

Lugha-Llama Cross-Lingual Alignment via Targeted Lexical Injection on XTREME-R

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does the cross-lingual alignment performance of Lugha-Llama with Targeted Lexical Injection compare to Lugha-Llama fine-tuned using other LoRA-based methods on the XTREME-R benchmark. 10 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Targeted Lexical Injection: Unlocking Latent Cross-Lingual Alignment in Lugha-Llama via Early-Layer LoRA Fine-Tuning. Research question: How does the cross-lingual alignment performance of Lugha-Llama with Targeted Lexical Injection compare to Lugha-Llama fine-tuned using other LoRA-based methods on the XTREME-R benchmark?.

2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.2/10.

3 Results

12 papers retrieved. 10 claims extracted; 2 independently verified. Quality review score: 4.2/10.

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
Layer 0 (input embeddings) showed a modest average cosine similarity of approximately 0.3153.	×	0.07
Layer 1 showed an average cosine similarity of 0.9808.	×	0.09
Layer 2 exhibited the peak average cosine similarity, reaching 0.99998.	×	0.08
Layer 31 showed an average similarity of 0.9876 in the pilot scan.	×	0.04
The baseline output similarity observed on the full evaluation set was approximately 0.32.	×	0.09
The average cosine similarity at the final output layer (Layer 31) of the base model was approximately 0.3211 for the tr	✓	0.16
The model uses Lughu-Llama-8B-wura as the base model.	×	0.09
Lughu-Llama-8B-wura is an open-source LLM specifically adapted for several African languages, including Swahili, built u	×	0.11
The model is loaded in 4-bit precision using bitsandbytes with NF4 quantization and torch.bfloat16 as the compute data t	×	0.02
The pilot study revealed that Lughu-Llama-8B-wura inherently achieves very high lexical alignment in its early layers, p	✓	0.19

References

- <http://arxiv.org/abs/2506.15415v1>
- <http://arxiv.org/abs/2106.09063v4>
- <http://arxiv.org/abs/2204.06487v3>