

Contrastive Multilingual Pre-training for Robust Zero-Shot Cross-Lingual Transfer

Assignee Research

July 8, 2026

Abstract

Multilingual pre-trained models have achieved remarkable performance on cross-lingual transfer learning. Some multilingual models such as mBERT, have been pre-trained on unlabeled corpora, therefore the embeddings of different languages in the models may not be aligned very well. In this paper, we aim to improve the zero-shot cross-lingual transfer performance by proposing a pre-training task named Word-Exchange Aligning Model (WEAM), which uses the statistical alignment information as the prior knowledge to guide cross-lingual word prediction. We evaluate our model on multilingual machine rea

1 Introduction

This paper examines: Bilingual Alignment Pre-Training for Zero-Shot Cross-Lingual Transfer. Research question: Does contrastive multilingual pre-training (e.g., SimCSE, LaBSE) improve zero-shot cross-lingual transfer robustness compared to standard MLM pre-training, as measured by average accuracy across XTREME-R benchmarks?.

2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.7/10.

3 Results

14 papers retrieved. 12 claims extracted; 10 independently verified. Quality review score: 7.7/10.

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 mBERT+TLM model outperforms mBERT by a large margin in the zero-shot setting.	✓	0.27
The mBERT+TLM model is not as good as the mBERT in the translate-train setting.	✓	0.24
The mBERT+WEAM model improves the scores in the zero-shot setting and also outperforms mBERT in the translate-train sett	✓	0.30
The mBERT+TLM and word-aligned mBERT achieved similar improvements on XNLI compared to mBERT.	✓	0.24
The mBERT+WEAM model significantly outperformed both mBERT+TLM and word-aligned mBERT on XNLI.	✓	0.18
The mBERT+WEAM result is slightly lower but close to the translate-train result on XNLI.	✓	0.24
The examples in XNLI have shorter input sequences and thus have fewer translation noises.	✓	0.17
The masking probability is empirically set to 0.3 for better performance.	×	0.14
The learning rate is set to 5e-5, the batch size to 32, the max sequence length to 128, and the number of pre-training e	✓	0.24
The value of λ is set to 1.	×	0.03
The WEAM model uses FastAlign to identify bilingual word pairs in parallel bilingual sentence pairs.	✓	0.23
The WEAM model performs multilingual prediction and cross-lingual prediction for each masked token.	✓	0.18

References

- <http://arxiv.org/abs/2104.08645v2>
- <http://arxiv.org/abs/2104.08821v4>
- <http://arxiv.org/abs/2106.01732v2>