

Tackling Sparse Data Issue in MT Evaluation



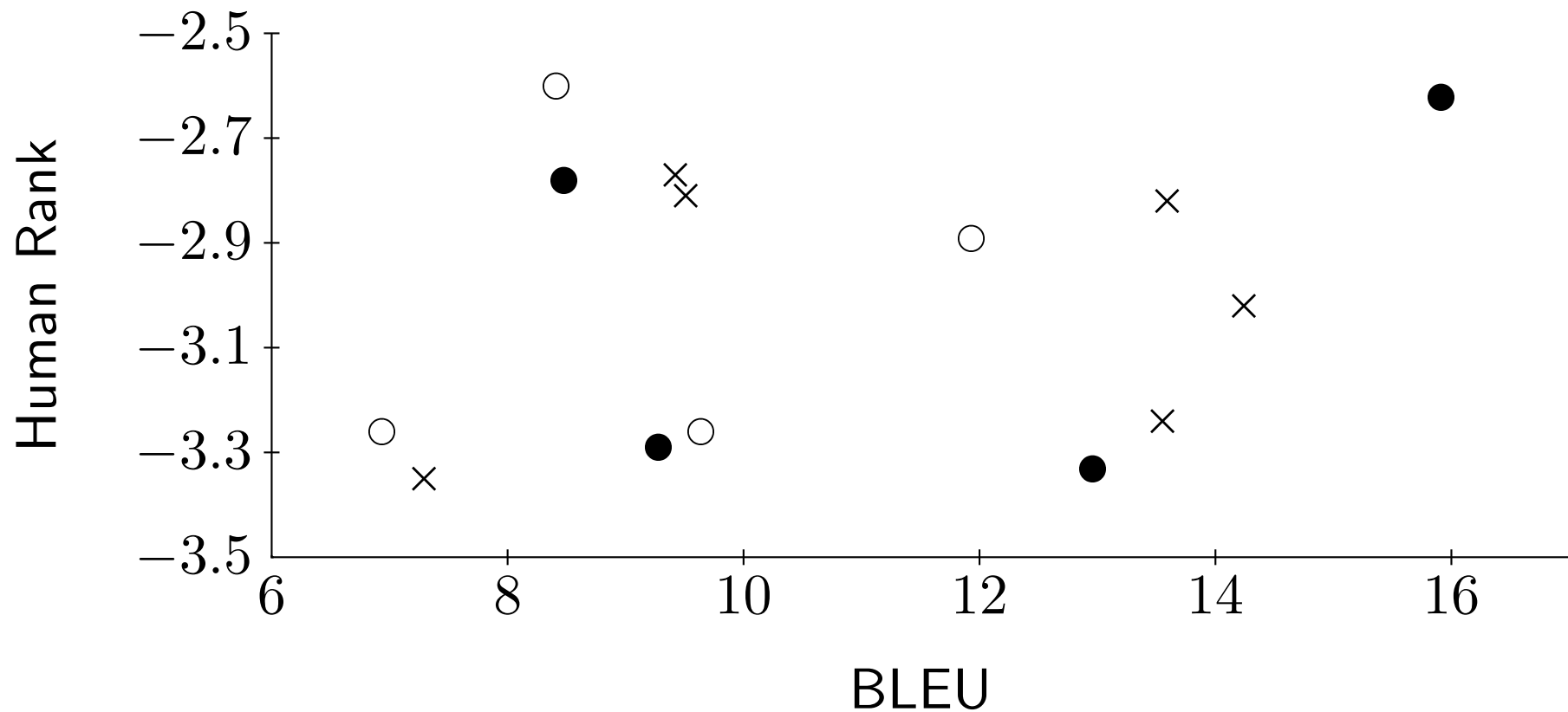
Ondřej Bojar, Kamil Kos, David Mareček
{bojar,marecek}@ufal.mff.cuni.cz, kamilkos@email.cz
Institute of Formal and Applied Linguistics
Faculty of Mathematics and Physics
Charles University in Prague

- Issues of BLEU score.
 - Especially for morphologically rich languages.
- SemPOS
 - Coarser representation of words.
 - Performs better for both Czech and English.



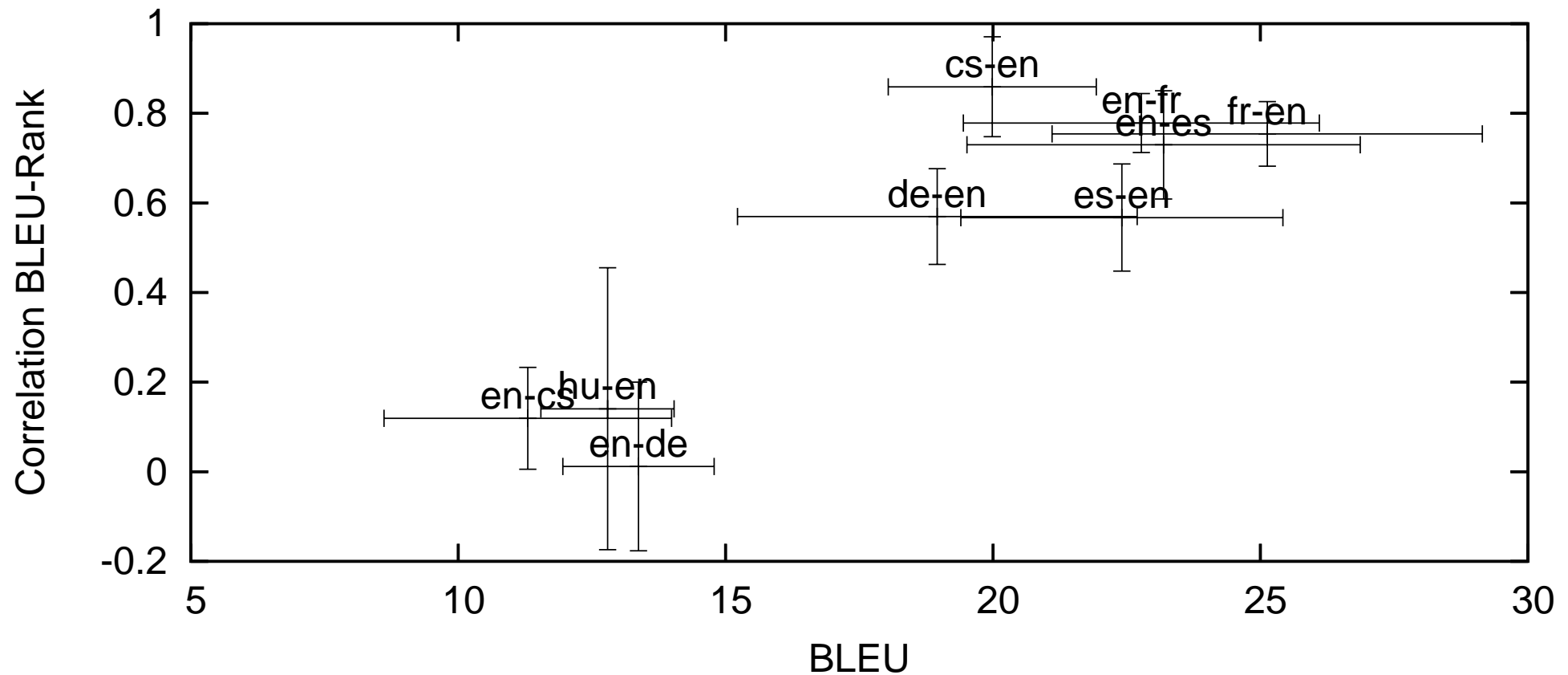
BLEU vs. Human Rank

- English → Czech Translation
 - Evaluation campaigns WMT08 and WMT09.



. . . across Languages

- Various Language Pairs
 - Evaluation campaigns WMT08, WMT09, MetricsMATR.



Reason 1: Focus on Forms



SRC	Prague Stock Market falls to minus by the end of the trading day
REF	pražská burza se ke konci obchodování propadla do minusu
cu-bojar	praha stock market klesne k minus na <u>konci</u> obchodního dne
pctrans	praha trh cenných papírů padá minus <u>do</u> konce obchodního dne

- Only a single unigram in each hyp. confirmed by the reference.
- Large chunks of hypotheses are not compared at all.

Confirmed by Reference	Yes	Yes	No	No
Contains Errors	Yes	No	Yes	No
Running words	6.34%	36.93%	22.33%	34.40%

Reason 2: Sequences \rightsquigarrow SemPOS



BLEU overly sensitive to sequences:

- Gives credit for 1, 3, 5 and 8 four-, three-, bi- and unigrams,
- Two of three serious errors not noticed,
 \Rightarrow Quality of cu-bojar overestimated.

SRC	Congress yields: US government can pump 700 billion dollars into banks							
REF	kongres ustoupil : vláda usa může do bank napumpovat 700 miliard dolarů							
cu-bojar	<u>kongres</u>	výnosy	<u>:</u>	<u>vláda usa může</u>	čerpadlo	<u>700 miliard dolarů</u>	<u>v</u>	bankách
pctrans	<u>kongres</u>	<u>vynáší</u>	<u>:</u>	<u>us vláda může</u>	<u>čerpat</u>	<u>700 miliardu dolarů</u>	<u>do bank</u>	

SemPOS (Kos and Bojar, 2009) gives credit for 8 lemmas:

REF	<u>kongres/n</u>	<u>ustoupit/v</u>	<u>:/n</u>	<u>vláda/n</u>	<u>usa/n</u>	<u>banka/n</u>	<u>napumpovat/v</u>	<u>700/n</u>	<u>miliarda/n</u>	<u>dolar/n</u>	
cu-bojar	<u>kongres/n</u>	<u>výnos/n</u>	<u>:</u>	<u>vláda/n</u>	<u>usa/n</u>	<u>moci/v</u>	<u>čerpadlo/n</u>	<u>700/n</u>	<u>miliarda/n</u>	<u>dolar/n</u>	<u>banka/n</u>
pctrans	<u>kongres/n</u>	<u>vynášet/v</u>	<u>:</u>	<u>us/n</u>	<u>vláda/n</u>	<u>čerpat/v</u>	<u>700/n</u>	<u>miliarda/n</u>	<u>dolar/n</u>	<u>banka/n</u>	

SemPOS Performs Well



- We evaluated:
 - several variants of SemPOS (e.g. Void, include synt. structure),
 - adaptation of SemPOS for English,
 - linear combination of n -grams and SemPOS.

Correlation with human judgments (selected results):

To English (10 test sets)		To Czech (3 test sets)	
Metric	Avg Correlation	Metric	Avg Correlation
Void _{par}	0.75	3·SemPOS+1·BLEU ₄	0.55
GTM	0.71	SemPOS	0.53
4·SemPOS+1·BLEU ₂	0.70	GTM	0.35
SemPOS	0.69	Void	0.33
BLEU	0.66	BLEU	0.33
TER	0.63	TER	0.07

- BLEU does not correlate with human ranks if below ~20.
 - Too much focus on exact forms and sequences.
 - $\Rightarrow > \frac{1}{3}$ of correct tokens not confirmed by reference.
- Suggested SemPOS:
 - Evaluates bags of lemmas instead of sequences of words.
- Evaluated several variants of SemPOS.
 - Improvement over several MT evaluation metrics.

MERT optimization towards SemPOS fails, see our poster at WMT10.

References



Kamil Kos and Ondřej Bojar. 2009. Evaluation of Machine Translation Metrics for Czech as the Target Language. Prague Bulletin of Mathematical Linguistics, 92.