

# Scaling Behavior of MMKD Across Teacher-Student Size Ratios in Low-Resource Settings on XTREME-R

Assignee Research

June 13, 2026

## Abstract

Zero-shot cross-lingual knowledge transfer enables the multilingual pretrained language model (mPLM), finetuned on a task in one language, make predictions for this task in other languages. While being broadly studied for natural language understanding tasks, the described setting is understudied for generation. Previous works notice a frequent problem of generation in a wrong language and propose approaches to address it, usually using mT5 as a backbone model. In this work, we test alternative mPLMs, such as mBART and NLLB-200, considering full finetuning and parameter-efficient finetuning wi

## 1 Introduction

This paper examines: Empirical study of pretrained multilingual language models for zero-shot cross-lingual knowledge transfer in generation. Research question: How does the performance of MMKD scale with varying teacher-student model size ratios in low-resource language settings on the XTREME-R benchmark?.

## 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 7.4/10.

## 3 Results

14 papers retrieved. 11 claims extracted; 8 independently verified. Quality review score: 7.4/10.

## 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 study evaluates pretrained multilingual language models for zero-shot cross-lingual knowledge transfer in generation	✓	0.22
The study tests alternative multilingual pretrained language models (mPLMs) including mBART and NLLB-200.	✓	0.17
The study considers both full finetuning and parameter-efficient finetuning with adapters.	✓	0.17
mBART with adapters performs similarly to mT5 of the same size.	✓	0.24
NLLB-200 can be competitive in some cases compared to other models tested.	×	0.14
Tuning the learning rate used for finetuning helps alleviate the problem of generation in the wrong language.	✓	0.37
Figure 1 shows validation curves for full finetuning of mT5-base and mBART with various learning rates for the Russian l	×	0.13
Full results demonstrating the effect of learning rate across all task-model-adaptation method-language combinations are	✓	0.22
Previous approaches to addressing wrong language generation often use mT5 as a backbone model.	✓	0.16
Drawbacks of translation-based cross-lingual generation approaches include high computational cost, lack of high-quality	✓	0.23
The study focuses on encoder-decoder multilingual pretrained language models.	×	0.10

## References

- <http://arxiv.org/abs/2004.08116v1>
- <http://arxiv.org/abs/2310.09917v3>
- <http://arxiv.org/abs/2106.09063v4>