

# Hybrid Training Data Ratios and Retrieval Robustness Across Diverse Language Pairs in MKQA

Assignee Research

June 20, 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 effect of varying the ratio of monolingual and cross-lingual data in the hybrid training approach on the robustness of retrieval performance across diverse language pairs in MKQA?.

## 2 Methodology

Systematic literature search across multiple databases yielded 10 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.1/10.

## 3 Results

10 papers retrieved. 17 claims extracted; 14 independently verified. Quality review score: 8.1/10.

## 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	✓	0.40
Experiments on XQuAD-R, MLQA-R, and MIRACL Datasets.	×	0.12
XQuAD-R and MLQA-R are question-answering datasets with parallel questions and passages in 11 languages and 7 languages,	✓	0.20
We report the mean average precision (mAP) for XQuAD-R and MLQA-R.	✓	0.16
The evaluation of the models is conducted on datasets that are completely separate and distinct from the ones used for training	✓	0.23
Hybrid batch sampling achieves the best performance in multilingual retrieval settings.	✓	0.28
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 fine-tunes multilingual language models using a balanced mix of monolingual and cross-lingual questions	✓	0.32
The models have not encountered any data samples, whether from the training or testing splits, of the evaluation dataset	✓	0.24
The results for the Recall metric are in Appendix A.3.1.	✓	0.17
The detailed monolingual retrieval effectiveness on MIRACL dev is reported in Table 12 and 13 in Appendix A.3.3.	✓	0.19
X-X sampling only performs well in monolingual retrieval settings.	×	0.08
X-Y sampling only performs well in cross-lingual retrieval settings.	×	0.12
Optimization for either monolingual or cross-lingual retrieval alone may come at the expense of the other.	✓	0.19
Hybrid batch sampling optimizes both monolingual and cross-lingual retrieval settings.	✓	0.20
The same conclusion holds when using XLM-R and LaBSE as initialization that hybrid batch sampling is better than the other	✓	0.29

## References

- <http://arxiv.org/abs/2305.05295v2>
- <http://arxiv.org/abs/2505.18673v1>
- <http://arxiv.org/abs/2408.10536v1>