

# SOVEREIGN: How does the Tree of Reviews framework’s F1 score on the MuSiQue benchmark vary when evaluated with Llama-3-8B

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 29, 2026

## Abstract

Large Language Models (LLMs) showcase impressive capabilities but encounter challenges like hallucination, outdated knowledge, and non-transparent, untraceable reasoning processes. Retrieval-Augmented Generation (RAG) has emerged as a promising solution by incorporating knowledge from external databases. This enhances the accuracy and credibility of the generation, particularly for knowledge-intensive tasks, and allows for continuous knowledge updates and integration of domain-specific information. RAG synergistically merges LLMs’ intrinsic knowledge with the vast, dynamic repositories of exte

## 1 Introduction

Analysis of: Retrieval-Augmented Generation for Large Language Models: A Survey. Research goal: How does the Tree of Reviews framework’s F1 score on the MuSiQue benchmark vary when evaluated with Llama-3-8B-128K under context lengths of 32K, 64K, and 128K, compared to the chain-based retrieval method’s F1 at each length?.

## 2 Methodology

Multi-query arXiv search (4 parallel queries, Relevance-sorted). TF-IDF cosine semantic verification (bigrams, threshold=0.15). NIM nv-embedqa-e5-v5 (dim=1024) for semantic indexing. Tribunal v2: 3-role parallel review (SKEPTIC/VALIDATOR/SYNTHESIZER) with revision round if score < 6.5.

### 3 Results

8 papers retrieved. 10 claims extracted, 10 verified. Tribunal: 8.8/10 → APPROVE (revision\_round=0). Policy: AUTO\_APPROVE.

### 4 Uncertainties

NIM free tier latency varies. TF-IDF verification is a weak signal. arXiv Relevance ranking is query-dependent. Tribunal consensus is LLM-based and prompt-sensitive.

### 5 Extracted Claims

Claim	Verified	Confidence
Large Language Models (LLMs) encounter challenges like hallucination, outdated knowledge, and non-transparent, untraceable	✓	0.35
Retrieval-Augmented Generation (RAG) has emerged as a promising solution by incorporating knowledge from external databases	✓	0.36
RAG enhances the accuracy and credibility of the generation, particularly for knowledge-intensive tasks.	✓	0.26
RAG allows for continuous knowledge updates and integration of domain-specific information.	✓	0.26
RAG synergistically merges LLMs' intrinsic knowledge with the vast, dynamic repositories of external databases.	✓	0.33
This comprehensive review paper offers a detailed examination of the progression of RAG paradigms, encompassing the Naive	✓	0.37
The paper meticulously scrutinizes the tripartite foundation of RAG frameworks, which includes the retrieval, the generation	✓	0.31
The paper highlights the state-of-the-art technologies embedded in each of these critical components, providing a profound	✓	0.34
This paper introduces up-to-date evaluation framework and benchmark.	✓	0.21
This article delineates the challenges currently faced and points out prospective avenues for research and development.	✓	0.29

## References

- <https://doi.org/10.48550/arxiv.2402.07927>
- <https://doi.org/10.48550/arxiv.2403.05530>
- <https://doi.org/10.48550/arxiv.2312.10997>