

Optimal Transport-Based Knowledge Distillation for Multilingual Retrieval Performance Gaps Under Training Data Noise

Assignee Research

June 19, 2026

Abstract

Benefiting from transformer-based pre-trained language models, neural ranking models have made significant progress. More recently, the advent of multilingual pre-trained language models provides great support for designing neural cross-lingual retrieval models. However, due to unbalanced pre-training data in different languages, multilingual language models have already shown a performance gap between high and low-resource languages in many downstream tasks. And cross-lingual retrieval models built on such pre-trained models can inherit language bias, leading to suboptimal result for low-reso

1 Introduction

This paper examines: Improving Cross-lingual Information Retrieval on Low-Resource Languages via Optimal Transport Distillation. Research question: How does the performance gap between high and low-resource languages in multilingual retrieval tasks change when using different optimal transport-based knowledge distillation techniques under varying degrees of training data noise?.

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 7.9/10.

3 Results

12 papers retrieved. 11 claims extracted; 9 independently verified. Quality review score: 7.9/10.

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
Sasaki et al. proposed a large cross-lingual retrieval collection named WikiCLIR based on linked foreign language articles	✓	0.25
The relevant judgments in the WikiCLIR dataset are synthetically generated based on mutual links across pages.	✓	0.17
Bonifacio et al. built the mMARCO multilingual passage ranking dataset by translating queries and passages from MS MARCO	✓	0.28
The MS MARCO dataset is generated from query logs.	×	0.09
OPTICAL is a knowledge distillation framework based on Optimal Transport designed for low-resource Cross-lingual Informa	✓	0.17
OPTICAL formulates the cross-lingual token alignment task as an optimal transport problem where the cost matrix is the c	✓	0.23
In OPTICAL, the loss is defined as the Frobenius inner product of the transportation plan and the cost matrix.	✓	0.25
OPTICAL requires only bitext data for distillation training.	×	0.14
Experiments were performed on seven language pairs, including four low-resource languages and three medium or high-resou	✓	0.18
OPTICAL achieved a 13.7% improvement in Mean Average Precision (MAP) over a method based on neural machine translation o	✓	0.20
OPTICAL significantly outperforms several strong baseline methods on low-resource languages in terms of Mean Average Pre	✓	0.22

References

- <http://arxiv.org/abs/2510.00908v1>
- <http://arxiv.org/abs/2301.12566v1>
- <http://arxiv.org/abs/1909.07342v1>