

# Multilingual vs. Monolingual Training Robustness in Low-Resource Hate Speech Detection

Assignee Research

June 9, 2026

## Abstract

This report synthesises findings from 7 peer-reviewed papers addressing the following research question: Does joint multilingual training for hate speech detection degrade per-language robustness against adversarial perturbations in low-resource settings compared to monolingual fine-tuning. 13 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Leveraging Multilingual Transformers for Hate Speech Detection. Research question: Does joint multilingual training for hate speech detection degrade per-language robustness against adversarial perturbations in low-resource settings compared to monolingual fine-tuning?.

## 2 Methodology

Systematic literature search across multiple databases yielded 7 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.2/10.

## 3 Results

7 papers retrieved. 13 claims extracted; 1 independently verified. Quality review score: 4.2/10.

## 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 evaluation metric used throughout the study is the macro F1-score.	×	0.07
Perspective API features combined with a multi-layer perceptron classifier provide respectable results on hate and offen	✓	0.16
In monolingual mode, identity activation is the most effective MLP hidden layer activation setting for the English task.	×	0.02
In monolingual mode, tanh activation is the most effective MLP hidden layer activation setting for the German task.	×	0.02
German Task 2 benefits from the multilingual mode due to additional data from English training examples allowing better	×	0.05
A drop in English results is observed in multilingual mode, potentially due to a reduction in the number of available fe	×	0.04
The study utilized the python library 'tweet-preprocessor' for tweet tokenization.	×	0.03
The study utilized the python library 'ekphrasis' for hashtag segmentation.	×	0.01
For Hindi tweets, tokenization was performed on whitespaces and symbols including colons, commas, and semicolons.	×	0.02
Preprocessing for Hindi tweets involved the removal of hashtags, smileys, emojis, URLs, mentions, numbers, and reserved	×	0.02
Hashtag text was segmented into meaningful tokens using the ekphrasis segmenter trained on the twitter corpus.	×	0.03
The study initially experimented with the 'emot5' python library to obtain textual descriptions of emojis.	×	0.01
The study ultimately chose to utilize 'emoji2vec' to obtain semantic vectors representing emojis instead of textual desc	×	0.02

## References

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

- <http://arxiv.org/abs/2101.03207v1>
- <http://arxiv.org/abs/2109.13711v1>