

Targeted Lexical Injection Effects on Lugha-Llama Zero-Shot Cross-Lingual Alignment in African Languages

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: What is the impact of Targeted Lexical Injection on the zero-shot cross-lingual alignment performance of Lugha-Llama when evaluated on the FLORES-200 benchmark for other low-resource African languages. 9 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. 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: What is the impact of Targeted Lexical Injection on the zero-shot cross-lingual alignment performance of Lugha-Llama when evaluated on the FLORES-200 benchmark for other low-resource African languages?.

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 4.5/10.

3 Results

12 papers retrieved. 9 claims extracted; 1 independently verified. Quality review score: 4.5/10.

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.	×	0.04
The baseline output similarity observed on the full evaluation set was approximately 0.32.	×	0.09
The average similarity at the final output layer (Layer 31) of the base model was approximately 0.3211 for the trained s	×	0.15
The base model used is Lugha-Llama-8B-wura, an open-source LLM adapted for several African languages, including Swahili,	×	0.10
The model is loaded in 4-bit precision using bitsandbytes with NF4 quantization and torch.bfloat16 as the compute data t	×	0.03
The pilot study revealed that Lugha-Llama-8B-wura inherently achieves very high lexical alignment in its early layers, p	✓	0.20

References

- <http://arxiv.org/abs/2605.17152v1>
- <http://arxiv.org/abs/2506.15415v1>
- <http://arxiv.org/abs/2603.19273v1>