

Cross-Lingual Transfer Performance of DPO-Aligned and SFT-Only OPT-350M in Multilingual Hate Speech Detection

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 5 peer-reviewed papers addressing the following research question: How does the cross-lingual transfer performance of DPO-aligned OPT-350M models compare to SFT-only variants on multilingual hate speech detection tasks within the XTREME-R benchmark. 9 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: One to rule them all: Towards Joint Indic Language Hate Speech Detection. Research question: How does the cross-lingual transfer performance of DPO-aligned OPT-350M models compare to SFT-only variants on multilingual hate speech detection tasks within the XTREME-R benchmark?.

2 Methodology

Systematic literature search across multiple databases yielded 5 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.0/10.

3 Results

5 papers retrieved. 9 claims extracted; 0 independently verified. Quality review score: 4.0/10.

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
Fine-tuned multilingual models beat the baselines by at least % on the hate speech detection task.	×	0.14
XLM-RoBERTa outperformed similar multilingual Transformer models such as mBERT and distilmBERT on the hate speech detect	×	0.13
XLM-RoBERTa secured the 1st position among 24 participants and the 5th position among 34 participants on the HASOC 2021	×	0.06
Using SOUP (Similarity-based Oversampling and Undersampling processing) resulted in a drop of 5% in accuracy compared to	×	0.05
Data augmentation using back-translation increased the dataset size by three times but did not result in performance gai	×	0.02
Applying machine learning algorithms like random forest and LightGBM resulted in an average drop of 5.3% in performance.	×	0.02
The tweet-preprocessor and ekphrasis libraries were used for preprocessing tweet data and hashtags.	×	0.02
NeuralSpace’s transliteration tool and langdetect library were used to extract pure Hindi and Marathi text within the tw	×	0.03
Ekphrasis segmenter was used to segment hashtag text into constituent and meaningful tokens for feature extraction.	×	0.05

References

- <http://arxiv.org/abs/2112.09986v1>

- <http://arxiv.org/abs/2201.04227v1>
- <http://arxiv.org/abs/2109.13711v1>