

Analyzing the effect of linguistic similarity on cross-lingual transfer: Tasks and experimental setups matter

Verena Blaschke,¹ Masha Fedzechkina,² Maartje ter Hoeve² ACL Findings | ¹LMU Munich & MCML, ²Apple | July 2025

Motivation

- 7000+ languages in the world
 - Only a few of which are focused on in NLP research (Joshi+, 2020)
 - Only a few of which have training data available
- \rightarrow Cross-lingual transfer



The State and Fate of Linguistic Diversity and Inclusion in the NLP World (Joshi et al., ACL 2020)



Motivation

- Given a target language, how do I select a good training language? Intuitively: pick a language/dataset that is in some way similar...
 - ... but in what way?
 - Prior work: either relatively few languages or NLP tasks (Philippy+, 2023)
 - Here: 263 languages, 3 NLP tasks, 10 similarity measures





Overview

Large-scale cross-lingual transfer experiments Correlations with similarity measures Practical takeaways for picking source languages

Overview

Large-scale cross-lingual transfer experiments Correlations with similarity measures Practical takeaways for picking source languages

NLP experiments

POS tagging & dep. parsing

Topic classification

NLP experiments — grammatical tasks

POS tagging & dep. parsing

- Universal Dependencies
- 70 × 153 languages
- UDPipe 2 (mono- and multilingual) char/word embeddings)



Topic classification

NLP experiments — grammatical tasks

POS tagging & dep. parsing

- Universal Dependencies
- 70 × 153 languages
- UDPipe 2 (mono- and multilingual char/word embeddings)



LAS = labelled attachment score

NLP Experiments — topic classification

POS tagging & dep. parsing

- Universal Dependencies
- 70 × 153 languages
- UDPipe 2 (mono- and multilingual char/word embeddings)

Topic classification

- SIB-200
- 194 × 194 languages
- MLPs → comparable & competitive
- Input representations
 - character n-grams
 - n-grams (transliterated text)
 - mBERT embeddings



NLP Experiments — topic classification



SIB-200: A Simple, Inclusive, and Big Evaluation Dataset for Topic Classification in 200+ Languages and Dialects (Adelani et al., EACL 2024)

Topic classification

- SIB-200
- 194 × 194 languages
- MLPs \rightarrow comparable & competitive
- Input representations
 - character n-grams
 - n-grams (transliterated text)
 - mBERT embeddings



Overview

Large-scale cross-lingual transfer experiments Correlations with similarity measures Practical takeaways for picking source languages

Similarity measures

- Linguistic measures

 - **Lex**ical similarity
 - Phylogenetic relatedness
 - **Geo**graphic proximity
- Dataset measures
 - Character overlap
 - Word* overlap (words, character trigrams, subword tokens)
 - **Size** of training split

Structural similarities: Grammar, syntax, phonology, phoneme inventory



Parsing (labelled attachment score)

Pearson's r

0.8

Mean correlation (*r*) over test languages



size pho inv geo syn gram gen lex char word*



POS tagging (accuracy)



Topic classification (accuracy) — MLP with n-grams

Pearson's r





Topic classification (accuracy) — MLP with n-grams (transliterated)

Pearson's r





Topic classification (accuracy) — MLP with mBERT embeddings



Overview

Large-scale cross-lingual transfer experiments Comparing transfer trends Practical takeaways for picking source languages

Picking source languages based on similarity measures

	size	pho	inv	geo	syn	gram	gen	lex	char	word*
Top-1 source candidate	(<i>=</i> mc	ost sin	nilar l	angu	age)					
POS	29	15	14	15	10	12	9	10	15	12
Parsing (LAS)	21	13	13	13	7	10	8	8	16	11
Topics (n-grams)	—	17	17	13	15	14	9	9	13	4
Topics (n-grams, translit	.) —	13	13	11	11	10	7	7	20	3
Topics (mBERT)	—	12	11	10	9	8	8	8	12	9

Mean performance loss in percentage points if picking the best training language according to one measure (instead of the overall best one)

Picking source languages based on similarity measures

	cizo	nho	inv	000		aram	aon		ohar	word	
	SIZE	μιο		geo	Syll	gram	gen		Chai	word	
Top-1 source candidate (= most similar language)											
POS	29	15	14	15	10	12	9	10	15	12	
Parsing (LAS)	21	13	13	13	7	10	8	8	16	11	
Topics (n-grams)	-	17	17	13	15	14	9	9	13	4	
Topics (n-grams, translit.) —	13	13	11	11	10	7	7	20	3	
Topics (mBERT)	-	12	11	10	9	8	8	8	12	9	
Top-3 source candidates											
POS	25	7	7	5	3	4	5	4	7	5	
Parsing (LAS)	18	8	7	5	3	4	4	3	8	6	
Topics (n-grams)	-	10	9	5	6	6	5	3	8	2	
Topics (n-grams, translit.) —	8	7	4	5	5	3	3	14	1	
Topics (mBERT)	—	6	5	4	4	4	4	4	5	6	

Mean performance loss in percentage points if picking the best training language according to one measure (instead of the overall best one)

Picking source languages based on other transfer experiments?

Mean performance loss in percentage points if picking the best training language according to the results of another transfer experiment (instead of the overall best one)

Conclusions

- Different tasks/input representations
 → different similarity measures matter
- Selecting training languages based on relevant similarity measures (or on similar experiments) works well
 - If possible: compare multiple promising training languages

More details in the paper :)

TM and © 2025 Apple Inc. All rights reserved.

