

# Dataset Size and Diversity in Zero-Shot Cross-Lingual Transfer Robustness

Assignee Research

June 29, 2026

## Abstract

Intermediate-task training—fine-tuning a pretrained model on an intermediate task before fine-tuning again on the target task—often improves model performance substantially on language understanding tasks in monolingual English settings. We investigate whether English intermediate-task training is still helpful on non-English target tasks. Using nine intermediate language-understanding tasks, we evaluate intermediate-task transfer in a zero-shot cross-lingual setting on the XTREME benchmark. We see large improvements from intermediate training on the BUCC and Tatoeba sentence retrieval tas

## 1 Introduction

This paper examines: English Intermediate-Task Training Improves Zero-Shot Cross-Lingual Transfer Too. Research question: How does the size and diversity of the intermediate-task training dataset influence the robustness of zero-shot cross-lingual transfer performance across different language families in the XTREME benchmark?.

## 2 Methodology

Systematic literature search across multiple databases yielded 11 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

11 papers retrieved. 18 claims extracted; 14 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
XLM-R Large model achieves state-of-the-art performance on many zero-shot cross-lingual transfer tasks.	✓	0.21
The XTREME benchmark evaluates zero-shot cross-lingual transfer performance across diverse target tasks across up to 40	✓	0.18
Intermediate-task training on SQuAD, MNLI, and HellaSwag yields large target-task improvements of 8.2, 7.5, and 7.0 poin	✓	0.27
Multi-task intermediate-task training on all 9 tasks performs best, improving by 8.7 points.	✓	0.24
Applying intermediate-task training to BUCC and Tatoeba yields dramatic improvements with almost every intermediate trai	✓	0.22
TyDiQA shows consistent improvements with many intermediate tasks, whereas XNLI does not see benefits from intermediate	✓	0.17
Evaluating the best performing models for each target task on the XTREME benchmark yields an average improvement of 5.4	✓	0.32
Training on English intermediate tasks outperforms continuing multilingual MLM during intermediate-task training.	✓	0.23
Training on English intermediate tasks outperforms using machine-translated intermediate-task data.	✓	0.23
The pretrained XLM-R model is used as a starting point for all experiments.	×	0.15
The baseline involves fine-tuning the pretrained XLM-R model on each target task’s English training data and evaluating	✓	0.20
The main approach includes an additional intermediate-task training phase before training and evaluating on the target t	✓	0.16
Multi-task training on all available intermediate tasks is also experimented with.	×	0.15
The three-phase approach to training includes MLM, intermediate-task training, and fine-tuning on English target-task tr	✓	0.21
Intermediate tasks have English input data, with an alternative of machine-translating intermediate-task data to other l	✓	0.17
Both single- and multi-task training <sup>4</sup> are experimented with for intermediate-task training.	×	0.14
Target tasks from the XTREME benchmark are used for zero-shot cross-lingual transfer.	✓	0.18
Nine different English intermediate tasks are studied, covering various task formats and sources.	×	0.10

## References

- <http://arxiv.org/abs/2003.11080v5>
- <http://arxiv.org/abs/2104.08645v2>
- <http://arxiv.org/abs/2005.13013v2>