国际人工智能会议 AAAI 2021论文北京预讲会



# Revisiting Iterative Back- Microsoft Translation from the Perspective of Compositional Generalization

郭一诺(

### **Our Team - Main Contributors**



Peking University Yinuo Guo



Microsoft Asia Bei Chen



Beihang University Hualei Zhu

Microsoft Asia Jian-Guang



Microsoft Asia Zeqi Lin



Microsoft Asia Dongmei Zhang

### **Compositional Generalization**

 The algebraic ability to understand and produce unseen combinations of seen atoms.
 — Chomsky

Natural LanguageProgramming LanguageTrainrun twice $\Rightarrow$ JUMP WALK

Test jump twice and run ⇒JUMP JUMP RUN

# Background: Seq2seq Tasks in NLP

- Machine Translation
- Semantic Parsing
- Summarization



# Semi-Supervised Learning

- Parallel data are limited and expensive
- Monolingual data are cheap and abundant, containing lots of unseen combinations
- Hypothesis: semi-supervised learning can enable models understand and produce much more combinations beyond labelled data, thus tackling the bottleneck of lacking compositione'



Unlabeled source-side corpus

#### **Iterative Back-Translation**

 We focus on Iterative Back-Translation (IBT), a simple yet effective semisupervised method that has been successfully applied in machine translation.



#### **Three Research Questions**

 RQ1:How does IBT affect compositional generalization of seq2seq models?

• Yes

- RQ2: What is the key that contributes to the success of IBT?
  Quality of pseudo parallel data & Perturbations
- RQ3: How to further improve the performance of IBT?
  - Curriculum Iterative Back-translation

### Evaluate on CFQ & SCAN

- Substantially improves the performance on CG benchmarks.
- Better monolingual data, better results.

	Models	MCD1	MCD2	MCD3
Which Swedish founder of [M0] produced [M2] ?	LSTM+Attn Transformer	$28.9 \pm 1.8$ 34.9 ± 1.1	$5.0 \pm 0.8$ 8 2 ± 0.3	$10.8 \pm 0.6$ 10.6 ± 1.1
<pre>SELECT DISTINCT ?x0 WHERE {</pre>	Uni-Transformer CGPS T5-11B	$37.4 \pm 2.2$ $13.2 \pm 3.9$ $61.4 \pm 4.8$	$8.1 \pm 1.6$ $1.6 \pm 0.8$ $30.1 \pm 2.2$	$\begin{array}{c} 11.3 \pm 0.3 \\ 6.6 \pm 0.6 \\ 31.2 \pm 5.7 \end{array}$
	GRU+Attn (Ours) +mono30	$\begin{array}{c} 32.6 \pm 0.22 \\ 64.8 \pm 4.4 \end{array}$	$\begin{array}{c} 6.0 \pm 0.25 \\ \textbf{57.8} \pm \textbf{4.9} \end{array}$	$\begin{array}{c} 9.5 \pm 0.25 \\ 64.6 \pm 4.9 \end{array}$
	+mono100 +transductive	$83.2 \pm 3.1$ $88.4 \pm 0.7$	$71.5 \pm 6.9 \\ 81.6 \pm 6.5$	$81.3 \pm 1.6$ $88.2 \pm 2.2$

#### **Quality of Pseudo Parallel Data**

 Iterative back-translation can increasingly correct errors in pseudo-parallel data



### Impact of Error-Prone Data & Perturbations

- Even noise pseudo-parallel data can bring gains!
  - · As they bring implicit knowledge of unseen combinations
- Perturbations brought by OTF (on-the-fly) is very important!
  - Pseudo-parallel data are generated dynamically, which prevent learning specific incorrect bias



#### **Curriculum Iterative Back-Translation**

- We want to help reduce errors more efficiently
- CIBT: during the training process:
  - start out with easy monolingual data,
  - then gradually increase the difficulty.



#### **Curriculum Iterative Back-Translation**

- Curriculum learning benefits iterative back-translation.
- Curriculum learning is more beneficial to difficult data than simple data.

0	IBT	CIBT with hyperparameter $c$ (steps in each stage)					
	IDT	2000	2500	3000	3500	4000	
MCD1	$64.8 \pm 4.4$	$66.1 \pm 5.0$	$66.0 \pm 4.8$	$66.6 \pm 5.4$	$65.9 \pm 3.7$	$65.4 \pm 3.8$	
MCD2	$57.8 \pm 4.9$	$68.6\pm2.6$	69.1 ± 3.1	$68.0 \pm 1.9$	$66.8 \pm 2.4$	$65.4 \pm 3.1$	
MCD3	$64.6 \pm 4.9$	$70.2 \pm 4.9$	$68.4 \pm 7.0$	$70.4 \pm 4.8$	$69.2 \pm 4.1$	$67.0\pm6.3$	
Mean	$62.4\pm6.1$	$68.3 \pm 4.1$	$67.8 \pm 4.7$	$68.3 \pm 4.1$	$67.3 \pm 3.4$	$65.9 \pm 4.1$	





Figure 6: Performance on different subsets. This figure indicates that curriculum learning is more beneficial to difficult data (larger k) than simple data (smaller k).

## Takeaways

- Iterative back-translation can significantly improve CG.
- Why IBT works well:
  - Unseen combinations
  - Increasingly improving the quality of pseudo-parallel data
  - Perturbations
- We propose curriculum iterative back-translation to further improving the performance.

国际人工智能会议 AAAI 2021论文北京预讲会

# THANKS

#### 2020.12.19

Related papers from our team (MSRA DKI): Hierarchical Poset Decoding for Compositional Generalization in Language (NeurIPS 2020) Compositional Generalization by Learning Analytical Expressions (NeurIPS 2020 Spotlight)

Iterating Utterance Segmentation for Neural Semantic Parsing (AAAI 2021)