

# A Study of Reinforcement Learning for Neural Machine Translation

Lijun Wu<sup>1\*</sup>, Fei Tian<sup>2</sup>, Tao Qin<sup>2</sup>, Jianhuang Lai<sup>1</sup> and Tie-Yan Liu<sup>2</sup>

<sup>1</sup>School of Data and Computer Science, Sun Yat-sen University

<sup>2</sup>Microsoft Research

wulijun3@mail2.sysu.edu.cn; stsljh@mail.sysu.edu.cn;

{fetia,taoqin,tyliu}@microsoft.com

## Abstract

Recent studies have shown that reinforcement learning (RL) is an effective approach for improving the performance of neural machine translation (NMT) system. However, due to its instability, successfully RL training is challenging, especially in real-world systems where deep models and large datasets are leveraged. In this paper, taking several large-scale translation tasks as testbeds, we conduct a systematic study on how to train better NMT models using reinforcement learning. We provide a comprehensive comparison of several important factors (e.g., baseline reward, reward shaping) in RL training. Furthermore, to fill in the gap that it remains unclear whether RL is still beneficial when monolingual data is used, we propose a new method to leverage RL to further boost the performance of NMT systems trained with source/target monolingual data. By integrating all our findings, we obtain competitive results on WMT14 English-German, WMT17 English-Chinese, and WMT17 Chinese-English translation tasks, especially setting a state-of-the-art performance on WMT17 Chinese-English translation task.

## 1 Introduction

Recently, neural machine translation (NMT) (Bahdanau et al., 2015; Hassan et al., 2018; Wu et al., 2016; He et al., 2017; Xia et al., 2016, 2017; Wu et al., 2018b,a) has become more and more popular given its superior performance without the demand of heavily hand-crafted engineering efforts. It is usually trained to maximize the likelihood of each token in the target sentence, by taking the source sentence and the preceding (ground-truth) target tokens as inputs. Such training approach is referred as maximum likelihood estimation (MLE) (Scholz, 1985). Although easy to implement, the **token-level**

objective function during training is inconsistent with sequence-level evaluation metrics such as BLEU (Papineni et al., 2002).

To address the inconsistency issue, **reinforcement learning (RL) methods have been adopted to optimize sequence-level objectives**. For example, policy optimization methods such as REINFORCE (Ranzato et al., 2016; Wu et al., 2017b) and actor-critic (Bahdanau et al., 2017) are leveraged for sequence generation tasks including NMT. In machine translation community, a similar method is proposed with the name ‘minimum risk training’ (Shen et al., 2016). All these works demonstrate the effectiveness of RL techniques for NMT models (Wu et al., 2016).

However, effectively applying RL to real-world NMT systems has not been fulfilled by previous works. First, most of, if not all, previous works verified their methods based on shallow recurrent neural network (RNN) models. However, to obtain state-of-the-art (SOTA) performance, it is essential to leverage recently derived deep models (Gehring et al., 2017; Vaswani et al., 2017), which are much more powerful.

Second, it is not easy to make RL practically effective given quite a few widely acknowledged limitations of RL method (Henderson et al., 2018) such as high variance of gradient estimation (Weaver and Tao, 2001), and objective instability (Mnih et al., 2013). Therefore, several tricks are proposed in previous works. However, it remains unclear, and no agreement is achieved on how to use these tricks in machine translation. For example, baseline reward method (Weaver and Tao, 2001) is suggested in (Ranzato et al., 2016; Nguyen et al., 2017; Wu et al., 2016) but not leveraged in (He and Deng, 2012; Shen et al., 2016).

Third, large-scale datasets, especially monolingual datasets are shown to significantly improve translation quality (Sennrich et al., 2015a; Xia et al.,

---

This work was conducted at Microsoft Research Asia.

2016) with MLE training, while it remains nearly empty on how to combine RL with monolingual data in NMT.

In this paper, we try to fulfill these gaps and study how to practically apply RL to obtain strong NMT systems with quite competitive, even state-of-the-art performance. Several comprehensive studies are conducted on different aspects of RL training to figure out how to: 1) set efficient rewards; 2) combine MLE and RL objectives with different weights, which aims to stabilize the training procedure; 3) reduce the variance of gradient estimation.

In addition, given the effectiveness of leveraging monolingual data in improving translation quality, we further propose a new method to **combine the strength of both RL training and source/target monolingual data**. To the best of our knowledge, this is the first work that tries to explore the power of monolingual data when training NMT model with RL method.

We obtain some useful findings through the experiments on WMT17 Chinese-English (Zh-En), WMT17 English-Chinese (En-Zh) and WMT14 English-German (En-De) translation tasks. For instance, multinomial sampling is better than beam search in reward computation, and the combination of RL and monolingual data significantly enhances the NMT model performance. Our main contributions are summarized as follows.

- We provide the first comprehensive study on different aspects of RL training, such as how to setup reward and baseline reward, on top of quite competitive NMT models.
- We propose a new method that effectively leverages large-scale monolingual data, from both the source and target side, when training NMT models with RL.
- Combined with several of our findings and method, we obtain the SOTA translation quality on WMT17 Zh-En translation task, surpassing strong baseline (Transformer big model + back translation) by nearly 1.5 BLEU points. Furthermore, on WMT14 En-De and WMT17 En-Zh translation tasks, we can also obtain strong competitive results.

We hope that our studies and findings will benefit the community to better understand and leverage reinforcement learning for developing strong

NMT models, especially in real-world scenarios faced with deep models and large amount of training data (including both parallel and monolingual data). Towards this end, we open source all our codes/dataset at <https://github.com/apeterswu/RL4NMT> to provide a clear recipe for performance reproduction.

## 2 Background

In this section, we first introduce the attention-based sequence-to-sequence learning framework for neural machine translation (NMT), and then introduce the basis of applying reinforcement learning to training NMT models.

### 2.1 Neural Machine Translation

Typical NMT models are based on the encoder-decoder framework with attention mechanism. The encoder first maps a source sentence  $x = (x_1, x_2, \dots, x_n)$  to a set of continuous representations  $z = (z_1, z_2, \dots, z_n)$ . Given  $z$ , the decoder then generates a target sentence  $y = (y_1, y_2, \dots, y_m)$  of word tokens one by one. At each decoding step  $t$  of model training, the probability of generating a token  $y_t$  is maximized conditioned on  $x$  and  $y_{<t} = (y_1, \dots, y_{t-1})$ . Given  $N$  training sentence pairs  $\{x^i, y^i\}_{i=1}^N$ , maximum likelihood estimation (MLE) is usually adopted to optimize the model, and the training objective is defined as:

$$\begin{aligned} L_{mle} &= \sum_{i=1}^N \log p(y^i | x^i) \\ &= \sum_{i=1}^N \sum_{t=1}^m \log p(y_t^i | y_1^i, \dots, y_{t-1}^i, x^i), \end{aligned} \quad (1)$$

where  $m$  is the length of sentence  $y^i$ .

Among all the encoder-decoder models, the recently proposed Transformer (Vaswani et al., 2017) architecture achieves the best translation quality so far. The main difference between Transformer and previous RNNSearch (Bahdanau et al., 2015) or ConvS2S (Gehring et al., 2017) is that Transformer relies entirely on self-attention (Lin et al., 2017) to compute representations of source and target side sentences, without using recurrent or convolutional operations.

### 2.2 Training NMT with Reinforcement Learning

As aforementioned, reinforcement learning (RL) is leveraged to bridge the gap between training and

inference of NMT, by directly optimizing the evaluation measure (e.g., BLEU) at training time. Specifically, NMT model can be viewed as an *agent*, which interacts with the *environment* (the previous words  $y_{<t}$  and the context vector  $z$  available at each step  $t$ ). The parameters of the agent define a *policy*, i.e., a conditional probability  $p(y_t|x, y_{<t})$ . The agent will pick an *action*, i.e., a candidate word out from the vocabulary, according to the policy. **A terminal reward is observed once the agent generates a complete sequence  $\hat{y}$ .** The reward for machine translation is the BLEU (Papineni et al., 2002) score, denoted as  $R(\hat{y}, y)$ , which is defined by comparing the generated  $\hat{y}$  with the ground-truth sentence  $y$ . Note that here the reward  $R(\hat{y}, y)$  is the sentence-level reward, i.e., a scalar for each complete sentence  $\hat{y}$ . The goal of the RL training is to maximize the expected reward:

$$\begin{aligned} L_{rl} &= \sum_{i=1}^N E_{\hat{y} \sim p(\hat{y}|x^i)} R(\hat{y}, y^i) \\ &= \sum_{i=1}^N \sum_{\hat{y} \in Y} p(\hat{y}|x^i) R(\hat{y}, y^i), \end{aligned} \quad (2)$$

where  $Y$  is the space of all candidate translation sentences, which is exponentially large due to the large vocabulary size, making it impossible to exactly maximize  $L_{rl}$ . In practice, REINFORCE (Williams, 1992) is usually leveraged to approximate the above expectation via sampling  $\hat{y}$  from the policy  $p(y|x)$ , leading to the objective as maximizing:

$$\hat{L}_{rl} = \sum_{i=1}^N R(\hat{y}^i, y^i), \hat{y}^i \sim p(y|x^i), \forall i \in [N]. \quad (3)$$

Throughout the paper we will use REINFORCE as our policy optimization method for RL training.

### 3 Strategies for RL Training

Although training NMT with RL can fill in the gap between training objectives and evaluation metrics, it is not easy to successfully put RL training into practice. A key challenge is that RL methods are highly unstable and inefficient, due to the noise in gradient estimation and reward computation. To our best knowledge, currently there is no consensus, or even a systematic study on how to configure different setups for RL training to avoid such problems, especially for training deep NMT models on large scale datasets. We therefore aim to shed light

on practical applications of RL for NMT training. For this purpose, we provide a comprehensive review of several important methods to stabilize RL training process in this section.

#### 3.1 Reward Computation

It is critical to set up appropriate rewards for RL training, i.e., the  $R(\hat{y}, y)$  in Eqn. (3). There are two important aspects to consider in configuring the reward  $R(\hat{y}, y)$ : how to sample training instance  $\hat{y}$  and whether to use reward shaping.

**Generate  $\hat{y}$**  There are two strategies to sample  $\hat{y}$  for computing the BLEU reward  $R(\hat{y}, y)$ . The first one is *beam search* (Sutskever et al., 2014), it is a breadth-first search method that maintains a “beam” of the top- $K$  scoring candidates (prefix hypothesis sentences) at each generation step. Then, for each candidate sentence in the beam,  $K$  most likely words are appended, resulting in a pool of  $K \times K$  new candidates. Out from this pool, the top- $K$  translations with largest probabilities are selected, and the beam search process continues. The second strategy is *multinomial sampling* (Chatterjee and Cancedda, 2010), which produces each word one by one through multinomial sampling over the model’s output distribution. Both sampling strategies terminate the expansion of a candidate sentence when an ‘end of sentence’ (<EOS>) token is met.

The choice of different sampling strategies reflects the *exploration-exploitation* dilemma. Beam search strategy generates more accurate  $\hat{y}$  by exploiting the probabilistic space output via current NMT model, while multinomial sampling pays more attention to explore more diverse candidates.

**Whether to Use Reward Shaping** From Eqn. (3) we can see that for the entire sequence  $\hat{y}$ , there is only one terminal reward  $R(\hat{y}, y)$  available for model training. Note that the agent needs to take tens of actions (with the number depending on the length of  $\hat{y}$ ) to generate a complete sentence  $\hat{y}$ , but only one reward is available for all those actions. Consequently, RL training is inefficient due to the sparsity of rewards, and the model updates each token in the training sentence with the same reward value without distinction. Reward shaping (Ng et al., 1999) is a strategy to overcome this shortcoming. In reward shaping, intermediate reward at each decoding step  $t$  is imposed and denoted as  $r_t(\hat{y}_t, y)$ . Bahdanau et al. (2017) sets up the intermediate reward as  $r_t(\hat{y}_t, y) = R(\hat{y}_{1..t}, y) - R(\hat{y}_{1..t-1}, y)$ , where  $R(\hat{y}_{1..t}, y)$  is defined as the BLEU score

of  $\hat{y}_{1...t}$  with respect to  $y$ . Note that we have  $R(\hat{y}, y) = \sum_{t=1}^m r_t(\hat{y}_t, y)$ , where  $m$  is the length of  $\hat{y}$ . During RL training, the cumulative reward  $\sum_{\tau=t}^m r_\tau(\hat{y}_\tau, y)$  is used to update the policy at time step  $t$ . It is verified that using the shaped reward  $r_t$  instead of awarding the whole score  $R(\hat{y}, y)$  does not change the optimal policy (Ng et al., 1999).

### 3.2 Variance Reduction of Gradient Estimation

As mentioned before, the REINFORCE algorithm suffers from high variance in gradient estimation, mainly caused by using single sample  $\hat{y}$  to estimate the expectation. To reduce the variance, Ranzato et al. (2016) subtracts an average reward from the returned reward at each time step  $t$ , and the actual reward used to update the policy is

$$R(\hat{y}, y) - \hat{r}_t, \quad (4)$$

where  $\hat{r}_t$  is the estimated average reward at step  $t$ , named as *baseline reward* (Weaver and Tao, 2001). Together with reward shaping, the updated reward becomes  $\sum_{\tau=t}^m r_\tau(\hat{y}_\tau, y) - \hat{r}_t$  at step  $t$ .

Intuitively speaking, a baseline reward  $\hat{r}_t$  is established, which either encourages a word choice  $\hat{y}_t$  if the induced reward  $R$  satisfies  $R > \hat{r}_t$ , or discourages it if  $R < \hat{r}_t$ . Here  $R$  is either the terminal reward  $R(\hat{y}, y)$  or the cumulative reward  $\sum_{\tau=t}^m r_\tau(\hat{y}_\tau, y)$ . Such estimated baseline reward  $\hat{r}_t$  is designed to decrease the high variance of the gradient estimator.

In practice, the baseline reward  $\hat{r}_t$  can be obtained through different approaches. For example, one may sample multiple sentences and use the mean terminal reward for these sentences as baseline reward. In our work, we adopt the function learning approach, using simple network (e.g., multi-layer perceptron) to build the learning function, which is the same as used in (Ranzato et al., 2016; Bahdanau et al., 2017).

### 3.3 Combine MLE and RL Objectives

The last important strategy we would like to mention is the combination of MLE training objective with RL objective, which is assumed to further stabilize RL training process (Wu et al., 2016; Li et al., 2017; Wu et al., 2017a).

A simple way is to linearly combine the MLE (Eqn. (1)) and RL (Eqn. (3)) objectives as follows:

$$L_{com} = \alpha * L_{mle} + (1 - \alpha) * \hat{L}_{rl}, \quad (5)$$

where  $\alpha$  is the hyperparameter controlling the trade-off between MLE and RL objectives. We will empirically evaluate how different values of  $\alpha$  impact the final translation accuracy.

## 4 RL Training with Monolingual Data

Previous works typically conduct RL training with only bilingual data for NMT. Monolingual data has been proved to be able to significantly improve the performance of NMT systems (Sennrich et al., 2015a; Xia et al., 2016; Cheng et al., 2016). It remains an open problem whether it is possible to combine the benefits of RL training and monolingual data such that even more competitive results can be obtained. In this section we provide several solutions for combination and will study them in next section. Note that all the settings discussed in this section are semi-supervised learning, i.e., both bilingual and monolingual data are available.

### 4.1 With Source-Side Monolingual Data

We first provide a solution to RL training with source-side monolingual data. As shown in Eqn. (3), in RL training we need to calculate the reward signal  $R(\hat{y}, y)$  for each generated sentence  $\hat{y}$ , and therefore the reference sentence  $y$  seems to be a must-have, which unfortunately is missing for source-side monolingual data.

We tackle this challenge via generating pseudo target reference  $y$  by bootstrapping with the model itself. Apparently, for the source-side monolingual data, the pseudo target reference  $y$  should have good translation quality. Therefore, for each source-side monolingual sentence, we use the NMT model trained from the bilingual data to beam search a target sentence and treat it as the pseudo target reference  $y$ . Afterwards  $\hat{y}$  is obtained via multinomial sampling to calculate the reward. Although multinomial sampling is usually not as good as sampling via beam search, the combination of beam search (to get the pseudo target reference sentence) and the multinomial sampling (to generate the action sequence of the agent) achieves good exploration-exploitation trade-off, since the pseudo target reference exploits the accuracy of current NMT model while  $\hat{y}$  achieves better exploration.

### 4.2 With Target-Side Monolingual Data

For a target-side monolingual sentence, its source sentence  $x$  is missing, and consequently  $\hat{y}$  is unavailable since it is sampled based on  $x$ . We tackle



this challenge via back translation (Sennrich et al., 2015a). We first train a reverse NMT model from the target language to the source language with bilingual data. For each target-side monolingual sentence, using the reverse NMT model, we back translate it to get its pseudo source sentence  $x$ . We then pair the target monolingual data and its back-translated sentence as a pseudo bilingual sentence pair, which can be used for RL training in the same way as the genuine bilingual sentence pairs.

### 4.3 With both Source-Side and Target-Side Monolingual Data

A natural extension of previous discussions is to combine both the source-side and target-side monolingual data for RL training. We consider two combinations, the *sequential* method and the *unified* method. The former one sequentially leverages the source-side and target-side monolingual data for RL training. Specifically, we first train an MLE model using the bilingual data and source-side (or target-side) monolingual data; based on this MLE model, we then use REINFORCE for training with target-side (or source-side) monolingual data. For *unified* approach, we pack the paired data out from three domains together: the genuine bilingual data, the source monolingual data with its pseudo target references (introduced in subsection 4.1), and the target monolingual data with its back-translated samples (introduced in subsection 4.2). Then we treat the combined data as normal bilingual data on which the NMT model is trained via MLE or RL principles. Our goal is to investigate the model performance with different training data and find the best recipe of how to use these data in RL training. More details are introduced in next section.

## 5 Experiments

In this section, we provide a systematic study on aforementioned RL training strategies and the solutions of leveraging monolingual data. The RL training strategies are evaluated on bilingual datasets from three translation tasks, WMT14 English-German (En-De), WMT17 English-Chinese (En-Zh) and WMT17 Chinese-English (Zh-En), and we further conduct the experiments to leverage monolingual data in WMT17 Zh-En translation.

### 5.1 Experimental Settings

For the bilingual datasets, WMT17 (Bojar et al., 2017) En-Zh<sup>1</sup> and WMT17 Zh-En use the same dataset, which contains about 24M sentences pairs, including CWMT Corpus 2017 and UN Parallel Corpus V1.0. The Jieba<sup>2</sup> segmenter is used to perform Chinese word segmentation. We use byte pair encoding (BPE) (Sennrich et al., 2015b) to preprocess the source and target sentences, forming source-side and target-side dictionary with 40,000 and 37,000 types, respectively. We use the *news-dev2017* as the dev set and *newstest2017* as the test set. For the WMT14 En-De dataset, it contains about 4.5M training pairs, *newstest2012* and *newstest2013* are concatenated as the dev set and *newstest2014* acts as test set. Same as (Vaswani et al., 2017), we also perform BPE to process the En-De dataset, the shared source-target vocabulary contains about 37,000 tokens.

For the monolingual dataset on Zh-En translation task, similar to (Sennrich et al., 2017), the Chinese monolingual data comes from LDC Chinese Gigaword (4th edition) and the English monolingual data comes from News Crawl 2016 articles. After preprocessing (e.g., language detection and filtering sentences with more than 80 words), we keep 4M Chinese sentences and 7M English sentences.

We adopt the Transformer model with *transformer.big* setting as defined in (Vaswani et al., 2017) for Zh-En and En-Zh translations, which achieves SOTA translation quality in several other datasets. For En-De translation, we utilize the *transformer\_base\_v1* setting. These settings are exactly same as used in the original paper, except we set the *layer\_prepostprocess\_dropout* for Zh-En and En-Zh translation to be 0.05. The optimizer used for MLE training is Adam (Kingma and Ba, 2015) with initial learning rate is 0.1, and we follow the same learning rate schedule in (Vaswani et al., 2017). During training, roughly 4,096 source tokens and 4,096 target tokens are paired in one mini batch. Each model is trained using 8 NVIDIA Tesla M40 GPUs. For RL training, the model is initialized with parameters of the MLE model (trained with only bilingual data), and we continue training it with learning rate 0.0001. Same as (Bahdanau et al., 2017), to calculate the BLEU reward, we start all n-gram counts from 1 instead of 0 and

<sup>1</sup><http://www.statmt.org/wmt17/translation-task.html>

<sup>2</sup><https://github.com/fxsjy/jieba>

Training Strategy	En-De	En-Zh	Zh-En
MLE	27.02	34.12	24.29
RL (beam + terminal)	27.06	34.25	24.42
RL (multinomial + terminal)	<b>27.22</b>	<b>34.46</b>	<b>24.70</b>
RL (beam + shaping)	27.04	34.28	24.47
RL (multinomial + shaping)	<b>27.23</b>	<b>34.47</b>	<b>24.72</b>

Table 1: Results of different strategies for reward computation. ‘beam’ refers to ‘beam search and ‘multinomial’ to ‘multinomial sampling’. While generating  $\hat{y}$  through beam search, we use width 4. ‘shaping’ refers to using reward shaping and ‘terminal’ refers not.

multiply the resulting score by the length of the target reference sentence. For inference, we use beam search with width 6. We run each setting for at least 5 times and report the averaged case sensitive BLEU scores<sup>3</sup> (Papineni et al., 2002) on test set. The test set BLEU is chosen via the best configuration based on the validation set.

## 5.2 Results of of RL Training Strategies

We first evaluate different strategies for RL training, based only on bilingual datasets from previously introduced three translation tasks.

**Reward Computation** As reviewed in subsection 3.1, for reward computation, we need to consider how to sample  $\hat{y}$  and whether to use reward shaping.

The results are shown in Table 1, where “RL” stands for RL training with the REINFORCE algorithm. We also report the performance of the pre-trained NMT model with the MLE loss. From the table, an interesting finding is that  $\hat{y}$  sampled via beam search strategy is worse than that by multinomial sampling, with a gap of roughly 0.2-0.3 BLEU points on the test set (with significant test score  $\rho < 0.05$ ). We therefore conjecture that exploration is more important than exploitation in reward computing: multinomial sampling brings more data diversity to the training of NMT model, while sentences generated by beam search are usually very similar to each other. Furthermore, we find that there is no big difference between the leverage of reward shaping or terminal reward, with only slightly better performance of reward shaping. We therefore use multinomial sampling and reward shaping in later experiments.

<sup>3</sup>Calculated by SacreBLEU toolkit, which produces exactly the same evaluation result as that in WMT17 Zh-En campaign. <https://github.com/aws-labs/sockeye/tree/master/contrib/sacrebleu>

Training Strategy	En-De	En-Zh	Zh-En
RL	27.23	34.47	24.72
RL (baseline function)	27.25	34.43	24.73

Table 2: Results of variance reduction of gradient estimation.

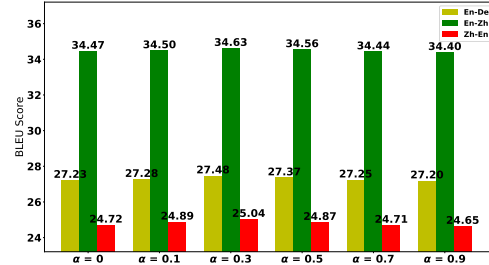


Figure 1: Results of different weights  $\alpha$  to combine MLE and RL objectives.

## Variance Reduction of Gradient Estimation

Next we evaluate the strategies for reducing variance of gradient estimation (see section 3.2). We want to know whether the *baseline reward* is necessary. To compute the baseline reward, similar to (Ranzato et al., 2016; Bahdanau et al., 2017), we build a two-layer MLP regressor with Relu (Nair and Hinton, 2010) activation units. The function takes the hidden states from decoder as input, and the parameters of the regressor are trained to minimize the mean squared loss of Eqn. (4). We first pre-train the baseline function for 20k steps/mini-batches, and then jointly train NMT model (with RL) and the baseline reward function.

Table 2 shows that the learning of baseline reward does not help RL training. This contradicts with previous observations (Ranzato et al., 2016), and seems to suggest that the variance of gradient estimation in NMT is not as large as we expected. The reason might be that the probability mass on the target-side language space induced by the NMT model is highly concentrated, making the sampled  $\hat{y}$  representative enough in terms of estimating the expectation. Therefore, for the economic perspective, it is not necessary to add the additional steps of using baseline reward on RL training for NMT.

**Combine MLE and RL Objectives** As shown in Eqn. (5), the hyperparameter  $\alpha$  controls the trade-off between MLE and RL objectives. For comparison, we set  $\alpha$  to be [0, 0.1, 0.3, 0.5, 0.7, 0.9] in our experiments. The results are presented in Figure 1.

[Data] (Objective)	Valid	Test
[B] (MLE)	22.32	24.29
[B] (MLE) + [B] (RL)	22.87	25.04
[B] (MLE) + [Ms] (RL)	<b>23.03</b>	<b>25.22</b>
[B & Ms] (MLE)	24.31	25.31
[B & Ms] (MLE) + [B & Ms] (RL)	24.58	25.60

Table 3: Results with source monolingual data. “B” denotes bilingual data, “Ms” denotes source-side monolingual data, “&” denotes data combination.

The results show that combining the MLE objective with the RL objective achieves better performance (27.48 for En-De, 34.63 for En-Zh and 25.04 for Zh-En with  $\alpha = 0.3$ ). This indicates that MLE objective is helpful to stabilize the training and improve the model performance, as we expected. However, further increasing  $\alpha$  does not bring more gain. The best trade-off between MLE and RL objectives in our experiment is  $\alpha = 0.3$ . Therefore, we set  $\alpha = 0.3$  in the following experiments.

### 5.3 Results of RL Training with Monolingual Data

In this subsection, we report the results on both valid and test set of RL training using bilingual and monolingual data in Zh-En translation. From Table 3 to Table 6, “RL” denotes the model trained with RL using multinomial sampling, reward shaping, no baseline reward, and combined objective, based on the observations in the last subsection. “B” denotes bilingual data, “Ms” denotes source-side monolingual data and “Mt” denotes target-side monolingual data, “&” denotes data combination.

**With Source-Side Monolingual Data** As discussed before, we use beam search with beam width 4 to sample the pseudo target sentence  $y$  for each monolingual sentence  $x$ . We consider several settings for RL training: 1) only source-side monolingual data; 2) the combination of bilingual and source-side monolingual data. We first train an MLE model using the augmented dataset combining the genuine bilingual data with the pseudo bilingual data generated from the monolingual data, and then perform RL training on this combined dataset. The results are shown in Table 3.

**With Target-Side Monolingual Data** For target-side monolingual data, we first pre-train a translation model from English to Chinese<sup>4</sup>, and use it to back translate target-side monolingual

<sup>4</sup>The BLEU score of the En-Zh model is 34.12.

[Data] (Objective)	Valid	Test
[B] (MLE)	22.32	24.29
[B] (MLE) + [B] (RL)	22.87	25.04
[B] (MLE) + [Mt] (RL)	<b>22.96</b>	<b>25.15</b>
[B & Mt] (MLE)	24.14	25.24
[B & Mt] (MLE) + [B & Mt] (RL)	24.41	25.58

Table 4: Results with target monolingual data. “B” denotes bilingual data, “Mt” denotes target-side monolingual data, “&” denotes data combination.

[Data] (Objective)	Valid	Test
[B & Ms] (MLE)	24.31	25.31
[B & Ms] (MLE) + [B & Ms] (RL)	24.58	25.60
[B & Ms] (MLE) + [Mt] (RL)	<b>24.61</b>	<b>25.72</b>
[B & Mt] (MLE)	24.14	25.24
[B & Mt] (MLE) + [B & Mt] (RL)	24.41	25.58
[B & Mt] (MLE) + [Ms] (RL)	<b>24.75</b>	<b>25.92</b>

Table 5: Results of *sequential* approach for monolingual data. “B” denotes bilingual data, “Ms” denotes source-side monolingual data and “Mt” denotes target-side monolingual data, “&” denotes data combination.

[Data] (Objective)	Valid	Test
[B & Ms & Mt] (MLE)	25.58	26.13
+ [B & Ms & Mt] (RL)	<b>25.90</b>	<b>26.73</b>

Table 6: Results of *unified* approach for monolingual data. “+” means to initialize the RL model using above MLE model, which is trained on the combination of bilingual data, source-side monolingual data and target-side monolingual data.

sentence  $y$  to get pseudo source sentence  $x$ . Similarly, we consider several settings for RL training: 1) only target-side monolingual data; 2) the combination of bilingual data and target-side monolingual data. We train an MLE model using both the genuine and the generated pseudo bilingual data, and then perform RL training on this data. The results are presented in Table 4.

From Table 3 and 4, we have several observations. First, monolingual data helps RL training, improving BLEU score from 25.04 to 25.22 ( $\rho < 0.05$ ) in Table 3. Second, when we only add monolingual data for RL training, the model achieves similar performance compared to MLE training with bilingual and monolingual data (e.g., 25.15 vs. 25.24 ( $\rho < 0.05$ ) in Table 4).

**With both Source-Side and Target-Side Monolingual Data** We have two approaches to use both source-side and target-side monolingual data, as described in subsection 4.3. The results are reported in Table 5 and Table 6.

From Table 5, we can observe that the *sequen-*

System	Architecture	BLEU
<i>Existing end-to-end NMT systems</i>		
Vaswani et al. (2017)	Transformer	24.29
Sennrich et al. (2015a)	Transformer + Target Monolingual Data (i.e., back translation)	25.24
SougouKnowing	Stacked LSTM model + Reranking	24.00
SougouKnowing-ensemble	Stacked LSTM model + Reranking + Ensemble	26.40
<i>Our end-to-end NMT</i>		
<i>this work</i>	Transformer + <i>RL</i>	25.04
	Transformer + Source Monolingual Data	25.31
	Transformer + Source Monolingual Data + <i>RL</i>	25.60
	Transformer + Target Monolingual Data	25.24
	Transformer + Target Monolingual Data + <i>RL</i>	25.58
	Transformer + Source & Target Monolingual Data	26.13
	Transformer + Source & Target Monolingual Data + <i>RL</i>	<b>26.73</b>

Table 7: Comparisons of different competitive end-to-end NMT systems. SougouKnowing results come from [http://matrix.statmt.org/matrix/systems\\_list/1878](http://matrix.statmt.org/matrix/systems_list/1878).

tial training of monolingual data can benefit the model performance. Taking the last three rows as an example, the BLEU score of the MLE model trained on the combination of bilingual data and target-side monolingual data is 25.24; based on this model, RL training using the source-side monolingual data further improves the model performance by 0.7 ( $\rho < 0.01$ ) BLEU points. From Table 6, we can observe on top of a quite strong MLE baseline (26.13), through the *unified* RL training, we can still improve the test set by 0.6 points to 26.73 ( $\rho < 0.01$ ), which shows the effectiveness of combining source/target monolingual data and reinforcement learning.

#### 5.4 Comparison with Other Models

At last, as a summary of our empirical results, we compare several representative end-to-end NMT systems to our work in Table 7, which includes the Transformer (Vaswani et al., 2017) model, with/without back-translation (Sennrich et al., 2015a) and the best NMT system in WMT17 Chinese-English translation challenge<sup>5</sup> (SougouKnowing-ensemble). The results clearly show that after combining both source-side and target-side monolingual data with RL training, we obtain the state-of-the-art BLEU score 26.73, even surpassing the best ensemble model in WMT17 Zh-En translation challenge.

## 6 Related Work

Our work is mainly related with the literature of using reinforcement learning to directly optimize the evaluation measure for neural machine translation. Several representative works are (Ranzato et al.,

2016; Shen et al., 2016; Bahdanau et al., 2017). In (Ranzato et al., 2016), the authors propose to train a neural translation model with the objective gradually shifting from maximizing token-level likelihood to optimizing the sentence-level BLEU score. Shen et al. (2016) proposes to adopt minimum risk training (Goel and Byrne, 2000) to minimize the task specific expected loss (i.e., induced by BLEU score) on NMT training data. Instead of the REINFORCE (Williams, 1992) algorithm used in the above two works, Bahdanau et al. (2017) further optimizes the policy by actor-critic algorithm. Wu et al. (2016) introduces a simple RL based method to optimize the stacked LSTM model for NMT, achieving better BLEU scores on English-French translation but not on English-German. Edunov et al. (2017) presents a comparative study of several classical structural prediction losses for NMT model, which also includes sequence-level loss but not exactly the same as RL.

Our work is also related with the research works that leverage monolingual data for improving NMT models (Zhang and Zong, 2016; Sennrich et al., 2015a; Wang et al., 2018; Xia et al., 2016; Cheng et al., 2016). Zhang and Zong (2016) exploits the source-side monolingual data in NMT. Sennrich et al. (2015a) proposes back-translation method to leverage target-side monolingual data for NMT. Xia et al. (2016) formulates the machine translation as a communication game, which leverages the power of two directional translation models and source/target monolingual data. Cheng et al. (2016) proposes a similar semi-supervised approach. However, none of these works have explored the power of monolingual data in the context of training NMT model with reinforcement learning.

<sup>5</sup>[http://matrix.statmt.org/matrix/systems\\_list/1878](http://matrix.statmt.org/matrix/systems_list/1878)



## 7 Conclusion

In this work, we presented a study of how to effectively train NMT models using reinforcement learning. Different RL strategies were evaluated in German-English, English-Chinese and Chinese-English translation tasks on large-scale bilingual datasets. We found that (1) multinomial sampling is better than beam search, (2) several previous tricks such as reward shaping and baseline reward does not make significant difference, and (3) the combination of the MLE and RL objectives is important. In addition, we explored the source/target monolingual data for RL training. By combining the power of RL and monolingual data, we achieve the state-of-the-art BLEU score on WMT17 Chinese-English translation task. We hope that our study and results can benefit the community and bring some insights on how to train deep NMT models with reinforcement learning and big data.

## References

- Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2017. An actor-critic algorithm for sequence prediction. *Fifth International Conference on Learning Representations*.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. *Third International Conference on Learning Representations*.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. 2017. Findings of the 2017 conference on machine translation. In *Proceedings of the Second Conference on Machine Translation, Volume 2: Shared Task Papers*, pages 169–214, Copenhagen, Denmark. Association for Computational Linguistics.
- Samidh Chatterjee and Nicola Cancedda. 2010. Minimum error rate training by sampling the translation lattice. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 606–615. Association for Computational Linguistics.
- Yong Cheng, Wei Xu, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. Semi-supervised learning for neural machine translation. *meeting of the association for computational linguistics*, pages 1965–1974.
- Sergey Edunov, Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2017. Classical structured prediction losses for sequence to sequence learning. *Proceedings of NAACL-HLT 2018*.
- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional sequence to sequence learning. *Proceedings of the 34th international conference on machine learning*.
- Vaibhava Goel and William J Byrne. 2000. Minimum bayes-risk automatic speech recognition. *Computer Speech & Language*, 14(2):115–135.
- Hany Hassan, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, et al. 2018. Achieving human parity on automatic chinese to english news translation. *arXiv preprint arXiv:1803.05567*.
- Di He, Hanqing Lu, Yingce Xia, Tao Qin, Liwei Wang, and Tieyan Liu. 2017. Decoding with value networks for neural machine translation. In *Advances in Neural Information Processing Systems*, pages 178–187.
- Xiaodong He and Li Deng. 2012. Maximum expected bleu training of phrase and lexicon translation models. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 292–301. Association for Computational Linguistics.
- Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. 2018. Deep reinforcement learning that matters. *Thirty-Second AAAI Conference On Artificial Intelligence*.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. *Third International Conference on Learning Representations*.
- Jiwei Li, Will Monroe, Tianlin Shi, Alan Ritter, and Dan Jurafsky. 2017. Adversarial learning for neural dialogue generation. *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*.
- Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017. A structured self-attentive sentence embedding. *Fifth International Conference on Learning Representations*.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing atari with deep reinforcement learning. *Advances in neural information processing systems, workshop*.
- Vinod Nair and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning*, pages 807–814.

- Andrew Y Ng, Daishi Harada, and Stuart Russell. 1999. Policy invariance under reward transformations: Theory and application to reward shaping. In *ICML*, volume 99, pages 278–287.
- Khanh Nguyen, Hal Daumé III, and Jordan Boyd-Graber. 2017. Reinforcement learning for bandit neural machine translation with simulated human feedback. *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2016. Sequence level training with recurrent neural networks. *Fourth International Conference on Learning Representations*.
- FW Scholz. 1985. Maximum likelihood estimation. *Encyclopedia of statistical sciences*.
- Rico Sennrich, Alexandra Birch, Anna Currey, Ulrich Germann, Barry Haddow, Kenneth Heafield, Antonio Valerio Miceli Barone, and Philip Williams. 2017. The university of edinburgh’s neural mt systems for wmt17. In *Proceedings of the Second Conference on Machine Translation*, pages 389–399. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015a. Improving neural machine translation models with monolingual data. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015b. Neural machine translation of rare words with subword units. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*.
- Shiqi Shen, Yong Cheng, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. Minimum risk training for neural machine translation. *meeting of the association for computational linguistics*, pages 1683–1692.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 6000–6010.
- Yijun Wang, Yingce Xia, Li Zhao, Jiang Bian, Tao Qin, Guiquan Liu, and T Liu. 2018. Dual transfer learning for neural machine translation with marginal distribution regularization. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*.
- Lex Weaver and Nigel Tao. 2001. The optimal reward baseline for gradient-based reinforcement learning. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, pages 538–545. Morgan Kaufmann Publishers Inc.
- Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. In *Reinforcement Learning*, pages 5–32. Springer.
- Lijun Wu, Xu Tan, Di He, Fei Tian, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2018a. Beyond error propagation in neural machine translation: Characteristics of language also matter. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Lijun Wu, Fei Tian, Li Zhao, Jianhuang Lai, and Tie-Yan Liu. 2018b. Word attention for sequence to sequence text understanding. In *AAAI*.
- Lijun Wu, Yingce Xia, Li Zhao, Fei Tian, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2017a. Adversarial neural machine translation. *arXiv preprint arXiv:1704.06933*.
- Lijun Wu, Li Zhao, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2017b. Sequence prediction with unlabeled data by reward function learning. *IJCAI-17*, pages 3098–3104.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Yingce Xia, Di He, Tao Qin, Liwei Wang, Nenghai Yu, Tieyan Liu, and Weiyang Ma. 2016. Dual learning for machine translation. *neural information processing systems*, pages 820–828.
- Yingce Xia, Fei Tian, Lijun Wu, Jianxin Lin, Tao Qin, Nenghai Yu, and Tie-Yan Liu. 2017. Deliberation networks: Sequence generation beyond one-pass decoding. In *Advances in Neural Information Processing Systems*, pages 1784–1794.
- Jiajun Zhang and Chengqing Zong. 2016. Exploiting source-side monolingual data in neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1535–1545.