

Zero-Shot Cross-Lingual Image-Text Retrieval Performance of Multimodal Models Trained on Artificially Code-Switched Data

Assignee Research

June 18, 2026

Abstract

Transferring information retrieval (IR) models from a high-resource language (typically English) to other languages in a zero-shot fashion has become a widely adopted approach. In this work, we show that the effectiveness of zero-shot rankers diminishes when queries and documents are present in different languages. Motivated by this, we propose to train ranking models on artificially code-switched data instead, which we generate by utilizing bilingual lexicons. To this end, we experiment with lexicons induced from (1) cross-lingual word embeddings and (2) parallel Wikipedia page titles. We use

1 Introduction

This paper examines: Boosting Zero-shot Cross-lingual Retrieval by Training on Artificially Code-Switched Data. Research question: How do multimodal models trained on artificially code-switched text and image pairs perform in zero-shot cross-lingual image-text retrieval tasks, as measured by recall@k on the M3IT benchmark?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.2/10.

3 Results

15 papers retrieved. 23 claims extracted; 17 independently verified. Quality review score: 7.2/10.

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
Code-switching improves cross-lingual and multilingual re-ranking performance.	×	0.14
Code-switching training does not impede monolingual (MoIR) setup performance.	×	0.06
The average zero-shot MoIR performance is 15.7 MRR@10.	✓	0.17
The average zero-shot MLIR performance is 16.6 MRR@10.	×	0.13
In CLIR, the performance drop when transferring from EN-EN is larger for typologically distant languages (AR-IT, AR-RU)	✓	0.22
The performance gap between zero-shot and fine-tuning on translated data in MoIR is +4 MRR@10.	✓	0.23
The performance gap between zero-shot and fine-tuning on translated data in CLIR is +11.1 MRR@10.	✓	0.25
The performance gap between zero-shot and fine-tuning on translated data in MLIR is +8.3 MRR@10.	✓	0.23
Training on code-switched data consistently outperforms zero-shot models in CLIR and MLIR.	✓	0.25
For the AR-IT language pair, code-switching training improved MRR@10 from 7.7 to 15.6.	×	0.14
For the AR-RU language pair, code-switching training improved MRR@10 from 7.1 to 14.1.	✓	0.16
The difference in MoIR performance between code-switching approaches (BL-CS, ML-CS) and Zero-shot is not statistically significant.	✓	0.21
Specializing one zero-shot model for multiple CLIR language pairs (ML-CS, Wiki-CS) performs almost on par with specializing multiple models.	✓	0.31
Wiki-CS results are slightly worse in MoIR compared to other approaches.	✓	0.18
In MoIR, Zero-shot Translate Test and ML-CS Translate Test underperform compared to other approaches.	✓	0.23
Zero-shot rankers perform better on clean monolingual data in the target language than on noisy monolingual data in English.	✓	0.22
In CLIR, Translate Test yields improvements of +0.2 and +2.2 MRR@10.	×	0.15
Translate Test consistently falls behind code-switching at training time in both MoIR and CLIR.	✓	0.20
Code-switching gains remain virtually unchanged when moving from six seen languages to fourteen languages including eight unseen languages.	✓	0.23
The gain for six seen languages was +4.1 MRR@10 and +3.8 MRR@10.	✓	0.16

References

- <http://arxiv.org/abs/2305.05295v2>
- <http://arxiv.org/abs/2603.23521v1>
- <http://arxiv.org/abs/2308.15273v1>