

Comparative Analysis of Zero-Shot Cross-Lingual Transfer in Prefix-Adapted Versus Instruction-Tuned 7B Models on XTREME-R

Assignee Research

June 26, 2026

Abstract

With the release of new large language models (LLMs) like Llama and Mistral, zero-shot cross-lingual transfer has become increasingly feasible due to their multilingual pretraining and strong generalization capabilities. However, adapting these decoder-only LLMs to new tasks across languages remains challenging. While parameter-efficient fine-tuning (PeFT) techniques like Low-Rank Adaptation (LoRA) are widely used, prefix-based techniques such as soft prompt tuning, prefix tuning, and Llama Adapter are less explored, especially for zero-shot transfer in decoder-only models. We present a compre

1 Introduction

This paper examines: Zero-Shot Cross-Lingual Transfer using Prefix-Based Adaptation. Research question: How does the zero-shot cross-lingual transfer performance of prefix-based adapted 7B-parameter models compare to that of instruction-tuned 7B-parameter models on XTREME-R benchmark tasks?.

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 7.4/10.

3 Results

11 papers retrieved. 26 claims extracted; 19 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
Zhao and Schtze (2021) systematically compared discrete prompting, soft prompting, and fine-tuning on the few-shot mult	✓	0.29
Tu et al. (2022) compared prompt tuning with fine-tuning across diverse NLU tasks on XLM-R and mBERT.	✓	0.27
Tu et al. (2022) evaluated prefix tuning on the encoder-only XLM-R model and showed its effectiveness over full fine-tun	✓	0.35
Tu et al. (2022) investigated the decoder-based multilingual model XGLM, limiting their analysis to a single small model	✓	0.26
Tu et al. (2022) showed that prompt tuning can sometimes surpass fine-tuning, particularly for low-resource languages.	✓	0.32
The experiments in this study were conducted on Llama 3.1 (8B), Mistral v0.3 (7B), Llama 3.2 (1B), and Mistral Small (24	✓	0.24
Llama 3.1 and 3.2 series are multilingual large language models developed by Meta.	×	0.12
Mistral v0.3 (7B) has an extended vocabulary compared to Mistral v0.1.	✓	0.20
Mistral Small (24B) is categorized as a 'small' LLM (under 70B) with improved multilingual capabilities and a larger voc	✓	0.22
The study limits experiments to base model variants only.	×	0.09
XQUAD (Artetxe et al., 2019) was used as a benchmark for cross-lingual question answering.	✓	0.19
XNLI (Conneau et al., 2018) was used as a benchmark for cross-lingual natural language inference.	✓	0.20
Belebele (Bandarkar et al., 2024) was used as a benchmark for cross-lingual machine reading comprehension.	✓	0.21
MGSM (Shi et al., 2023) was used to assess reasoning capabilities in multilingual settings.	✓	0.20
Prefix-based adaptation methods and LoRA (rank 4) were fine-tuned using the English SQuAD training set containing 87.6K	✓	0.22
A subset of English NLI training data containing 100K samples was used for XNLI evaluations.	✓	0.22
The suggested Belebele training set containing 67.5K English samples was used for Belebele evaluations.	✓	0.21
The GSM8K English training dataset with 7.47K samples was used for MGSM evaluations.	✓	0.18
Learning rates of 3e-3, 1e-3, and 3e-4 were experimented with.	✓	0.18
On the Kyrgyz language benchmark, the Base	×	0.10

References

- <http://arxiv.org/abs/2402.14778v2>
- <http://arxiv.org/abs/2310.09917v3>
- <http://arxiv.org/abs/2510.24619v1>