

---

# Appendix for: Cross-model Back-translated Distillation for Unsupervised Machine Translation

---

Xuan-Phi Nguyen<sup>1,2</sup> Shafiq Joty<sup>1,3</sup> Thanh-Tung Nguyen<sup>1,2</sup> Wu Kui<sup>2</sup> Ai Ti Aw<sup>2</sup>

## 1. Appendix

In the following supplementary material, we first provide the full mathematical derivations of the loss function  $\mathcal{L}$  presented in the paper (§1.1). Then, we provide the generalized version of our method cross-model back-translated distillation, or GCBD, and measure its effectiveness in the IWSLT English-German, German-English, English-French and French-English unsupervised tasks (§1.2). In addition, we investigate why ensemble knowledge distillation (Freitag et al., 2017), which boosts the performance in a supervised setup, fails to do so in an unsupervised setup where we replace the supervised agents used in the method with the UMT agents (§1.3). Finally, in §1.5, we provide a comparison between unsupervised models and supervised counterparts to provide a perspective of how far unsupervised machine translation research has progressed.

### 1.1. Derivations of negative log likelihood $\mathcal{N}(\theta_\alpha, \theta_\beta)$

In this section, we provide the complete mathematical derivations of the loss function  $\mathcal{L}$  in the paper. Recalling that we are supposed to maximize the log probabilities of the variables  $x_s, y_t, z_s, x_t, y_s$  and  $z_t$  according to the sampling process in Figure 1 and the graphical model in Figure 2. Otherwise speaking, we seek to minimize the following negative log likelihood:

$$\mathcal{J}(\theta) = -\log P_\theta(x_s, y_t, z_s) - \log P_\theta(x_t, y_s, z_t) \quad (1)$$

Then we can expand the first term as follows:

---

<sup>1</sup>Nanyang Technological University <sup>2</sup>Institute for Infocomm Research (I<sup>2</sup>R), A\*STAR <sup>3</sup>Salesforce Research Asia. Correspondence to: Xuan-Phi Nguyen <nguyenxu002@e.ntu.edu.sg>.

$$\begin{aligned} \log P_\theta(x_s, y_t, z_s) &= \log \frac{P_\theta(x_s, y_t, z_s)}{P_\theta(x_s, y_t)} P_\theta(x_s, y_t) \\ &= \log P_\theta(z_s|x_s, y_t) + \log P_\theta(x_s, y_t) \\ &= \log P_\theta(z_s|x_s, y_t) + \log \frac{P_\theta(x_s, y_t)}{P_\theta(y_t)} P_\theta(y_t) \\ &= \log P_\theta(z_s|x_s, y_t) + \log P_\theta(x_s|y_t) + \log P_\theta(y_t) \end{aligned} \quad (2)$$

Since  $z_s$  is independent from  $x_s$  given  $y_t$  according to the graphical model (fig. 1), we have  $P_\theta(z_s|x_s, y_t) = P_\theta(z_s|y_t)$ , then Eq. 2 can be reduced to:

$$\log P_\theta(x_s, y_t, z_s) = \log P_\theta(z_s|y_t) + \log P_\theta(x_s|y_t) + \log P_\theta(y_t) \quad (3)$$

Alternatively, the first term can also be express as follows:

$$\begin{aligned} \log P_\theta(x_s, y_t, z_s) &= \log P_\theta(z_s|x_s, y_t) + \log P_\theta(x_s, y_t) \\ &= \log P_\theta(z_s|y_t) + \log P_\theta(y_t, x_s) \\ &= \log \frac{P_\theta(y_t|z_s)P_\theta(z_s)}{P_\theta(y_t)} + \log \frac{P_\theta(y_t, x_s)}{P_\theta(x_s)} P_\theta(x_s) \\ &= \log P_\theta(y_t|z_s) + \log P_\theta(z_s) - \log P_\theta(y_t) \\ &\quad + \log P_\theta(y_t|x_s) + \log P_\theta(x_s) \end{aligned} \quad (4)$$

After that, we expand the second term in similar fashion, which we yield:

$$\log P_\theta(x_t, y_s, z_t) = \log P_\theta(z_t|y_s) + \log P_\theta(x_t|y_s) + \log P_\theta(y_s) \quad (5)$$

$$\begin{aligned} \log P_\theta(x_t, y_s, z_t) &= \log P_\theta(y_s|z_t) + \log P_\theta(y_s|x_t) \\ &\quad + \log P_\theta(z_t) + \log P_\theta(x_t) - \log P_\theta(y_s) \end{aligned} \quad (6)$$

Then, by adding up Eq. 3, 4, 5 and 6 together, and then divide it by 2, we will derive the negative log likelihood of Eq. 1 as:

$$\begin{aligned} \mathcal{J}(\theta) &= \frac{1}{2} [-\log P_\theta(y_t|z_s) - \log P_\theta(y_t|x_s) \\ &\quad - \log P_\theta(z_s|y_t) - \log P_\theta(x_s|y_t) - \log P_\theta(y_s|z_t) \\ &\quad - \log P_\theta(y_s|x_t) - \log P_\theta(z_t|y_s) - \log P_\theta(y_s|x_t) \\ &\quad - \log P_\theta(x_s) - \log P_\theta(z_s) - \log P_\theta(x_t) - \log P_\theta(z_t)] \end{aligned} \quad (7)$$



Table 2: Percentage of tri-gram repetitions in the synthetic data generated by ensemble knowledge distillation (Freitag et al., 2017), compared to those created by CBD; and the respective test BLEU scores in the *base* WMT’14 En-Fr, WMT’16 En-De and En-Ro unsupervised tasks.

Method	En-Fr	Fr-En	En-De	De-En	En-Ro	Ro-En
<b>% tri-gram repetition</b>						
Ens-Distil	30.3%	34%	73%	76%	43%	86%
CBD	$10^{-3}\%$	$10^{-2}\%$	$10^{-2}\%$	$10^{-1}\%$	$10^{-2}\%$	$10^{-2}\%$
<b>BLEU on test set</b>						
Ens-Distil	17.3	20.0	3.5	3.7	1.2	1.1
CBD	26.6	25.7	16.6	20.5	18.1	17.8

### 1.3. Analysis of degeneration in ensemble knowledge distillation

Ensemble knowledge distillation (Freitag et al., 2017) has been used to enhance supervised machine translation. It uses multiple strong (supervised) teachers to generate synthetic parallel data from both sides of the parallel corpora by averaging the decoding probabilities of the teachers at each step. The synthetic data are then used to train the student model. Having seen its effectiveness in the supervised setup, we apply this same tactic to unsupervised MT tasks by replacing the supervised teachers with unsupervised MT agents. However, the method surprisingly causes drastic performance drop in the WMT’14 En-Fr, WMT’16 En-De and En-Ro unsupervised MT tasks.

By manual inspection, we found that many instances of the synthetic data are incomprehensible and contain repetitions, which is a degeneration behavior. We then quantitatively measure the percentage of sentences in the synthetic data containing tri-gram repetitions by counting the number of sentences where a word/sub-word is repeated at least three consecutive times. As reported in the main paper, from 30% to 86% of the synthetic data generated by the ensemble knowledge distillation (Ens-Distil) method are incomprehensible and contain repetitions. Relative to the performance of CBD, the performance drop in ensemble distillation is also more dramatic for language pairs with higher percentage of degeneration (En-Ro and En-De). This explains why the downstream student model fails to learn from these corrupted data. The results indicate that UMT agents are unable to jointly translate through ensembling strategy the monolingual data that they were trained on. This phenomenon may require further research to be fully understood. On the other hand, with less than 0.1% tri-gram repetitions, CBD generates little to no repetitions, which partly explains why it is able to improve the performance.

Convergence curve with En-Fr BLEU vs Updates

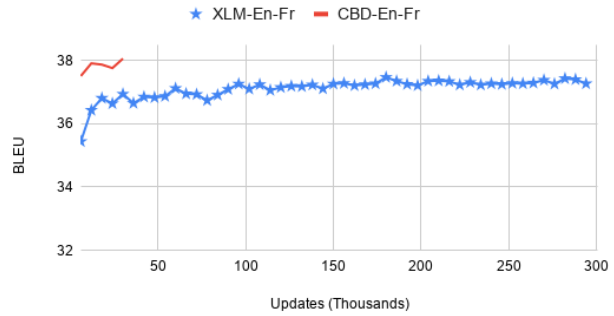


Figure 3: Convergence speed of CBD in comparison with baseline XLM, represented by the test BLEU score of the WMT En-Fr task after a given number of training updates.

Convergence curve with Fr-En BLEU vs Updates

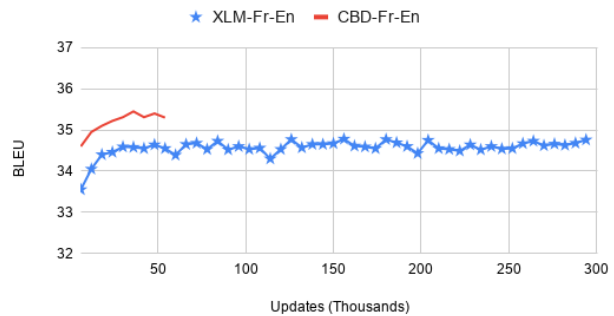


Figure 4: Convergence speed of CBD in comparison with baseline XLM, represented by the test BLEU score of the WMT Fr-En task after a given number of training updates.

### 1.4. Convergence curves of CBD compared with the baselines

This section provides extra convergence curve charts for all 6 of the language pairs in the large scale WMT English-French (Figure 3 & Figure 4), WMT English-German (Figure 5 & Figure 6) and WMT English-Romanian (Figure 7 & Figure 8) tasks. As it can be seen from the charts, CBD converges rapidly and outperforms the baselines with little additional resources, given the pretrained models provided by Conneau & Lample (2019) and Song et al. (2019).

### 1.5. Comparison with supervised MT

In this section, we compare the performances of the CBD method, along with previous SOTA unsupervised models, with the standard supervised Transformer model (Ott et al., 2018) to present a perspective of how much progress the field of unsupervised machine translation has made. More specifically, we use the provided Transformer models pretrained on the parallel WMT’14 English-French and

Convergence curve with En-De BLEU vs Updates

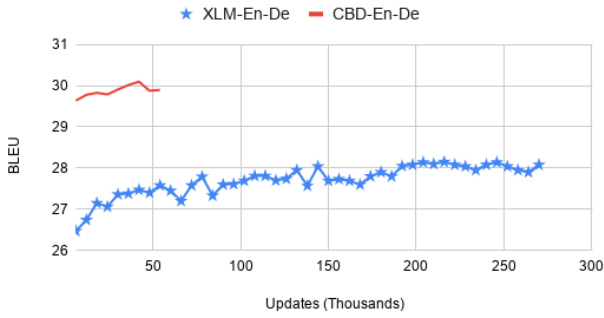


Figure 5: Convergence speed of CBD in comparison with baseline XLM, represented by the test BLEU score of the WMT En-De task after a given number of training updates.

Convergence curve with De-En BLEU vs Updates

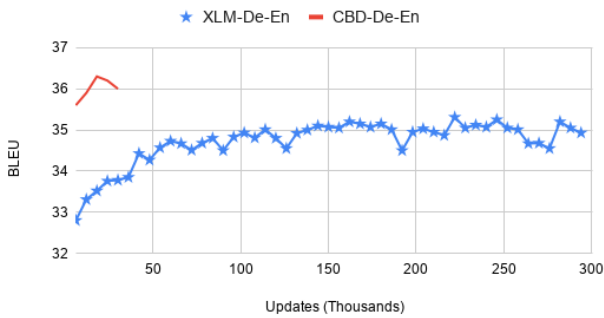


Figure 6: Convergence speed of CBD in comparison with baseline XLM, represented by the test BLEU score of the WMT De-en task after a given number of training updates.

Convergence curve with En-Ro BLEU vs Updates

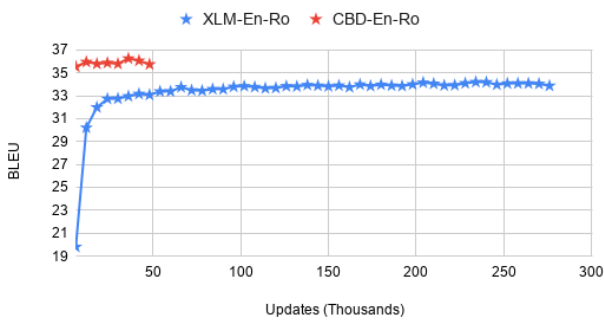


Figure 7: Convergence speed of CBD in comparison with baseline MASS, represented by the test BLEU score of the WMT En-Ro task after a given number of training updates.

English-German datasets and evaluate them on the WMT’14 En-Fr and WMT’16 En-De test sets, as similarly done for unsupervised counterparts. The results are presented in Table 3. As it can be seen, unsupervised MT models have made significant advancement throughout multiple

Convergence curve with Ro-En BLEU vs Updates

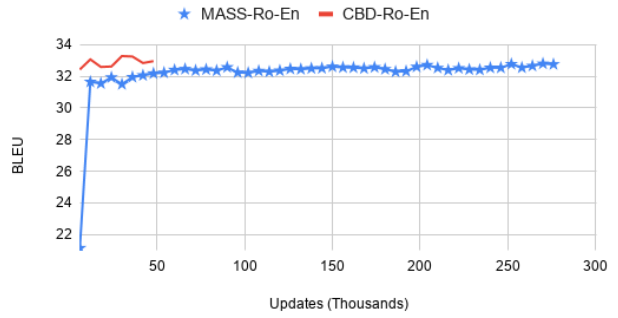


Figure 8: Convergence speed of CBD in comparison with baseline MASS, represented by the test BLEU score of the WMT Ro-En task after a given number of training updates.

Table 3: BLEU scores on the WMT’14 English-French (En-Fr) and WMT’16 English-German (En-De) tasks of unsupervised MT methods (MASS and CBD), in comparison to supervised MT method (Ott et al., 2018).

Method	En-Fr	En-De
Unsupervised MT		
XLM (Conneau & Lample, 2019)	33.4	26.4
MASS (Song et al., 2019)	37.5	28.3
CBD	38.2	30.1
Supervised MT		
Transformer (Ott et al., 2018)	43.2	33.0

iterations and refinement (Conneau & Lample, 2019; Song et al., 2019). However, while the CBD method further improve the performance, it still lags behind the supervised MT model (Ott et al., 2018) by around 3 to 5 BLEU points.

## References

Conneau, A. and Lample, G. Cross-lingual language model pretraining. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32*, pp. 7059–7069. Curran Associates, Inc., 2019.

Freitag, M., Al-Onaizan, Y., and Sankaran, B. Ensemble distillation for neural machine translation. *arXiv preprint arXiv:1702.01802*, 2017.

Ott, M., Edunov, S., Grangier, D., and Auli, M. Scaling neural machine translation. In *Proceedings of the Third Conference on Machine Translation (WMT)*, 2018.

Song, K., Tan, X., Qin, T., Lu, J., and Liu, T.-Y. Mass: Masked sequence to sequence pre-training for language generation. *arXiv preprint arXiv:1905.02450*, 2019.