

Minimum Risk Training for Neural Machine Translation

Shiqi Shen[†], Yong Cheng[#], Zhongjun He⁺, Wei He⁺, Hua Wu⁺, Maosong Sun[†], Yang Liu^{†*}

[†]State Key Laboratory of Intelligent Technology and Systems

Tsinghua National Laboratory for Information Science and Technology

Department of Computer Science and Technology, Tsinghua University, Beijing, China

[#]Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing, China

⁺Baidu Inc., Beijing, China

{vicapple22, chengyong3001}@gmail.com, {hezongjun, hewei06, wu_hua}@baidu.com, {sms, liuyang2011}@tsinghua.edu.cn

Abstract

We propose minimum risk training for end-to-end neural machine translation. Unlike conventional maximum likelihood estimation, minimum risk training is capable of optimizing model parameters directly with respect to arbitrary evaluation metrics, which are not necessarily differentiable. Experiments show that our approach achieves significant improvements over maximum likelihood estimation on a state-of-the-art neural machine translation system across various languages pairs. Transparent to architectures, our approach can be applied to more neural networks and potentially benefit more NLP tasks.

1 Introduction

Recently, end-to-end neural machine translation (NMT) (Kalchbrenner and Blunsom, 2013; Sutskever et al., 2014; Bahdanau et al., 2015) has attracted increasing attention from the community. Providing a new paradigm for machine translation, NMT aims at training a single, large neural network that directly transforms a source-language sentence to a target-language sentence without explicitly modeling latent structures (e.g., word alignment, phrase segmentation, phrase re-ordering, and SCFG derivation) that are vital in conventional statistical machine translation (SMT) (Brown et al., 1993; Koehn et al., 2003; Chiang, 2005).

Current NMT models are based on the *encoder-decoder* framework (Cho et al., 2014; Sutskever et al., 2014), with an encoder to read and encode a source-language sentence into a vector, from which a decoder generates a target-language sentence. While early efforts encode the input into a

fixed-length vector, Bahdanau et al. (2015) advocate the attention mechanism to dynamically generate a context vector for a target word being generated.

Although NMT models have achieved results on par with or better than conventional SMT, they still suffer from a major drawback: the models are optimized to maximize the likelihood of training data instead of evaluation metrics that actually quantify translation quality. Ranzato et al. (2015) indicate two drawbacks of *maximum likelihood estimation* (MLE) for NMT. First, the models are only exposed to the training distribution instead of model predictions. Second, the loss function is defined at the word level instead of the sentence level.

In this work, we introduce *minimum risk training* (MRT) for neural machine translation. The new training objective is to minimize the expected loss (i.e., risk) on the training data. MRT has the following advantages over MLE:

1. *Direct optimization with respect to evaluation metrics*: MRT introduces evaluation metrics as loss functions and aims to minimize expected loss on the training data.
2. *Applicable to arbitrary loss functions*: our approach allows arbitrary sentence-level loss functions, which are not necessarily differentiable.
3. *Transparent to architectures*: MRT does not assume the specific architectures of NMT and can be applied to any end-to-end NMT systems.

While MRT has been widely used in conventional SMT (Och, 2003; Smith and Eisner, 2006; He and Deng, 2012) and deep learning based MT (Gao et al., 2014), to the best of our knowledge, this work is the first effort to introduce MRT

*Corresponding author: Yang Liu.

into end-to-end NMT. Experiments on a variety of language pairs (Chinese-English, English-French, and English-German) show that MRT leads to significant improvements over MLE on a state-of-the-art NMT system (Bahdanau et al., 2015).

2 Background

Given a source sentence $\mathbf{x} = \mathbf{x}_1, \dots, \mathbf{x}_m, \dots, \mathbf{x}_M$ and a target sentence $\mathbf{y} = \mathbf{y}_1, \dots, \mathbf{y}_n, \dots, \mathbf{y}_N$, end-to-end NMT directly models the translation probability:

$$P(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) = \prod_{n=1}^N P(\mathbf{y}_n|\mathbf{x}, \mathbf{y}_{<n}; \boldsymbol{\theta}), \quad (1)$$

where $\boldsymbol{\theta}$ is a set of model parameters and $\mathbf{y}_{<n} = \mathbf{y}_1, \dots, \mathbf{y}_{n-1}$ is a partial translation.

Predicting the n -th target word can be modeled by using a recurrent neural network:

$$P(\mathbf{y}_n|\mathbf{x}, \mathbf{y}_{<n}; \boldsymbol{\theta}) \propto \exp \left\{ q(\mathbf{y}_{n-1}, \mathbf{z}_n, \mathbf{c}_n, \boldsymbol{\theta}) \right\}, \quad (2)$$

where \mathbf{z}_n is the n -th hidden state on the target side, \mathbf{c}_n is the context for generating the n -th target word, and $q(\cdot)$ is a non-linear function. Current NMT approaches differ in calculating \mathbf{z}_n and \mathbf{c}_n and defining $q(\cdot)$. Please refer to (Sutskever et al., 2014; Bahdanau et al., 2015) for more details.

Given a set of training examples $D = \{(\mathbf{x}^{(s)}, \mathbf{y}^{(s)})\}_{s=1}^S$, the standard training objective is to maximize the log-likelihood of the training data:

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left\{ \mathcal{L}(\boldsymbol{\theta}) \right\}, \quad (3)$$

where

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{s=1}^S \log P(\mathbf{y}^{(s)}|\mathbf{x}^{(s)}; \boldsymbol{\theta}) \quad (4)$$

$$= \sum_{s=1}^S \sum_{n=1}^{N^{(s)}} \log P(\mathbf{y}_n^{(s)}|\mathbf{x}^{(s)}, \mathbf{y}_{<n}^{(s)}; \boldsymbol{\theta}). \quad (5)$$

We use $N^{(s)}$ to denote the length of the s -th target sentence $\mathbf{y}^{(s)}$.

The partial derivative with respect to a model parameter θ_i is calculated as

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \theta_i} = \sum_{s=1}^S \sum_{n=1}^{N^{(s)}} \frac{\partial P(\mathbf{y}_n^{(s)}|\mathbf{x}^{(s)}, \mathbf{y}_{<n}^{(s)}; \boldsymbol{\theta}) / \partial \theta_i}{P(\mathbf{y}_n^{(s)}|\mathbf{x}^{(s)}, \mathbf{y}_{<n}^{(s)}; \boldsymbol{\theta})}. \quad (6)$$

Ranzato et al. (2015) point out that MLE for end-to-end NMT suffers from two drawbacks.

First, while the models are trained only on the training data distribution, they are used to generate target words on previous model predictions, which can be erroneous, at test time. This is referred to as *exposure bias* (Ranzato et al., 2015). Second, MLE usually uses the cross-entropy loss focusing on word-level errors to maximize the probability of the next correct word, which might hardly correlate well with corpus-level and sentence-level evaluation metrics such as BLEU (Papineni et al., 2002) and TER (Snover et al., 2006).

As a result, it is important to introduce new training algorithms for end-to-end NMT to include model predictions during training and optimize model parameters directly with respect to evaluation metrics.

3 Minimum Risk Training for Neural Machine Translation

Minimum risk training (MRT), which aims to minimize the expected loss on the training data, has been widely used in conventional SMT (Och, 2003; Smith and Eisner, 2006; He and Deng, 2012) and deep learning based MT (Gao et al., 2014). The basic idea is to introduce evaluation metrics as loss functions and assume that the optimal set of model parameters should minimize the expected loss on the training data.

Let $(\mathbf{x}^{(s)}, \mathbf{y}^{(s)})$ be the s -th sentence pair in the training data and \mathbf{y} be a model prediction. We use a *loss function* $\Delta(\mathbf{y}, \mathbf{y}^{(s)})$ to measure the discrepancy between the model prediction \mathbf{y} and the gold-standard translation $\mathbf{y}^{(s)}$. Such a loss function can be negative smoothed sentence-level evaluation metrics such as BLEU (Papineni et al., 2002), NIST (Doddington, 2002), TER (Snover et al., 2006), or METEOR (Lavie and Denkowski, 2009) that have been widely used in machine translation evaluation. Note that a loss function is not parameterized and thus not differentiable.

In MRT, the *risk* is defined as the expected loss with respect to the posterior distribution:

$$\mathcal{R}(\boldsymbol{\theta}) = \sum_{s=1}^S \mathbb{E}_{\mathbf{y}|\mathbf{x}^{(s)}; \boldsymbol{\theta}} \left[\Delta(\mathbf{y}, \mathbf{y}^{(s)}) \right] \quad (7)$$

$$= \sum_{s=1}^S \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x}^{(s)})} P(\mathbf{y}|\mathbf{x}^{(s)}; \boldsymbol{\theta}) \Delta(\mathbf{y}, \mathbf{y}^{(s)}), \quad (8)$$

where $\mathcal{Y}(\mathbf{x}^{(s)})$ is a set of all possible candidate translations for $\mathbf{x}^{(s)}$.

	$\Delta(\mathbf{y}, \mathbf{y}^{(s)})$	$P(\mathbf{y} \mathbf{x}^{(s)}; \boldsymbol{\theta})$			
\mathbf{y}_1	-1.0	0.2	0.3	0.5	0.7
\mathbf{y}_2	-0.3	0.5	0.2	0.2	0.1
\mathbf{y}_3	-0.5	0.3	0.5	0.3	0.2
$\mathbb{E}_{\mathbf{y} \mathbf{x}^{(s)}; \boldsymbol{\theta}}[\Delta(\mathbf{y}, \mathbf{y}^{(s)})]$		-0.50	-0.61	-0.71	-0.83

Table 1: Example of minimum risk training. $\mathbf{x}^{(s)}$ is an observed source sentence, $\mathbf{y}^{(s)}$ is its corresponding gold-standard translation, and \mathbf{y}_1 , \mathbf{y}_2 , and \mathbf{y}_3 are model predictions. For simplicity, we suppose that the full search space contains only three candidates. The loss function $\Delta(\mathbf{y}, \mathbf{y}^{(s)})$ measures the difference between model prediction and gold-standard. The goal of MRT is to find a distribution (the last column) that correlates well with the gold-standard by minimizing the expected loss.

The training objective of MRT is to minimize the risk on the training data:

$$\hat{\boldsymbol{\theta}}_{\text{MRT}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left\{ \mathcal{R}(\boldsymbol{\theta}) \right\}. \quad (9)$$

Intuitively, while MLE aims to maximize the likelihood of training data, our training objective is to discriminate between candidates. For example, in Table 1, suppose the candidate set $\mathcal{Y}(\mathbf{x}^{(s)})$ contains only three candidates: \mathbf{y}_1 , \mathbf{y}_2 , and \mathbf{y}_3 . According to the losses calculated by comparing with the gold-standard translation $\mathbf{y}^{(s)}$, it is clear that \mathbf{y}_1 is the best candidate, \mathbf{y}_3 is the second best, and \mathbf{y}_2 is the worst: $\mathbf{y}_1 > \mathbf{y}_3 > \mathbf{y}_2$. The right half of Table 1 shows four models. As model 1 (column 3) ranks the candidates in a reverse order as compared with the gold-standard (i.e., $\mathbf{y}_2 > \mathbf{y}_3 > \mathbf{y}_1$), it obtains the highest risk of -0.50 . Achieving a better correlation with the gold-standard than model 1 by predicting $\mathbf{y}_3 > \mathbf{y}_1 > \mathbf{y}_2$, model 2 (column 4) reduces the risk to -0.61 . As model 3 (column 5) ranks the candidates in the same order with the gold-standard, the risk goes down to -0.71 . The risk can be further reduced by concentrating the probability mass on \mathbf{y}_1 (column 6). As a result, by minimizing the risk on the training data, we expect to obtain a model that correlates well with the gold-standard.

In MRT, the partial derivative with respect to a model parameter θ_i is given by

$$\begin{aligned} & \frac{\partial \mathcal{R}(\boldsymbol{\theta})}{\partial \theta_i} \\ &= \sum_{s=1}^S \mathbb{E}_{\mathbf{y}|\mathbf{x}^{(s)}; \boldsymbol{\theta}} \left[\Delta(\mathbf{y}, \mathbf{y}^{(s)}) \times \right. \\ & \quad \left. \sum_{n=1}^{N^{(s)}} \frac{\partial P(\mathbf{y}_n|\mathbf{x}^{(s)}, \mathbf{y}_{<n}; \boldsymbol{\theta}) / \partial \theta_i}{P(\mathbf{y}_n|\mathbf{x}^{(s)}, \mathbf{y}_{<n}; \boldsymbol{\theta})} \right]. \quad (10) \end{aligned}$$

Since Eq. (10) suggests there is no need to differentiate $\Delta(\mathbf{y}, \mathbf{y}^{(s)})$, MRT allows arbitrary non-differentiable loss functions. In addition, our approach is transparent to architectures and can be applied to arbitrary end-to-end NMT models.

Despite these advantages, MRT faces a major challenge: the expectations in Eq. (10) are usually intractable to calculate due to the exponential search space of $\mathcal{Y}(\mathbf{x}^{(s)})$, the non-decomposability of the loss function $\Delta(\mathbf{y}, \mathbf{y}^{(s)})$, and the context sensitiveness of NMT.

To alleviate this problem, we propose to only use a subset of the full search space to approximate the posterior distribution and introduce a new training objective:

$$\begin{aligned} \tilde{\mathcal{R}}(\boldsymbol{\theta}) &= \sum_{s=1}^S \mathbb{E}_{\mathbf{y}|\mathbf{x}^{(s)}; \boldsymbol{\theta}, \alpha} \left[\Delta(\mathbf{y}, \mathbf{y}^{(s)}) \right] \quad (11) \\ &= \sum_{s=1}^S \sum_{\mathbf{y} \in \mathcal{S}(\mathbf{x}^{(s)})} Q(\mathbf{y}|\mathbf{x}^{(s)}; \boldsymbol{\theta}, \alpha) \Delta(\mathbf{y}, \mathbf{y}^{(s)}), \quad (12) \end{aligned}$$

where $\mathcal{S}(\mathbf{x}^{(s)}) \subset \mathcal{Y}(\mathbf{x}^{(s)})$ is a sampled subset of the full search space, and $Q(\mathbf{y}|\mathbf{x}^{(s)}; \boldsymbol{\theta}, \alpha)$ is a distribution defined on the subspace $\mathcal{S}(\mathbf{x}^{(s)})$:

$$Q(\mathbf{y}|\mathbf{x}^{(s)}; \boldsymbol{\theta}, \alpha) = \frac{P(\mathbf{y}|\mathbf{x}^{(s)}; \boldsymbol{\theta})^\alpha}{\sum_{\mathbf{y}' \in \mathcal{S}(\mathbf{x}^{(s)})} P(\mathbf{y}'|\mathbf{x}^{(s)}; \boldsymbol{\theta})^\alpha}. \quad (13)$$

Note that α is a hyper-parameter that controls the sharpness of the Q distribution (Och, 2003).

Algorithm 1 shows how to build $\mathcal{S}(\mathbf{x}^{(s)})$ by sampling the full search space. The sampled subset initializes with the gold-standard translation (line 1). Then, the algorithm keeps sampling a target word given the source sentence and the partial translation until reaching the end of sentence (lines 3-16). Note that sampling might produce duplicate candidates, which are removed when building

Input: the s -th source sentence in the training data $\mathbf{x}^{(s)}$, the s -th target sentence in the training data $\mathbf{y}^{(s)}$, the set of model parameters θ , the limit on the length of a candidate translation l , and the limit on the size of sampled space k .

Output: sampled space $\mathcal{S}(\mathbf{x}^{(s)})$.

```

1  $\mathcal{S}(\mathbf{x}^{(s)}) \leftarrow \{\mathbf{y}^{(s)}\}$ ; // the gold-standard translation is included
2  $i \leftarrow 1$ ;
3 while  $i \leq k$  do
4    $\mathbf{y} \leftarrow \emptyset$ ; // an empty candidate translation
5    $n \leftarrow 1$ ;
6   while  $n \leq l$  do
7      $y \sim P(\mathbf{y}_n | \mathbf{x}^{(s)}, \mathbf{y}_{<n}; \theta)$ ; // sample the  $n$ -th target word
8      $\mathbf{y} \leftarrow \mathbf{y} \cup \{y\}$ ;
9     if  $y = \text{EOS}$  then
10      | break; // terminate if reach the end of sentence
11     end
12      $n \leftarrow n + 1$ ;
13   end
14    $\mathcal{S}(\mathbf{x}^{(s)}) \leftarrow \mathcal{S}(\mathbf{x}^{(s)}) \cup \{\mathbf{y}\}$ ;
15    $i \leftarrow i + 1$ ;
16 end

```

Algorithm 1: Sampling the full search space.

the subspace. We find that it is inefficient to force the algorithm to generate exactly k distinct candidates because high-probability candidates can be sampled repeatedly, especially when the probability mass highly concentrates on a few candidates. In practice, we take advantage of GPU’s parallel architectures to speed up the sampling.¹

Given the sampled space, the partial derivative with respect to a model parameter θ_i of $\tilde{\mathcal{R}}(\theta)$ is given by

$$\begin{aligned}
& \frac{\partial \tilde{\mathcal{R}}(\theta)}{\partial \theta_i} \\
= & \alpha \sum_{s=1}^S \mathbb{E}_{\mathbf{y} | \mathbf{x}^{(s)}; \theta, \alpha} \left[\frac{\partial P(\mathbf{y} | \mathbf{x}^{(s)}; \theta) / \partial \theta_i}{P(\mathbf{y} | \mathbf{x}^{(s)}; \theta)} \times \right. \\
& \left. \left(\Delta(\mathbf{y}, \mathbf{y}^{(s)}) - \mathbb{E}_{\mathbf{y}' | \mathbf{x}^{(s)}; \theta, \alpha} [\Delta(\mathbf{y}', \mathbf{y}^{(s)})] \right) \right]. \quad (14)
\end{aligned}$$

Since $|\mathcal{S}(\mathbf{x}^{(s)})| \ll |\mathcal{Y}(\mathbf{x}^{(s)})|$, the expectations in Eq. (14) can be efficiently calculated by explicitly enumerating all candidates in $\mathcal{S}(\mathbf{x}^{(s)})$. In our experiments, we find that approximating the full space with 100 samples (i.e., $k = 100$) works very well and further increasing sample size does not lead to significant improvements (see Section 4.3).

¹To build the subset, an alternative to sampling is computing top- k translations. We prefer sampling to computing top- k translations for efficiency: sampling is more efficient and easy-to-implement than calculating k -best lists, especially given the extremely parallel architectures of GPUs.

4 Experiments

4.1 Setup

We evaluated our approach on three translation tasks: Chinese-English, English-French, and English-German. The evaluation metric is BLEU (Papineni et al., 2002) as calculated by the `multi-bleu.perl` script.

For Chinese-English, the training data consists of 2.56M pairs of sentences with 67.5M Chinese words and 74.8M English words, respectively. We used the NIST 2006 dataset as the validation set (hyper-parameter optimization and model selection) and the NIST 2002, 2003, 2004, 2005, and 2008 datasets as test sets.

For English-French, to compare with the results reported by previous work on end-to-end NMT (Sutskever et al., 2014; Bahdanau et al., 2015; Jean et al., 2015; Luong et al., 2015b), we used the same subset of the WMT 2014 training corpus that contains 12M sentence pairs with 304M English words and 348M French words. The concatenation of news-test 2012 and news-test 2013 serves as the validation set and news-test 2014 as the test set.

For English-German, to compare with the results reported by previous work (Jean et al., 2015; Luong et al., 2015a), we used the same subset of the WMT 2014 training corpus that contains 4M sentence pairs with 91M English words and 87M German words. The concatenation of news-test 2012 and news-test 2013 is used as the validation set and news-test 2014 as the test set.

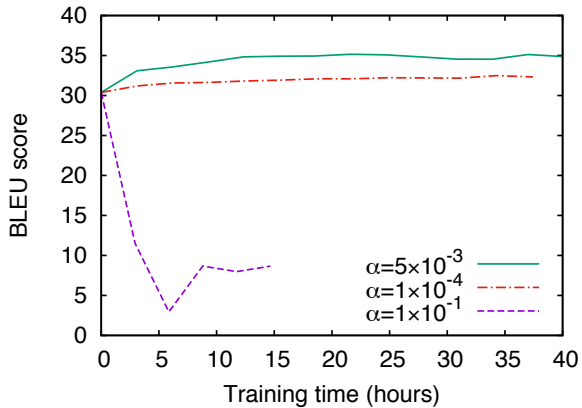


Figure 1: Effect of α on the Chinese-English validation set.

We compare our approach with two state-of-the-art SMT and NMT systems:

1. MOSES (Koehn and Hoang, 2007): a phrase-based SMT system using minimum error rate training (Och, 2003).
2. RNNSEARCH (Bahdanau et al., 2015): an attention-based NMT system using maximum likelihood estimation.

MOSES uses the parallel corpus to train a phrase-based translation model and the target part to train a 4-gram language model using the SRILM toolkit (Stolcke, 2002).² The log-linear model Moses uses is trained by the minimum error rate training (MERT) algorithm (Och, 2003) that directly optimizes model parameters with respect to evaluation metrics.

RNNSEARCH uses the parallel corpus to train an attention-based neural translation model using the maximum likelihood criterion.

On top of RNNSEARCH, our approach replaces MLE with MRT. We initialize our model with the RNNsearch50 model (Bahdanau et al., 2015). We set the vocabulary size to 30K for Chinese-English and English-French and 50K for English-German. The beam size for decoding is 10. The default loss function is negative smoothed sentence-level BLEU.

4.2 Effect of α

The hyper-parameter α controls the smoothness of the Q distribution (see Eq. (13)). As shown in

²It is possible to exploit larger monolingual corpora for both MOSES and RNNSEARCH (Gulcehre et al., 2015; Senrich et al., 2015). We leave this for future work.

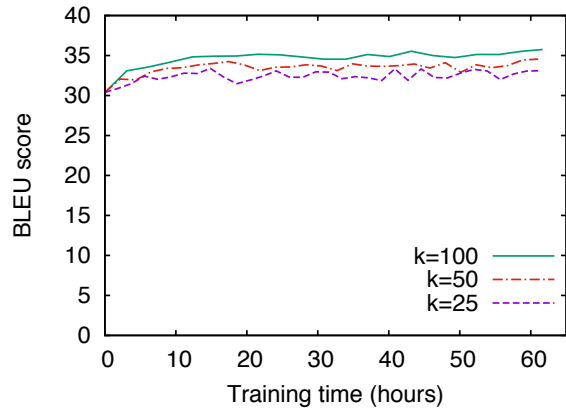


Figure 2: Effect of sample size on the Chinese-English validation set.

critierion	loss	BLEU	TER	NIST
MLE	N/A	30.48	60.85	8.26
MRT	-sBLEU	36.71	53.48	8.90
	sTER	30.14	53.83	6.02
	-sNIST	32.32	56.85	8.90

Table 2: Effect of loss function on the Chinese-English validation set.

Figure 1, we find that α has a critical effect on BLEU scores on the Chinese-English validation set. While $\alpha = 1 \times 10^{-1}$ decreases BLEU scores dramatically, $\alpha = 5 \times 10^{-3}$ improves translation quality significantly and consistently. Reducing α further to 1×10^{-4} , however, results in lower BLEU scores. Therefore, we set $\alpha = 5 \times 10^{-3}$ in the following experiments.

4.3 Effect of Sample Size

For efficiency, we sample k candidate translations from the full search space $\mathcal{Y}(\mathbf{x}^{(s)})$ to build an approximate posterior distribution Q (Section 3). Figure 2 shows the effect of sample size k on the Chinese-English validation set. It is clear that BLEU scores consistently rise with the increase of k . However, we find that a sample size larger than 100 (e.g., $k = 200$) usually does not lead to significant improvements and increases the GPU memory requirement. Therefore, we set $k = 100$ in the following experiments.

4.4 Effect of Loss Function

As our approach is capable of incorporating evaluation metrics as loss functions, we investigate the effect of different loss functions on BLEU, TER

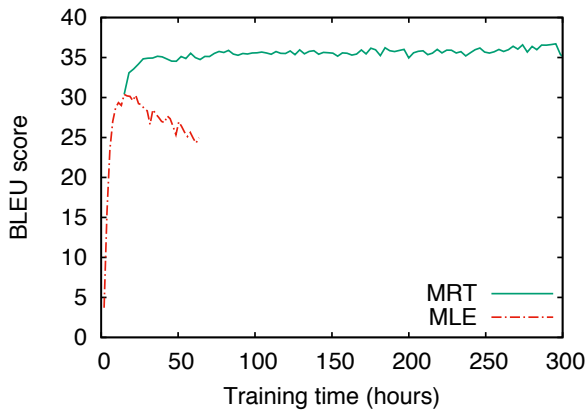


Figure 3: Comparison of training time on the Chinese-English validation set.

and NIST scores on the Chinese-English validation set. As shown in Table 2, negative smoothed sentence-level BLEU (i.e., $-sBLEU$) leads to statistically significant improvements over MLE ($p < 0.01$). Note that the loss functions are all defined at the sentence level while evaluation metrics are calculated at the corpus level. This discrepancy might explain why optimizing with respect to $sTER$ does not result in the lowest TER on the validation set. As $-sBLEU$ consistently improves all evaluation metrics, we use it as the default loss function in our experiments.

4.5 Comparison of Training Time

We used a cluster with 20 Tesla K40 GPUs to train the NMT model. For MLE, it takes the cluster about one hour to train 20,000 mini-batches, each of which contains 80 sentences. The training time for MRT is longer than MLE: 13,000 mini-batches can be processed in one hour on the same cluster.

Figure 3 shows the learning curves of MLE and MRT on the validation set. For MLE, the BLEU score reaches its peak after about 20 hours and then keeps going down dramatically. Initializing with the best MLE model, MRT increases BLEU scores dramatically within about 30 hours.³ Afterwards, the BLEU score keeps improving gradually but there are slight oscillations.

4.6 Results on Chinese-English Translation

4.6.1 Comparison of BLEU Scores

Table 3 shows BLEU scores on Chinese-English datasets. For RNNSEARCH, we follow Luong

³Although it is possible to initialize with a randomized model, it takes much longer time to converge.

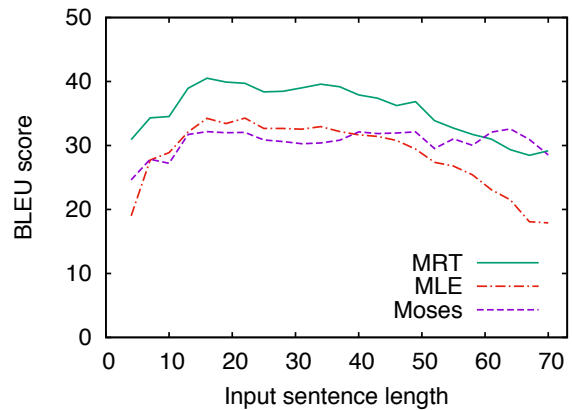


Figure 4: BLEU scores on the Chinese-English test set over various input sentence lengths.

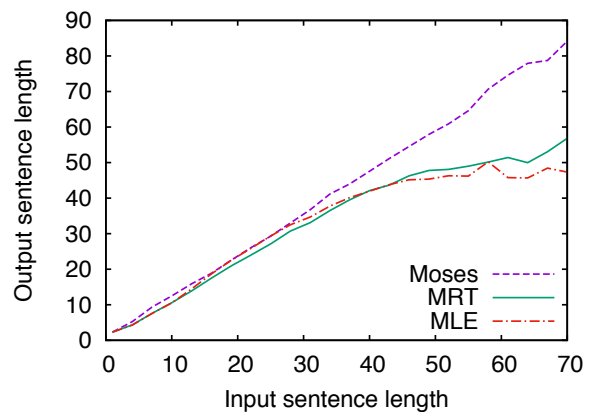


Figure 5: Comparison of output sentences lengths on the Chinese-English test set.

et al. (2015b) to handle rare words. We find that introducing minimum risk training into neural MT leads to surprisingly substantial improvements over MOSES and RNNSEARCH with MLE as the training criterion (up to +8.61 and +7.20 BLEU points, respectively) across all test sets. All the improvements are statistically significant.

4.6.2 Comparison of TER Scores

Table 4 gives TER scores on Chinese-English datasets. The loss function used in MRT is $-sBLEU$. MRT still obtains dramatic improvements over MOSES and RNNSEARCH with MLE as the training criterion (up to -10.27 and -8.32 TER points, respectively) across all test sets. All the improvements are statistically significant.

4.6.3 BLEU Scores over Sentence Lengths

Figure 4 shows the BLEU scores of translations generated by MOSES, RNNSEARCH with MLE,

System	Training	MT06	MT02	MT03	MT04	MT05	MT08
MOSES	MERT	32.74	32.49	32.40	33.38	30.20	25.28
RNNSEARCH	MLE	30.70	35.13	33.73	34.58	31.76	23.57
	MRT	37.34	40.36	40.93	41.37	38.81	29.23

Table 3: Case-insensitive BLEU scores on Chinese-English translation.

System	Training	MT06	MT02	MT03	MT04	MT05	MT08
MOSES	MERT	59.22	62.97	62.44	61.20	63.44	62.36
RNNSEARCH	MLE	60.74	58.94	60.10	58.91	61.74	64.52
	MRT	52.86	52.87	52.17	51.49	53.42	57.21

Table 4: Case-insensitive TER scores on Chinese-English translation.

	MLE vs. MRT		
	<	=	>
evaluator 1	54%	24%	22%
evaluator 2	53%	22%	25%

Table 5: Subjective evaluation of MLE and MRT on Chinese-English translation.

and RNNSEARCH with MRT on the Chinese-English test set with respect to input sentence lengths. While MRT consistently improves over MLE for all lengths, it achieves worse translation performance for sentences longer than 60 words.

One reason is that RNNSEARCH tends to produce short translations for long sentences. As shown in Figure 5, both MLE and MRE generate much shorter translations than MOSES. This results from the length limit imposed by RNNSEARCH for efficiency reasons: a sentence in the training set is no longer than 50 words. This limit deteriorates translation performance because the sentences in the test set are usually longer than 50 words.

4.6.4 Subjective Evaluation

We also conducted a subjective evaluation to validate the benefit of replacing MLE with MRT. Two human evaluators were asked to compare MLE and MRT translations of 100 source sentences randomly sampled from the test sets without knowing from which system a candidate translation was generated.

Table 5 shows the results of subjective evaluation. The two human evaluators made close judgments: around 54% of MLE translations are worse than MRE, 23% are equal, and 23% are better.

4.6.5 Example Translations

Table 6 shows some example translations. We find that MOSES translates a Chinese string “*yi wei fuze yu pingrang dangju da jiaodao de qian guowuyuan guanyuan*” that requires long-distance reordering in a wrong way, which is a notorious challenge for statistical machine translation. In contrast, RNNSEARCH-MLE seems to overcome this problem in this example thanks to the capability of gated RNNs to capture long-distance dependencies. However, as MLE uses a loss function defined only at the word level, its translation lacks sentence-level consistency: “chinese” occurs twice while “two senate” is missing. By optimizing model parameters directly with respect to sentence-level BLEU, RNNSEARCH-MRT seems to be able to generate translations more consistently at the sentence level.

4.7 Results on English-French Translation

Table 7 shows the results on English-French translation. We list existing end-to-end NMT systems that are comparable to our system. All these systems use the same subset of the WMT 2014 training corpus and adopt MLE as the training criterion. They differ in network architectures and vocabulary sizes. Our RNNSEARCH-MLE system achieves a BLEU score comparable to that of Jean et al. (2015). RNNSEARCH-MRT achieves the highest BLEU score in this setting even with a vocabulary size smaller than Luong et al. (2015b) and Sutskever et al. (2014). Note that our approach does not assume specific architectures and can in principle be applied to any NMT systems.

4.8 Results on English-German Translation

Table 8 shows the results on English-German translation. Our approach still significantly out-

Source	<i>meiguo daibiao tuan baokuo laizi shidanfu daxue de yi wei zhongguo zhuanjia , liang ming canyuan waijiao zhengce zhuli yiji yi wei fuze yu pingrang dangju da jiaodao de qian guowuyuan guanyuan .</i>
Reference	the us delegation consists of a chinese expert from the stanford university , two senate foreign affairs policy assistants and a former state department official who was in charge of dealing with pyongyang authority .
MOSES	the united states to members of the delegation include representatives from the stanford university , a chinese expert , two assistant senate foreign policy and a responsible for dealing with pyongyang before the officials of the state council .
RNNSEARCH-MLE	the us delegation comprises a chinese expert from stanford university , a chinese foreign office assistant policy assistant and a former official who is responsible for dealing with the pyongyang authorities .
RNNSEARCH-MRT	the us delegation included a chinese expert from the stanford university , two senate foreign policy assistants , and a former state department official who had dealings with the pyongyang authorities .

Table 6: Example Chinese-English translations. “Source” is a romanized Chinese sentence, “Reference” is a gold-standard translation. “MOSES” and “RNNSEARCH-MLE” are baseline SMT and NMT systems. “RNNSEARCH-MRT” is our system.

System	Architecture	Training	Vocab	BLEU
<i>Existing end-to-end NMT systems</i>				
Bahdanau et al. (2015)	gated RNN with search	MLE	30K	28.45
Jean et al. (2015)	gated RNN with search		30K	29.97
Jean et al. (2015)	gated RNN with search + PosUnk		30K	33.08
Luong et al. (2015b)	LSTM with 4 layers		40K	29.50
Luong et al. (2015b)	LSTM with 4 layers + PosUnk		40K	31.80
Luong et al. (2015b)	LSTM with 6 layers		40K	30.40
Luong et al. (2015b)	LSTM with 6 layers + PosUnk		40K	32.70
Sutskever et al. (2014)	LSTM with 4 layers		80K	30.59
<i>Our end-to-end NMT systems</i>				
<i>this work</i>	gated RNN with search	MLE	30K	29.88
	gated RNN with search	MRT	30K	31.30
	gated RNN with search + PosUnk	MRT	30K	34.23

Table 7: Comparison with previous work on English-French translation. The BLEU scores are case-sensitive. “PosUnk” denotes Luong et al. (2015b)’s technique of handling rare words.

System	Architecture	Training	BLEU
<i>Existing end-to-end NMT systems</i>			
Jean et al. (2015)	gated RNN with search	MLE	16.46
Jean et al. (2015)	gated RNN with search + PosUnk		18.97
Jean et al. (2015)	gated RNN with search + LV + PosUnk		19.40
Luong et al. (2015a)	LSTM with 4 layers + dropout + local att. + PosUnk		20.90
<i>Our end-to-end NMT systems</i>			
<i>this work</i>	gated RNN with search	MLE	16.45
	gated RNN with search	MRT	18.02
	gated RNN with search + PosUnk	MRT	20.45

Table 8: Comparison with previous work on English-German translation. The BLEU scores are case-sensitive.

performs MLE and achieves comparable results with state-of-the-art systems even though Luong et al. (2015a) used a much deeper neural network. We believe that our work can be applied to their architecture easily.

Despite these significant improvements, the margins on English-German and English-French datasets are much smaller than Chinese-English. We conjecture that there are two possible reasons. First, the Chinese-English datasets contain four reference translations for each sentence while both English-French and English-German datasets only have single references. Second, Chinese and English are more distantly related than English, French and German and thus benefit more from MRT that incorporates evaluation metrics into optimization to capture structural divergence.

5 Related Work

Our work originated from the minimum risk training algorithms in conventional statistical machine translation (Och, 2003; Smith and Eisner, 2006; He and Deng, 2012). Och (2003) describes a smoothed error count to allow calculating gradients, which directly inspires us to use a parameter α to adjust the smoothness of the objective function. As neural networks are non-linear, our approach has to minimize the expected loss on the sentence level rather than the loss of 1-best translations on the corpus level. Smith and Eisner (2006) introduce minimum risk annealing for training log-linear models that is capable of gradually annealing to focus on the 1-best hypothesis. He et al. (2012) apply minimum risk training to learning phrase translation probabilities. Gao et al. (2014) leverage MRT for learning continuous phrase representations for statistical machine translation. The difference is that they use MRT to optimize a sub-model of SMT while we are interested in directly optimizing end-to-end neural translation models.

The Mixed Incremental Cross-Entropy Reinforce (MIXER) algorithm (Ranzato et al., 2015) is in spirit closest to our work. Building on the REINFORCE algorithm proposed by Williams (1992), MIXER allows incremental learning and the use of hybrid loss function that combines both REINFORCE and cross-entropy. The major difference is that Ranzato et al. (2015) leverage reinforcement learning while our work resorts to minimum risk training. In addition, MIXER only sam-

ples one candidate to calculate reinforcement reward while MRT generates multiple samples to calculate the expected risk. Figure 2 indicates that multiple samples potentially increases MRT’s capability of discriminating between diverse candidates and thus benefit translation quality. Our experiments confirm their finding that taking evaluation metrics into account when optimizing model parameters does help to improve sentence-level text generation.

More recently, our approach has been successfully applied to summarization (Ayana et al., 2016). They optimize neural networks for headline generation with respect to ROUGE (Lin, 2004) and also achieve significant improvements, confirming the effectiveness and applicability of our approach.

6 Conclusion

In this paper, we have presented a framework for minimum risk training in end-to-end neural machine translation. The basic idea is to minimize the expected loss in terms of evaluation metrics on the training data. We sample the full search space to approximate the posterior distribution to improve efficiency. Experiments show that MRT leads to significant improvements over maximum likelihood estimation for neural machine translation, especially for distantly-related languages such as Chinese and English.

In the future, we plan to test our approach on more language pairs and more end-to-end neural MT systems. It is also interesting to extend minimum risk training to minimum risk annealing following Smith and Eisner (2006). As our approach is transparent to loss functions and architectures, we believe that it will also benefit more end-to-end neural architectures for other NLP tasks.

Acknowledgments

This work was done while Shiqi Shen and Yong Cheng were visiting Baidu. Maosong Sun and Hua Wu are supported by the 973 Program (2014CB340501 & 2014CB34505). Yang Liu is supported by the National Natural Science Foundation of China (No.61522204 and No.61432013) and the 863 Program (2015AA011808). This research is also supported by the Singapore National Research Foundation under its International Research Centre@Singapore Funding Initiative and administered by the IDM Programme.

References

- Ayana, Shiqi Shen, Zhiyuan Liu, and Maosong Sun. 2016. Neural headline generation with minimum risk training. arXiv:1604.01904.
- Dzmitry Bahdanau, KyungHyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proceedings of ICLR*.
- Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, and Robert L. Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*.
- David Chiang. 2005. A hierarchical phrase-based model for statistical machine translation. In *Proceedings of ACL*.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of EMNLP*.
- George Doddington. 2002. Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. In *Proceedings of HLT*.
- Jianfeng Gao, Xiaodong He, Wen tao Yih, and Li Deng. 2014. Learning continuous phrase representations for translation modeling. In *Proceedings of ACL*.
- Caglar Gulcehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Loïc Barrault, Hui-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2015. On using monolingual corpora in neural machine translation. arXiv:1503.03535.
- Xiaodong He and Li Deng. 2012. Maximum expected bleu training of phrase and lexicon translation models. In *Proceedings of ACL*.
- Sebastien Jean, Kyunghyun Cho, Roland Memisevic, and Yoshua Bengio. 2015. On using very large target vocabulary for neural machine translation. In *Proceedings of ACL*.
- Nal Kalchbrenner and Phil Blunsom. 2013. Recurrent continuous translation models. In *Proceedings of EMNLP*.
- Philipp Koehn and Hieu Hoang. 2007. Factored translation models. In *Proceedings of EMNLP*.
- Philipp Koehn, Franz J. Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proceedings of HLT-NAACL*.
- Alon Lavie and Michael Denkowski. 2009. The mereor metric for automatic evaluation of machine translation. *Machine Translation*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Proceedings of ACL*.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015a. Effective approaches to attention-based neural machine translation. In *Proceedings of EMNLP*.
- Minh-Thang Luong, Ilya Sutskever, Quoc V. Le, Oriol Vinyals, and Wojciech Zaremba. 2015b. Addressing the rare word problem in neural machine translation. In *Proceedings of ACL*.
- Franz J. Och. 2003. Minimum error rate training in statistical machine translation. In *Proceedings of ACL*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of ACL*.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. arXiv:1511.06732v1.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Improving neural machine translation models with monolingual data. arXiv:1511.06709.
- David A. Smith and Jason Eisner. 2006. Minimum risk annealing for training log-linear models. In *Proceedings of ACL*.
- Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of AMTA*.
- Andreas Stolcke. 2002. Srlm - an extensible language modeling toolkit. In *Proceedings of ICSLP*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Proceedings of NIPS*.
- Ronald J. Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*.