

Block-Sparse FlashAttention Latency Reductions in Multimodal Reasoning on ChartQA

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: What is the latency reduction of Block-Sparse FlashAttention versus sliding window attention during multimodal reasoning tasks on the ChartQA dataset. 17 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Block Sparse Flash Attention. Research question: What is the latency reduction of Block-Sparse FlashAttention versus sliding window attention during multimodal reasoning tasks on the ChartQA dataset?.

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 5.8/10.

3 Results

13 papers retrieved. 17 claims extracted; 4 independently verified. Quality review score: 5.8/10.

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.

5 Extracted Claims

Claim	Verified	Confidence
Block Sparse Flash Attention achieves up to 1.10 \times speedup on real-world reasoning tasks while maintaining 99% of baselin	✓	0.26
Block Sparse Flash Attention achieves up to 1.24 \times speedup for needle-in-a-haystack retrieval tasks.	✓	0.19
Block Sparse Flash Attention substantially outperforms methods that approximate attention scores.	×	0.14
Block Sparse Flash Attention provides a CUDA kernel implementation that extends FlashAttention-2.	✓	0.18
Block Sparse Flash Attention is a production-ready solution that can be immediately deployed in real-world applications.	×	0.10
Transformers use multi-head scaled dot-product attention to process sequences of tokens.	×	0.03
The computational costs of standard attention implementations include linear projections, score computation, softmax nor	×	0.03
For long sequences where $N \gg d_{model}$, the operations quadratic in N dominate: both the QK score computation and PV aggreg	×	0.03
In Llama-3.1-8B with $d_{model} = 4096$ ($d = 128$, $H = 32$), processing a sequence of $N = 128K$ tokens requires $N^2 d_{model} \approx 6.7 \times$	×	0.04
Block Sparse Flash Attention partitions the query sequence into $MQ = N/BM$ blocks of size BM and the key/value sequence	×	0.08
Block Sparse Flash Attention uses online softmax with incremental updates instead of computing and storing the full atte	×	0.09
Block Sparse Flash Attention maintains running statistics (maximum values and normalizers) that are updated incrementall	×	0.07
Block-Sparse FlashAttention (BSFA) computes all query-key scores exactly to determine importance, then uses these scores	✓	0.18
BSFA skips loading and processing value blocks whose maximum scores fall below calibrated thresholds.	×	0.12
BSFA exploits the observation that blocks with uniformly low scores contribute negligibly after softmax normalization.	×	0.04
If no query-key pair in a block has high enough similarity, the entire block can be safely skipped without meaningful im	×	0.04
By preserving exact score computation while selectively skipping value operations, BSFA achieves significant speedups wi	×	0.08

References

- <http://arxiv.org/abs/2601.15305v1>
- <http://arxiv.org/abs/2512.07011v1>
- <http://arxiv.org/abs/2407.04973v1>