

Cross-Lingual Retrieval Performance of Block-Sparse FlashAttention in Multilingual Models on MLQA

Assignee Research

June 11, 2026

Abstract

Information retrieval across different languages is an increasingly important challenge in natural language processing. Recent approaches based on multilingual pre-trained language models have achieved remarkable success, yet they often optimize for either monolingual, cross-lingual, or multilingual retrieval performance at the expense of others. This paper proposes a novel hybrid batch training strategy to simultaneously improve zero-shot retrieval performance across monolingual, cross-lingual, and multilingual settings while mitigating language bias. The approach fine-tunes multilingual lang

1 Introduction

This paper examines: Synergistic Approach for Simultaneous Optimization of Monolingual, Cross-lingual, and Multilingual Information Retrieval. Research question: To what extent does the cross-lingual retrieval performance of Block-Sparse FlashAttention degrade when scaling to multilingual models with 10+ languages on the MLQA benchmark compared to local attention mechanisms?.

2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.3/10.

3 Results

14 papers retrieved. 14 claims extracted; 11 independently verified. Quality review score: 7.3/10.

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
The approach fine-tunes multilingual language models using a mix of monolingual and cross-lingual question-answer pair b	✓	0.41
Experiments are conducted on XQuAD-R, MLQA-R, and MIRACL Datasets.	×	0.08
XQuAD-R and MLQA-R are question-answering datasets with parallel questions and passages in 11 languages and 7 languages,	✓	0.20
The evaluation of the models is conducted on datasets that are completely separate and distinct from the ones used for t	✓	0.23
The models have not encountered any data samples, whether from the training or testing splits, of the evaluation dataset	✓	0.24
The mean average precision (mAP) is reported for XQuAD-R and MLQA-R.	×	0.11
The hybrid batch sampling achieves the best performance in multilingual retrieval settings.	✓	0.29
Hybrid batch training substantially reduces language bias in multilingual retrieval compared to monolingual training.	✓	0.36
Hybrid batch training enables strong zero-shot retrieval performance across diverse languages.	✓	0.30
The proposed approach learns language-agnostic representations.	×	0.13
The hybrid batch training strategy simultaneously optimizes retrieval performance across monolingual, cross-lingual, and	✓	0.29
The hybrid batch training strategy mitigates language bias.	✓	0.15
The hybrid batch training strategy uses a balanced mix of monolingual and cross-lingual question-answer pair batches.	✓	0.24
The hybrid batch training strategy collects a diverse set of English question-answer datasets and uses machine translati	✓	0.26

References

- <http://arxiv.org/abs/2107.11976v2>

- <http://arxiv.org/abs/2408.10536v1>
- <http://arxiv.org/abs/2509.22472v1>