

Block-Sparse FlashAttention Scaling in Multi-Document Retrieval on PG-19

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does the memory efficiency of Block-Sparse FlashAttention scale relative to sliding window attention and linear attention models during multi-document retrieval on the PG-19 dataset. 17 claims were extracted from source literature; 5 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.3/10. 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: How does the memory efficiency of Block-Sparse FlashAttention scale relative to sliding window attention and linear attention models during multi-document retrieval on the PG-19 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 6.3/10.

3 Results

13 papers retrieved. 17 claims extracted; 5 independently verified. Quality review score: 6.3/10.

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.	✓	0.23
Block Sparse Flash Attention maintains 99% of baseline accuracy on real-world reasoning tasks.	✓	0.25
Block Sparse Flash Attention achieves up to 1.24 \times speedup for needle-in-a-haystack retrieval tasks.	✓	0.21
Block Sparse Flash Attention substantially outperforms methods that approximate attention scores.	✓	0.16
The authors provide a CUDA kernel implementation that extends FlashAttention-2.	×	0.13
Transformers use multi-head scaled dot-product attention to process sequences of tokens.	×	0.04
In standard implementations, linear projections require $O(Nd^2_{\text{model}})$ FLOPs total for all Q, K, V projections across all	×	0.02
In standard implementations, score computation (QK) requires $O(N^2d)$ FLOPs per head.	×	0.04
In standard implementations, softmax normalization requires $O(N^2)$ operations per head.	×	0.06
In standard implementations, value aggregation (PV) requires $O(N^2d)$ FLOPs per head.	×	0.03
For long sequences where $N \gg d_{\text{model}}$, QK score computation and PV aggregation scale as $O(N^2d_{\text{model}})$.	×	0.04
In Llama-3.1-8B, d_{model} is 4096, d is 128, and H is 32.	×	0.04
Processing a sequence of $N = 128\text{K}$ tokens in Llama-3.1-8B requires approximately 6.7×10^{13} operations for QK and PV each.	×	0.04
The ratio of operations for QK/PV versus linear projections in the Llama-3.1-8B example is approximately 32:1.	×	0.04
Block-Sparse FlashAttention (BSFA) computes all query-key scores exactly to determine importance before deciding which v	✓	0.20
BSFA skips loading and processing value blocks whose maximum scores fall below calibrated thresholds.	×	0.13
Blocks with uniformly low scores contribute negligibly after softmax normalization.	×	0.04

References

- <http://arxiv.org/abs/2205.14135v2>
- <http://arxiv.org/abs/2512.07011v1>
- <http://arxiv.org/abs/2105.02358v2>