

LongNet Dilated Attention vs. FlashAttention Throughput on PG-19 at Scale

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 11 peer-reviewed papers addressing the following research question: How does the training throughput of LongNet’s dilated attention compare to standard FlashAttention on the PG-19 dataset when scaling context lengths from 32k to 100k tokens. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: LongNet: Scaling Transformers to 1,000,000,000 Tokens. Research question: How does the training throughput of LongNet’s dilated attention compare to standard FlashAttention on the PG-19 dataset when scaling context lengths from 32k to 100k tokens?.

2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.0/10.

3 Results

11 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 5.0/10.

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/2601.15305v1>
- <http://arxiv.org/abs/2307.02486v2>
- <http://arxiv.org/abs/2205.14135v2>