

# Morphological and Syntactic Drivers of Performance Drops in Zero-Shot Cross-Lingual Transfer on XTREME-R

Assignee Research

June 26, 2026

## Abstract

We present a method for discovering morphological features in low-resource Bantu languages by combining cross-lingual transfer learning with unsupervised clustering. Applied to Giriama (nyf), a language with only 91 labeled paradigms, our pipeline discovers noun class assignments for 2,455 words and identifies two previously undocumented morphological patterns: an a- prefix variant for Class 2 (vowel coalescence - the merger of two adjacent vowels - of wa-, 95.1% consistency) and a contracted k'- prefix (98.5% consistency). External validation on 444 known Giriama verb paradigms confirms 78.2%

## 1 Introduction

This paper examines: Zero-Shot Morphological Discovery in Low-Resource Bantu Languages via Cross-Lingual Transfer and Unsupervised Clustering. Research question: Which specific morphological or syntactic features contribute most to performance drops in zero-shot cross-lingual transfer for low-resource languages on the XTREME-R suite?.

## 2 Methodology

Systematic literature search across multiple databases yielded 16 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

16 papers retrieved. 19 claims extracted; 15 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
The ensemble pipeline discovers noun class assignments for 2,455 words in the Giriama corpus (7,812 sentences), represent	✓	0.21
Transfer learning contributes 8,698 predictions with a mean confidence of 0.71.	✓	0.19
Unsupervised clustering contributes 18,508 predictions.	×	0.14
The high-confidence ensemble retains 5,279 predictions.	×	0.14
Transfer-clustering agreement on Giriama is 36.7%.	✓	0.17
Agreement is highest for morphologically transparent features (non-finite forms 78.3%, present tense 61.2%) and lowest for	✓	0.27
Unsupervised clustering identified Cluster 1 (266 words, 95.1% consistency) using the a- prefix variant in Giriama.	✓	0.24
This a- prefix variant pattern accounts for 19.6% of all Class 2 words in the Giriama corpus.	✓	0.20
Clustering identified Cluster 8 (206 words, 98.5% consistency) with k'- (apostrophe = elision) in Giriama.	✓	0.21
The BantuMorph model achieves 78.2% lemmatization accuracy on 444 known Giriama verb paradigms (347/444).	✓	0.23
The BantuMorph model achieves 54.3% inflection (completion) accuracy on 444 known Giriama verb paradigms (241/444).	✓	0.21
The model was not trained on the 444 known Giriama verb paradigms used for external evaluation.	✓	0.19
Bantu languages exhibit rich agglutinative morphology with noun class systems.	✓	0.20
Each noun in Bantu languages belongs to one of 15–20 classes marked by prefixes that trigger agreement on verbs, adjecti	✓	0.20
Most of the 500+ Bantu languages remain understudied computationally.	✓	0.15
Contextualized embeddings capture morphosyntactic information.	×	0.10
Character-level models handle morphological variation naturally.	✓	0.16
Cross-lingual embeddings enable transfer.	×	0.09
ByT5 operates at the character level and has shown strong cross-lingual transfer for morphologically rich languages.	✓	0.23

## References

- <http://arxiv.org/abs/2406.09549v2>
- <http://arxiv.org/abs/2604.22723v1>
- <http://arxiv.org/abs/2303.02357v1>