

Cross-lingual Calibration in XTREME: Intermediate-Task Training Effects

Assignee Research

June 26, 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: What is the impact of intermediate-task training strategies on the calibration of confidence scores in zero-shot cross-lingual settings across diverse typological groups in XTREME?.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.3/10.

3 Results

13 papers retrieved. 20 claims extracted; 12 independently verified. Quality review score: 7.3/10.

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 uses the pretrained XLM-R Large model as the starting point for all experiments.	×	0.15
Experiments are performed on 9 target tasks from the XTREME benchmark.	×	0.10
The XTREME benchmark evaluates zero-shot cross-lingual transfer performance across up to 40 languages for each task.	×	0.14
Intermediate-task training on SQuAD yields a target-task improvement of 8.2 points on the development set.	✓	0.19
Intermediate-task training on MNLI yields a target-task improvement of 7.5 points on the development set.	✓	0.18
Intermediate-task training on HellaSwag yields a target-task improvement of 7.0 points on the development set.	✓	0.19
Multi-task intermediate-task training on all 9 tasks improves performance by 8.7 points.	✓	0.16
Applying intermediate-task training to BUCC and Tatoeba yields dramatic improvements with almost every intermediate tra	✓	0.22
BUCC and Tatoeba are sentence retrieval target tasks that have no training data of their own.	✓	0.20
TyDiQA shows consistent improvements with many intermediate tasks.	✓	0.16
XNLI does not see benefits from intermediate training.	×	0.09
Evaluating the best performing models for each target task on the XTREME benchmark yields an average improvement of 5.4	✓	0.28
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 study investigates nine different English intermediate tasks including question answering, sentence tagging, sentenc	✓	0.22
The ANLI+ intermediate task contains 1,104,934 training examples.	×	0.07
The MNLI intermediate task contains 392,702 training examples.	×	0.09
The SQuAD v2.0 intermediate task is based on Wikipedia data.	×	0.10
The HellaSwag intermediate task is based on Video captions and Wikihow data.	×	0.10
The methodology follows a three-phase approach: MLM pre-training, intermediate-task training on English data, and fine-t	✓	0.24

References

- <http://arxiv.org/abs/2003.11080v5>
- <http://arxiv.org/abs/2009.05166v3>
- <http://arxiv.org/abs/2005.13013v2>