

Vocabulary Augmentation for Robust Cross-Lingual Universal Dependency Parsing Under Domain Shift in Low-Resource Settings

Assignee Research

June 16, 2026

Abstract

Pretrained multilingual language models have become a common tool in transferring NLP capabilities to low-resource languages, often with adaptations. In this work, we study the performance, extensibility, and interaction of two such adaptations: vocabulary augmentation and script transliteration. Our evaluations on part-of-speech tagging, universal dependency parsing, and named entity recognition in nine diverse low-resource languages uphold the viability of these approaches while raising new questions around how to optimally adapt multilingual models to low-resource settings.

1 Introduction

This paper examines: Specializing Multilingual Language Models: An Empirical Study. Research question: Does vocabulary augmentation improve cross-lingual transfer robustness for Universal Dependency Parsing under domain shift conditions in low-resource settings?.

2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.6/10.

3 Results

14 papers retrieved. 14 claims extracted; 11 independently verified. Quality review score: 7.6/10.

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
Chau et al. (2020) augment the model’s vocabulary to more effectively tokenize text and then pretrain on a small amount	✓	0.28
Chau et al. (2020) report significant performance improvements on a small set of low-resource languages.	✓	0.20
Muller et al. (2021) propose transliterating text in the target language to Latin script to be better tokenized by the e	✓	0.24
Muller et al. (2021) observe mixed results and note that transliteration quality may be a confounding factor.	✓	0.20
The study verifies the performance of vocabulary augmentation on three tasks across nine low-resource languages using th	✓	0.19
Performance gains from vocabulary augmentation are associated with improved vocabulary coverage of the target language.	✓	0.20
There is a negative interaction between vocabulary augmentation and transliteration methods.	×	0.14
Vocabulary augmentation offers an appealing balance of performance and cost.	✓	0.21
The study expands on Chau et al. (2020) by evaluating named entity recognition and part-of-speech tagging in addition to	✓	0.20
The study computes CWR for each token as a weighted sum of the activations at each MBERT layer.	✓	0.18
In the benchmark tables, the VA method achieved a score of 95.28 \pm 0.51 on one metric compared to LAPT’s 95.74 \pm 0.44.	×	0.12
In the benchmark tables, the VA method achieved a score of 73.22 \pm 1.23, outperforming MBERT (71.83 \pm 0.90) and LAPT (72	✓	0.16
In the benchmark tables, the VA method achieved a score of 68.93 \pm 3.30 on a specific metric, outperforming BERT (54.64	✓	0.15
In the benchmark tables, the VA method achieved an overall score of 83.74, compared to 81.72 for LAPT and 78.46 for MBER	×	0.14

References

- <http://arxiv.org/abs/2106.09063v4>
- <http://arxiv.org/abs/2511.20872v1>
- <http://arxiv.org/abs/2204.08143v2>