

Optimizing Sequence Composition for Causal Masking in Zero-Shot Logical Deduction Tasks

Assignee Research

June 9, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: Does optimizing sequence composition for causal masking improve zero-shot performance on the Big-Bench Hard logical deduction tasks relative to standard concatenation strategies. 9 claims were extracted from source literature; 7 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Enriching Location Representation with Detailed Semantic Information. Research question: Does optimizing sequence composition for causal masking improve zero-shot performance on the Big-Bench Hard logical deduction tasks relative to standard concatenation strategies?.

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 6.8/10.

3 Results

13 papers retrieved. 9 claims extracted; 7 independently verified. Quality review score: 6.8/10.

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
The study evaluates open-source Large Language Models (LLMs), including Mistral 7B and Llama3.1:8b-instruct-fp16, for an	✓	0.31
The methodology employs retrieval-augmented generation (RAG) techniques.	×	0.15
The methodology utilizes a two-step process where LLMs first infer operational rules from normal behavior before applyin	✓	0.31
The original prompt design yielded strong results for the battery dataset.	✓	0.27
The original prompt design required modification for the powertrain dataset to improve performance.	✓	0.26
An adjusted prompt emphasizing rule inference significantly improved anomaly detection for the powertrain dataset.	✓	0.29
Mistral 7B achieved F1-scores up to 0.99 in the experiments.	✓	0.25
Llama3.1:8b-instruct-fp16 reached an F1-score of 1.0 in complex scenarios.	✓	0.22
Gemma 2 reached an F1-score of 1.0 in complex scenarios.	×	0.12

References

- <https://doi.org/10.1145/3649506>
- <https://doi.org/10.48550/arxiv.2310.14735>
- <https://doi.org/10.4230/lipics.giscience.2025.3>