

Tree of Reviews vs. Chain-Based Retrieval in F1 Stability for Llama-3-8B-128K Context Scaling

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How does the Tree of Reviews framework compare to the chain-based retrieval method in terms of F1 score stability when scaling Llama-3-8B-128K’s context length from 4K to 128K on the MuSiQue benchmark. Multi-hop question answering is a knowledge-intensive complex problem. Large Language Models (LLMs) use their Chain of Thoughts (CoT) capability to reason complex problems step by step, and retrieval-augmentation can effectively alleviate factual errors caused by outdated and. 5 claims were extracted from source literature; 3 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 7.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Tree of Reviews: A Tree-based Dynamic Iterative Retrieval Framework for Multi-hop Question Answering. Research question: How does the Tree of Reviews framework compare to the chain-based retrieval method in terms of F1 score stability when scaling Llama-3-8B-128K’s context length from 4K to 128K on the MuSiQue benchmark?.

2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.0/10.

3 Results

14 papers retrieved. 5 claims extracted; 3 independently verified. Quality review score: 7.0/10.

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
Tree of Reviews (TOR) is a dynamic retrieval framework that integrates a tree structure into the iterative retrieval pro	✓	0.21
The TOR framework mitigates the negative impact associated with the inherent vulnerabilities of chain-like retrieval met	×	0.07
The proposed framework dynamically decides whether to initiate a new search, reject, or accept based on paragraphs on th	✓	0.29
Two tree-based search optimization strategies, pruning and effective expansion, are proposed to reduce time overhead and	×	0.08
Experiments conducted on three different multi-hop question answering datasets show that TOR achieves state-of-the-art p	✓	0.31

References

- <http://arxiv.org/abs/2510.22344v1>
- <http://arxiv.org/abs/2404.14464v1>
- <http://arxiv.org/abs/2407.13739v1>