

Sliding Window Attention in Mistral-7B vs. Llama Models on Long-Context Reasoning Benchmarks

Assignee Research

May 31, 2026

Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does the sliding window attention mechanism in Mistral-7B affect its performance on long-context reasoning benchmarks compared to Llama-3-8B-128K under memory-constrained inference conditions. We introduce Mistral 7B v0.1, a 7-billion-parameter language model engineered for superior performance and efficiency. Mistral 7B outperforms Llama 2 13B across all evaluated benchmarks, and Llama 1 34B in reasoning, mathematics, and code generation. 15 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Mistral 7B. Research question: How does the sliding window attention mechanism in Mistral-7B affect its performance on long-context reasoning benchmarks compared to Llama-3-8B-128K under memory-constrained inference conditions?.

2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.7/10.

3 Results

12 papers retrieved. 15 claims extracted; 2 independently verified. Quality review score: 5.7/10.

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.

5 Extracted Claims

Claim	Verified	Confidence
Mistral 7B outperforms Llama 2 13B on all benchmarks.	✓	0.29
Mistral 7B outperforms Llama 1 34B on most benchmarks.	✓	0.25
Mistral 7B achieves 3.9% accuracy on GSM8K benchmark.	×	0.00
Mistral 7B achieves 13.1% accuracy on MATH benchmark.	×	0.05
Mistral 7B achieves 30.5% accuracy on HumanEval benchmark.	×	0.05
Mistral 7B achieves 47.5% accuracy on MBPP benchmark.	×	0.05
Mistral 7B achieves 60.1% accuracy on MMLU benchmark.	×	0.05
Mistral 7B achieves 81.3% accuracy on HellaSwag benchmark.	×	0.05
Mistral 7B achieves 75.3% accuracy on Winograndemark benchmark.	×	0.05
Mistral 7B achieves 83.0% accuracy on PIQA benchmark.	×	0.01
Mistral 7B achieves 55.5% accuracy on ARCEasy benchmark.	×	0.04
Mistral 7B achieves 28.8% accuracy on NaturalQuestions benchmark.	×	0.05
Mistral 7B achieves performance equivalent to a Llama 2 model with more than 3x its size on MMLU benchmark.	×	0.10
Mistral 7B has a compression rate of 1.9x on Knowledge benchmarks due to limited parameter count.	×	0.07

References

- <http://arxiv.org/abs/2307.13365v3>
- <http://arxiv.org/abs/2310.06825v1>
- <http://arxiv.org/abs/2501.15089v3>