

XTREME-R Zero-Shot Cross-Lingual Transfer Robustness Across Language Families via English Intermediate-Task Fine-Tuning

Assignee Research

June 25, 2026

Abstract

Transfer learning from large language models (LLMs) has emerged as a powerful technique to enable knowledge-based fine-tuning for a number of tasks, adaptation of models for different domains and even languages. However, it remains an open question, if and when transfer learning will work, i.e. leading to positive or negative transfer. In this paper, we analyze the knowledge transfer across three natural language processing (NLP) tasks - text classification, sentimental analysis, and sentence similarity, using three LLMs - BERT, RoBERTa, and XLNet - and analyzing their performance, by fine-tun

1 Introduction

This paper examines: The (In)Effectiveness of Intermediate Task Training For Domain Adaptation and Cross-Lingual Transfer Learning. Research question: How does the robustness of zero-shot cross-lingual transfer performance on XTREME-R vary across different language families (e.g., Romance, Germanic, Semitic) when using English intermediate-task fine-tuning versus monolingual fine-tuning?.

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 7.6/10.

3 Results

13 papers retrieved. 11 claims extracted; 10 independently verified. Quality review score: 7.6/10.

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
RoBERTa and BERT with intermediate task training are the best models, depending on the task	✓	0.25
RoBERTa outperformed in three out of six tasks, across three NLP tasks (text classification, sentiment analysis, sentenc	✓	0.27
BERT outperformed others in the rest three tasks	×	0.11
XLNet was consistently the worst performing model in all of our experiments	✓	0.20
Similar trends for transfer learning using LLMs, where RoBERTa and BERT have similar performance, and both outperform XL	✓	0.44
In target tasks per NLP task, the first task is for domain adaptation, and the next one is for cross-lingual adaptation	✓	0.19
For text classification, we performed intermediate task training using the IMDB movie reviews dataset	✓	0.24
For the intermediate task training, each pre-trained LLM was trained for 100 epochs using the large dataset	✓	0.28
For fine-tuning after and without intermediate task training, transfer learning to the target dataset was performed by t	✓	0.28
In both cases of transfer learning, all the model weights were updated, or none of the layers were frozen	✓	0.22
In each of the NLP tasks, the dataset used for the intermediate task training from the LLM is at least an order of magni	✓	0.25

References

- <http://arxiv.org/abs/2212.01757v1>
- <http://arxiv.org/abs/2210.01091v2>
- <http://arxiv.org/abs/2503.19979v1>