

Hybrid Batch Training for Generalization to Unseen Languages in the MIRACL Dataset Versus Multilingual Pre-trained Models

Assignee Research

June 18, 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 impact of the hybrid batch training method on generalization to unseen languages within the MIRACL dataset compared to existing multilingual pre-trained language models?.

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 8.3/10.

3 Results

10 papers retrieved. 7 claims extracted; 7 independently verified. Quality review score: 8.3/10.

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
Recent approaches based on multilingual pre-trained language models have achieved remarkable success in information retr	✓	0.34
Recent approaches often optimize for either monolingual, cross-lingual, or multilingual retrieval performance at the exp	✓	0.33
The paper proposes a novel hybrid batch training strategy to simultaneously improve zero-shot retrieval performance acro	✓	0.51
The approach fine-tunes multilingual language models using a mix of monolingual and cross-lingual question-answer pair b	✓	0.43
Experiments on XQuAD-R, MLQA-R, and MIRACL benchmark datasets show that the proposed method consistently achieves compar	✓	0.51
Hybrid batch training substantially reduces language bias in multilingual retrieval compared to monolingual training.	✓	0.38
The proposed approach enables learning language-agnostic representations that enable strong zero-shot retrieval performa	✓	0.35

References

- <https://doi.org/10.48550/arxiv.2402.03216>
- <https://doi.org/10.18653/v1/2024.findings-acl.137>
- <https://doi.org/10.48550/arxiv.2408.10536>