

# Bridging the Data Gap between Training and Inference for Unsupervised Neural Machine Translation

Zhiwei He<sup>1</sup>, Xing Wang<sup>2</sup>, Rui Wang<sup>1</sup>, Shuming Shi<sup>2</sup>, Zhaopeng Tu<sup>2</sup>

<sup>1</sup>Shanghai Jiao Tong University, <sup>2</sup>Tencent AI Lab

hez.w.tkcw@gmail.com

## 1. Summary

- Findings
  - Unsupervised neural machine translation (UNMT) suffers from the data gap between training and inference. This data discrepancy results in the overestimation of UNMT on the previous benchmark, which is reflected in the performance gap between UNMT and supervised neural machine translation (SNMT) on the source-original test sets.
  - We identify two representative characteristics of the data gap: style gap and content gap.
- Solution
  - We propose an online self-training approach, which simultaneously uses the pseudo parallel data {natural source, translated target} to mimic the inference scenario.
- Results
  - Our method achieves significant improvement on the source-original test sets.
  - Better natural-to-natural and named entities translation (more details in the paper).

## 2. Data Gap

Existing methods adopt online back-translation causing data gap between training and inference.

- The model is trained with **translated source** ( $\mathcal{X}^*$ ).
- But it translates **natural source** ( $\mathcal{X}$ ) sentences in inference.

Types of training and inference data. \* stands for translated sentences

	Source	Target
Train	$\mathcal{X}^*$	$\mathcal{Y}$
Inference	$\mathcal{X}$	$\mathcal{Y}^*$

## 3. The Overestimated UNMT

- Full set: SNMT  $\approx$  UNMT (previous works)
- Tgt-Ori: SNMT < UNMT
- Src-Ori: SNMT > UNMT (what we need)

Model	En-Fr		En-De		En-Ro		Avg.
	$\Rightarrow$	$\Leftarrow$	$\Rightarrow$	$\Leftarrow$	$\Rightarrow$	$\Leftarrow$	

### Full Test Set

SNMT	38.4	33.6	29.5	33.9	33.7	32.5	33.6
UNMT	37.8	34.9	27.1	35.2	35.1	33.4	33.9

### Target-Original Test Set / Translated Input

SNMT	37.4	32.4	25.6	37.1	38.2	28.2	33.2
UNMT	<b>39.2</b>	<b>37.6</b>	<b>27.0</b>	<b>42.9</b>	<b>43.1</b>	<b>35.6</b>	<b>37.6</b>

### Source-Original Test Set / Natural Input

SNMT	<b>38.2</b>	<b>34.1</b>	<b>32.3</b>	<b>28.8</b>	<b>29.4</b>	<b>35.9</b>	<b>33.1</b>
UNMT	35.2	30.2	26.1	23.6	27.4	30.8	28.9

## 4. Two Factors of Data Gap

### 1 Style Gap

When training, the input is in translated style; while in inference, it's in the natural style.

UNMT has a lower perplexity on the translated input than on natural input

Inference Input	Perplexity
Natural	242
Translated	219

UNMT improves significantly when the input style switches from natural to translated

Model	Natural In.		Translated In.	
	BLEU	$\Delta$	BLEU	$\Delta$
SNMT	28.8	-	44.9	-
UNMT	22.5	-6.3	42.1	-2.8

### 2 Content Gap

The content of input in training is biased towards the target language. While the input during inference is more biased towards the source language.

10 most frequent entities in the source sentences of **De-En** translation. The training data of UNMT has more entities biased towards the target language **English**

Data	Most Frequent Name Entities
Src-Ori Test	Deutschland, Stadt, CDU, deutschen, Zeit SPD, USA, deutsche, China, Mittwoch Großbritannien, London, Trump, USA,
Tgt-Ori Test	Russland, Vereinigten Staaten, Europa Mexiko, Amerikaner, Obama
UNMT Train	Deutschland, dpa, USA, China, Obama, Stadt Hause, Europa, Großbritannien, Russland

UNMT model outputs the hallucinated translation "U.S." which is biased towards target language **English**

Input	Die <b>deutschen</b> Kohlekraftwerke ... in <b>Deutschland</b> emittierten ...
Ref	<b>German</b> coal plants , ..., total amount emitted in <b>Germany</b> .
SNMT	..., <b>German</b> coal-fired power stations ..., emissions in <b>Germany</b> .
UNMT	<b>U.S.</b> coal-fired power plants ... amount emitted in the <b>U.S.</b> ...

## 5. Our Approach

We incorporate the self-training method into UNMT framework to remedy the data gap between the training and inference.

Given translation task  $X \rightarrow Y$ , for each batch:

- $x^* = \arg \max_x P_{Y \rightarrow X}(x | y; \tilde{\theta})$
- construct sample  $(x^*, y)$
- reverse the sample and get  $(y, x^*)$
- train the model using  $(x^*, y)$  and  $(y, x^*)$ <sup>a</sup>

<sup>a</sup>UNMT models are typically bi-directional.

## 6. Experiments

Our method achieves significant improvement on the source-original test sets

Testset	Model	Approach	En-Fr		En-De		En-Ro		Avg.	$\Delta$
			$\Rightarrow$	$\Leftarrow$	$\Rightarrow$	$\Leftarrow$	$\Rightarrow$	$\Leftarrow$		
<i>Our Implementation</i>										
Full set	XLM	UNMT	37.4	34.5	27.2	34.3	34.6	32.7	33.5	-
		+Self-training	<b>37.8</b>	<b>35.1</b>	<b>28.1</b>	<b>34.8</b>	<b>36.2</b>	<b>33.9</b>	<b>34.3</b>	+0.8
	MASS	UNMT	37.8	34.9	27.1	35.2	35.1	33.4	33.9	-
		+Self-training	<b>38.0</b>	<b>35.2</b>	<b>28.9</b>	<b>35.6</b>	<b>36.5</b>	<b>34.0</b>	<b>34.7</b>	+0.8
Src-Ori	XLM	UNMT	34.7	<b>30.4</b>	26.6	22.5	27.4	30.6	28.7	-
		+Self-training	<b>35.4</b> <sup>↑</sup>	30.2	<b>28.0</b> <sup>↑</sup>	<b>23.1</b> <sup>↑</sup>	<b>29.6</b> <sup>↑</sup>	<b>32.7</b> <sup>↑</sup>	<b>29.8</b>	+1.1
	MASS	UNMT	35.2	30.2	26.1	23.6	27.4	30.8	28.9	-
		+Self-training	<b>35.9</b> <sup>↑</sup>	<b>30.9</b> <sup>↑</sup>	<b>28.7</b> <sup>↑</sup>	<b>24.9</b> <sup>↑</sup>	<b>30.1</b> <sup>↑</sup>	<b>31.9</b> <sup>↑</sup>	<b>30.4</b>	+1.5

Natural-to-natural and named entities translation

Note: HQ(R) and HQ(all 4) are natural-to-natural test sets provided by Google (detailed in the paper).

Model	HQ(R)	HQ(all 4)
XLM+UNMT	24.5	19.6
+Self-training	<b>25.9</b>	<b>20.7</b>
MASS+UNMT	24.3	19.6
+Self-training	<b>26.0</b>	<b>20.8</b>

  

Model	Approach	NE Acc.
XLM	UNMT	0.46
	+Self-training	<b>0.53</b>
MASS	UNMT	0.44
	+Self-training	<b>0.52</b>