

Enough is Enough! A Case Study on the Effect of Data Size for Evaluation Using Universal Dependencies

Rob van der Goot, Zoey Liu, Max Müller-Eberstein

Motivation

Too large evaluation datasets:

- ▶ less training data
- ▶ less coverage
- ▶ more costs

Setup

- ▶ Compare 33 multi-lingual language models for parsing
- ▶ Use 50k for training, rest for obtaining the “gold ranking”
- ▶ Get ranking for smaller sizes, use Kendall's Tau for comparing them

Splitting strategies

SEQ

T	T	T				D	D
---	---	---	--	--	--	---	---

RAND

T	T	T	D		D		
---	---	---	---	--	---	--	--

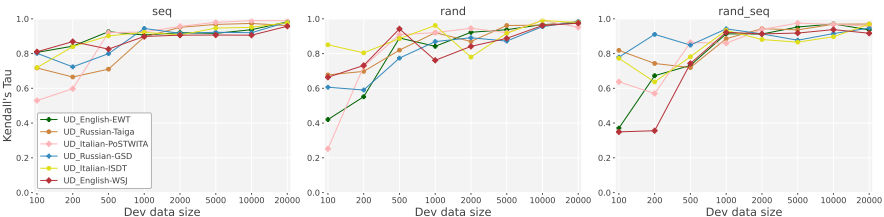
RAND_SEQ

T	T	T		D	D		
---	---	---	--	---	---	--	--

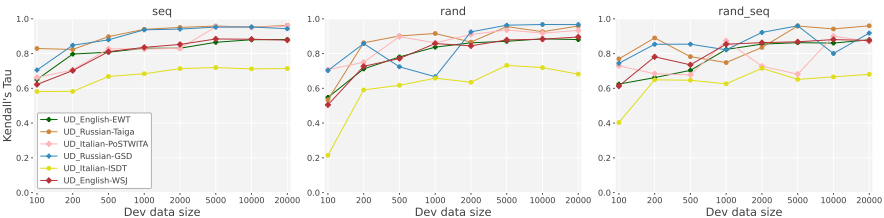
Data

Domain/Language	English	Italian	Russian
News	WSJ	ISDT	GSD
Web	EWT	PoSTWITA	Taiga

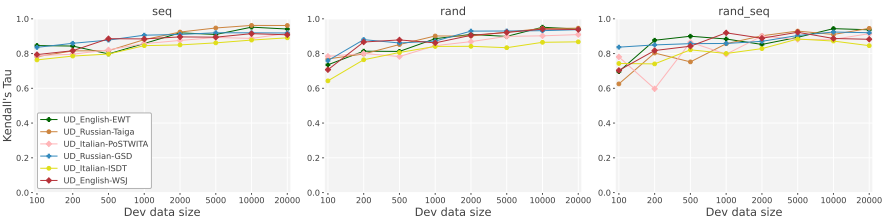
Results



Results Cross-domain



Results Cross-lingual



Thanks!

Code and predictions are available:

https://bitbucket.org/robvanderger/data_size/