

Past-Token Prediction Training and Retrieval Accuracy Decay in Llama-3 Long-Context Tasks

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: How does past-token prediction training affect Needle-in-a-Haystack retrieval accuracy decay rates in Llama-3 compared to next-token prediction across 100K to 500K token contexts. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: To Memorize or to Retrieve: Scaling Laws for RAG-Considerate Pretraining. Research question: How does past-token prediction training affect Needle-in-a-Haystack retrieval accuracy decay rates in Llama-3 compared to next-token prediction across 100K to 500K token contexts?.

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 3.8/10.

3 Results

15 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 3.8/10.

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.

References

- <http://arxiv.org/abs/2412.18619v2>
- <http://arxiv.org/abs/2505.09561v2>
- <http://arxiv.org/abs/2604.00715v1>