

Sequence Generation with Target Attention

Yingce Xia^{1,*}, Fei Tian², Tao Qin², Nenghai Yu¹, and Tie-Yan Liu²

¹ University of Science and Technology of China, Hefei, Anhui, P.R.China
yingce.xia@gmail.com, ynh@ustc.edu.cn

² Microsoft Research, Beijing, P.R.China
{fetia,taoqin,Tie-Yan.Liu}@microsoft.com

Abstract. Source-target attention mechanism (briefly, source attention) has become one of the key components in a wide range of sequence generation tasks, such as neural machine translation, image caption, and open-domain dialogue generation. In these tasks, the attention mechanism, typically in control of information flow from the encoder to the decoder, enables to generate every component in the target sequence relying on different source components. While source attention mechanism has attracted many research interests, few of them turn eyes to if the generation of target sequence can additionally benefit from attending back to itself, which however is intuitively motivated by the nature of attention. To investigate the question, in this paper, we propose a new *target-target* attention mechanism (briefly, target attention). Along the progress of generating target sequence, target attention mechanism takes into account the relationship between the component to generate and its preceding context within the target sequence, such that it can better keep the coherent consistency and improve the readability of the generated sequence. Furthermore, it complements the information from source attention so as to further enhance semantic adequacy. After designing an effective approach to incorporate target attention in encoder-decoder framework, we conduct extensive experiments on both neural machine translation and image caption. Experimental results clearly demonstrate the effectiveness of our design of integrating both source and target attention for sequence generation tasks.

Keywords: Sequence generation, target-target attention model, neural machine translation, image captioning

1 Introduction

Recurrent Neural Network (RNN) based sequence generation, which aims to decode a target sequence y given source input x , has been widely adopted in real-world applications such as machine translation [1, 3], image caption [28, 30], and document summarization [19]. Although some RNN variants which include multiplicative gating mechanisms, such as LSTM [7] and GRU [1], can help

* This work was done when Yingce Xia was an intern at Microsoft Research Asia.

smooth the flow of historical information, it is not guaranteed to be sufficient, especially when faced with long sequences. To overcome this difficulty, the attention mechanism [1, 30, 15] is introduced to RNN based neural models. Inspired by human cognitive process, attention mechanism assumes that the generation of each target component (i.e., a word) can rely on different contexts, either in a “soft” form that depends on a weighted combination of all contextual components [1, 30], or in a “hard” way that assumes only one contextual component can affect the target one [15]. Accordingly, attention mechanism can enable to discover semantic dependencies between the source and the target sequences in an end-to-end way.

Typical attention mechanisms model the dependency of target sequence on the source input, which implies that the context of the attention only comes from the source side input. Taking the “soft” attention mechanism as an example, to generate the j -th component y_j in y , an attentive weight associated with each source side component x_i is employed to describe how important each x_i is for y_j . The attentive weights in fact follow the multinomial distribution, derived from the matching degree between decoder’s hidden state s_{j-1} and every hidden state of the source encoder, represented as $\{h_1, \dots, h_{T_x}\}$ with T_x as the sequence length of x . In other words, such attention is a source-target mechanism, in which attentive information used in decoding only depends on the source-side encoder.

While the source-target attention (briefly, *source* attention) enables to pass important semantics from source input to target sequence, it overlooks the important coherent information hidden in the target sequence itself. Intuitively, beyond the selected source side parts, the generation of each target side component can be affected by certain preceding components in target side as well, especially when the target sequence is comparatively long. Apparently, such attentive dependency cannot be captured by source attention alone.

Therefore, to build a more comprehensive mechanism for sequence generation, in this paper, we propose the *target-target* attention (briefly, *target* attention), as a powerful complement to the source attention. In the proposed approach, the generation of each target side component depends not only on certain components in source side, but also on its prefix (i.e, the preceding components) in the target sequence. Acting in this way, more accurate probabilistic modelling for target sequence is achieved based on the better characterization for dependency within target side. Furthermore, we observe that in the decoding phase, even the semantics contained in source side could be enhanced due to the stimulation brought by attending to target-side prefix. As a result, compared with the source attention, our new approach can generate the target sequence with a couple of advantages:

1. The coherent consistency gets improved;
2. Eliminated repeated and redundant textual fragments [26];
3. Adequate semantics reflecting source-side information.

Examples in Table 2 clearly demonstrate all these improvements.

We conduct extensive experiments on both neural machine translation (English-to-French, English-to-Germany and Chinese-to-English translations) and image

caption. The results show that incorporating the target attention mechanism effectively improves the quality of generated sequences.

The rest of the paper is organized as follows: the mathematical details of target attention are introduced in Section 2. We report and analyze the experiments on neural machine translation in Section 3, and image caption in Section 4. Background related works are summarized in Section 5. The paper is concluded in Section 6 together with perspectives on future works.

2 Target Attention Framework

In this section, we introduce the proposed target attention mechanism for sequence generation. The overall framework, together with several mathematical notations used in this section, is illustrated in Fig. 1. As a preliminary, we incorporate the target-target attention mechanism into the RNN based sequence-to-sequence network with source-target attention, which is briefly introduced in the next subsection.

2.1 Source Attention based Sequence Generation

The sequence-to-sequence networks³ typically include two components: the encoder network which encodes the source-side sequence and the decoder network which decodes and generates the target-side sequence. The attention mechanism acts as a bridge effectively transmitting information between the encoder and decoder.

Concretely speaking, the encoder network reads the source input x and processes it into a source-side memory buffer $M_{\text{src}} = \{h_1, h_2, \dots, h_{T_x}\}$ with size T_x . For each $i \in [T_x]$,⁴ the vector h_i acts as a representation for a particular part of x . Here are several examples.

- In neural machine translation [1] and neural dialogue generation [21], the encoder networks are RNNs with LSTM/GRU units, which sequentially process each word in x and generate a sequence of hidden states. The source-side memory M_{src} is composed of RNN hidden states at each time-step.
- In image captioning [30, 28], the encoder network is a convolution neural network (CNN) working on an image. In this task, M_{src} contains low level local feature map vectors extracted by the CNN, representing different parts of the input image.

The decoder network is typically implemented using LSTM/GRU RNN together with a softmax layer. Specifically the decoder consumes every component (i.e, word) $y_j, j \in [T_y]$ in the target sequence y , meanwhile selectively reads from

³ In some scenarios such as image caption, the source-side input is not in a typical sequential form. For the ease of statement but with a little inaccuracy, we still use “sequence-to-sequence” as a general name even for these scenarios.

⁴ For ease of reference, $[T_x]$ denotes the set $\{1, 2, \dots, T_x\}$.

the source-side memory M_{src} to form attentive contextual vectors c_j^e (Eqn.(1)), and finally generates each RNN hidden state s_j for any $j \in [T_y]$ (Eqn.(2)). All these signals are then fed into the softmax-layer to generate the next-step component, i.e., y_j (Eqn.(3)):

$$c_j^e = q(s_{j-1}, M_{\text{src}}), \quad (1)$$

$$s_j = g(s_{j-1}, y_{j-1}, c_j^e), \quad (2)$$

$$P(y_j | y_{<j}, x) \propto \exp(y_j; s_j, c_j^e). \quad (3)$$

The source attention plays an important role in the generation of c_j^e , i.e., the function $q(\cdot, \cdot)$. Intuitively, the attention mechanism grants different weights on source-side memory vectors $h_i \in M_{\text{src}}$ in generating each c_j^e :

$$\alpha_{ij} = \frac{\exp(A_e(s_{j-1}, h_i))}{\sum_{k=1}^{T_x} \exp(A_e(s_{j-1}, h_k))}, \quad (4)$$

$$c_j^e = \sum_{i=1}^{T_x} \alpha_{ij} h_i. \quad (5)$$

In Eqn.(4), $A_e(\cdot, \cdot)$ acts as the key component in source attention, typically implemented as a feed-forward neural network.

2.2 Target Attention based Sequence Generation

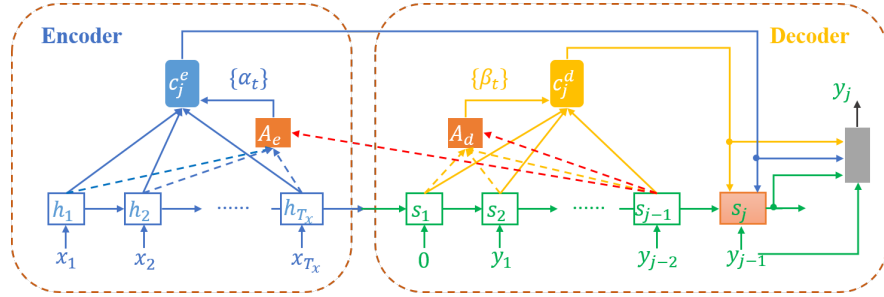


Fig. 1. The structure of sequence-to-sequence learning with target attention.

From Eqn.(1) and Eqn.(4), it is not difficult to observe that the attention weights are associated with the source-side memory M_{src} . Apart from that, as we have argued before, in sequence generation, better control over target-side contexts, i.e., the components that have been generated so far in the target sequence, is important as well. To add such target attention, we augment the memory space read by decoder RNN by adding an extra target-side memory M_{tgt} to original source-side one M_{src} . The j -th step slice of such a target-side

memory is defined as $M_{\text{tgt}}^j = \{s_1, \dots, s_{j-1}\}$, i.e., the hidden states before time step⁵ j .

Afterwards, to decode the word at j -th timestep, M_{tgt}^j is read and weighted averaged to form an extra contextual representation c_j^d (Eqn.(7)), where the weights are computed using a new attentive function $A_d(\cdot, \cdot)$ (Eqn.(6)). Intuitively speaking, β_{tj} represents how important the t -th (already) generated word is for current decoding at time step j . Such attentive signals rising from target-side memory M_{tgt}^j are integrated into c_j^d :

$$\beta_{tj} = \frac{\exp(A_d(s_{j-1}, s_t))}{\sum_{k=1}^{j-1} \exp(A_d(s_{j-1}, s_k))}, \quad (6)$$

$$c_j^d = \sum_{t=1}^{j-1} \beta_{tj} s_t. \quad (7)$$

Finally, c_j^d is provided as an addition input to derive the hidden state s_j (Eqn.(8)) and softmax distribution $P(y_j)$ (Eqn.(9)), from which the j -th component is chosen. The target-side memory also consequently gets updated as $M_{\text{tgt}}^{j+1} = M_{\text{tgt}}^j \cup \{s_j\}$.

$$s_j = g(s_{j-1}, y_{j-1}, c_j^e, \mathbf{c}_j^d) \quad (8)$$

$$P(y_j | y_{<j}, x) \propto \exp(y_j; s_j, c_j^e, \mathbf{c}_j^d). \quad (9)$$

To make it more clear, we further give the mathematical details of Eqn.(8) and (9) in Eqn.(10) and (11) respectively. (We take GRU as an example here and the mathematical formulation of LSTM can be similarly defined.)

$$\begin{aligned} s_j &= (1 - z_j) \circ s_{j-1} + z_j \circ \tilde{s}_j; \\ \tilde{s}_j &= \tanh(W y_{j-1} + U[r_j \circ s_{j-1}] + C^e c_j^e + \mathbf{C}^d \mathbf{c}_j^d); \\ z_j &= \sigma(W_z y_{j-1} + U_z s_{j-1} + C_z^e c_j^e + \mathbf{C}_z^d \mathbf{c}_j^d); \\ r_j &= \sigma(W_r y_{j-1} + U_r s_{j-1} + C_r^e c_j^e + \mathbf{C}_r^d \mathbf{c}_j^d). \end{aligned} \quad (10)$$

$$P(y_j | y_{<j}, x) \propto \exp(W_y y_{j-1} + W s_j + W^e c_j^e + \mathbf{W}^d \mathbf{c}_j^d). \quad (11)$$

In Eqn.(8) ~ (11), bold symbols are what make our target attention enhanced model different from conventional sequence-to-sequence model. The W 's, U 's and C 's are the parameters of the network unit, $\sigma(\cdot)$ denotes the sigmoid activation function, and \circ indicates element-wise product.

3 Application to Neural Machine Translation

We evaluate our proposed algorithm on three translation tasks: English→French (En→Fr), English→Germany (En→De) and Chinese→English (Zh→En).

⁵ For $j < 2$, let $M_{\text{tgt}}^1 = \emptyset$. Target attention starts to work when $j \geq 2$.

3.1 Settings

For En→Fr and En→De, we conduct experiments on the same datasets as those used in [8]. To be more specific, part of data in WMT’14 is used as the bilingual training data, which consists of 12M and 4.5M sentence pairs for En→Fr and En→De respectively. We remove the sentences with more than 50 words. *newstest2012* and *newstest2013* are concatenated as the validation set and *newstest2014* acts as test set. For En→Fr, we limit the source and target vocabularies as the most frequent 30k words, while for En→De, such vocabulary size is set as 50k. The out-of-vocabulary words will be replaced with a special token “UNK”⁶. For Zh→En, we use 1.25M bilingual sentence pairs from LDC dataset as training set and NIST2003 as validation set. Furthermore, NIST2004, NIST2005 and NIST2006 all act as the test sets. Both source and target vocabulary sizes are set as 30k for Zh→En.

The basic NMT model is the sequence-to-sequence GRU model widely adopted in previous works [1, 8]. In such a model, the word embedding size and GRU hidden layer size are respectively set as 620 and 1000. For Zh→En, dropout with a ratio 0.5 is applied to the last layer before softmax.

We use *beam search* algorithm [10] with beam width 12 to generate translations. The most common metric to evaluate translation quality is the BLEU score [16], which is defined as follows:

$$\text{BLEU} = \text{BP} \cdot \exp\left(\sum_{n=1}^4 \frac{1}{4} \log p_n\right), \quad \text{BP} = \begin{cases} 1, & c > r \\ e^{1-r/c}, & c \leq r \end{cases}, \quad (12)$$

where c is the total number of candidate sentence pairs in the corpus, r is the sum of the lengths for maximum perfect aligned translation subsequence in every translation pair, and p_n is a measure for n -gram translation precision.

Following the common practice in NMT, for En→Fr and En→De, the translation quality is evaluated by tokenized case-sensitive BLEU score⁷. For Zh→En, the BLEU scores are case-insensitive. Furthermore, to demonstrate better language modelling brought by target attention, we apply the perplexity [13] of target sentences conditioned on source sentences as another evaluation metric, which measures not only the translative correspondence of (source, target) pair, but also the naturalness of the target sentences. The perplexity is defined as follows:

$$\text{ppl}(D) = \exp\left\{-\frac{\sum_{i=1}^M \log P(y(i)|x(i))}{\sum_{i=1}^M N_i}\right\}, \quad (13)$$

where $D = \{(x(i), y(i))\}_{i=1}^M$ is the test set containing M bilingual language sentence pairs and N_i is the number of words in target sentence $y(i)$.

⁶ We focus on the word-level translations, instead of subword-level ones like BPE [20].

The reason is that BPE cannot be applied to some languages like Chinese, although it works well for other languages like German and Czech.

⁷ The script is from <https://github.com/moses-smt/mosesdecoder/blob/master/scripts/generic/multi-bleu.perl>

To reduce training time and stabilize the training process, following the common practice [23, 25], for all the three translation tasks, we initialize the target attention accompanied NMT models with the basic NMT models without target attention, i.e., the RNNSearch models [1]. (Note that training from scratch would lead to similar results as those obtained by training from warm-start models.) The three RNNSearch models used for warm start are all trained by Adadelta [33] with mini-batch size as 80, and these initialized models are able to reproduce the public reported BLEU scores in previous works.

After model initialization, we adopt vanilla SGD with minibatch size 80 to continue model training. According to the validation performance, the initial learning rates are set as 0.4, 0.4 and 0.1 for En→Fr, En→De and Zh→En respectively. During the training process, we halve the learning rates once the BLEU on validation set drops. We clip the gradients [17] with clipping threshold 1.0, 5.0 and 1.0 for En→Fr, En→De and Zh→En respectively. When halving learning rates cannot improve validation BLEU scores, we freeze the word embedding and continue model training for additional several days [8], leading to total training time of roughly one week for all the three translation tasks.

3.2 Quantitative Results

The experimental results are shown in Table 1. In this table, “RNNsearch” refers to the warm-start model, i.e., the sequence-to-sequence neural machine translation model with only source attention; “Target Attn” refers to our target attention enhanced NMT model.

Table 1. BLEU scores and perplexities of different Neural Machine Translation Models

	BLEU					Perplexities				
	En→Fr	En→De	MT04	MT05	MT06	En→Fr	En→De	MT04	MT05	MT06
RNNSearch	29.93	16.47	34.96	34.57	32.74	4.71	7.36	13.05	11.85	15.03
Target Attn	31.63	17.67	36.71	35.62	33.78	4.19	6.87	12.34	11.24	14.27

One can see that by introducing target attention into conventional NMT model, we can achieve significant improvements on all the three translation tasks. For En→Fr, the gain of BLEU brought by target attention is 1.7; for En→De, we improve BLEU by 1.2; furthermore, the average improvement for Zh→En is 1.28 BLEU point. By applying the statistical significance test for machine translation [11], we get that the results of our target attention mechanism is significantly better than RNNSearch with p -values $p < 0.01$. These results well demonstrate that our target attention algorithm is an effective approach that consistently and significantly improves the performances of different NMT tasks.

After obtaining the translation results, we further process them by the widely used post-processing technique proposed by [8], that is able to replace the “UNK” token in translation sentences. The steps include:

1. Get a word-level translation table \mathcal{T} that maps the source language words to target language words; we use the *fastAlign*⁸ [5];
2. Given any translated sentence \hat{y} , for each of its word \hat{y}_j , if \hat{y}_j is UNK, find the corresponding source side word x_i according to the attention weights learnt in NMT model, i.e., $i = \operatorname{argmax}_k \alpha_{kj}$ (refer to the definition in Eqn.(4));
3. Look up the table \mathcal{T} to get the corresponding translation for source word x_i .

By applying this technique to En→Fr translation, we can further improve the BLEU to **34.49**, which is a new best result for En→Fr translation conditioned on (i) the model is a single-layer NMT model; (ii) the model is trained with only bilingual data.

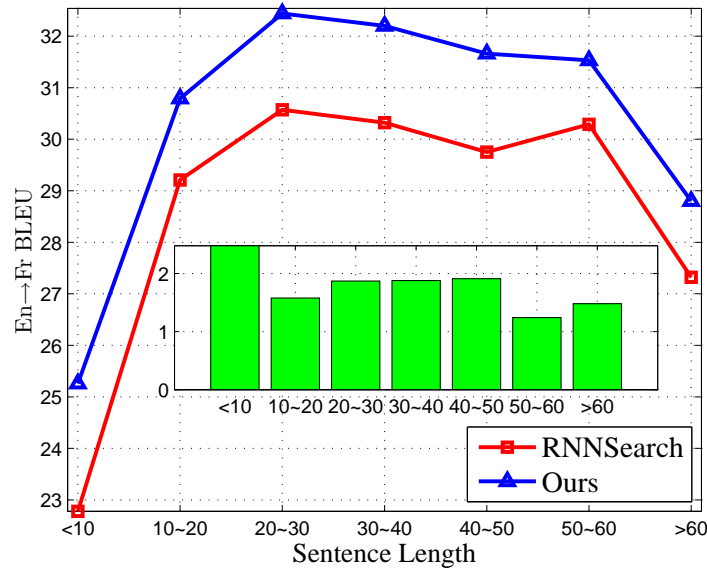


Fig. 2. BLEU w.r.t. input sentence length

We further record the average BLEU scores in En→Fr task for different bilingual sentence pair buckets that are determined by source sentence length and visualize them in Fig. 2. The chart inside Fig. 2 shows the improvements of BLEU score by adding target attention. From this figure we can see that with target attention, the BLEU scores for all buckets clearly increase. Specially,

1. For sentences with $[10, 50)$ words, the longer the sentence is, the more improvement is brought by the target attention. This clearly verifies the effectiveness of target attention in capturing long-range dependency.

⁸ The script is from https://github.com/clab/fast_align

2. For sentences with fewer than 10 words, our target attention also brings significant improvements. Note that although the calculation of BLEU is quite sensitive for very short sentences (e.g., p_4 in Eqn.(12) is very likely to be zero for sentence-level BLEU), we can still generate better sentences with target attention.
3. Even we remove the sentences with more than 50 words from the training corpus, target attention can achieve improvements on such un-trained subsets containing extremely long sentences (i.e., > 50 words), although improvements brought by target attention are slightly less compared to those in the region $[10, 50)$.

In the right part of Table 1 we additionally report the performances of different models using another evaluation measure, i.e., the perplexity, to represent how smooth and natural the translation is⁹. It can be observed that by incorporating target attention, the perplexity scores are decreased by 0.52, 0.49 and 0.69 points for En→Fr, En→De and Zh→En respectively. Such improvements clearly demonstrate the advantages of target attention in generating more smooth translations, mainly out from better dependency modeling within target translation sentence.

3.3 Subjective Evaluation

To better understand and visualize the benefits that target attention brings to NMT, we give three concrete Zh→En examples in Table 2. For each example, we highlight some important parts by bold words to demonstrate either the limitations of baseline model, or the improvements led by incorporating target attention. As discussed in Section 1, the three examples respectively show that target attention:

1. Improves long range semantic consistency such as matching the subjects perfectly and avoiding such pattern as “*The founding was the company*”;
2. Eliminates repeated translations such as “*economy could slow down*”¹⁰;
3. Enhances semantic integrity such as successfully translates “*bimen*” and “*xijieweijian gongbu*” in the last example.

Our intuitive explanation for these improvements is that target attention enhances the ability to model a natural sentence, due to which the inconsistent and repeated translations would be assigned with low probabilities and thus not selected. The punishment towards such wrong translation patterns will comparatively improve the possibility of right translation that has not been translated yet, thereby alleviating the semantic inadequacy issue.

⁹ For Zh→En, each source sentence x has four references $y(j)$ $j \in \{1, 2, 3, 4\}$. To calculate the perplexities, we simply regard them as four individual sentence pairs $(x, y(j))$.

¹⁰ Although the coverage model in [26] can eliminate repeated translations, it is actually based on source-target attention but not target-target attention. Therefore, [26] can be further combined with our proposed target-target attention. We leave it as a future work.

Table 2. Translation examples. Source, Base, Ours, Ref denote the source sentence, the RNNSearch, target attention and the reference sentence respectively.

Source	<i>e tianranqi gongye gufengongsi chengli yu 1993nian2yue , shi shijieshang zuida de tianranqi kaicai gongsi .</i>
Base	<i>the founding of the russian gas industry in February 1993 was the world 's largest natural gas mining company .</i>
Ours	<i>founded in February 1993 , the russian gas industrial corporation is the world 's largest producer of natural gas mining .</i>
Ref	<i>Established in February 1993 , Gazprom is the largest natural gas exploitation company in the world .</i>
Source	<i>youyu riyuan shengzhi he pinfu chaju rijian kuoda keneng pohuai jinnian shangbannian xiangyou de nazhong hexie qifen, riben jingji keneng fanghuan sudu , mairu 2005 nian</i>
Base	<i>japan 's economy could slow down as the japanese economy could slow down as the yen appreciated and the disparity between the rich and the poor as a result of a growing gap between the rich and the poor .</i>
Ours	<i>japan 's economy may slow down in 2005 as the yen 's appreciation and the growing gap between the rich and the poor may damage the harmonious atmosphere in the first half of the year .</i>
Ref	<i>Japan 's economy may slow down towards 2005 as yen appreciation and a widening gap between the rich and poor could break the harmonious atmosphere it enjoyed in the first half of this year .</i>
Source	<i>Xizang liuwang jingshen lingxiu dalai lama de teshi zhengzai beijing he zhongguo guanyuan jinxing bimen huiyi, xijie weijian gongbu .</i>
Base	<i>the special envoy of the dalai lama , tibet 's exiled spiritual leader , was scheduled to meet with chinese officials and officials from the tibetan spiritual leader in beijing and chinese officials .</i>
Ours	<i>the special envoy of tibet 's exiled spiritual leader , the dalai lama , is holding a closed-door meeting with chinese officials , and the details were not disclosed.</i>
Ref	<i>The special envoy of the Dalai Lama , exiled Tibetan spiritual leader , is currently in Beijing carrying out closed meetings with Chinese officials, but the details have not been released .</i>

3.4 Discussion

In this subsection, we carry out some discussions about our proposed target-target attention framework:

1. For decoding speed, our approach indeed takes 17% more time for decoding than previous approach without target attention, considering we have an additional target-target attention model. Such a cost is acceptable considering the BLEU score improvements (i.e., 1.7pts for En→Fr, 1.2pts on En→De and 1.28pts for Zh→En.)
2. A degenerated version of our target-target attention is to use delayed memory. Mathematically, when predicting the t 'th word, not only the hidden state s_{t-1} at step $t-1$, but also the one $s_{t-\tau}$ at step $t-\tau$ are used where τ is a fixed number. Note such the degenerated model might not work well.

Take $\tau = 2$ as an example. For the third example in Table 2, “Details were not disclosed” should be attended to “closed-door”, which is far from $t - 2$. Delayed memory cannot handle this case since it does not know how many steps to delay. Our model works well due to its adaptive and dynamic attention mechanism.

3. The improvements of target-target attention are not caused by better optimization properties of the new RNN, since its architecture is more complex and more difficult to optimize.

4 Application to Image Captioning

To further verify the effectiveness of incorporating target attention into sequence generation models, we then apply the proposed model to image caption, which targets at describing the content of an image with a natural sentence.

4.1 Settings

We choose a public benchmark dataset, Flickr30k [32], for image caption. In this dataset, every image is provided with 5 reference sentences. We follow the data splitting rule as that used in [9, 30] such that the dataset is separated into 28000 training samples, 1000 validation samples and 1000 test samples.

We follow the same model structure as that proposed in [30]. The Oxford VGGnet [24] pretrained on ImageNet is used as the image encoder, which could eventually output 196 features per image, with each feature as a 512 dimensional vector. The decoder is a 512×512 LSTM RNN. Dropout with drop ratio 0.5 is applied to the layer before softmax. The vocabulary is set as the most 10k frequent words in the dataset. Soft source attention is chosen due to better performance in our implementation. In implementation, we base our codes on the open-source project provided by [30]¹¹.

The captions for all images are generated using beam search with beam size 10. To comprehensively evaluate their quality, we adopt several different measures. Following [30], we report the BLEU-3, BLEU-4 without a brevity penalty, which respectively indicates tri-gram and four-gram matching degree with groundtruth caption. Besides, we report the CIDEr, another widely used metric for image caption [27, 28]. CIDEr also measures n-gram matching degree, in which each n-gram will be assigned a TF-IDF weighting [18]. For all these metrics, the larger, the better. In addition, we report test perplexities (the smaller, the better) to evaluate whether target attention improve sentence smoothness in image caption task. (Refer to Eqn.(13) for the definition of perplexity.)

At the beginning of training process, we warm start our model with a pre-trained model with only source attention, which is previously optimized with Adadelta for about one day. Then in the training process, we incorporate target attention mechanism into the initialized captioning model and continue to train it using plain stochastic gradient descent for another one day, with a fixed learning rate as 0.05 selected by validation performance.

¹¹ <https://github.com/kelvinxu/arctic-captions>

4.2 Results

We present our results in Table 3. The second row, labeled with “source-attn” represents the performance of baseline captioning model, which approximately matches the numbers reported in [30] (e.g, BLEU-3 as 28.8 and BLEU-4 as 19.1), while the third row labeled with “Ours” records the results after adding target attention into caption generation.

Table 3. Results of image caption with/without target attention. For perplexity, the lower, the better, while for other measures, the higher, the better.

	BLEU-3	BLEU-4	CIDEr	perplexity
source-attn	28.3	19.0	37.6	28.56
Ours	29.6	20.5	40.4	23.06

From Table 3, we can clearly see that target attention achieves significant improvements over baseline in terms of all evaluation measures. In particular, the decrease of perplexity further demonstrates the better probabilistic modelling for captions by adding target attention.

Table 4. Two examples showing different captioning result



source-attn	<i>A group of people are playing instruments on stage</i>
Ours	<i>A band is playing on stage in front of a crowd</i>
Ref	<i>(1) Two men , one sitting , one standing , are playing their guitars on stage while the audience is looking on .</i> <i>(2) band doing a concert for people</i>
source-attn	<i>A black dog is running through the grass</i>
Ours	<i>A black dog is running through the grass with a toy in its mouth</i>
Ref	<i>(1) A curly brown dog runs across the lawn carrying a toy in its mouth</i> <i>(2) The black dog is running on the grass with a toy in its mouth .</i>

We also list two examples of image caption in Table 4, including the image, its two referenced captions (marked by “Ref”) and the captioning results

with/without target attention (respectively marked by “source-attn” and “Ours”). It is clearly shown that target attention mechanism generates better captions for both images (see the highlighted underlined words), mainly owing to the enhanced stimulation from already generated words in target side, right brought by target attention.

To better demonstrate how target attention helps to improve the image caption quality, we take the right figure in Table 4 as an example and analyze the target attention weights in the decoder. To be concrete, in the internal green bar chart of Fig. 3, we show the weights of the previously generated words when generating the last word *mouth* in our caption (i.e, *A black dog is running through the grass with a toy in its mouth*). Here we remove the most common words like *a*, *is* to make the illustration more compact and clearer. Simultaneously, the outer bar chart in Fig. 3 shows the words co-occurrence statistics for the word *mouth*, which is calculated by

$$\text{co-occurrence}(\text{word}_i) = \frac{\text{The number of sentences containing word}_i \text{ and } \textit{mouth}}{\text{The number of sentences containing } \textit{mouth}}.$$

From the outer chart, it is clearly observed that *mouth* is highly correlated to relevant words such as *dog* and *black*. The target attention thereby stimulates the generation of word *mouth*, by learning from the significant weights assigned to these two previously decoded words. This clearly shows that target attention mechanism accurately characterizes the semantic relation among the decoded words, thus improves both the coherence and completeness of the caption by attending to the past.

5 Related Work

The attempt of applying attention mechanism with deep learning dates back to several works in computer vision, represented by [4], [15] and [30]. Particularly [15] and [30] leverage a “hard” attention mechanism, which in every step attend to only one part of the image. For NMT, [1] first incorporates the “soft” attention mechanism that automatically assign attentive weights to every source side words in decoding target sentence word. Since then many varieties have been designed to improve such a source-target attention mechanism for neural machine translation. For example, [12] proposes a local attention model, which is effective and computationally efficient. In this model, for each decoding step t , the attention weights are only assigned over a subset of source hidden states within a window $[p_t - D, p_t + D]$, where D is the window size, and p_t is the window center determined by a selective mechanism.

There are some other works that target improving the network structure for attention model in neural machine translation. For example, [14] propose an interactive attention model for NMT named as NMT_{IA} . NMT_{IA} can keep track of the interaction between the decoder and the representation of source sentence during translation by both reading and writing operations, which are helpful to

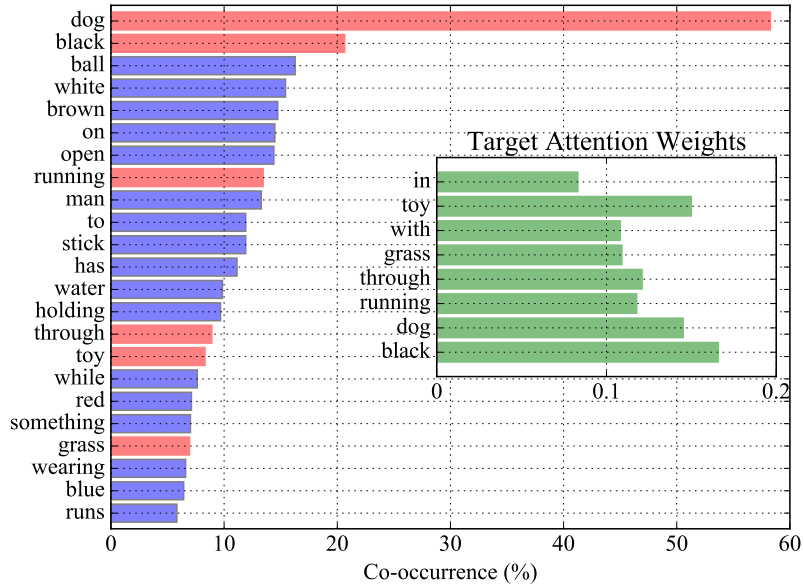


Fig. 3. Analysis of co-occurrence and target attention weights

improve translation quality. Similarly, [31] design a recurrent attention model, in which the attention weights are tracked by an LSTM.

As to decoder side, [22] proposes a self-attention model to maintain coherence in longer responses for neural conversation model: the decoded words would be concatenated with words in the source side encoder by their embeddings, and the self-attention model will generate contexts by these “faked” words. Such a proposal does not fit to other sequence-to-sequence tasks when the words in the encoder and decoder are not in the same language, like NMT and image caption. Therefore, such a model based on source-target words concatenation is limited and cannot be generalized to more general scenarios. [2] proposed a similar model like ours but [2] did not focus on sequence generation tasks.

6 Conclusion

In this work, motivated from tailored observations and analysis, we design a target attention model to enhance the dependency within decoder side components and thereby improve performances of sequence generation tasks. We conduct extensive evaluations and analysis on neural machine translation and image caption. Significant better results with the proposed model are observed on these two tasks, which illustrate the effectiveness of target attention.

There are many interesting directions left as future works. First, we aim to adapt target attention into different model structures [29] and different training techniques [23, 6]. Second, we plan to study how to make the target attention model more effective by combining it with the advanced source attention mechanisms discussed in related work. Last, target attention will be tested on more tasks such as document summarization, neural dialogue generation and question-answering.

References

1. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: International Conference on Learning Representations (2015)
2. Cheng, J., Dong, L., Lapata, M.: Long short-term memory-networks for machine reading. In: EMNLP (2016)
3. Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder–decoder for statistical machine translation. In: EMNLP. pp. 1724–1734 (2014)
4. Denil, M., Bazzani, L., Larochelle, H., de Freitas, N.: Learning where to attend with deep architectures for image tracking. *Neural computation* 24(8), 2151–2184 (2012)
5. Dyer, C., Chahuneau, V., Smith, N.A.: A simple, fast, and effective reparameterization of ibm model 2. *Association for Computational Linguistics* (2013)
6. He, D., Xia, Y., Qin, T., Wang, L., Yu, N., Liu, T., Ma, W.Y.: Dual learning for machine translation. In: *Advances In Neural Information Processing Systems*. pp. 820–828 (2016)
7. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Comput.* 9(8), 1735–1780 (Nov 1997), <http://dx.doi.org/10.1162/neco.1997.9.8.1735>
8. Jean, S., Cho, K., Memisevic, R., Bengio, Y.: On using very large target vocabulary for neural machine translation. In: the annual meeting of the Association for Computational Linguistics (2015)
9. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3128–3137 (2015)
10. Koehn, P.: Pharaoh: a beam search decoder for phrase-based statistical machine translation models. *Machine translation: From real users to research* pp. 115–124 (2004)
11. Koehn, P.: Statistical significance tests for machine translation evaluation. In: *EMNLP*. pp. 388–395. Citeseer (2004)
12. Luong, M.T., Pham, H., Manning, C.D.: Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025* (2015)
13. Luong, T., Cho, K., D. Manning, C.: Tutorial: Neural machine translation. *ACL 2016 tutorial* (2016)
14. Meng, F., Lu, Z., Li, H., Liu, Q.: Interactive attention for neural machine translation. *COLING* (2016)
15. Mnih, V., Heess, N., Graves, A., et al.: Recurrent models of visual attention. In: *Advances in neural information processing systems*. pp. 2204–2212 (2014)

16. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: the annual meeting of the Association for Computational Linguistics. pp. 311–318 (2002)
17. Pascanu, R., Mikolov, T., Bengio, Y.: On the difficulty of training recurrent neural networks. *International Conference on Machine Learning* 28, 1310–1318 (2013)
18. Robertson, S.: Understanding inverse document frequency: on theoretical arguments for idf. *Journal of documentation* 60(5), 503–520 (2004)
19. Rush, A.M., Chopra, S., Weston, J.: A neural attention model for abstractive sentence summarization. In: *EMNLP*. pp. 379–389 (2015)
20. Sennrich, R., Haddow, B., Birch, A.: Neural machine translation of rare words with subword units. the annual meeting of the Association for Computational Linguistics (2016)
21. Shang, L., Lu, Z., Li, H.: Neural responding machine for short-text conversation. arXiv preprint arXiv:1503.02364 (2015)
22. Shao, L., Gouws, S., Britz, D., Goldie, A., Strope, B., Kurzweil, R.: Generating long and diverse responses with neural conversation models. arXiv preprint arXiv:1701.03185 (2017)
23. Shen, S., Cheng, Y., He, Z., He, W., Wu, H., Sun, M., Liu, Y.: Minimum risk training for neural machine translation. the annual meeting of the Association for Computational Linguistics (2016)
24. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations* (2015)
25. Tu, Z., Liu, Y., Shang, L., Liu, X., Li, H.: Neural machine translation with reconstruction. In: *AAAI* (2017)
26. Tu, Z., Lu, Z., Liu, Y., Liu, X., Li, H.: Modeling coverage for neural machine translation. In: the annual meeting of the Association for Computational Linguistics. pp. 76–85 (2016)
27. Vedantam, R., Lawrence Zitnick, C., Parikh, D.: Cider: Consensus-based image description evaluation. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4566–4575 (2015)
28. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3156–3164 (2015)
29. Wu, Y., Schuster, M., Chen, Z., Le, Q.V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al.: Google’s neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint arXiv:1609.08144 (2016)
30. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A.C., Salakhutdinov, R., Zemel, R.S., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: *International Conference on Machine Learning*. vol. 14, pp. 77–81 (2015)
31. Yang, Z., Hu, Z., Deng, Y., Dyer, C., Smola, A.: Neural machine translation with recurrent attention modeling. arXiv preprint arXiv:1607.05108 (2016)
32. Young, P., Lai, A., Hodosh, M., Hockenmaier, J.: From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics* 2, 67–78 (2014)
33. Zeiler, M.D.: Adadelta: an adaptive learning rate method. arXiv preprint arXiv:1212.5701 (2012)