

ReST-KV Spatial-Temporal Smoothing Overhead and Throughput in Llama-3-70B Inference

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: What is the computational overhead and throughput degradation of ReST-KV's spatial-temporal smoothing mechanism relative to top-K retention methods during inference on Llama-3-70B. 14 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: What is the computational overhead and throughput degradation of ReST-KV's spatial-temporal smoothing mechanism relative to top-K retention methods during inference on Llama-3-70B?.

2 Methodology

Systematic literature search across multiple databases yielded 10 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

10 papers retrieved. 14 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.	×	0.05
ReST-KV is evaluated on four benchmarks: LongBench, RULER, Needle-in-a-Haystack, and InniteBench.	×	0.11
ReST-KV achieves the best performance in most cases on the LongBench benchmark across 16 datasets.	×	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.	×	0.11
ReST-KV is compatible with prell sparse attention approaches, yielding a Time-To-First-Token (TTFT) speedup of up to 3.	×	0.06
ReST-KV's computational complexity is comparable to that of SnapKV.	×	0.05
LLMs typically decode text in an auto-regressive manner, which is computationally expensive.	×	0.04
KV cache reduces redundant computation by storing previously computed keys and values.	×	0.05
At each decoding step t , the KV cache stores previously computed keys and values $K_{1:t-1}$, $V_{1:t-1}$ for $X[1 : t - 1]$.	×	0.05
The model requires only the current token x_t to generate x_{t+1} , rather than the full sequence $X = [x_1, \dots, x_t]$.	×	0.02
The query q_t , key k_t , and value v_t are computed as $q_t = x_t W_Q$, $k_t = x_t W_K$, $v_t = x_t W_V$.	×	0.02
The currently computed k_t and v_t will be concatenated with the previously cached keys and values, and used in the attent	×	0.03

References

- <http://arxiv.org/abs/2603.21389v1>
- <http://arxiv.org/abs/2601.11584v1>
- <http://arxiv.org/abs/2605.08840v1>