

Trade-off Between Model Size Scaling and Inference Efficiency in Zero-Shot Cross-Lingual Retrieval via Hybrid Batch Training

Assignee Research

June 19, 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: What is the trade-off between model size scaling and inference efficiency in zero-shot cross-lingual retrieval when using the hybrid batch training strategy, as evaluated on throughput and latency metrics (e.g., tokens/sec)?.

2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.7/10.

3 Results

8 papers retrieved. 11 claims extracted; 9 independently verified. Quality review score: 7.7/10.

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
The approach fine-tunes multilingual language models using a mix of monolingual and cross-lingual question-answer pairs	✓	0.41
Experiments were 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 training	✓	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.	×	0.10
For XQuAD-R (MLQA-R), there are 11 and 7 parallel languages; thus, there are 110 (42) and 11 (7) cross-lingual and monolingual	✓	0.23
Hybrid batch sampling achieves the best performance in multilingual retrieval settings.	✓	0.28
Hybrid batch sampling is better than the other two baseline batch sampling methods when using XLM-R and LaBSE as initial	✓	0.25
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.	✓	0.29

References

- <http://arxiv.org/abs/2305.05295v2>
- <http://arxiv.org/abs/2410.21676v4>
- <http://arxiv.org/abs/2408.10536v1>