

Pre-training Data Size Variation in mT5 and Cross-Lingual Transfer Performance on the XTREME-R Benchmark

Assignee Research

July 6, 2026

Abstract

Multi-lingual language models (LM), such as mBERT, XLM-R, mT5, mBART, have been remarkably successful in enabling natural language tasks in low-resource languages through cross-lingual transfer from high-resource ones. In this work, we try to better understand how such models, specifically mT5, transfer *any* linguistic and semantic knowledge across languages, even though no explicit cross-lingual signals are provided during pre-training. Rather, only unannotated texts from each language are presented to the model separately and independently of one another, and the model appears to implicitly

1 Introduction

This paper examines: Languages You Know Influence Those You Learn: Impact of Language Characteristics on Multi-Lingual Text-to-Text Transfer. Research question: What is the impact of varying the pre-training data size for each language in mT5 on cross-lingual transfer performance in the XTREME-R benchmark, measured by accuracy and language-specific F1 scores?.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.2/10.

3 Results

13 papers retrieved. 10 claims extracted; 9 independently verified. Quality review score: 8.2/10.

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 Exact-Match accuracy metric LMEM(L) is defined as the average of the indicator function $1(x_i = \hat{x}_i)$ over all masked	✓	0.16
The mT5 framework is a multi-lingual adaptation of T5, which formulates any NLP task as sequence generation.	✓	0.18
The T5 framework abstracts away the output feature engineering from meaningless indexes to meaningful language tokens.	✓	0.25
The architecture of T5 is a Transformer encoder-decoder, pre-trained with a span-masking objective inspired by the BERT	✓	0.21
The cross-lingual analysis is conducted on the base version of mT5.	✓	0.16
The analysis includes languages such as Arabic, Bengali, English, Finnish, Indonesian, Russian, Swahili, Spanish, German	✓	0.22
Each task (XNLI, NER, QA) gets at least 7 languages, with a detailed list provided in the Appendix in Table 4.	✓	0.18
The XNLI dataset is used for Natural Language Inference (NLI), the PANX dataset for Name-Entity Recognition (NER), and t	✓	0.26
The Pearson correlation between features and cross-lingual transfer performance is reported in Table 1 for XNLI, QA, and	✓	0.20
The benchmark tables show performance metrics for different language pairs and tasks, with values indicating accuracy or	×	0.05

References

- <http://arxiv.org/abs/2410.21676v4>
- <http://arxiv.org/abs/2212.01757v1>
- <http://arxiv.org/abs/2102.12407v1>