

Phi-3-Mini and Mistral-7B Hallucination Rates in Long-Context Religious RAG Benchmarks

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How do Phi-3-mini and Mistral-7B-v0.1 compare in hallucination rates on long-context RAG benchmarks for specialized religious domains. 8 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.1/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Hallucination Detection with Small Language Models. Research question: How do Phi-3-mini and Mistral-7B-v0.1 compare in hallucination rates on long-context RAG benchmarks for specialized religious domains?.

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 3.1/10.

3 Results

13 papers retrieved. 8 claims extracted; 0 independently verified. Quality review score: 3.1/10.

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
Small language models (SLMs) can produce accurate results for specific tasks.	×	0.11
The proposed framework improves verification accuracy by 10% over the baseline.	×	0.07
The proposed method is superior to both P(yes) and ChatGPT in detecting correct responses from partial responses.	×	0.15
The 'max' method achieves the highest F1 score of 0.99 in Fig. 5 (a).	×	0.04
The 'harmonic' method achieves the highest F1 score of 0.81 in Fig. 5 (b).	×	0.04
The geometric mean is calculated using the formula: $si(S, m = 'geo') = \exp(1/ S(ri) * \sum \log(s_{i,j}))$, where $s_{i,j} > 0$.	×	0.01
The minimum value in the dataset is found using the formula: $si(S, m = 'min') = \min(S(ri))$.	×	0.02
The maximum value in the dataset is found using the formula: $si(S, m = 'max') = \max(S(ri))$.	×	0.02

References

- <http://arxiv.org/abs/2402.12317v2>
- <http://arxiv.org/abs/2506.22486v1>
- <http://arxiv.org/abs/2409.03708v2>