

UdS-DFKI@WMT20: Unsupervised MT and Very Low Resource Supervised MT for German↔Upper Sorbian

Sourav Dutta¹ Jesujoba O. Alabi^{1,3} Saptarashmi Bandyopadhyay²
Dana Ruiter¹ Josef van Genabith^{1,3}

¹Saarland University (UdS)

²University of Maryland

³DFKI GmbH

WMT 2020

Corpus	# Sentences
<i>Unsupervised</i>	
Sorbian Institute _{hsb}	339k
Witaj Sprachzentrum _{hsb}	222k
Web-scraped data _{hsb}	134k
Europarl _{de}	2,107k
News Commentary _{de}	422k
Wikipedia _{de}	833k
Wikipedia _{hsb}	76k
<i>Supervised</i>	
Parallel	60k
Europarl	568k
JW300	1,179k

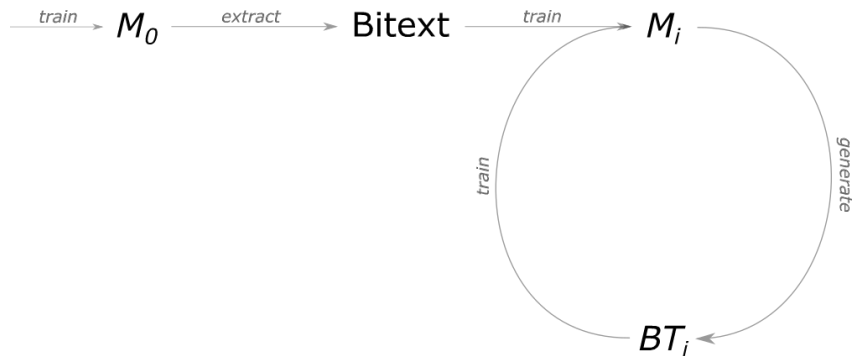
- ① Normalization, Deduplication, Tokenization, Truecasing
- ② **Factorization:** word token | lemma | POS
... hsb: *UDPipe* lemmatizer and POS-tagger
... de: *Snowball* stemmer + *spacy* POS-tagger
- ③ Byte-Pair Encoding (BPE) of 5k (default)
- ④ Quality token prefix (<low>, <medium>, <high>)
per sentence per corpus

- 1 German (de) \rightarrow Upper Sorbian (hsb)
- 2 Factorization of source text (German)

Model	BLEU		
	Dev20	DevTest20	Test20
de(fac.) \rightarrow hsb	44.34	41.15	40.90

Unsupervised Task

Method



Unsupervised Task

Results

Model	BLEU		
	Dev20	DevTest20	Test20
de \rightarrow hsb	13.60	9.90	10.30
hsb \rightarrow de	11.30	8.10	8.90
hsb(fac.) \rightarrow de	9.80	8.70	9.00

Thank you

Sourav Dutta

Saarland University, Germany
souravd@coli.uni-saarland.de