*switch-*GLAT: Multilingual Parallel Machine Translation Via Code-Switch Decoder

Zhenqiao Song, Hao Zhou, Lihua Qian, Jingjing Xu, Shanbo Cheng, Mingxuan Wang, Lei Li



Reporter: Xiachong Feng









Zhenqiao Song*Hao Zhou*Lihua QianJingjing XuFudan UniversityResearch AssociateByteDance AI Labpost-doc UCSBBytedance AI LabProfessor

Professor Insititute for AI Industry Research, Tsinghua University







heng Shanbo Cheng research scientist/manager Pro working at ByteDance AI Lab

Lei Li Assistant Professor UCSB

Shanbo Cheng Shand res scientis

2

Non-Autoregressive Decoding



http://xcfeng.net/res/presentation/Non-Autoregressive%20Decoding.pdf

Non-Autoregressive Decoding



https://docs.google.com/presentation/d/1lnXSxPL3hZQ_OHyV_Lksz3UcOmPyi DqC1eg29Q5IAr0/edit#slide=id.g84c46a9f82_0_0





字节跳动火山翻译团 队的并行生成翻译系 统 GLAT 拿下 了 WMT2021 De-En/En-De 的双料冠军。

Glancing Transformer



Glancing Transformer for Non-Autoregressive Neural Machine Translation ACL 2021

Glancing Transformer: Training



Note that only the second decoding will update the model parameters

Glancing Sampling Strategy

- Selected tokens provide "correct" information from the ground-truth target
- Adaptive sampling strategy guides the model to first learn the generation of fragments and then gradually turn to the whole sentences.
- As the model gets better progressively, the sampling strategy will sample fewer words to enable the model to learn the parallel generation of the whole sentence.



Glancing Sampling Strategy

- 1. Deciding a sampling number *S*
- 2. Randomly selecting *S* words from the reference.
 - 1. The random reference word selection is simple and yields good performance empirically.

$$\begin{split} \mathbb{GS}(Y,\hat{Y}) &= \mathsf{Random}(Y,S(Y,\hat{Y})) \\ S(Y,\hat{Y}) &= \lambda \cdot d(Y,\hat{Y}) \\ \lambda: \text{ sampling ratio } & d: \mathsf{metric for measuring the} \\ d: \mathsf{metric for measuring the} \\ \mathsf{differences between } Y \text{ and } \hat{Y} \\ & \underset{\mathsf{Hamming distance}}{\mathsf{Hamming distance}} \\ d(Y,\hat{Y}) &= \sum_{t=1}^{T} (y_t \neq \hat{y}_t) \end{split}$$

 $d(Y, \widehat{Y})$: the sampling number can be decided adaptively considering the current trained model's prediction capability.



Glancing Transformer: Inference

- We need to decide the **output lengths** before decoding.
- An additional [LENGTH] token is added to the source input, and the encoder output for the [LENGTH] token is used to predict the length.

switch-GLAT: Code-Switch Decoder



Switch-GLAT can translate between different languages using the indicative language tags, of which the decoder is called code-switch decoder.

$$\begin{split} L_{\text{multi}} &= \sum_{\mathbb{D}^l \in \mathbb{D}} \sum_{(X_j^l, Y_j^l) \in \mathbb{D}^l} \{ L_{\text{tag}}(Y_j^l | X_j^l; \theta_M) + L_{\text{len}}^l(j) \} \\ L_{\text{len}}^l(j) &= -P(L_j^l) \log \hat{P}(L_j^l | [F_e(X_j^l; \theta_M); E_{\text{src}}; E_{\text{tgt}}]; \theta_M) \end{split}$$

$$\begin{aligned} L_{\text{tag}}(Y_j^l | X_j^l; \theta_M) &= -\sum_{y_t \in \overline{\mathbb{GS}(Y_j^l, \hat{Y}_j^l)}} \log P(y_t | \mathbb{GS}(Y_j^l, \hat{Y}_j^l), X_j^l, \text{src}, \text{tgt}; \theta_M) \\ \hat{Y}_j^l &= F_d(\tilde{H}_d^0, F_e(X_j^l, \text{src}; \theta), \text{tgt}; \theta_M) \end{aligned}$$

$$\tilde{f}_i^0 = f_i^0 + E_{\text{src}}; \tilde{f}_i^K = f_i^K + E_{\text{src}}$$
$$\tilde{h}_j^0 = h_j^0 + E_{\text{tgt}}; \tilde{h}_j^K = h_j^K + E_{\text{tgt}}$$

switch-GLAT: Code-Switch Back-Translation

Through this process, abundant code-switched sentences can be generated, which helps to learn better-aligned cross-lingual representations.



Loss



Experiments

- (1) WMT-EDF: We collect 4 language pairs from WMT-14 English (En)
 German (De) and English (En)
 French (Fr). All three languages belong to Indo-European language family and are relatively close on linguistics.
- (2) WMT-EFZ: We also collect 4 language pairs from WMT-14 English (En) French (Fr) and WMT-17 English (En) Chinese (Zh), which are distant languages on linguistics and their relationships are more difficult to learn.
- (3) WMT-many: We also gather 10 language pairs from WMT-14 English (En)
 German (De), English (En)
 French (Fr), WMT-16 English (En)
 Russian (Ru), English (En)
 Romanian (Ro) and WMT-17 English (En)
 Chinese (Zh) to test switch-GLAT on more diverse language pairs.

Result

| Models | WMT-EDF | | | | | |) | WMT-EFZ | | | |
|--|--|--|--|--|--|--|--|---|--|---|---|
| | En-De | De-En | En-Fr | Fr-En | Avg | Speed | En-Fr | Fr-En | En-Zh | Zh-En | Avg |
| Bilingual models | | | | | | | | | | | |
| Transformer GLAT | $\begin{array}{c} 27.77\\ 26.09 \end{array}$ | $\begin{array}{c} 31.55\\ 30.53 \end{array}$ | $\begin{array}{c} 38.80\\ 38.62 \end{array}$ | $\begin{array}{c} 37.35\\ 34.44\end{array}$ | 33.86 32.42 | 1.3 	imes 6.3 	imes | $\begin{array}{c} 38.80\\ 38.62 \end{array}$ | $\begin{array}{c} 37.35\\ 34.44\end{array}$ | $\begin{array}{c} 23.60\\ 21.05 \end{array}$ | $\begin{array}{c} 24.05\\ 22.89 \end{array}$ | 30.95 29.25 |
| Multilingual models | | | | | | | | | | | |
| M-Transformer CLSR Adapter - MNAT <i>switch</i> -GLAT - w/o glancing - w/o CSBT - w/o CCS | $\begin{array}{r} 25.84\\ 23.51\\ \underline{22.19}\\ 1\overline{3.82}\\ 25.27\\ 17.76\\ \underline{24.29}\\ 24.72\end{array}$ | $31.5731.2929.56-21.\overline{89}31.2923.2829.0329.51$ | $\begin{array}{r} 38.52\\ 38.58\\ 40.72\\ \hline 24.31\\ 40.81\\ 25.91\\ 36.45\\ 37.32\end{array}$ | $\begin{array}{r} 36.03\\ 34.67\\ 35.88\\ \hline 25.28\\ 36.00\\ 29.15\\ 34.09\\ 34.17\end{array}$ | 32.9932.0132.0821.3233.3424.0230.9731.43 | $1.0 \times \\ 0.9 \times \\ -5.9 \times \\ 6.2 \times \\ 6.0 \times \\ 6.2 \times \\ 6.1 \times $ | $\begin{array}{r} 38.06\\ 37.39\\ 40.13\\ \hline 19.46\\ 40.54\\ 21.28\\ 35.17\\ 35.96\end{array}$ | $\begin{array}{r} 35.07\\ 35.62\\ 35.02\\ \hline 20.18\\ 36.48\\ 21.97\\ 33.65\\ 33.61\\ \end{array}$ | $20.76 \\ 20.23 \\ 19.87 \\ -\overline{7.89} \\ 19.47 \\ 8.62 \\ 18.32 \\ 18.45$ | $\begin{array}{r} 22.19\\ 21.13\\ 21.29\\ -\overline{7.35}\\ 22.55\\ 8.02\\ 20.37\\ 20.83\end{array}$ | 29.02 28.59 29.08 13.72 29.76 14.97 26.87 27.21 |

Table 1: Translation performance (BLEU) on WMT-EDF/EFZ¹. Avg means the average BLEU score.

Result

| Models | WMT-many | | | | | | | | | | | |
|--|---|--|--|--|--|--|--|---|---|--|--|--|
| | En-De | De-En | En-Fr | Fr-En | En-Ro | Ro-En | En-Ru | Ru-En | En-Zh | Zh-En | Avg | |
| Bilingual models | | | | | | | | | | | | |
| Transformer GLAT | $27.77 \\ 26.09$ | $\begin{array}{c} 31.55\\ 30.53 \end{array}$ | $\begin{array}{c} 38.80\\ 38.62 \end{array}$ | $\begin{array}{c} 37.35\\ 34.44 \end{array}$ | $\begin{array}{c} 33.01\\ 31.83 \end{array}$ | $33.59 \\ 32.59$ | $\begin{array}{c} 28.22\\ 25.42 \end{array}$ | $\begin{array}{c} 29.89\\ 28.13 \end{array}$ | $\begin{array}{c} 23.60\\ 21.05 \end{array}$ | $\begin{array}{c} 24.05\\ 22.89 \end{array}$ | 30.78 29.20 | |
| Multilingual models | | | | | | | | | | | | |
| M-Transformer CLSR Adapter MNAT <i>switch</i> -GLAT – w/o glancing – w/o CSBT – w/o CCS | $\begin{array}{r} 23.14\\ 21.57\\ -\underline{23.26}\\ \overline{8.33}\\ 24.18\\ 6.83\\ 22.48\\ 22.25\end{array}$ | $29.38 \\ 30.01 \\ -29.87 \\ -13.86 \\ 30.49 \\ 16.29 \\ 27.91 \\ 28.13$ | $\begin{array}{r} 35.19\\ 34.07\\ 38.79\\ 12.07\\ 39.47\\ 11.08\\ 31.62\\ 34.01 \end{array}$ | $\begin{array}{r} 34.07\\ 32.72\\ 40.03\\ \overline{19.21}\\ 36.30\\ 20.33\\ 33.55\\ 33.62\end{array}$ | $\begin{array}{r} 34.36\\ 30.61\\ 30.42\\ \overline{12.89}\\ 31.93\\ 13.57\\ 29.89\\ 30.16\end{array}$ | $\begin{array}{r} 35.26\\ 35.29\\ -\underline{32.01}\\ -\underline{21.36}\\ 32.40\\ 23.34\\ 32.41\\ 33.07 \end{array}$ | $\begin{array}{r} 24.63 \\ 19.13 \\ 22.87 \\ \hline 7.66 \\ 24.16 \\ 7.57 \\ 22.43 \\ 21.68 \end{array}$ | $\begin{array}{r} 29.14\\ 30.27\\ \underline{26.38}\\ 15.43\\ 28.33\\ 18.25\\ 26.48\\ 26.53\end{array}$ | $16.63 \\ 16.28 \\ 18.29 \\ \overline{4.85} \\ 16.25 \\ 7.31 \\ 15.33 \\ 14.92$ | $\begin{array}{r} 20.17\\ 20.19\\ -\underline{20.31}\\ -6.0\overline{1}\\ 21.23\\ 6.74\\ 18.32\\ 19.07\end{array}$ | $\begin{array}{r} 28.19\\ 27.01\\ 28.22\\ \hline 12.16\\ \textbf{28.47}\\ 13.13\\ 26.04\\ 26.34\\ \end{array}$ | |

Table 2: Translation performance (BLEU) on WMT-many². Avg means the average score.





Figure 3: Inference speed evaluated on throughputs.





Figure 4: Representations learned by (a) Transformer (Vaswani et al., 2017), (b) GLAT (Qian et al., 2020), (c) M-Transformer and (d) *switch*-GLAT, projected to 2D.

Gold word pairs from Open Multilingual WordNet (OMW) datasets English words are displayed in blue color and German words in red.

Result

- Cross-lingual word induction performance to see how well the similar words from different languages are close to each other in the learned vector space
- Word Induction: golden word pairs are extracted from the OMW datasets
- Sentence Retrieval: Tatoeba dataset, Cosine similarity is leveraged to search the nearest neighbour

| Models | 1. | Wor | d Inductio | on | Sentence Retrieval | | | | | |
|------------------------------------|--------------|--------------|-------------------|---------------------|--------------------|-------|--------------|-------|------------------|-------------------|
| | En-De | De-En | En-Fr | Fr-En | Avg | En-De | De-En | En-Fr | Fr-En | Avg |
| M-Transformer | 30.2 | 31.7 | 32.7 | 38.9 | 33.3 | 23.6 | 22.3 | 24.5 | 23.8 | 23.5 |
| CLSR | 24.9 | 25.8 | 37.5 | 36.2 | 31.1 | 19.6 | 18.3 | 26.8 | 27.9 | 23.2 |
| Adapter | 32.8 | 34.5 | 40.1 | 42.8 | 37.5 | 34.6 | 28.7 | 26.9 | 35.2 | 31.4 |
| MNAT | $\bar{24.1}$ | $\bar{23.3}$ | $\overline{32.4}$ | $-3\bar{3}.\bar{6}$ | 28.3 | 13.8 | $17.\bar{6}$ | -18.2 | $\bar{16.3}^{-}$ | $\overline{16.5}$ |
| switch-GLAT | 33.8 | 36.2 | 41.9 | 46.3 | 39.5 | 34.8 | 36.3 | 36.9 | 37.2 | 36.3 |
| – w/o glancing | 24.3 | 24.6 | 32.9 | 33.7 | 28.9 | 18.2 | 22.8 | 19.2 | 20.3 | 20.1 |
| – w/o CSBT | 28.2 | 30.1 | 31.9 | 37.8 | 32.0 | 17.6 | 21.4 | 23.3 | 25.3 | 21.8 |
| – w/o CCS | 30.5 | 29.7 | 33.6 | 35.8 | 32.4 | 18.8 | 22.9 | 22.7 | 21.6 | 21.5 |

Table 3: Results of quality analyses. Avg means the average accuracy.

Conclusion

- switch-GLAT, a non-autoregressive multilingual neural machine translation model
- The multilingual translation performance and cross-lingual representations can both be improved
- Parallel decoder enables a highly efficient inference

Thanks~