

cases	doc_1		doc_2		decision	id
	authors	<ul style="list-style-type: none"><li>Wei Zhao</li><li>Goran GlavaÅ¡</li><li>Maxime Peyrard</li><li>Yang Gao</li><li>Robert West</li><li>Steffen Eger</li></ul>	authors	<ul style="list-style-type: none"><li>Steffen Eger</li><li>Goran GlavaÅ¡</li><li>Yang Gao</li><li>Robert West</li><li>Wei Zhao</li><li>Maxime Peyrard</li></ul>	DUPLICATES	254
	title	On the Limitations of Cross-lingual Encoders as Exposed by Reference-Free Machine Translation Evaluation	title	On the Limitations of Cross-lingual Encoders as Exposed by Reference-Free Machine Translation Evaluation		
	publication_date	2020-05-01 00:00:00	publication_date	2020-05-03 00:00:00		
	source	SupportedSources.OPENALEX	source	SupportedSources.PAPERS_WITH_CODE		
	journal	Meeting of the Association for Computational Linguistics	journal			
	volume		volume			
	doi	10.18653/v1/2020.acl-main.151	doi			
	urls	<ul style="list-style-type: none"><li>https://openalex.org/W3035459196</li><li>https://doi.org/10.18653/v1/2020.acl-main.151</li><li>https://www.aclweb.org/anthology/2020.acl-main.151.pdf</li></ul>	urls	<ul style="list-style-type: none"><li>https://arxiv.org/pdf/2005.01196v3.pdf</li><li>https://github.com/AIPHES/ACL20-Reference-Free-MT-Evaluation</li><li>https://aclanthology.org/2020.acl-main.151.pdf</li></ul>		
	id	id-4849408777373638810	id	id8263332987419513617		
	abstract		abstract	Evaluation of cross-lingual encoders is usually performed either via zero-shot cross-lingual transfer in supervised downstream tasks or via unsupervised cross-lingual textual similarity. In this paper, we concern ourselves with reference-free machine translation (MT) evaluation where we directly compare source texts to (sometimes low-quality) system translations, which represents a natural adversarial setup for multilingual encoders. Reference-free evaluation holds the promise of web-scale comparison of MT systems. We systematically investigate a range of metrics based on state-of-the-art cross-lingual semantic representations obtained with pretrained M-BERT and LASER. We find that they perform poorly as semantic encoders for reference-free MT evaluation and identify their two key limitations, namely, (a) a semantic mismatch between representations of mutual translations and, more prominently, (b) the inability to punish "translationese", i.e., low-quality literal translations. We propose two partial remedies: (1) post-hoc re-alignment of the vector spaces and (2) coupling of semantic-similarity based metrics with target-side language modeling. In segment-level MT evaluation, our best metric surpasses reference-based BLEU by 5.7 correlation points.		
	versions		versions			