

# Frontier Language Model Failures in Abstract Mathematical Reasoning

Assignee Research

June 6, 2026

## **Abstract**

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: What are the failure modes of frontier language models on abstract mathematical reasoning v11. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

## **1 Introduction**

This paper examines: Solving Challenging Math Word Problems Using GPT-4 Code Interpreter with Code-based Self-Verification. Research question: What are the failure modes of frontier language models on abstract mathematical reasoning v11.

## **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 3.8/10.

## **3 Results**

13 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 3.8/10.

## **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.

## References

- <http://arxiv.org/abs/2601.01982v1>
- <http://arxiv.org/abs/2310.03731v1>
- <http://arxiv.org/abs/2308.07921v1>