

Does manipulating tokenization aid cross-lingual transfer?

A study on POS tagging for non-standardized languages

Verena Blaschke, Hinrich Schütze & Barbara Plank

VarDial @ EACL 2023

Center for Information and Language Processing, LMU Munich
Munich Center for Machine Learning



“We speak Alemannic dialects.”

Wir sprechen alemannische
Wir sprechen al, #emann,
##ische

Mundarten .
Mund, ##arten .

“We speak Alemannic dialects.”

M'r	redd	alemànnschi	Mundàrte	.
M, ', r	red, ##d	al, ##em, ##à,	Mund,	.
		##nn, ##isch, ##i	##à, ##rte	
Wir	sprechen	alemannische	Mundarten	.
Wir	sprechen	al, #emann,	Mund, ##arten	.
		##ische		

“We speak Alemannic dialects.”

M'r	redd	alemànnischi	Mundàrte	.
M, ', r	red, ##d	al, ##em, ##à,	Mund,	.
		##nn, ##isch, ##i	##à, ##rte	

Wir	sprechen	alemannische	Mundarten	.
Wir	sprechen	al, #emann,	Mund, ##arten	.
		##ische		

W(r	sprechen	alemaInische	Mundarten	.
W, (, r	sprechen	al, ##ema, ##In,	Mund, ##arten	.
		##ische		

Improving Zero-Shot Cross-lingual Transfer Between Closely Related Languages by Injecting Character-Level Noise

Noëmi Aepli¹ and Rico Sennrich^{1,2}

Do the previous findings on noise injection hold for new data/models?

Do different languages and PLMs have different ideal noise levels?

How does noise injection affect the subword tokenization differences between source and target data?

Zero-shot transfer for POS tagging

Aepli & Sennrich This work

Alemannic German → 2 Alemannic varieties

Dutch, German → Low Saxon

Faroese Nynorsk, Bokmål → 3 Norwegian dialects

Old French French → Picard

Spanish, French → Occitan

2 Karelian lects Finnish → 6 Finnish dialects

Modern Standard Arabic → 4 Arabic varieties

Zero-shot transfer for POS tagging

Aepli & Sennrich This work

Alemannic German → 2 Alemannic varieties

Dutch, German → Low Saxon

Faroese Nynorsk, Bokmål → 3 Norwegian dialects

Old French French → Picard

Spanish, French → Occitan

2 Karelian lects Finnish → 6 Finnish dialects

Modern Standard Arabic → 4 Arabic varieties

Monolingual BERTs/RobERTas vs. mBERT vs. XLM-R

Zero-shot transfer for POS tagging

Aepli & Sennrich This work

Alemannic German → 2 Alemannic varieties

Dutch, German → Low Saxon

Faroese Nynorsk, Bokmål → 3 Norwegian dialects

Old French French → Picard

Spanish, French → Occitan

2 Karelian lects Finnish → 6 Finnish dialects

Modern Standard Arabic → 4 Arabic varieties

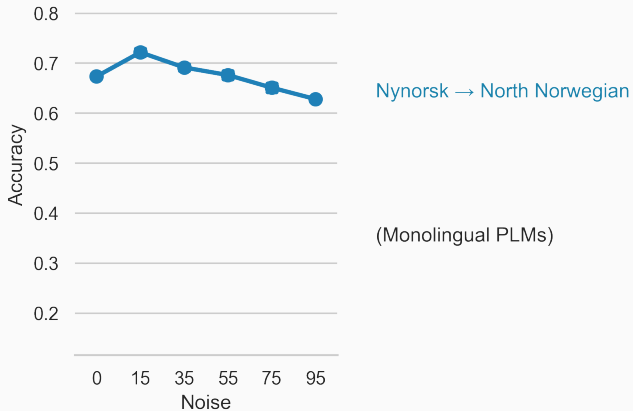
Monolingual BERTs/RoBERTas vs. mBERT vs. XLM-R

Noise levels: 0 / 15 / 35 / 55 / 75 / 95 %

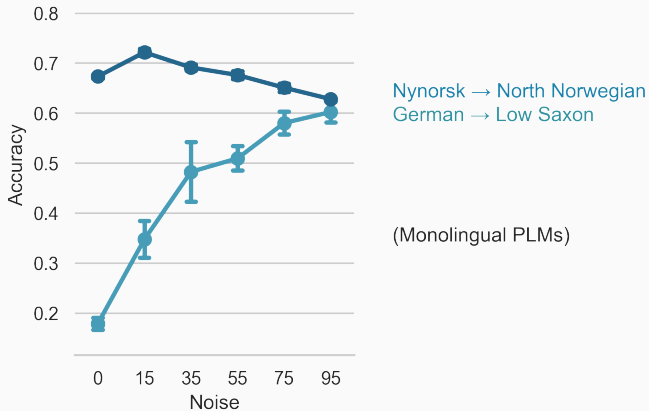
- Cross-lingual performance drop (language-dependent)

- Cross-lingual performance drop (language-dependent)
- PLM choice matters

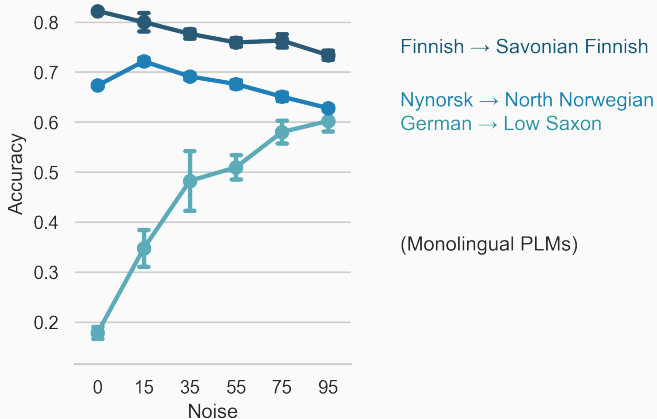
- One best noise level per set-up



- One best noise level per set-up
- No universal best noise level



- One best noise level per set-up
- No universal best noise level



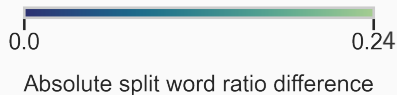
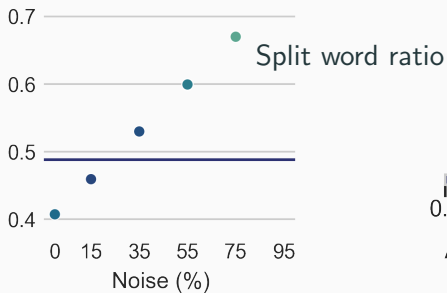
How does noise injection affect the subword tokenization differences between source and target data?

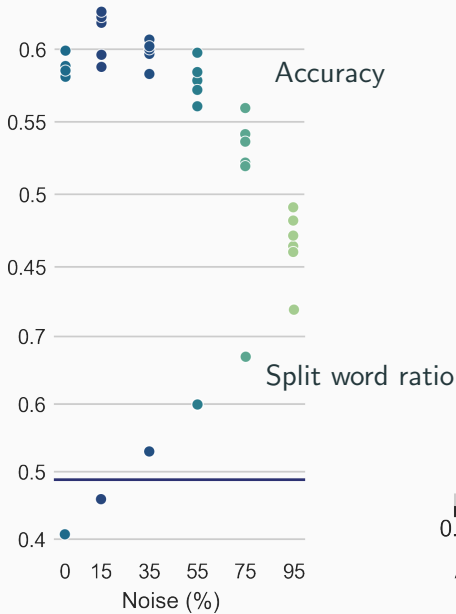
(Seen subwords, ...)

Split word ratio difference:

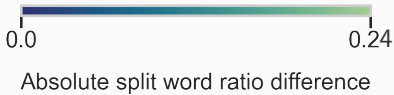
$$\left| \frac{\# \text{ words split in source data}}{\# \text{ words in source data}} - \frac{\# \text{ words split in target data}}{\# \text{ words in target data}} \right|$$

MSA → Egyptian Arabic
(mBERT)






MSA → Egyptian Arabic
(mBERT)



Recommendation:

- Don't want to tune the noise level as a hyperparameter?
Cheaply calculate the *split word ratio differences* for different noise levels and pick the noise level with the lowest difference
- Otherwise, find the best noise level by starting low and increasing the noise until the dev accuracy starts dropping

-  N. Aepli & R. Sennrich (2022). Improving zero-shot cross-lingual transfer between closely related languages by injecting character-level noise. *Findings of the Association for Computational Linguistics: ACL 2022*.

This research is supported by the European Research Council (ERC) Consolidator Grant DIALECT 101043235 and in parts Advanced Grant NonSequeToR 740516.