MetaMath exhibit substantial similarity to the original examples containe	×	0.01

References

- <http://arxiv.org/abs/2403.02884v1>
- <http://arxiv.org/abs/2511.11017v1>
- <http://arxiv.org/abs/2606.04326v1>