

# Test Set Influence on BLEU

Havlíček (2007) evaluates the influence of:

- number of reference translations,
- translation direction.

on *human-produced* text (1 human translation against 4 others).

Refs	cs→en, Professionals				en→cs, Math Students			
	Indiv. Results			Avg	Indiv. Results			Avg
1	41.15	32.66	34.03	<b>35.95</b>	3.66	8.62	5.79	<b>6.02</b>
2	49.09	49.78	41.26	<b>46.71</b>	9.82	8.26	9.36	<b>9.15</b>
3	52.63			<b>52.63</b>	13.06			<b>13.06</b>

⇒ heavy dependence on the number of references.

More references allow to match more n-grams of MT output.

⇒ heavy dependence on the translation direction and quality.