

Language Models and Multi-Hop Reasoning in Scientific Question Answering

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: How do language models handle multi-hop reasoning chains in scientific question answering v17. 12 claims were extracted from source literature; 12 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 8.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Can Small Agent Collaboration Beat a Single Big LLM?. Research question: How do language models handle multi-hop reasoning chains in scientific question answering v17.

2 Methodology

Systematic literature search across multiple databases yielded 4 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.7/10.

3 Results

4 papers retrieved. 12 claims extracted; 12 independently verified. Quality review score: 8.7/10.

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
Recent progress in language modeling has largely relied on scaling model size.	✓	0.27
Larger models do not reliably improve performance on tasks requiring multi-step reasoning and tool use.	✓	0.33
Multi-agent collaboration offers a potential alternative to scaling model size.	✓	0.28
The paper addresses whether well-organized systems built from smaller models can outperform much larger language models.	✓	0.24
The multi-agent system used in the paper has a single orchestrator and a small set of specialized sub-agents with restri	✓	0.30
The benchmarks used in the paper span factual retrieval, multi-hop reasoning, scientific question answering, and mathema	✓	0.23
The paper conducts controlled comparisons between small multi-agent systems and large single-agent models.	✓	0.31
Small multi-agent systems can outperform substantially larger single-agent models, even when the latter have direct acce	✓	0.40
Reasoning at the orchestrator yields the largest gains in performance.	✓	0.20
Enabling reasoning in sub-agents provides limited or negative benefits.	✓	0.27
Overall system performance is driven primarily by orchestrator capacity rather than sub-agent capacity.	✓	0.29
Improved agentic performance depends more on architectural orchestration than on raw model scaling.	✓	0.27

References

- <https://doi.org/10.48550/arxiv.2511.16283>

- <https://openalex.org/W7124817582>
- <https://openalex.org/W7106475840>