

# Scalability of WEAM’s Alignment Mechanism with Multilingual Pre-training and Model Size Trade-offs in XTREME-R Performance

Assignee Research

June 23, 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 the effectiveness of WEAM’s alignment mechanism scale with the number of languages included in pre-training, and what is the trade-off between model size and zero-shot transfer performance on XTREME-R?.

## 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 7.5/10.

## 3 Results

14 papers retrieved. 35 claims extracted; 25 independently verified. Quality review score: 7.5/10.

## 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 model was pre-trained on unlabeled corpora.	✓	0.21
The embeddings of different languages in mBERT may not be aligned very well due to pre-training on unlabeled corpora.	✓	0.20
The proposed pre-training task is named Word-Exchange Aligning Model (WEAM).	✓	0.29
WEAM uses statistical alignment information as prior knowledge to guide cross-lingual word prediction.	✓	0.41
The model was evaluated on the MLQA multilingual machine reading comprehension task.	✓	0.22
The model was evaluated on the XNLI natural language interface task.	✓	0.17
WEAM significantly improves zero-shot performance on MLQA and XNLI tasks.	✓	0.17
Three parallel corpora were used with English as the source language and Chinese, German, and Spanish as target language	×	0.13
The mBERT model was initialized with weights released by Google.	×	0.10
Three models were pre-trained separately for the three target languages to avoid alignment interference.	×	0.09
The masking probability was empirically set to 0.3 during pre-training.	×	0.15
A masking probability of 0.3 yielded better performance than other tested values.	×	0.08
The learning rate was set to 5e-5 for all three models.	×	0.08
The batch size was set to 32 for all three models.	×	0.11
The maximum sequence length was set to 128 for all three models.	×	0.07
The number of pre-training epochs was set to 2 for all three models.	✓	0.15
The hyper-parameter $\lambda$ was set to 1.	×	0.12
The mBERT+TLM model outperforms mBERT by a large margin in the zero-shot setting on MLQA.	✓	0.28
The mBERT+TLM model performs worse than mBERT in the translate-train setting on MLQA.	✓	0.19
The mBERT+WEAM model improves scores in the zero-shot setting on MLQA compared to mBERT.	✓	0.22
The mBERT+WEAM model outperforms mBERT in the translate-train setting on MLQA.	✓	0.27
A properly aligned pre-training model can exceed the performance of translate-train even with zero-shot training on MLQA.	✓	0.28

## References

- <http://arxiv.org/abs/2203.09982v1>
- <http://arxiv.org/abs/2310.10378v5>
- <http://arxiv.org/abs/2106.01732v2>