

Scaling XGLM Performance on Indonesian Language Tasks with Model Size

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How does the accuracy of XGLM-564M on Indonesian language understanding tasks (e.g., XNLI, PAWS-X) scale with increasing model size compared to English. 19 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Improving Indonesian Text Classification Using Multilingual Language Model. Research question: How does the accuracy of XGLM-564M on Indonesian language understanding tasks (e.g., XNLI, PAWS-X) scale with increasing model size compared to English?.

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 5.0/10.

3 Results

14 papers retrieved. 19 claims extracted; 2 independently verified. Quality review score: 5.0/10.

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
LSTM with word embedding enhanced with paragraph vector was used by Crisdayanti and Purwarianti in their experiments on	×	0.04
Random forest with unigram as its feature was used by Ibrohim and Budi in their experiments on multi-label hate speech a	×	0.06
The largest multilingual model showed competitive performance and even outperformed monolingual language models in vario	×	0.11
There is no research in Indonesian text classification using multilingual language models, especially in sentiment analy	✓	0.38
The combination of word embedding with RNN, LSTM, or GRU as the model has been very successful but still possesses some	×	0.03
BERT uses Transformer architecture and is pre-trained on millions of texts with masked language model (MLM) objective.	×	0.04
BERT has proven massively successful, topping various benchmarks and significantly improving performance from previously	×	0.06
Labeled Indonesian text data is scarce in comparison to English text data.	✓	0.16
Cross-lingual representation of text has enabled models to do transfer learning across languages.	×	0.12
There are two main methods of producing cross-lingual representation on shared vector space: alignment and joint optimiz	×	0.07
Farhan and Khodra used reviews crawled from TripAdvisor for their sentiment analysis dataset.	×	0.07
Crisdayanti and Purwarianti used text from Twitter, Zomato, TripAdvisor, Facebook, Instagram, and Qraved for their senti	×	0.05
Yelp Review dataset contains 299000 positive and 299000 negative samples.	×	0.05
Ibrohim and Budi crawled tweets from Twitter and annotated the text with the help of 30 diverse annotators for their hat	×	0.04
Jigsaw Toxic Comment dataset contains 152111 hate speech and 1750083 normal samples.	×	0.03
XLNet achieved an average gain of 0.176221 with 500 data points.	×	0.03
mBERT achieved an average gain of 0.129394 with 500 data points.	×	0.03
XLNet achieved an average gain of 0.077875 with MAX data points.	×	0.03
mBERT achieved an average gain of 0.020184 with MAX data points.	×	0.03

References

- <http://arxiv.org/abs/2009.05713v1>
- <http://arxiv.org/abs/2310.09917v3>
- <http://arxiv.org/abs/2207.08179v1>