

# Impact of TLI Early-Layer LoRA Fine-Tuning on Lugha-Llama Robustness Against Adversarial Lexical Perturbations in Low-Resource

Assignee Research

June 15, 2026

## Abstract

Large Language Models (LLMs) have demonstrated remarkable capabilities, yet their performance in low-resource languages (LRLs), such as Swahili, often lags due to data scarcity and underrepresentation in pre-training. A key challenge is achieving robust cross-lingual lexical alignment, crucial for tasks like translation and cross-lingual information retrieval. This paper introduces Targeted Lexical Injection (TLI), a novel and efficient fine-tuning approach. We first demonstrate that Lugha-Llama-8B-wura, a Swahili-centric LLM, exhibits strong, near-perfect lexical alignment for Swahili-English

## 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 TLI early-layer LoRA fine-tuning on the robustness of Lugha-Llama against adversarial lexical perturbations in low-resource Bantu languages compared to standard fine-tuning?.

## 2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.5/10.

## 3 Results

8 papers retrieved. 9 claims extracted; 8 independently verified. Quality review score: 8.5/10.

## 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.	✓	0.25
Layer 1 exhibited an average cosine similarity of 0.9808.	✓	0.16
Layer 2 exhibited the peak average cosine similarity, reaching 0.99998.	✓	0.21
Layer 31 showed an average similarity of 0.9876 in the pilot scan.	✓	0.19
The baseline output similarity observed on the full evaluation set was approximately 0.32.	×	0.13
The average cosine similarity at the final output layer (Layer 31) of the base model was approximately 0.3211 for the tr	✓	0.33
Lugha-Llama-8B-wura is an open-source LLM specifically adapted for several African languages, including Swahili, built u	✓	0.28
The model is loaded in 4-bit precision using bitsandbytes with NF4 quantization and torch.bfloat16 as the compute data t	✓	0.25
The pilot study revealed that Lugha-Llama-8B-wura inherently achieves very high lexical alignment in its early layers, p	✓	0.33

## References

- <http://arxiv.org/abs/2504.15610v3>
- <http://arxiv.org/abs/2504.07887v2>
- <http://arxiv.org/abs/2506.15415v1>