

Scaling of Source Languages and Cross-Lingual Transfer in Multi-Level Contrastive Models

Assignee Research

June 19, 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 scaling relationship between the number of source languages in the training set and the cross-lingual transfer performance of multi-level contrastive models on low-resource slots in X-TREME-SLU?.

2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.5/10.

3 Results

11 papers retrieved. 17 claims extracted; 13 independently verified. Quality review score: 7.5/10.

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 $1/ X * \sum_{x,s} \prod_{i \in s} 1(x_i = s_i)$, where x_i	✓	0.29
The span masking procedure follows the pre-training span masking procedure defined in [XCR+20].	✓	0.17
The statistics LML(L) and LMEM(L) are estimated on the training dataset of each task.	×	0.12
The analysis focuses on the mT5 framework, a multi-lingual adaptation of T5 [RSR+19].	✓	0.18
T5 formulates any NLP tasks as sequence generation.	×	0.15
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 [VSP+17] encoder-decoder, pre-trained with a span-masking objective closely insp	✓	0.28
The cross-lingual analysis is run on the base version of mT5.	×	0.15
The analysis is conducted on Arabic, Bengali, English, Finnish, Indonesian, Russian, Swahili, Spanish, German, and Hindi	✓	0.27
Not all languages have training data for all the three tasks, but each task gets at least 7 languages.	✓	0.16
Each language is used both as a source language (S) and as a target language (T), leading to up to 90 language pairs.	✓	0.22
The tasks focused on are Natural Language Inference (NLI), Name-Entity Recognition (NER), and Question Answering (QA).	✓	0.25
For NLI, the XNLI dataset [CRL+18] is used.	✓	0.17
For NER, the PANX dataset [GL17] is used.	✓	0.16
For QA, the TyDiQA dataset [CCC+20] is used.	✓	0.15
The Pearson correlation between the features introduced and the cross-lingual transfer performance in the zero-shot sett	✓	0.17
The benchmark tables show performance metrics for various language pairs and tasks.	×	0.06

References

- <http://arxiv.org/abs/2212.01757v1>
- <http://arxiv.org/abs/1810.12836v4>
- <http://arxiv.org/abs/2304.01492v5>