

The TALP-UPC Machine Translation System for the WMT19 News Translation Task: Pivoting Techniques for Low Resource MT

Noe Casas José A. R. Fonollosa Carlos Escolano
noe.casas@upc.edu jose.fonollosa@upc.edu carlos.escolano@upc.edu

Christine Basta Marta R. Costa-jussà
christine.raouf.saad.basta@tsc.upc.edu marta.ruiz@upc.edu

Task Description

- News translation: English \longleftrightarrow Kazakh
- Low resources ($\sim 100K$ parallel segments)
- Large amount of Russian-English data ($\sim 6M$).
- Large amount of Russian-Kazakh data ($\sim 4M$).

Corpora (after cleaning)

- Kazakh-English training data:

	sents.	words	vocab.	L_{max}	L_{mean}
Kazakh	99.6K	1.2M	139.6K	85	11.7
English		1.5M	85.3K	102	14.9

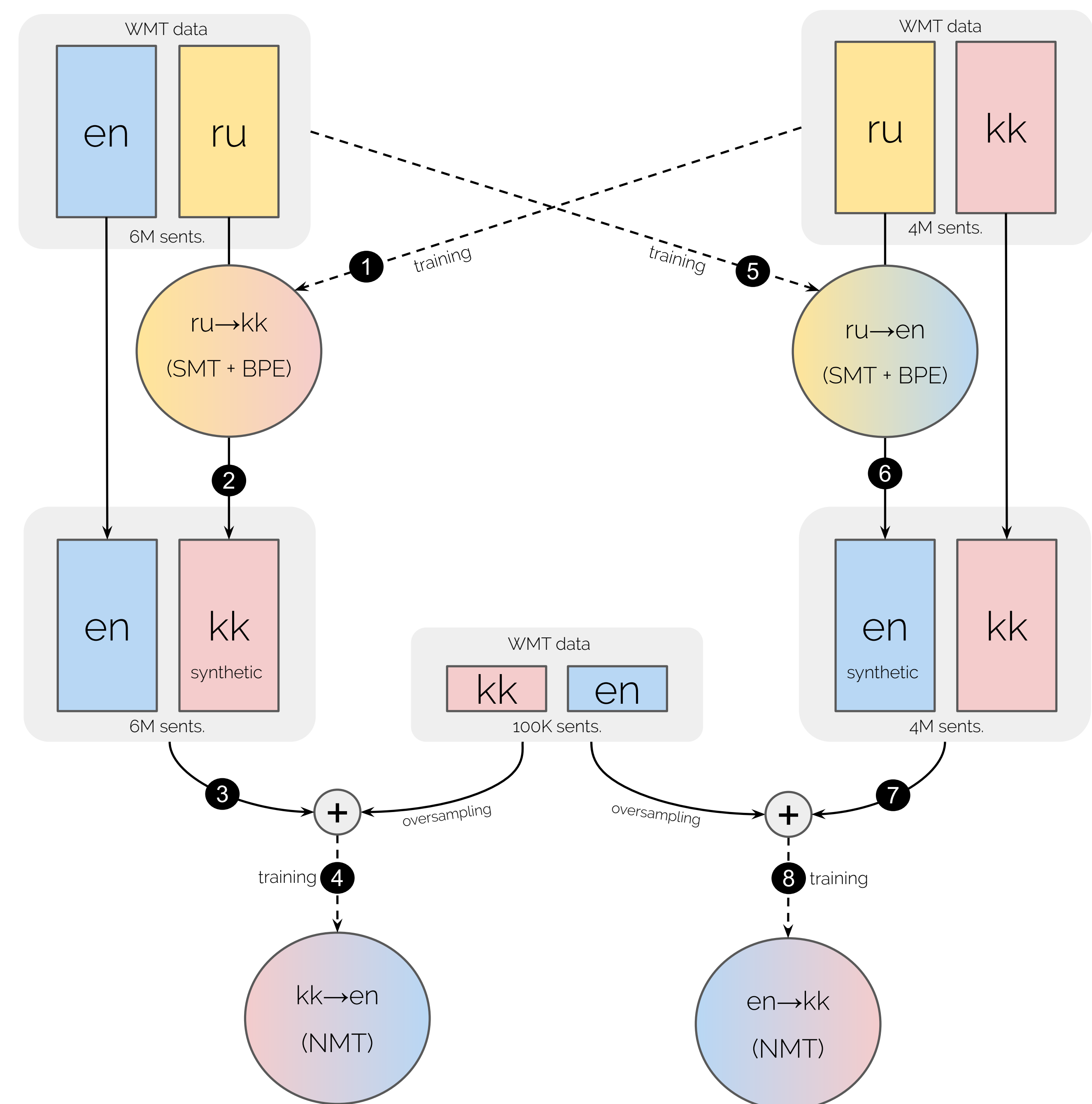
- English-Russian training data:

	sents.	words	vocab.	L_{max}	L_{mean}
Russian	6.1M	125.6M	3.2M	80	20.7
English		144.9M	2.0M	80	23.9

- Russian-Kazakh training data:

	sents.	words	vocab.	L_{max}	L_{mean}
Russian	4.2M	78.8M	1.4M	96	18.9
Kazakh		75.3M	1.6M	70	18.0

Strategy: pivoting through Russian



- 1 Train ru \rightarrow kk SMT engine (BPE vocab).
- 2 Create synthetic Kazakh-English corpus: use ru \rightarrow kk engine to translate into Kazakh the Russian side of the Russian-Kazakh data.
- 3 Combine synthetic and parallel data.
- 4 Train NMT system on the combined data.
- 5 Train ru \rightarrow en SMT engine (BPE vocab).
- 6 Create synthetic English-Kazakh corpus: use ru \rightarrow en engine to translate into English the Russian side of the Russian-Kazakh data.
- 7 Combine synthetic and parallel data.
- 8 Train NMT system on the combined data.

SMT Pivot Systems

- BPE to cope with morphologically-rich langs (32K merge operations).
- 5-gram order LM to compensate for subword segmentation.
- Phrase-based SMT with Moses.

NMT Systems

- BPE vocabulary with 32K merge operations.
- Joint source-target attention model (He et al., 2018; Fonollosa et al., 2019).
- Hyperparams: 14 layers, embedding dimensionality of 1024, feedforward expansion of dimensionality 4096 and 16 attention heads.

Evaluation (BLEU on held-out subset of dev)

	RBMT	SMT	SMT+BPE	NMT	NMT
Kazakh \rightarrow English	1.51	6.34	7.48	2.32	21.00
English \rightarrow Kazakh	1.46	3.53	3.82	1.42	15.47

parallel training data augmented data

Competition Results

Kazakh \rightarrow English			English \rightarrow Kazakh		
Ave.	Ave. z	System	Ave.	Ave. z	System
72.2	0.270	online-B	81.5	0.746	HUMAN
70.1	0.218	NEU	67.6	0.262	UAlacant-NMT
69.7	0.189	rug-morfessor	63.8	0.243	online-B
68.1	0.133	online-G	63.8	0.222	UAlacant-NM
67.1	0.113	talp-upc-2019	63.8	0.222	RBMT
67.0	0.092	NRC-CNRC	63.3	0.126	NEU
65.8	0.066	Frank-s-MT	63.3	0.108	MSRA-CrossBERT
65.6	0.064	NICT	60.4	0.097	CUNI-T2T-transfer
64.5	0.003	CUNI-T2T-transfer	61.7	0.078	online-G
48.9	-0.477	UMD	55.2	-0.049	rug-bpe
32.1	-1.058	DBMS-KU	49.0	-0.328	talp-upc-2019
			41.4	-0.493	NICT
			11.6	-1.395	DBMS-KU

Linguistic Background

Kazakh:

- Turkic language
- Cyrillic script.
- Agglutinative.

English:

- West Germanic lang.
- Latin script.
- Simple morphology.

Russian:

- East Slavic language.
- Cyrillic script.
- Fusional morphology.