

Hybrid Batch Training for Multimodal Alignment in Zero-Shot Cross-Lingual Image-Text Retrieval

Assignee Research

June 16, 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: Does the hybrid batch training strategy proposed for information retrieval improve multimodal alignment accuracy on zero-shot cross-lingual image-text retrieval benchmarks like XM3600 compared to separate task fine-tuning?.

2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.5/10.

3 Results

12 papers retrieved. 16 claims extracted; 15 independently verified. Quality review score: 8.5/10.

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	✓	0.42
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.	✓	0.15
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 also substantially reduces language bias in multilingual retrieval compared to monolingual training	✓	0.36
The proposed approach enables strong zero-shot retrieval performance across diverse languages.	✓	0.27
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.	✓	0.16
We report the detailed monolingual retrieval effectiveness on MIRACL dev (Zhang et al., 2022) in Table 12 and 13 in Appendix	✓	0.28
In XQuAD-R (MLQA-R), we have 11 and 7 parallel languages; thus, there are 110 (42) and 11 (7) cross-lingual and monolingual	✓	0.23
The same conclusion holds when using XLM-R and LaBSE as initialization that hybrid batch sampling is better than the other	✓	0.31
Hybrid batch sampling achieves the best performance in multilingual retrieval settings, where the ability of the models	✓	0.38
Hybrid batch training also substantially reduces language bias in multilingual retrieval compared to monolingual training	✓	0.36
The proposed approach enables strong zero-shot retrieval performance across diverse languages.	✓	0.27

References

- <http://arxiv.org/abs/2503.06380v1>
- <http://arxiv.org/abs/2408.10536v1>
- <http://arxiv.org/abs/2506.15415v1>