

Or

## **Can Synthetic Translations Improve Bitext Quality?**

Eleftheria Briakou & Marine Carpuat

{ebriakou, marine}@cs.umd.edu

Synthetic Translations improve MT quality		Revising Bitext w/ Selective Replacement of Synthetic Translations	
ST: Translations generated by MT		procedure REVISE( $D = (S, T), \mathcal{R}$ )	Given semantic equivalence classifier*
Prior work primary uses synthetic translations to reliably improve MT translation quality in two ways:		$M_{S \to T} \leftarrow TRAIN\_MT(D = (S, T))$ $M_{T \to S} \leftarrow TRAIN\_MT(D = (T, S))$ $\widetilde{D} \leftarrow \emptyset$	Train NMT models to translate in opposite directions
Augmenting original translations:	<b>Replacing</b> original translations:	$ \begin{array}{c c} D \leftarrow Q \\ \hline \mathbf{for} \ i \in  D  \ \mathbf{do} \\ \hline (S_i, \hat{T}_i) \leftarrow (S_i, M_{S \to T}(S_i)) \end{array} \end{array} $	Generate synthetic bitext by pairing
<ul> <li>[1] Forward translation</li> <li>[2] Backward translation</li> <li>[3] Data Diversification</li> </ul>	<ul><li>[4] Sequence Distillation</li><li>[5] Data Rejuvenation</li></ul>	$(S_i, T_i) \leftarrow (M_{T \to S}(T_i), T_i)$ $d_F \leftarrow \mathcal{R}(S_i, \hat{T}_i) - \mathcal{R}(S_i, T_i)$ $d_R \leftarrow \mathcal{R}(\hat{S}_i, T_i) - \mathcal{R}(S_i, T_i)$	original references w/ synthetic transl. Compute equivalence scores for original & synthetic pairs
Yet, it remains unclear: Where does this improvement come from?		$if \max(d_F, d_B) > t \text{ then}$ $if \max(d_F, d_B) = d_F \text{ then}$ $\tilde{D} \leftarrow \tilde{D} \cup \{(S_i, \hat{T}_i)\}$ Perplace the original with a synthetic	
WE HYPOTHESIZE: Synthetic translations are of higher quality (i.e., preserve		else $\tilde{D} \leftarrow \tilde{D} \cup \{(\hat{S}_i, T_i)\}$ end if	translation only if it yields a more equivalent translation
translation equivalence) better than naturally occurring bitext WE CONTRIBUTE:		else $\tilde{D} \leftarrow \tilde{D} \cup \{(S_i, T_i)\}$ end if	otherwise keep the original
An extensive empirical evaluation of the quality of bitext revised with synthetic translations		end for return $\tilde{D}$ end procedure Fine-tuned mBE granurality base	alence Classifier: ERT of synthetic divergences of varying ed on our previous work [7]
Driginal vs. Rovisod Bitoxt Eva	Luation Exporimontal Sotti	ings Intrinsic Ev	aluation Regults
Intrinsic:       Human Assessments of Equivalence       Medium Resource Focus:         (a)       Sufficient MT Quality         (b)       Bitext Improvement needer		"Which sentence (A vs.B) conveys the meaning of the source better?" <u>Source (original)</u>	
renormance on downstream N	DATA:		

Ένας από τους οικισμούς που δημιούργησαν ήταν ο Καραβάς.



of low- & medium-frequency words, which we are more sensitive to noisy misalignments that result from poor quality bitext.

**<u>Findings</u>** Revised bitexts yields better translation quality than training on the original for both MT settings (training from scratch & continued training), which further confirms that it yields more reliable training signal due to the reduced noise in the synthetic samples.

## REFERENCES

[1] Jiajun Zhang and Chengqing Zong. 2016. Exploiting source-side monolingual data in neural machine translation. In EMNLP [2] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. Improving neural machine translation models with monolingual data. In ACL [3] Xuan-Phi Nguyen, Shafiq Joty, Kui Wu, and Ai Ti Aw. 2020. Data diversification: A simple strategy for neural machine translation. In. NEURIPS [4] Jiatao Gu, James Bradbury, Caiming Xiong, Victor O.K. Li, and Richard Socher. 2018. Non-autoregressive neural machine translation. In ICLR [5] Wenxiang Jiao, Xipg Wang, Shilin He, Irwin King, Michael Lyu, and Zhaopeng Tu. 2020. Data Rejuvenation: Exploiting Inactive Training Examples for Neural Machine Translation. In EMNLP

[6] Haoyue Shi, Luke Zettlemoyer, and Sida Wang. 2021 Bilingual lexicon induction via unsupervised bitext construction and word alignment. In ACL

[7] Eleftheria Briakou, Marine Carpuat. 2020. Detecting Fine-Grained Cross-Lingual Semantic Divergences without Supervision by Learning to Rank. In EMNLP

## Can Synthetic Translations Improve Bitext Quality?

they selectively replacing imperfect Yes, when... translations in naturally occurring bitexts under a semantic equivalence condition

According to... intrinsic evaluations of semantic equivalence and extrinsic evaluations on **BLI and MT tasks** 

**Data:** https://github.com/Elbria/xling-SemDiv-Equivalize.