

# Multimodal Model Performance in Zero-Shot Cross-Lingual Image-Text Retrieval on MMMU

Assignee Research

July 8, 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: Do multimodal models trained on artificially code-switched text-to-image pairs outperform text-only models in zero-shot cross-lingual image-text retrieval on the MMMU benchmark?.

## 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 8.2/10.

## 3 Results

8 papers retrieved. 15 claims extracted; 13 independently verified. Quality review score: 8.2/10.

## 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 without impeding monolingual setups.                   | ×        | 0.14       |
| The average MoIR performance is substantially higher than CLIR with 15.7 MRR@10 and MLIR with 16.6 MRR@10.               | ✓        | 0.24       |
| The performance drop is larger for setups involving typologically distant languages (AR-IT, AR-RU) in CLIR and to a less | ✓        | 0.31       |
| The performance gap to Fine-tuning on translated data is much smaller in MoIR (+4 MRR@10) than in CLIR (+11.1 MRR@10) an | ✓        | 0.31       |
| Training on code-switched data consistently outperforms zero-shot models in CLIR and MLIR.                               | ✓        | 0.24       |
| In AR-IT and AR-RU, improvements from 7.7 and 7.1 MRR@10 up to 15.6 and 14.1 MRR@10 are observed.                        | ✓        | 0.22       |
| The differences between both CS approaches (BL-CS and ML-CS) versus Zero-shot is not statistically significant.          | ✓        | 0.26       |
| Specializing one zero-shot model for multiple CLIR language pairs (ML-CS, Wiki-CS) performs almost on par with specializ | ✓        | 0.30       |
| Zero-shotTranslate Test and ML-CSTranslate Test underperform compared to other approaches in MoIR.                       | ✓        | 0.21       |
| Translate Test shows slight improvements of +0.2 and +2.2 MRR@10 in CLIR.  | ✓        | 0.15       |
| Translate Test consistently falls behind code-switching at training time in both MoIR and CLIR.                          | ✓        | 0.20       |
| The gains remain virtually unchanged when moving from six seen (+4.1 MRR@10 / +3.8 MRR@10) to fourteen languages includi | ✓        | 0.30       |
| Training on code-switched data is a cheap and effective way of generalizing zero-shot rankers for cross-lingual and mult | ✓        | 0.34       |
| Fine-tuning CEs on monolingual data biases the encoder towards encoding features that are only useful when the target se | ✓        | 0.25       |
| Artificial code-switching is a method to modify monolingual training data by borrowing words from other languages.       | ×        | 0.12       |

## References

- <http://arxiv.org/abs/2406.13361v1>
- <http://arxiv.org/abs/2109.07622v1>
- <http://arxiv.org/abs/2305.05295v2>