



# Language Models Are Polyglots

Language Similarity Predicts Cross-Lingual  
Transfer Learning Performance

---

Juuso Eronen, **Michal Ptaszynski**, Tomasz Wicherkiewicz, Robert Borges, Katarzyna Janic, Zhenzhen Liu, Tanjim Mahmud & Fumito Masui

**7,000+ languages.**

**Most NLP tools work for ~20.**

### **High-resource**

English, Chinese, etc.  
Abundant labeled data

### **Low-resource**

Most of the world's  
languages lack data

### **The gap**

Cross-lingual transfer  
is the bridge

# Everyone just picks English.

But is it always the best source language for transfer?

## Current approach

- Pick English by default
- Or choose source intuitively
- No principled method
- Large performance gaps ignored

## Our approach

- Measure linguistic similarity
- Use it to predict transfer
- Principled source selection
- Better results, less guessing



## Core Hypothesis

**The more similar two languages are,  
the better one transfers to the other.**

---

If true → we can pick optimal source languages  
systematically, not by pure intuition

# Three independent lines of evidence

We triangulate using causal inference principles



## Similarity Metrics

qWALS, lang2vec,  
eLinguistics, EzGlot



## Transfer Performance

mBERT & XLM-R  
4 tasks, 8 languages



## Expert Survey

3 professional  
linguists, 769 Qs

# The World Atlas of Language Structures

2,662

Languages

192

Features

~12%

Populated

## Covers phonology, morphology, syntax & lexicon

- Problem: Most data is missing — English has ~150 features, Danish only 58
- Previous metrics use only a handful of features from WALS (or none at all)
- **Our solution: Use ALL shared features for each language pair separately**

# Quantified WALS (qWALS)

A typology-based similarity metric from nearly 200 features



Select shared features  
for each language pair



Map values to 0–1  
(ordinal encoding)



Compute Manhattan  
distance, normalized



Symmetric distance  
metric for 2,000+ langs

## Key advantages over lang2vec

- Transparent: every feature is inspectable & interpretable
- Task-tunable: features can be re-weighted per downstream task
- Ordinal: preserves gradation (unlike one-hot encoding)

# Four similarity metrics

Metric	Focus	Scope	Limitation
eLinguistics	Consonant use	Phonetic only	Accuracy drops for distant languages
EzGlot	Lexical overlap	Vocabulary only	Many missing values, asymmetric
lang2vec	Multiple sources	Heterogeneous	Opaque weighting, not task-tunable
<b>qWALS (ours)</b>	~192 WALS features	Multi-domain	Sparse coverage for some pairs



**qWALS is the only metric that is both multi-domain AND task-tunable**

# 8 languages, 3 families

## Germanic

USGB English  
DE German  
DK Danish

## Slavic

PL Polish  
RU Russian  
HR Croatian

## Koreano-Japonic

JP Japanese  
KR Korean

*Within-family and cross-family transfer tested in zero-shot setting*

# 4 NLP tasks × 2 models

**DEP**

**Dependency Parsing**

Syntax-heavy

**NER**

**Named Entity Recognition**

Mixed syntax + semantics

**SA**

**Sentiment Analysis**

Semantics-heavy

**ALD**

**Abusive Language Identification**

Mixed + pragmatics

## Models

**mBERT**

**XLNet**

= 64 fine-tuned models, each tested on all 8 languages  
(4 tasks \* 2 models \* trained on 8 languages \* tested on 8 languages)



## Quick Quiz!

You want to build a cyberbullying detector for Polish.

You have no Polish training data.

Which language should you train on?

**A**

**English**

(most data available)

**B**

**German**

(also Germanic, like English)

**C**

**Russian**

(also Slavic, like Polish)

# Answer: C — Russian

Polish ← Russian transfer: F1 = 0.654 (XLM-R)

Polish ← English transfer: F1 = 0.468

F1 Scores for Abusive Language ID (source → Polish target)



# Transfer performance patterns



## Family matters

Same-family languages transfer best across most tasks



## XLM-R > mBERT

XLM-R outperforms on most tasks, except sentiment analysis



## DEP is hardest

Japanese/Korean → Indo-European: near-zero for dependency parsing

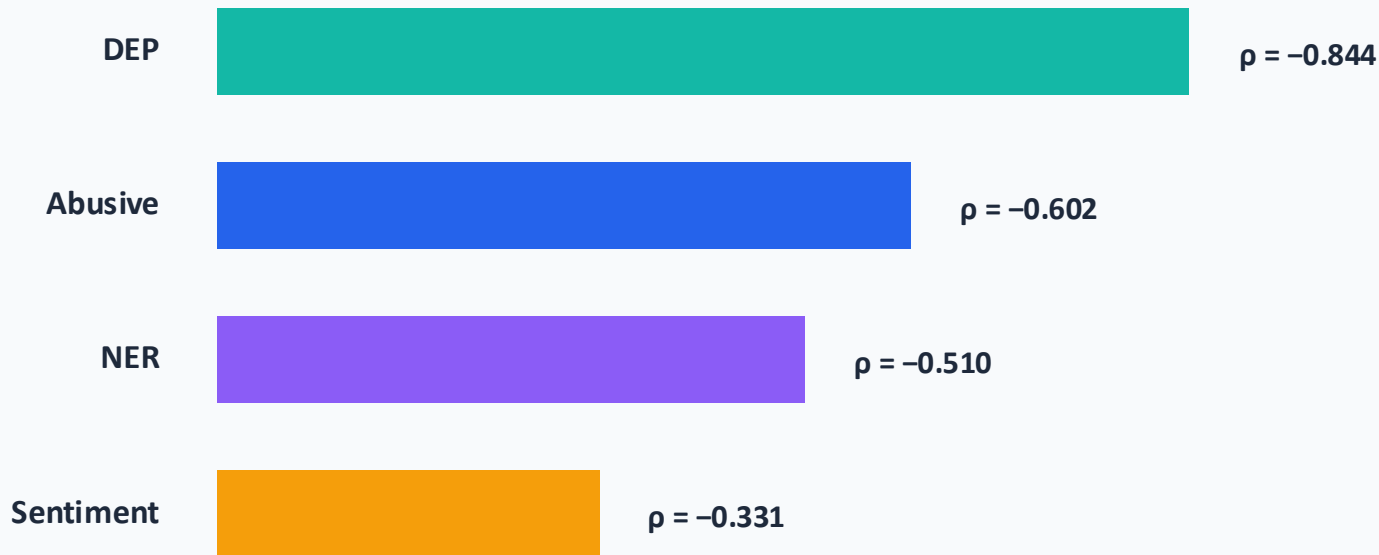


## SA is different

Sentiment analysis shows high transfer even across families

# Similarity strongly correlates with transfer

Spearman's  $\rho$  between qWALS (ordinal) and model performance (XLM-R)



All  $p$ -values  $< 0.001$  | Stronger correlation = more syntax-dependent task

Dependency parsing correlation with eLinguistics:  $\rho = -0.897$  (near-perfect!)

# Removing the monolingual "cheat codes"

Same-language pairs inflate correlations. What happens after removing them?

Task	$\rho$ (all pairs)	$\rho$ (zero-shot only)	Change
DEP	-0.844	<b>-0.769</b>	Still strong ✓
Abusive	-0.602	-0.415	Moderate
NER	-0.510	-0.316	Moderate
Sentiment	-0.331	-0.168	Weak

DEP stays strong — syntax-driven tasks genuinely depend on linguistic similarity.

Sentiment relies more on semantic/contextual cues that cross family boundaries.

# Optimizing qWALS per task

Leave-one-feature-out method: drop features that don't help

1. Remove one feature at a time → check correlation change
2. Drop the feature whose removal improves correlation most
3. Repeat until no further improvement

DEP	
Before	-0.771
After	<b>-0.990</b>
169 → 75 features	

Abusive	
Before	-0.646
After	<b>-0.822</b>
169 → 53 features	

NER	
Before	-0.561
After	<b>-0.808</b>
169 → 63 features	

Sentiment	
Before	-0.361
After	<b>-0.803</b>
169 → 21 features	

# -0.990

Pearson correlation for dependency parsing

With just 75 optimized WALS features,  
qWALS can near-perfectly predict  
which source language will transfer best.

**Syntax-heavy tasks → similarity matters most**

# Validating with human experts

3

Expert linguists  
(20-40 yrs experience)

769

Binary-choice  
questions

30

Languages  
covered

Question format:

*"Polish is more similar to..."*

**(A) Korean**      **(B) Latvian**

Simple binary choice — experts use ALL their knowledge  
(grammar, phonology, history, contact, geography)

# Experts agree with qWALS ~77% of the time

Comparison	Agreement	Cohen's $\kappa$
qWALS ↔ Linguist 1	76.5%	0.755
qWALS ↔ Linguist 2	73.0%	0.719
qWALS ↔ Linguist 3	71.3%	0.701
lang2vec ↔ Linguist 1	71.9%	0.707



qWALS consistently  
outperforms lang2vec  
in expert alignment

## Agreement by task (experts vs. best transfer language):

Abusive

**100%**

agreement

DEP

**76.9%**

agreement

NER

**69.2%**

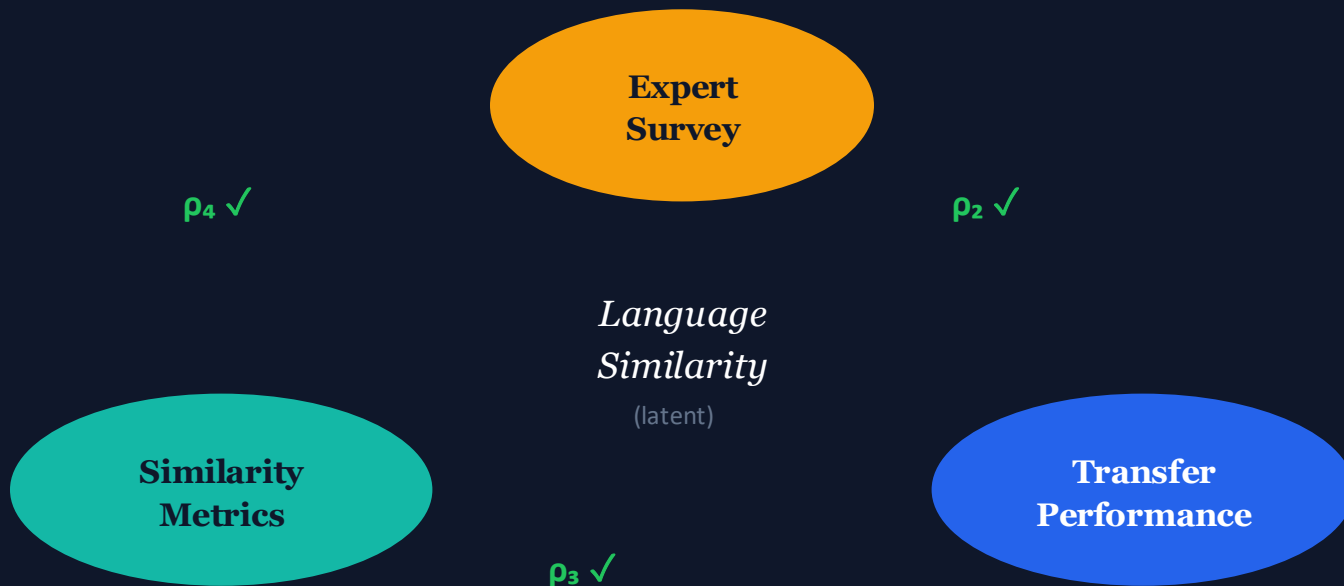
agreement

Sentiment

**61.5%**

agreement

# Triangulation: All three signals align



All three pairwise correlations are positive and significant  $\rightarrow$  supports causal role of similarity

# Don't default to English.

Use linguistic similarity to pick a smarter source language.

- ✓ qWALS: open-source, transparent, task-tunable similarity metric
- ✓ Covers 2,662 languages — including many low-resource ones
- ✓ Optimized feature sets per task improve prediction dramatically
- ✓ Expert validation confirms computational metrics match human judgment
- ✓ A similar source language usually outperforms English

**qWALS is released as open source — use it for your own research!**

# Honest caveats & next steps

## Limitations

- 8 languages, 3 families — need broader coverage (follow-up paper already in review)
- WALs sparsity: some pairs share <30% features
- No Romance languages tested yet
- Only encoder models (mBERT, XLM-R)
- Correlation  $\neq$  Causation (despite triangulation)

## Future directions

- Add Romance cluster (FR, ES, PT, IT)
- Integrate Grambank + PHOIBLE data
- Test on decoder-only LLMs (GPT-4, LLaMA)
- Use KANs for non-linear feature learning
- Continuous pretraining & probing studies
- Add online DEMO



# Thank You!

Questions, comments, ... DEMO!

---

[doi.org/10.3390/make8030065](https://doi.org/10.3390/make8030065)

[michal@mail.kitami-it.ac.jp](mailto:michal@mail.kitami-it.ac.jp) | [jeronen@mail.kitami-it.ac.jp](mailto:jeronen@mail.kitami-it.ac.jp)

**qWALS is open source — try it on your language pair!**



<https://pypi.org/project/qwals/>