PAPER

# Document-level Neural Machine Translation with Associated Memory Network

**Shu JIANG**[†,††], **Rui WANG**[†,††], **Zuchao LI**[†,††], **Masao UTIYAMA**[†††], **Kehai CHEN**[†††], **Eiichiro SUMITA**[†††], **Hai ZHAO**[†,††a)], *and* **Bao-liang LU**[†,††b)],

**SUMMARY**    Standard neural machine translation (NMT) is on the assumption that the document-level context is independent. Most existing document-level NMT approaches are satisfied with a smattering sense of global document-level information, while this work focuses on exploiting detailed document-level context in terms of a memory network. The capacity of the memory network that detecting the most relevant part of the current sentence from memory renders a natural solution to model the rich document-level context. In this work, the proposed document-aware memory network is implemented to enhance the Transformer NMT baseline. Experiments on several tasks show that the proposed method significantly improves the NMT performance over strong Transformer baselines and other related studies.

*key words:*  *memory network, neural machine translation, document-level context*

## 1. Introduction

Neural Machine Translation (NMT) [1]–[5] established on the encoder-decoder framework, where the encoder takes a source sentence as input and encodes it into a fixed-length embedding vector, and the decoder generates the translation sentence according to the encoder embedding, has achieved advanced translation performance in recent years. So far, most models take a standard assumption to translate every sentence independently, ignoring the document-level contextual clues during translation.

However, document-level information can improve the translation performance from multiple aspects: consistency, disambiguation, and coherence [6]. If translating every sentence is independent of the document-level context, it will be challenging to keep every sentence translation across the entire text consistent with each other. Moreover, the document-level context can also assist the model to disambiguate words with multiple senses, and the global context is of great benefit to translation in a coherent way.

There have been few recent attempts to introduce the document-level information into the existing standard NMT models. Various existing methods [7]–[11] focus on modeling the context from the surrounding text in addition to the source sentence. For the more high-level context, Miculicich *et al.* [12] propose a multi-head hierarchical attention machine translation model to capture the word-level and sentence-level information. The cache-based model raised by Kuang *et al.* [6] uses the dynamic cache and topic cache to capture the inter-sentence connection. Tan *et al.* [13] propose a hierarchical model of global document context to improve document-level translation. In addition, many studies [9]–[11] all add the contextual information to the NMT model by applying the gating mechanism [14] to dynamically control the auxiliary global context information at each decoding step.

However, most of the existing document-level NMT methods focus on briefly introducing the global document-level information but fail to consider selecting the most related part inside the document context.

Inspired by the observation that human and document-level machine translation models always refer to the source sentence's context during the translation, like query in their memory, we propose to utilize the document-level sentences associated with the source sentences to help predict the target sentence. To reach such a goal, we adopt a Memory Network component [15]–[17] which provides a natural solution for the requirement of modeling document-level context in document-level NMT. In fact, Maruf and Haffari [18] have already presented a document-level NMT model which projects the document contexts into the tiny dense hidden state space for RNN model using memory networks and updates word by word, and their model is effective in exploiting both source and target document context.

Different from any previous work, this paper presents a Transformer NMT model with document-level Memory Network enhancement [15], [16] which concludes contextual clues into the encoder of the source sentence by the Memory Network. Not like Maruf and Haffari [18] which memorizes the whole document information into a tiny dense hidden state, the memory in our work calculates the associated document-level contextualized information in the memory with the current source sentence using the attention mechanism. In this way, our proposed model is able to focus on the most relevant part of the concerned translation from the memory, which precisely encodes the concerned document-level context.

The empirical results indicate that our proposed method significantly improves the BLEU score compared with a strong Transformer baseline and performs better than other

---

[†]The authors are with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, 200240, China.

[††]The authors are with the Key Laboratory of Shanghai Education Commission for Intelligent Interaction and Cognitive Engineering, Shanghai Jiao Tong University, Shanghai, 200240, China.

[†††]The authors are with National Institute of Information and Communications Technology, Kyoto-shi, 6190289 Japan.

   a) E-mail: zhaohai@cs.sjtu.edu.cn (Corresponding author.)
   b) E-mail: bllu@sjtu.edu.cn (Corresponding author.)

related models for document-level machine translation on multiple language pairs with multiple domains.

## 2. Background

### 2.1 Neural Machine Translation

Given a source sentence with $S$ tokens $\mathbf{x} = \{x_1, ..., x_i, ..., x_S\}$ in the document to be translated and a target sentence with $T$ tokens $\mathbf{y} = \{y_1, ..., y_i, ..., y_T\}$, NMT model computes the probability of translation from the source sentence to the target sentence word by word:

$$P(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{T} P(y_i|y_{1:i-1}, \mathbf{x}), \tag{1}$$

where $y_{1:i-1}$ is a substring containing words $y_1, ..., y_{i-1}$. Generally, with an RNN, the probability of generating the $i$-th word $y_i$ is modeled as:

$$P(y_i|y_{1:i-1}, \mathbf{x}) = \text{softmax}(g(y_{i-1}, \mathbf{s}_{i-1}, \mathbf{c}_i)), \tag{2}$$

where $g(\cdot)$ is a nonlinear function that outputs the probability of previously generated word $y_i$, and $\mathbf{c}_i$ is the $i$-th source representation. Then $i$-th decoding hidden state $\mathbf{s}_i$ is computed as

$$\mathbf{s}_i = f(\mathbf{s}_{i-1}, y_{i-1}, \mathbf{c}_i). \tag{3}$$

For NMT models with an encoder-decoder framework, the encoder maps an input sequence of symbol representations $\mathbf{x}$ to a sequence of continuous representations $\mathbf{z} = \{z_1, ..., z_i, ..., z_S\}$. Then, the decoder generates the corresponding target sequence of symbols $\mathbf{y}$ one element at a time.

### 2.2 Transformer Architecture

Only based on the attention mechanism, a network architecture called Transformer [5] for NMT uses stacked self-attention and point-wise, fully connected layers for both encoder and decoder.

A stack of $N$ (usually equals to 6) identical layers constitutes the encoder, and each layer has two sub-layers: (1) multi-head self-attention mechanism, and (2) a simple, position-wise fully connected feed-forward network.

Multi-head attention in the Transformer allows the model to process information jointly from different representation spaces at different positions. It linearly projects the queries $Q$, keys $K$, and values $V$ $h$ times with different, learned linear projections to $d_k$, $d_k$, and $d_v$ dimensions respectively, and then the attention function is performed in parallel, generating $d_v$-dimensional output values, and yielding the final results by concatenating and once again projecting them. The core of multi-head attention is Scaled Dot-Product Attention and calculated as:

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V. \tag{4}$$

The second sub-layer is a feed-forward network containing two linear transformations with a ReLU activation in between.

Similar to the encoder, the decoder is also composed of a stack of $N$ identical layers, but it inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. The Transformer also employs residual connections around each of the sub-layers, followed by layer normalization. Thus, the Transformer is more parallelizable and faster for translating than earlier RNN methods.

### 2.3 Memory Network

Memory networks [15] utilize the external memories as inference components based on long-range dependencies, which can be categorized into a sort of lazy machine learning [19]. Using the similar memorizing mechanism, memory-based learning methods have been also applied in multiple traditional models [20]–[25]. A memory network [15] is a set of vectors $\mathcal{M} = \{\mathbf{m}_1, ..., \mathbf{m}_K\}$ and the memory cell $\mathbf{m}_k$ is potentially relevant to a discrete object (for example, a word) $x_k$. The memory is equipped with a *read* and optionally a *write* operation. Given a query vector $\mathbf{q}$, the output vector produced by reading from the memory is $\sum_{i=1}^{K} p_i\mathbf{m}_i$, where $p_i = \text{softmax}(\mathbf{q}^T \cdot \mathcal{M})$ scores the match between the query vector $\mathbf{q}$ and the $i$-th memory cell $\mathbf{m}_i$.

## 3. Model

### 3.1 Framework

Our NMT model consists of two components: *Contextual Associated Memory Network* and a Transformer model. For the *Contextual Associated Memory Network*, the core part is a neural controller, which acts as a "processor" to read memory from the contextual storage "RAM" according to the input before sending it to other components. The controller calculates the correlation between the input and memory data, i.e., "memory addressing".

### 3.2 Encoders

Our model requires two encoders: the source encoder for translation from input sentence representation and the context encoder for the *Contextual Associated Memory Network* from context sentence representation. The source encoder is composed of a stack of $N$ layers, the same as the source encoder in the original Transformer [5]. The proposed *Contextual Associated Memory Network* consists of four parts: context selection, inter-sentence attention, embedding merging, and context gating.

### 3.3 Contextual Associated Memory Network

For each source sentence $\mathbf{x}$ at each training step, we assume the $m$ context sentences $\{\mathbf{c}_j\}_{j=1}^{m}$ related with the current sentence $\mathbf{x}$ as the *contextual memory* with the memory size
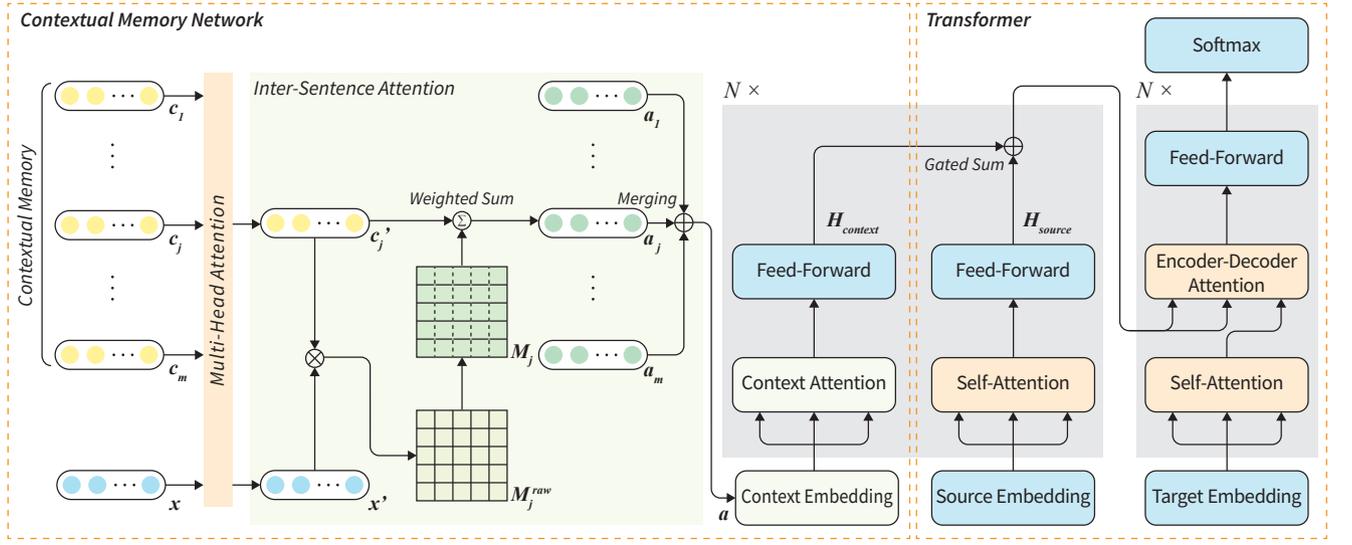
**Fig. 1**    The framework of our model.

$m$.

### 3.3.1   Context Selection

For the sake of fairness, we can treat all sentences in the document as our memory. However, it is impossible to attend all the sentences in the training dataset because of the extremely high computing and memorizing cost. We aim to utilize the context sentences and their representations to help our model predict the target sentences. There are three common ways to select the context sentences: *previous* sentences of the current sentence, *next* sentences of the current sentence, and *context* sentences randomly selected from the training corpus [10].

### 3.3.2   Inter-Sentence Attention

This part aims to attain the inter-sentence attention matrix, which can also be regarded as the core part of the *Contextual Associated Memory Network*. The input sentence $\mathbf{x}$ and the context sentences $\{\mathbf{c}_j\}_{j=1}^m$ in the *contextual memory* first go through a multi-head attention layer to encode the word representation:

$$\mathbf{x}' = \text{MultiHead}(\mathbf{x}, \mathbf{x}, \mathbf{x}), \tag{5}$$

and

$$\mathbf{c}'_j = \text{MultiHead}(\mathbf{c}_j, \mathbf{c}_j, \mathbf{c}_j) \;\; j \in \{1, 2, ..., m\}, \tag{6}$$

The lists of new word representations are denoted as follows:

$$\mathbf{x}' = \{x'_1, ..., x'_i, ..., x'_S\}, \tag{7}$$

and

$$\mathbf{c}'_j = \{c'_{1,j}, ..., c'_{k,j}, ..., c'_{K_j,j}\} \;\; j \in \{1, 2, ..., m\}. \tag{8}$$

Each word representation is as a vector $x \in \mathbb{R}^d$, where $d$ is the size of hidden state in MultiHead function.

Then, for each context sentence representation $\mathbf{c}'_j$, we apply the multi-head attention by treating the input sentence representation $\mathbf{x}'$ as the query sequence, on them and get the attention matrix $\mathcal{M}_j^{raw}$:

$$\mathcal{M}_j^{raw} = \mathbf{x}' \otimes \mathbf{c}'^T_j \;\; j \in \{1, 2, ..., m\}. \tag{9}$$

Every element $\mathcal{M}_{raw}(i, k) = x'_i \cdot c'^T_{k,j}$ can be regarded as an indicator of similarity between the $i$-th word in input sentence representation $\mathbf{x}'$ and the $k$-th word in memory sentence representation $\mathbf{c}'_j$.

Finally, we perform a softmax operation on every column in $\mathcal{M}_j^{raw}$ to normalize the value so that it can be considered as the probability from input sentence representation $\mathbf{x}'$ to memory sentence representation $\mathbf{c}'_j$:

$$\alpha_{i,j} = \text{softmax}([\mathcal{M}_j^{raw}(i, 1), ..., [\mathcal{M}_j^{raw}(i, K_j)]), \tag{10}$$

and

$$\mathbf{M}_j = [\alpha_{1,j}, ..., \alpha_{i,j}, ..., \alpha_{S,j}]. \tag{11}$$

We treat the probability vector $\alpha_{i,j}$ as a set of weights to sum all the representations in $\mathbf{c}'_j$ and get the memory-sentence-specified argument embedding $\mathbf{a}_j$:

$$\mathbf{a}_j = [a_{1,j}, ..., a_{i,j}, ..., a_{K_j,j}], \tag{12}$$

where

$$a_{i,j} = \sum_{k=1}^{K_j} \alpha_{i,j} c'_{k,j}. \tag{13}$$

### 3.3.3 Embedding Merging

To utilize the contextual embeddings $\mathbf{a}_j$ of the context sentences during training, embedding merging needs to be done.

Because the context sentences are different, the overall contributions of these word representations should be different. We let the model itself learn how to make use of these contextual word representations. Following the attention combination mechanism [17], [26], we consider four ways to merge the label information.

(1) Concatenation

All the contextual argument embedding are concatenated as the final attention embeddings.

$$\mathbf{a} = [\mathbf{a}_1, ..., \mathbf{a}_j, ..., \mathbf{a}_m]. \tag{14}$$

(2) Average

The average value of all the contextual argument embeddings is used as the final attention embedding.

$$\mathbf{a} = \frac{1}{m} \sum_{j=1}^{m} \mathbf{a}_j. \tag{15}$$

(3) Weighted Average

The weighted average of all the contextual argument embedding is used as the final attention embedding. We calculate the mean value of every raw similarity matrix $\mathcal{M}_j^{raw}$ to indicate the similarity between input sentence $\mathbf{x}$ and context sentence $\mathbf{c}_j$, and we use the softmax function to normalize them to get a probability vector $\beta$ indicating the similarity of input sentence $\mathbf{x}$ towards all the context sentences $\{\mathbf{c}_j\}_{j=1}^{m}$:

$$\begin{aligned} \beta &= \text{softmax}([g(\mathcal{M}_1^{raw}), ..., g(\mathcal{M}_m^{raw})]) \\ &= [\beta_1, ..., \beta_j, ...\beta_m], \end{aligned} \tag{16}$$

where $g(\cdot)$ represents the mean function.

Then, we use the probability vector $\beta$ as weight to sum all the contextual attention embedding $a_{i,j}$ for the final contextual attention embedding $\mathbf{a}$ of the $i$-th word $x_i$ in input sentence $\mathbf{x}$:

$$\mathbf{a} = \sum_{j=1}^{m} \beta_j \mathbf{a}_j. \tag{17}$$

(4) Flat

This method does not use $\mathbf{a}_j$. First, we concatenate all the raw similarity matrix $\mathcal{M}_j^{raw}$ along the row.

$$\mathcal{M}^{raw} = [\mathcal{M}_1^{raw}, ..., \mathcal{M}_j^{raw}, ..., \mathcal{M}_m^{raw}] \tag{18}$$

Then, we perform softmax operation on every row in $M^{raw}$ to normalize the value so that it can be considered as probability from input sentence $\mathbf{x}$ to all context sentences $\mathbf{c}_j$.

$$\gamma = f([\mathcal{M}_1^{raw}, ..., \mathcal{M}_k^{raw}, ..., \mathcal{M}_{K_{all}}^{raw}]), \tag{19}$$

where $f(\cdot)$ stands for softmax operation and $K_{all}$ is the total length of all context sentences, i.e.

$$K_{all} = \sum_{j=1}^{m} K_j, \tag{20}$$

and $K_j$ is the length of context sentences $\mathbf{c}_j$.

We also concatenate the contextual information $\mathbf{c} = [\mathbf{c}_1, ..., \mathbf{c}_j, ..., \mathbf{c}_m]$ and use $\gamma$ as weight to sum the concatenated contextual argument embedding as final contextual attention embedding.

$$\mathbf{a} = \gamma \cdot \mathbf{c}^T. \tag{21}$$

(5) Contextual RNN

We first pad and concatenate the contextual embeddings $\mathbf{a}_j$ of the context sentences by columns.

$$\mathbf{a}' = [\mathbf{a}_1; ...; \mathbf{a}_j; ...; \mathbf{a}_m]. \tag{22}$$

Inspired by the Document RNN method [9] to summarize the cross-sentence context information, we use RNN by column in the contextual argument embedding to generate the contextual attention embedding, and the hidden state at each time step can represent the relation from the first word embedding to the current word embedding.

As shown in the Figure 2, the RNN output of the $k$-th word embedding $a_{j,k}$ in the contextual argument embedding $a_j$ is

$$h_{j,k} = f(h_{j-1,k}, a_{j,k}) \tag{23}$$

where $f(\cdot)$ is an activation function, and $h_{j,k}$ is the hidden state at time $j$ of the $k$-th word embedding in the contextual argument embedding $a_j$.

Then we use the hidden state $h_{m,k}$ at the last time $m$, and the final contextual attention embedding $\mathbf{a}$ is concatenated by $h_{m,k}$.

$$\mathbf{a} = [h_{m,1}, ..., h_{m,k}, ..., h_{m,K_{max}}] \tag{24}$$

where $K_{max}$ is the max length of context sentences $\mathbf{c}_j$, i.e.

$$K_{max} = \max_{j \in \{1,2,...,m\}}\{K_j\}, \tag{25}$$

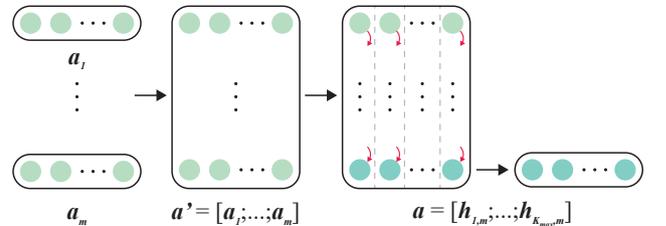and $K_j$ is the length of context sentences $\mathbf{c}_j$.



**Fig. 2** Contextual RNN.

### 3.3.4 Context Gating

Since the acquisition of contextual attention embedding **a**, we operate the MultiHead attention and feed-forward on the contextual attention embedding $a$ and source embedding **x** simultaneously like the original Transformer encoding steps, then we annotate the source attention embedding and contextual attention embedding after the above operations as $H_{source}$ and $H_{context}$. To control and analyze the flow of information from the extended context to the translation model, we use a context gate [14] to integrate the source and context attentions and control the flow from the source side and the context side. The gate $g$ is calculated by

$$g = \sigma(W_g[H_{source}, H_{context}] + b_g). \qquad (26)$$

Their gated sum $H$ is

$$H = g \otimes H_{source} + (1 - g) \otimes H_{context}, \qquad (27)$$

where $\sigma$ is the logistic sigmoid function, $\otimes$ is the point-wise multiplication and $W_g$ is trained by the model. As illustrated in Figure 1, the output of the gate $H$ is integrated into the encoder-decoder attention part at decoding step.

## 4. Experimental Setup

### 4.1 Data

The proposed document-level NMT model will be evaluated on multiple language pairs, i.e., Chinese-to-English (Zh-En), Spanish-to-English (Es-En), French-to-English (Fr-En), and English-to-French (Ja-En) on three domains: talks, subtitles, and news. Table 1 lists the statistics of all the concerned datasets.

(1)  TED Talks

The Zh-En TED talk documents are the parts of the IWSLT2015 Evaluation Campaign Machine Translation task[†]. We use *dev2010* as the development set and combine the *tst2010-2013* as the test set. The Es-En corpus is a subset of the IWSLT2014. We use *dev2010* for development set and *test2010-2012* as the test set. The Fr-En corpus is also a subset of the IWSLT2012, where *dev2010* is for development, and *test2010* is the test set. The Ja-En corpus is from IWSLT2017 and WMT, *dev2010* is for development and *test2010-2015* is for test.

(2)  Subtitles

The Es-En corpus is a subset of OpenSubtitles2018[††] [27][†††]. We randomly select 1,000 continuous sentences for each development set and test set.

(3)  News

The Es-En News-Commentaries11 corpus[††††] has document-level delimitation. We evaluate on the WMT sets [28]: *newstest2008* for development, and *newstest2009-2013* for testing.

### 4.2 Data Preprocessing

The English and Spanish datasets are tokenized by *tokenizer.perl* and truecased by *truecase.perl* provided by MOSES[†††††], a statistical machine translation system proposed by [29]. The Chinese corpus is tokenized by *Jieba* Chinese text segmentation[††††††]. Words in sentences are segmented into subwords by Byte-Pair Encoding (BPE) [30] with 32k BPE operations.

### 4.3 Model Configuration

We use the Transformer proposed by Vaswani *et al.* [5] as our baseline and implement our work using the THUMT, an open-source toolkit for NMT developed by the Natural Language Processing Group at Tsinghua University [31][*]. We follow the configuration of the Transformer "base model" described in the original paper [5]. Both encoder and decoder consist of 6 hidden layers each. All hidden states have 512 dimensions, eight heads for multi-head attention. The training batch contains about 6,520 source tokens, and we train the model about 200,000 training bathes. We use the original regularization and optimizer in Transformer [5]. Finally, we evaluate the performance of the model by BLEU score [32] using *multi-bleu.perl* on the *tokenized* text.

## 5. Results and Analysis

### 5.1 Translation Performance

We choose the previous $m = 3$ sentences as the contextual memory and using the Contextual RNN method to merge the embeddings. Table 2 demonstrates the BLEU scores for different models on multiple corpora. The baseline is a re-implemented attention-based NMT system RNNSearch* [33] and Transformer [5] using THUMT kit. We also employ the Context-aware model [10] on these datasets, when we set the contextual memory size $m = 1$ and without the inter-sentence attention. The results of RNN with Memory Network [18] and HAN model [12] are reported by the authors.

The results in Table 2 demonstrate that our proposed model significantly outperforms all the comparing models, especially, our model is significantly better than the baseline Transformer at significance level $p$-value<0.05. Our

---

[†]`https://wit3.fbk.eu`
[††]`http://www.opensubtitles.org/`
[†††]`http://opus.nlpl.eu/OpenSubtitles2018.php`

[††††]`https://opus.nlpl.eu/News-Commentary-v11.php`
[†††††]`https://github.com/moses-smt/mosesdecoder`
[††††††]`https://github.com/fxsjy/jieba`
[*]`https://github.com/thumt/THUMT`

**Table 1**  Data statistics of sentences.

| Dataset | TED Talks | | | | Subtitles | News |
| | Zh-En | Es-En | Fr-En | En-Ja | Es-En | Es-En |
|---|---|---|---|---|---|---|
| Training | 209,941 | 180,853 | 145,503 | 228,697 | 48,301,352 | 238,872 |
| Tuning | 887 | 887 | 934 | 871 | 1,000 | 2,000 |
| Test | 5,473 | 4,706 | 1,664 | 8,469 | 1,000 | 14,522 |

**Table 2**  BLEU scores on the different datasets. The scores in bold indicate the best ones on the same dataset. The last column indicates the time cost of different models on the News dataset.

| Model | TED Talks | | | | Subtitles | News | |
| | Zh-En | Es-En | Fr-En | En-Ja | Es-En | Es-En | Cost. (hours) |
|---|---|---|---|---|---|---|---|
| RNNSearch* | 16.09 | 36.55 | 30.79 | 10.41 | 39.90 | 22.95 | 23.60 |
| Transformer | 17.76 | 38.53 | 30.92 | 11.73 | 39.96 | 23.71 | 28.59 |
| RNN with Memory Network [18] | - | - | 22.00 | - | - | - | - |
| Context-aware Transformer [10] | 18.24 | 38.74 | 31.20 | 11.87 | 40.19 | 23.76 | 42.09 |
| Transformer with HAN [12] | 17.79 | 37.24 | - | - | 36.23 | 22.76 | - |
| Our model | **18.69** | **39.20** | **31.97** | **12.01** | **40.74** | **24.40** | 42.32 |

proposed model outperforms the RNNSearch* baseline by 2.60 BLEU point on the TED Talks (Zh-En) dataset, 2.65 BLEU point on the TED Talks (Es-En) dataset, 1.18 BLEU point on the TED Talks (Fr-En) dataset, 1.60 BLEU point on the TED Talks (En-Ja) dataset, 0.84 BLEU point on the OpenSubtitles (Es-En) dataset, and 1.45 BLEU point on the WMT dataset (Es-En).

Furthermore, our proposed model achieves the gains of 0.93, 0.67, 1.05, 0.28, 0.78, and 0.69 BLEU points on these four datasets individually over the Transformer baseline. Compared with the Context-aware Transformer proposed by [10], our proposed approach also raises the average 0.50 BLEU score on these different datasets. Moreover, the average increase of the BLEU score over the Transformer with HAN [12] is 2.25 points.

We note that the results on TED Fr-En are much higher than the result reported by Maruf and Haffai, and we deduce that it may be aroused by different prepossessing methods or BLEU styles.

The last column in Table 1 indicates the time cost of different models on the News dataset (Es-En) under the same model setting mentioned in Section 4.3. We can figure out that with the complexity of the model, the performance improves at the cost of running speed.

## 5.2  Translation Analysis

To reflect the improvements of our proposed model more exactly, we will analyze the overall performance from three aspects mentioned above: consistency, disambiguation, and coherence. The translation of HAN model [12] for comparison is downloaded from Miculicich's GitHub[†].

## 5.2.1  Consistency

In our work, the contextual memory is able to store the contextual sentences and help the model refine the translation.

Thus, we follow the previous work [6], and calculate the average number of words in generated translations which are also in the contextual sentences fed into the contextual memory. During our calculating process, punctuations, stop words, and *UNK* are removed from the contextual sentences and translations. Table 3 shows the results of consistency on TED datasets with the memory size $m = 3$. As shown in Table 3, HAN and our memory method can improve translation consistency compared to the baseline, confirming the claim that document translation can improve consistency between sentences. Our method is clearly closer to the reference than HAN and the baseline, demonstrating that our memory method is a more powerful approach for enhancing translation consistency.

**Table 3**  Consistency test on TED Zh↔En test sets.

| Model | TED Zh→En | | TED Zh→En | |
| | *pre-3* | *next-3* | *pre-3* | *next-3* |
|---|---|---|---|---|
| Reference | 1.22 | 1.23 | 1.21 | 1.21 |
| Transformer model | 1.04 | 1.05 | 1.02 | 1.04 |
| HAN model | 1.04 | 1.06 | 1.03 | 1.04 |
| **Our model** | **1.12** | **1.15** | **1.11** | **1.13** |

## 5.2.2  Disambiguation

We also want to investigate the ability of the word disambiguation of our model. We download the English-to-Chinese dictionary from free dictionary project[††], and select the words with multiple translation words in the source language to build a new dict $dict = \{word^{src} : trans_1^{tgt}, \ldots, trans_n^{tgt}\}$. When the token $w$ in the source sentence and $w \in \{word^{src}\}$, we count the appearance of the translation words $\{trans^{tgt}\}$ of $w$ in the corresponding translation sentence. We argue that if a model is weak at disambiguation, to translate an ambiguous word with multiple word senses, the model would prefer one of the

senses with the highest probability. The other corresponding candidate words' appearance will decrease accordingly. Thus, we use the Standard Deviation to evaluate the disambiguation ability.

**Table 4** Disambiguation ability test on TED Zh→En and Es→En test sets.

| Model | TED Zh→En | TED Es→En |
|---|---|---|
| Reference | 652.15 | 541.16 |
| Transformer model | 2691.81 | 2143.84 |
| HAN model | 2059.89 | 1457.06 |
| **Our model (*pre-3*)** | **1380.87** | **1060.50** |

Table 4 illustrates the results of the disambiguation ability of different models. First, comparing the Transformer baseline and reference, it can be seen that the lower the standard deviation, the better the disambiguation will be. Second, compared with the baseline transformer, HAN decrease this metric on the two datasets. At the same time, our model achieved the lowest deviation value, indicating that the introduction of document information can alleviate the translation variety, i.e., have a disambiguation effect.

### 5.2.3 Coherence

To further study how our proposed context-aware neural model improves the coherence in document translation, we follow the work of Lapata and Barzilay [34] to measure coherence as sentence similarity. We represent each sentence as the mean of the distributed vectors of its words. Then, the similarity between the two sentences is determined by the cosine of their means. For a fair comparison, we use the pre-trained language model BERT [35] to get the distributed vectors of words.

**Table 5** Coherence test on TED Es→En, Es→En, and Subtitles Zh→En test sets.

| Model | TED Talks | | Subtitles |
|---|---|---|---|
| | Zh→En | Es→En | Es→En |
| Reference | 0.67 | 0.67 | 0.60 |
| Transformer model | 0.62 | 0.61 | 0.54 |
| HAN model | 0.65 | 0.65 | 0.60 |
| **Our model** | **0.66** | **0.66** | **0.60** |

Table 5 summarizes the comparison results. The Transformer baseline without document information has the lowest coherence score, while our system outperforms the HAN model slightly. On the one side, it demonstrates that both HAN and our model can improve the translation coherence, which leverages document features; on the other hand, it shows that our approach has certain advantages over HAN.

### 5.2.4 Case Study

#### (1)   Example on Chinese-to-English

We extract the 4,123-th parallel lines from TED Talks (Zh-En) and the contextual memory consists of three previous sentences before the source sentence. The the final contextual attention embedding **a** is merged by Contextual RNN method. Table 7 shows an example from the TED Talks (Zh-En), on which the translation of our model is compared to other methods. This example shows that our proposed model can recognize the tense and even discourse relation from the document-level context and enhance the translation to more consistent and coherent.

#### (2)   Example on English-to-Japanese

We select the 7,160-th parallel lines from TED Talks (En-Ja) test set and list the three previous sentences in the contextual memory. Compared with the baseline, for the word *figure* with multiple word senses, our proposed model could recognize the correct word sense *person* instead *number*, and the attribute *tragic* is also translated correctly. We infer that the word *guy* in the context sentence $c_3$ provides the translation clue, and it verifies that the contextual memory network enhances the disambiguation ability of our model.

## 6.   Ablation Study

### 6.1   Effect of Recurrent Core of Contextual RNN

In our proposed model, we use Contextual RNN to integrate the contextual information. Its recurrent core can also be replaced by GRU and LSTM with different styles. Table 8 illustrates the results when we change the recurrent core. We can observe that the recurrent core alteration influences our model slightly, and forward RNN is most efficient from the results.

### 6.2   Effect of Embedding Merging

We choose the different embedding merging ways introduced in Section.3.3.3 to produce the final contextual attention embedding and compare the model performance on the different datasets with contextual memory size $m = 3$.

Table 9 demonstrates the BLEU scores of the different embedding merging methods, and the model performs best on these datasets by contextual RNN merging method.

### 6.3   Effect of Context Gating

We also investigate the impact of context gate $g$ by using the different given constants and compare the results with the context gate trained by the model. For instance, if the context gate $g$ equals 0, the model is the vanilla Transformer model, and context gate $g = 1$ means the model only encodes the final contextual attention embedding **a** from the context sentence(s) without the source attention. Fig. 3 illustrates the performance of the different context gate values when the contextual memory size $m = 3$. Of course, the context gate obtained from the model performs better than the fixed context gate, and meanwhile, both source information and context information are essential to the model.

**Table 6**   Example of the translation result. The context sentences are three previous sentences before the source sentence and we use the Contextual RNN method to merge contextual argument embedding. The words in blue from context indicate the heuristic clues for better translation and the sentences in Chinese have been provided with English translation.

| | |
|---|---|
| Context sentence $c_3$ | 它 去年 秋天 遭遇 破产 因为 他们 遭到 入侵 。 |
| | *(it was running into bankruptcy last fall because they were hacked into .)* |
| Context sentence $c_2$ | 有人 闯进去 彻底 毁 了 它 |
| | *(somebody broke in and they hacked it thoroughly .)* |
| Context sentence $c_1$ | 我 上周 在 与 荷兰政府 代表 开会 时 问过 ， 我 问 一位 领导 是否 他 发现 有 可能 有人 会 因为 Diginotar 攻击 而 死亡 。 |
| | *(and I asked last week in a meeting with Dutch government representatives, I asked one of the leaders of the team whether he found plausible that people died because of the DigiNotar hack .)* |
| Source sentence | 他的回答是肯定的。 |
| Reference sentence | *and his answer was yes .* |
| Transformer model | *his answer is yes .* |
| HAN model | *his answer and yes .* |
| **Our model** | *and his answer was yes .* |

**Table 7**   We select the 7,160-th parallel lines from TED Talks (En-Ja) test set and list the three previous sentences in the contextual memory. Compared with the baseline, for the word *figure* with multiple word senses, our proposed model could recognize the correct word sense *person* instead *number*, and the attribute *tragic* is also translated correctly. We infer that the word *guy* in the context sentence $c_3$ provides the translation clue, and it verifies that the contextual memory network enhances the disambiguation ability of our model.

| | |
|---|---|
| Context sentence $c_3$ | this is a guy called e.p . |
| Context sentence $c_2$ | the worst memory in the world . |
| Context sentence $c_1$ | his memory was so bad , that he didn &apos;t even remember he had a memory problem , which is amazing . |
| Source sentence | and he was this incredibly *tragic figure* , but he was a window into the extent to which our memories make us who we are . |
| Reference sentence | とても 悲劇 的 な 人物 です が どの 程度 記憶 が 我々 を 形作っ ている か を 知る 手がかり と なる 存在 です |
| Transformer model | 彼 は 本当 に 悲惨 な 数字 *(disastrous number)* でした が 私 たち の 記憶 が 私 たち を どう 解釈 する か に 窓 を つけ てい ました |
| **Our model** | 彼 は 非常 に 悲劇 的 な 人物 *(tragic person)* でした が 私 たち の 記憶 が 私 たち を どの ように する か に ついて の 窓 だった のです |

**Table 8**   The results on TED Talks with the different recurrent cores of Contextual RNN

| Core | TED-Zh-En | | | TED-Es-En | | |
|---|---|---|---|---|---|---|
| | forward | backward | bi-directional | forward | backward | bi-directional |
| RNN | 18.67 | 18.55 | 18.57 | 39.20 | 39.28 | 39.24 |
| LSTM | 18.58 | 18.67 | 18.56 | 39.21 | 39.27 | 39.21 |
| GRU | 18.54 | 18.63 | 18.61 | 39.19 | 39.23 | 39.20 |

**Table 9**   BLEU scores on the different datasets with various embedding merging ways.

| Embedding Merging | TED Talks | | Subtitles | News | |
|---|---|---|---|---|---|
| | Zh-En | Es-En | Es-En | Es-En | Cost. (hours) |
| Concatenation | 18.19 | 38.9 | 40.14 | 23.71 | 30.01 |
| Average | 18.23 | 38.95 | 40.34 | 23.82 | 37.37 |
| Weighted Average | 18.44 | 39.16 | 40.68 | 24.37 | 39.58 |
| Flat | 18.48 | 39.15 | 40.59 | 24.33 | 32.59 |
| Contextual RNN | **18.67** | **39.2** | **40.74** | **24.4** | 42.32 |

## 6.4   Effect of Contextual Information

### (1)   Different context sentence definition

The context sentences in our work are the previous three sentences of the current sentence. We investigate the effect of the different context sentence definition on the TED Talks (Zh-En) dataset. Following the work of Context-aware Transformer [10], we use the previous sentence(s), next sentence(s) and the random selected context sentence(s) form the document as the context sentence(s). As shown in Fig.4,
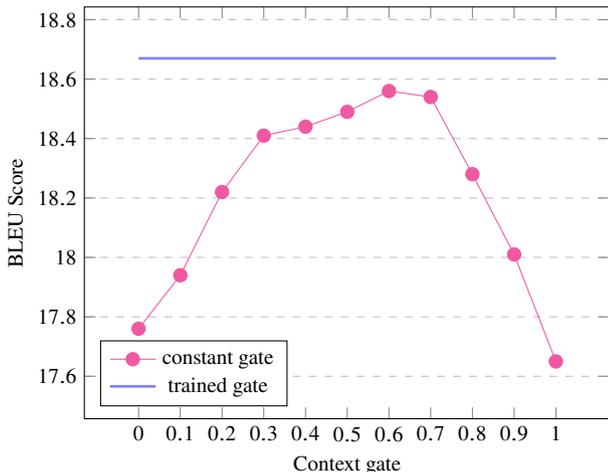
**Fig. 3** Results on TED Talks (Zh-En) dataset with different context gating ways.



**Fig. 5** The comparison of time consumption with different contextual memory size $m$ and different context selection.

the model which uses the previous sentence(s) as the context sentence(s) could get the best performance on the TED Talks (Zh-En) dataset, and it is in agreement with the work of the Context-aware Transformer [10].

(2)  Different contextual memory size

We also compare the effect with the different contextual memory size $m$ on the TED Talks (Zh-En) dataset. If the contextual memory size $m$ equals 0, the model is the original Transformer model. As shown in Fig.4, more contextual information appears beneficial to model translation, and the BLEU score gets better with more context sentences. However, it changes slightly when the contextual memory size $m$ greater than 4.
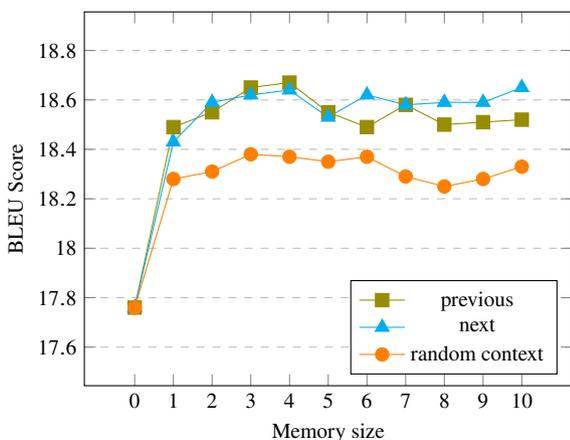


**Fig. 4** Results on TED Talks (Zh-En) dataset with different contextual memory size $m$ and different context selection.

### 6.5  Time Consumption

We also compare time consumption with different contextual memory sizes and context definitions mentioned in the
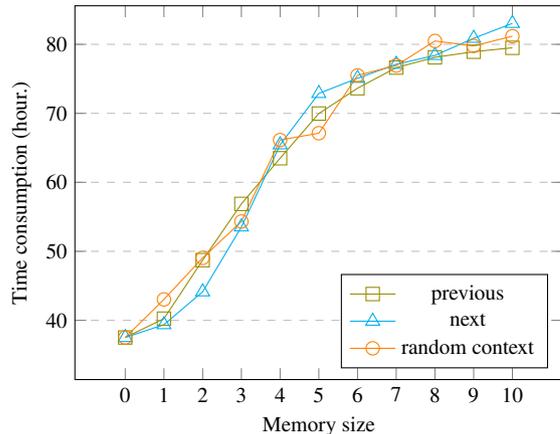
above section, and we illustrate Figure 5 according to the statistics. The original model needs 17.76 hours, and obviously, the proposed method needs more hours for 200 thousand steps on the TITAN RTX GPU device. Our proposed model needs a bit more time because the contextual associated memory network is more complex and has to train more hyper-parameters during the training process.

## 7.  Related Work

The existing work about NMT on the document-level can be divided into two parts: one is how to obtain the document-level information in NMT, and the other is how to integrate the document-level information.

### 7.1  Mining Document-level Information

Tiedemann *et al.* [8] merely concatenate the context in two ways: (1) extending the source sentence, which includes the context from the previous sentences in the source language, and (2) extending translation units, which increase the segments to translate.

Michel *et al.* [36] propose a simple yet parameter-efficient adaption method that only requires adapting the *Specific Vocabulary Bias* of output softmax to each particular use of the NMT system and allows the model to reflect distinct linguistic variations through translation better.

Mac *et al.* [37] present a *Word Embedding Average* method to add source context that captures the whole document with accurate boundaries, taking every word into account by an averaging method.

Kang *et al.* [38] propose to select dynamic context so that the document-level translation model can utilize the more useful selected context sentences to produce better translations via reinforcement learning.

## 7.2 Integrating Document-level Information

### (1) Gating Context

The context gate can automatically control the ratios of source and context representations contributions to the generation of target words [14]. Wang *et al.* [9] introduce this mechanism in their work to dynamically control the information flowing from the global text at each decoding step. Kuang *et al.* [11] propose an inter-sentence gate model, which is based on the attention-based NMT and uses the same encoder to encode two adjacent sentences and controls the amount of information flowing from the preceding sentence to the translation of the current sentence with an inter-sentence gate.

### (2) Document RNN

Wang *et al.* [9] propose a cross-sentence context-aware RNN approach to produce a global context representation called Document RNN. Given a source sentence in the document to be translated and its *m* previous sentences, they can obtain all sentence-level representations after processing each sentence. The last hidden state represents the summary of the whole sentence as it stores order-sensitive information. The last hidden state represents the summary of the global context over the sequence of the above sentence-level representations.

### (3) Cache-based Neural Model

Tu *et al.* [39] propose to augment the NMT models with an external cache to exploit translation history. At each decoding step, the probability distribution over generated words is updated online depending on the translation history retrieved from the cache with a query of the current attention vector, which helps NMT models adapt over time dynamically.

### (4) Context-Aware Transformer Model

Voita et al. [10] introduce the context information into the Transformer [5] and leave the Transformer's decoder intact while processing the context information on the encoder side. The model calculates the gate from the source sentence attention and the context sentence attention, exploiting their gated sum as the encoder output.

Zhang *et al.* [40] also extend the Transformer with a new context encoder to represent document-level context while incorporating it into both the original encoder and decoder by multi-head attention.

Miculicich *et al.* [12] propose a *Hierarchical Attention Networks (HAN) NMT* model to capture the context in a structured and dynamic pattern. Each predicted word uses word-level and sentence-level abstractions and selectively focuses on different words and sentences.

Tan *et al.* [13] propose a *hierarchical modeling of global document context model* to improve document-level translation, which is hierarchically extracted from the entire global text with a sentence encoder to model intra-sentence information and a document encoder to model document-level inter-sentence context representation.

Ma *et al.* [41] propose a *Flat-Transformer model* with a simple and effective unified encoder that model the bi-directional relationship between the contexts and the source sentences.

Chen *et al.* [42] propose to improve document-level NMT by the means of discourse structure information, and the encoder is based on a HAN [12]. They parse the document to obtain its discourse structure, then introduce a Transformer-based path encoder to embed the discourse structure information of each word and combine the discourse structure information with the word embedding.

Most of the previous works only focus on integrating context embedding or considering the context selection, but our work can mine the most related part among the contextual memory at each step.

## 8. Conclusion and Future Work

We propose a memory network enhancement over Transformer-based NMT, which provides a natural solution for modeling the detailed document-level context. Experiments show that our model performs better on the datasets of multiple domains and language pairs and can capture salient document-level contextual clues, select the most relevant part related to the input sequence from the contextual memory, and effectively enhance strong NMT baselines.

We will consider better context selection in our future work, like using discourse information to enhance our model. On the one hand, the discourse information will provide the heuristic, but on the other hand, it will bring much noise, and the internal structure may be incredibly complicated. Therefore, it is necessary to abstract its critical feature information effectively.

**References**

[1] N. Kalchbrenner and P. Blunsom, "Recurrent continuous translation models," Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, Seattle, Washington, USA, pp.1700–1709, Association for Computational Linguistics, 2013.

[2] I. Sutskever, O. Vinyals, and Q.V. Le, "Sequence to sequence learning with neural networks," in Advances in Neural Information Processing Systems 27, ed. Z. Ghahramani, M. Welling, C. Cortes, N.D. Lawrence, and K.Q. Weinberger, pp.3104–3112, Curran Associates, Inc., 2014.

[3] K. Cho, B. van Merrienboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder–decoder for statistical machine translation," Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar, pp.1724–1734, Association for Computational Linguistics, 2014.

[4] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," Proceedings of the 3rd International Conference on Learning Representations, San Diego, USA, pp.1–15, 2015.

[5] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A.N. Gomez, L.u. Kaiser, and I. Polosukhin, "Attention is all you need," in Advances in Neural Information Processing Systems 30, ed. I. Guyon, U.V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,

and R. Garnett, pp.5998–6008, Curran Associates, Inc., 2017.

[6] S. Kuang, D. Xiong, W. Luo, and G. Zhou, "Modeling coherence for neural machine translation with dynamic and topic caches," Proceedings of the 27th International Conference on Computational Linguistics, Santa Fe, New Mexico, USA, pp.596–606, Association for Computational Linguistics, 2018.

[7] S. Jean, S. Lauly, O. Firat, and K. Cho, "Does Neural Machine Translation Benefit from Larger Context?," arXiv e-prints, p.arXiv:1704.05135, April 2017.

[8] J. Tiedemann and Y. Scherrer, "Neural machine translation with extended context," Proceedings of the Third Workshop on Discourse in Machine Translation, Copenhagen, Denmark, pp.82–92, Association for Computational Linguistics, 2017.

[9] L. Wang, Z. Tu, A. Way, and Q. Liu, "Exploiting cross-sentence context for neural machine translation," Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Copenhagen, Denmark, pp.2826–2831, Association for Computational Linguistics, 2017.

[10] E. Voita, P. Serdyukov, R. Sennrich, and I. Titov, "Context-aware neural machine translation learns anaphora resolution," Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Melbourne, Australia, pp.1264–1274, Association for Computational Linguistics, 2018.

[11] S. Kuang and D. Xiong, "Fusing recency into neural machine translation with an inter-sentence gate model," Proceedings of the 27th International Conference on Computational Linguistics, Santa Fe, New Mexico, USA, pp.607–617, Association for Computational Linguistics, 2018.

[12] L. Miculicich, D. Ram, N. Pappas, and J. Henderson, "Document-level neural machine translation with hierarchical attention networks," Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, pp.2947–2954, Association for Computational Linguistics, 2018.

[13] X. Tan, L. Zhang, D. Xiong, and G. Zhou, "Hierarchical modeling of global context for document-level neural machine translation," Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), Hong Kong, China, pp.1576–1585, Association for Computational Linguistics, Nov. 2019.

[14] Z. Tu, Y. Liu, Z. Lu, X. Liu, and H. Li, "Context gates for neural machine translation," Transactions of the Association for Computational Linguistics, vol.5, pp.87–99, 2017.

[15] J. Weston, S. Chopra, and A. Bordes, "Memory networks," arXiv preprint arXiv:1410.3916, 2015.

[16] S. Sukhbaatar, J. Weston, R. Fergus, et al., "End-to-end memory networks," Advances in neural information processing systems, pp.2440–2448, 2015.

[17] C. Guan, Y. Cheng, and H. Zhao, "Semantic role labeling with associated memory network," Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Minneapolis, Minnesota, pp.3361–3371, Association for Computational Linguistics, June 2019.

[18] S. Maruf and G. Haffari, "Document context neural machine translation with memory networks," Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Melbourne, Australia, pp.1275–1284, Association for Computational Linguistics, July 2018.

[19] D.W. Aha, Lazy learning, Springer Science & Business Media, 2013.

[20] W. Daelemans, "Introduction to the special issue on memory-based language processing," Journal of Experimental & Theoretical Artificial Intelligence, vol.11, no.3, pp.287–296, 1999.

[21] E. Fix and J.L. Hodges Jr, "Discriminatory analysis-nonparametric discrimination: consistency properties," tech. rep., California Univ Berkeley, 1951.

[22] R. Skousen, Analogical modeling of language, Springer Science &

[23] R. Skousen, Analogy and structure, Springer Science & Business Media, 2013.

[24] M. Lebowitz, "Memory-based parsing," Artificial Intelligence, vol.21, no.4, pp.363–404, 1983.

[25] J. Nivre, J. Hall, and J. Nilsson, "Memory-based dependency parsing," Proceedings of the Eighth Conference on Computational Natural Language Learning (CoNLL-2004) at HLT-NAACL 2004, 2004.

[26] J. Libovický and J. Helcl, "Attention strategies for multi-source sequence-to-sequence learning," Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Vancouver, Canada, pp.196–202, Association for Computational Linguistics, July 2017.

[27] P. Lison and J. Tiedemann, "Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles," Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016), ed. N.C.C. Chair), K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, and S. Piperidis, Portoro, Slovenia, European Language Resources Association (ELRA), may 2016.

[28] O. Bojar, C. Buck, C. Callison-Burch, C. Federmann, B. Haddow, P. Koehn, C. Monz, M. Post, R. Soricut, and L. Specia, "Findings of the 2013 Workshop on Statistical Machine Translation," Proceedings of the Eighth Workshop on Statistical Machine Translation, Sofia, Bulgaria, pp.1–44, Association for Computational Linguistics, August 2013.

[29] P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst, "Moses: Open source toolkit for statistical machine translation," Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions, Prague, Czech Republic, pp.177–180, Association for Computational Linguistics, 2007.

[30] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Berlin, Germany, pp.1715–1725, Association for Computational Linguistics, 2016.

[31] J. Zhang, Y. Ding, S. Shen, Y. Cheng, M. Sun, H. Luan, and Y. Liu, "THUMT: an open source toolkit for neural machine translation," CoRR, vol.abs/1706.06415, 2017.

[32] K. Papineni, S. Roukos, T. Ward, and W.J. Zhu, "Bleu: a method for automatic evaluation of machine translation," Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, 2002.

[33] G.E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R.R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," 2012.

[34] M. Lapata and R. Barzilay, "Automatic evaluation of text coherence: Models and representations," IJCAI, 2005.

[35] J. Devlin, M.W. Chang, K. Lee, and K. Toutanova, "BERT: Pretraining of deep bidirectional transformers for language understanding," Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Minneapolis, Minnesota, pp.4171–4186, Association for Computational Linguistics, June 2019.

[36] P. Michel and G. Neubig, "Extreme adaptation for personalized neural machine translation," Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Melbourne, Australia, pp.312–318, Association for Computational Linguistics, 2018.

[37] V. Macé and C. Servan, "Using whole document context in neural machine translation," ArXiv, vol.abs/1910.07481, 2019.

[38] X. Kang, Y. Zhao, J. Zhang, and C. Zong, "Dynamic context selection for document-level neural machine translation via reinforcement

Business Media, 1989.

learning," Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp.2242–2254, 2020.

[39] Z. Tu, Y. Liu, S. Shi, and T. Zhang, "Learning to remember translation history with a continuous cache," Transactions of the Association for Computational Linguistics, vol.6, pp.407–420, 2018.

[40] J. Zhang, H. Luan, M. Sun, F. Zhai, J. Xu, M. Zhang, and Y. Liu, "Improving the transformer translation model with document-level context," Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, pp.533–542, Association for Computational Linguistics, Oct.-Nov. 2018.

[41] S. Ma, D. Zhang, and M. Zhou, "A simple and effective unified encoder for document-level machine translation," Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Online, pp.3505–3511, Association for Computational Linguistics, July 2020.

[42] J. Chen, X. Li, J. Zhang, C. Zhou, J. Cui, B. Wang, and J. Su, "Modeling discourse structure for document-level neural machine translation," Proceedings of the First Workshop on Automatic Simultaneous Translation, Seattle, Washington, pp.30–36, Association for Computational Linguistics, July 2020.