

# Noise Injection in mT5 Pre-training for Robust Zero-Shot Cross-Lingual Transfer on XQuAD and MLQA

Assignee Research

June 26, 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 incorporating noise injection during pre-training improve the robustness and zero-shot cross-lingual transfer performance of multilingual models like mT5 on the XQuAD and MLQA benchmarks?.

## 2 Methodology

Systematic literature search across multiple databases yielded 4 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

4 papers retrieved. 9 claims extracted; 7 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
mBERT+TLM model outperforms mBERT by a large margin in the zero-shot setting.	✓	0.26
mBERT+WEAM improves the scores in the zero-shot setting and also outperforms mBERT in the translate-train setting.	✓	0.32
mBERT+WEAM has 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 masking probability is empirically set to 0.3 for better performance.	×	0.15
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	✓	0.23
The value of $\lambda$ is set to 1.	×	0.03
WEAM performs two kinds of predictions: a multilingual prediction and a cross-lingual prediction.	✓	0.21
FastAlign is used to identify bilingual word pairs in parallel bilingual sentence pairs.	✓	0.22

## References

- <http://arxiv.org/abs/2103.08849v3>
- <http://arxiv.org/abs/2104.08645v2>
- <http://arxiv.org/abs/2106.01732v2>