

# FlowKV and Sliding Window Eviction: Throughput and Latency in LLaMA-3 Long-Context Inference

Assignee Research

June 7, 2026

## Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: What is the comparative throughput and latency overhead of FlowKV versus standard sliding window eviction strategies during inference on LLaMA-3 models with 200K 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 4.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Make Each Token Count: Towards Improving Long-Context Performance with KV Cache Eviction. Research question: What is the comparative throughput and latency overhead of FlowKV versus standard sliding window eviction strategies during inference on LLaMA-3 models with 200K token contexts?.

## 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 4.7/10.

## 3 Results

10 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 4.7/10.

## 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/2605.09649v1>
- <http://arxiv.org/abs/2504.03775v1>
- <http://arxiv.org/abs/2505.15347v2>