

# Context-Aware Chunking and Exact Match Performance in Multi-Hop QA with Extended Attention

Assignee Research

June 1, 2026

## Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: Does context-aware chunking improve answer exact match scores on multi-hop QA datasets compared to fixed-size segmentation for transformer models with extended attention spans. 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to reason complex problems step by step, and retrieval-augmentation can effectively alleviate factual errors caused by outdated and. 8 claims were extracted from source literature; 3 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.8/10. 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: Does context-aware chunking improve answer exact match scores on multi-hop QA datasets compared to fixed-size segmentation for transformer models with extended attention spans?.

## 2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.8/10.

### 3 Results

15 papers retrieved. 8 claims extracted; 3 independently verified. Quality review score: 5.8/10.

### 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.

### 5 Extracted Claims

Claim	Verified	Confidence
TOR achieves state-of-the-art performance in both retrieval and response generation on three different multi-hop questions	✓	0.32
Tree of Thought (ToT) enhances the problem-solving capabilities of Large Language Models (LLMs) by introducing a tree-like structure	×	0.08
The tree is an efficient structure for solving complex reasoning problems.	×	0.09
TOR is the first retrieval framework that uses a tree-like structure to dynamically initiate requests based on external information	×	0.13
TOR introduces a tree structure to handle each retrieved paragraph separately, alleviating the misleading effect of irrelevant information	✓	0.31
The diversity of reasoning path extension in TOR reduces the impact of a single reasoning error on the whole.	✓	0.22
TOR proposes two tree-based search optimization strategies: pruning and effective expansion.	×	0.06
Pruning and effective expansion strategies in TOR demonstrate significant improvements in reducing time overhead and increasing accuracy	×	0.04

### References

- <http://arxiv.org/abs/2208.10297v1>

- <http://arxiv.org/abs/2504.19754v1>
- <http://arxiv.org/abs/2404.14464v1>