

FlowKV and Sliding Window KV Eviction in Long-Context LLaMA-3 Models

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does the learned KV eviction strategy in FlowKV compare to sliding window eviction in terms of downstream task accuracy (e.g., RoME, HELM) when applied to LLaMA-3 models with 200K token contexts. 10 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: ReST-KV: Robust KV Cache Eviction with Layer-wise Output Reconstruction and Spatial-Temporal Smoothing. Research question: How does the learned KV eviction strategy in FlowKV compare to sliding window eviction in terms of downstream task accuracy (e.g., RoME, HELM) when applied to LLaMA-3 models with 200K token contexts?.

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 3.5/10.

3 Results

13 papers retrieved. 10 claims extracted; 0 independently verified. Quality review score: 3.5/10.

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 |
|--|----------|------------|
| ReST-KV is evaluated on five open-source LLMs: Llama2-Chat, Gemma-Instruct, Llama3-Instruct, Mistral-Instruct-v0.3, and | × | 0.03 |
| ReST-KV is compared with five baseline methods: StreamingLLM, H2O, TOVA, SnapKV, and LaCache. | × | 0.04 |
| ReST-KV is evaluated on four benchmarks: LongBench, RULER, Needle-in-a-Haystack, and InniteBench. | × | 0.09 |
| ReST-KV achieves the best performance in most cases on the LongBench benchmark. | × | 0.05 |
| ReST-KV reduces peak memory usage by approximately 36.0% compared to full cache at a context length of 128k. | × | 0.06 |
| ReST-KV achieves an approximate 10.61 \times speedup over the full cache method at a 128K context length. | × | 0.10 |
| ReST-KV is compatible with prell sparse attention approaches, yielding a Time-To-First-Token (TTFT) speedup of up to 3. | × | 0.08 |
| ReST-KV only requires computing attention outputs within a small query window, resulting in a computational complexity c | × | 0.04 |
| LLMs typically decode text in an auto-regressive manner, which allows them to generate high-quality, contextually cohere | × | 0.07 |
| KV cache reduces redundant computation by storing previously computed keys and values. | × | 0.05 |

References

- <http://arxiv.org/abs/2605.09649v1>

- <http://arxiv.org/abs/2602.10238v1>
- <http://arxiv.org/abs/2605.08840v1>