

P-Transformer: Towards Better Document-to-Document Neural Machine Translation

Yachao Li, Junhui Li, Jing Jiang, Shimin Tao, Hao Yang, Min Zhang

Abstract—Directly training a document-to-document (Doc2Doc) neural machine translation (NMT) via Transformer from scratch, especially on small datasets usually fails to converge. Our dedicated probing tasks show that 1) both the absolute position and relative position information gets gradually weakened or even vanished once it reaches the upper encoder layers, and 2) the vanishing of absolute position information in encoder output causes the training failure of Doc2Doc NMT. To alleviate this problem, we propose a position-aware Transformer (P-Transformer) to enhance both the absolute and relative position information in both self-attention and cross-attention. Specifically, we integrate absolute positional information, i.e., position embeddings, into the query-key pairs both in self-attention and cross-attention through a simple yet effective addition operation. Moreover, we also integrate relative position encoding in self-attention. The proposed P-Transformer utilizes sinusoidal position encoding and does not require any task-specified position embedding, segment embedding, or attention mechanism. Through the above methods, we build a Doc2Doc NMT model with P-Transformer, which ingests the source document and completely generates the target document in a sequence-to-sequence (seq2seq) way. In addition, P-Transformer can be applied to seq2seq-based document-to-sentence (Doc2Sent) and sentence-to-sentence (Sent2Sent) translation. Extensive experimental results of Doc2Doc NMT show that P-Transformer significantly outperforms strong baselines on widely-used 9 document-level datasets in 7 language pairs, covering small-, middle-, and large-scales, and achieves a new state-of-the-art. Experimentation on discourse phenomena shows that our Doc2Doc NMT models improve the translation quality in both BLEU and discourse coherence. We make our code available on Github.¹

Index Terms—neural machine translation, document-level NMT, document-to-document translation, position information, sequence-to-sequence.

I. INTRODUCTION

Document-level neural machine translation, which aims to make the translation more coherent and fluent between sentences, has recently received increasing attention. However, due to the challenge of modeling/generating an input/output

document as a single sequence, most related studies still translate a document sentence by sentence and propose various context-aware document-to-sentence models² ([2–5], to name a few). In these Doc2Sent models, document-level context is usually encoded with an extra-encoder and is properly integrated into a sentence-level translation model. Therefore, such context-aware Doc2Sent models are strictly sentence-level and suffer from limitations caused by not encoding current sentences and their context in the same encoding module [1, 6]. Moreover, independently generating target sentences in a document hinders the usage of target-side document-level context. Different from these studies, in this paper our goal is to build a document-to-document translation model which encodes a document as a single unit and generates its translation as well as a single unit.

Although extending the translation unit from a single sentence to multiple sentences (e.g., 4) achieves good performance [7–9], directly training a seq2seq-based Doc2Doc model, i.e., Transformer [10] from scratch tends to fail [6, 11], especially on small training data.³ For example, Bao et al. [6] show that the failure is due to the local minima during training. As a result, the attention weights in the cross-attention between the encoder and the decoder are flat, with large entropy values. Based on the observation of attention weight distribution, Bao et al. [6] propose G-Transformer, the first Doc2Doc NMT model by introducing local bias to attention. Later, Sun et al. [1] observe that vanilla Transformer [10] with appropriate data augmentation techniques can achieve good performance for Doc2Doc translation, which alleviates the issues of limited training data. Although both the studies [1, 6] successfully adapt Transformer to properly cater to long-sequence input and output, they have not answered the behind reasons for the failure when directly training a seq2seq-based Doc2Doc model on vanilla Transformer. Moreover, strictly forcing sentence-to-sentence alignment between the input and the output prohibits G-Transformer [6] from being applied to other seq2seq tasks, like text summarization while MR Doc2Doc [1] may not be able to recover the sentence-to-sentence alignment between a source document and its translation.⁴

Yachao Li and Jing Jiang are with the Key Laboratory of China’s Ethnic Languages and Information Technology of Ministry of Education, Northwest Minzu University, Lanzhou 730030, China (e-mail: harry_lyc@foxmail.com, jiangj0723@163.com)

Junhui Li and Min Zhang are with the School of Computer Science and Technology, Soochow University, Suzhou 215006, China (e-mail: {lijunhui, minzhang}@suda.edu.cn)

Shimin Tao and Hao Yang are with the Huawei Translation Services Center, Beijing, China (e-mail: {taoshimin, yanghao30}@huawei.com)

Yachao Li and Jiang Jiang are supported by the National Natural Science Foundation of China (Grant No. 62266038), Junhui Li and Min Zhang are supported by the NSFC (Grant No. 6203600).

This manuscript is under review.

¹<https://github.com/liyc7711/doc2doc>

²In this paper, we follow the notations in Sun et al. [1], where Doc2Sent models take partial or entire source document as context and generate target-side sentences independently while Doc2Doc models view the source document and target document as long sequences.

³We note that a Doc2Doc Transformer can be successfully trained upon a pre-trained model (e.g., mBART) or on large-scale training data. However, this requires additional large-scale datasets.

⁴In our re-implementation, 6 out of 115 documents of the English-to-German TED dataset have the different number of sentences in the source and the translation.

In this paper, we take a step further to explore the reason that causes the training failure of Doc2Doc NMT with vanilla Transformer model [10]. Similar to Bao et al. [6], our Doc2Doc Transformer fails to converge when training on small training data, e.g., English-to-German TED dataset with 10.9K documents and 210K sentences. We conjecture that even with residual connections, the Transformer model could not preserve sufficient position information at the top layer once the input sequence becomes long. As a result, the cross-attention between the encoder and the decoder has little knowledge about the position information of the source-side hidden states and thus results in flat attention weights with large entropy values. To verify the conjecture, we take three dedicated probing tasks to access the position information encoded in the encoder layers (Section II). As shown in Figure 1, the experimental results reveal that the accuracy of absolute position prediction for the top layer (i.e., Layer 6) of the Doc2Doc model is extremely as low as 0.5% while the accuracy for the Sent2Sent model is 62.2%. Moreover, as shown in Figure 2 and Figure 3, there is a visible tendency that with the increase of input sequence, the relative position information degrades significantly for both Sent2Sent and Doc2Doc models. That is to say, the position information has almost vanished when encoding a document as a single unit.

Motivated by the results of position information probing tasks, it is essential to enhance the position information of the Transformer model when building Doc2Doc NMT. To this end, we propose an incredibly simple yet effective approach and construct P-Transformer, a Transformer-based model with position-aware attention, in which the query-key pairs are explicitly equipped with their corresponding (absolute) position information. We then apply position-aware attention to both self-attention and cross-attention modules in Transformer. Moreover, we also integrate relative position information into self-attention both on the source and target side. The position-aware attention, in turn, will enhance the position information embedded in the hidden states of the top encoder layer, which is helpful for Doc2Doc translation. Experimental results on 9 popular document-level translation datasets in 7 language pairs show that our proposed P-Transformer significantly outperforms the strong baselines and achieves a new state-of-the-art performance.

Overall, the contributions of this paper are three-fold.

- For both Sent2Sent and Doc2Doc translation, we provide a probing study to investigate whether the position information gets weakened or even vanished at the top Transformer encoder layer. Through the probing tasks, we find the main reason causing the training failure of Doc2Doc translation with vanilla Transformer.
- We propose P-Transformer with position-aware attention which could be successfully trained for Doc2Doc translation. Moreover, P-Transformer could also be used to boost the performance of Doc2Sent and Sent2Sent translation.
- Our extensive experimental results show that the proposed P-Transformer significantly outperforms the strong baselines on 9 document-level datasets in 7 language pairs covering small-, middle- and large-scales. In addition, the P-Transformer achieves state-of-the-art performance.

II. POSITION INFORMATION PROBING TASKS

Positional encoding plays a crucial role in Transformer to make use of the token order of a sequence. Specifically, the word embeddings (WEs) and position embeddings (PEs) are summed, i.e., $\mathbf{WE} + \mathbf{PE}$ as the word representation, which is fed to Transformer. Specifically, vanilla Transformer [10] uses sinusoidal functions to parameterize PEs in a fixed ad hoc way. Theoretically, the position information from the input can efficiently be propagated to the upper layers through residual connections. To the best of our knowledge, however, few relevant studies have looked into whether the hidden states, especially in upper layers, can preserve appropriate position information. Therefore, we propose three probing tasks from different views to properly measure how much position information is encoded in Transformer encoder layers. We implement these probing tasks upon the linguistic features probing toolkit of Conneau et al. [12] with default parameter settings.⁵ The training set, validation set, and test set are from the popular English-to-German TED document-level dataset (See Section IV-A for data description).⁶

A. Absolute Position Probing Task

We formalize the absolute position probing task as follows: given the hidden state $\mathbf{h}_i \in \mathbb{R}^D$ of the i -th source-side word, our goal is to predict its absolute position value, i.e., i through a neural network. Inspired by the linguistic feature probing tasks proposed in Conneau et al. [12], we treat it as K -label classification problem and train a classifier through a fully-connected network to map a hidden state \mathbf{h} into a K -class probability distribution via:

$$P_A(\mathbf{k}|\mathbf{h}) = \text{Softmax}(\mathbf{W}_A \mathbf{h}), \quad (1)$$

where $\mathbf{W}_A \in \mathbb{R}^{K \times D}$ is a trainable matrix, $\mathbf{k} \in \{1, 2, \dots, K\}$ is an absolute position, K is the maximum input length, and D is the hidden state size.

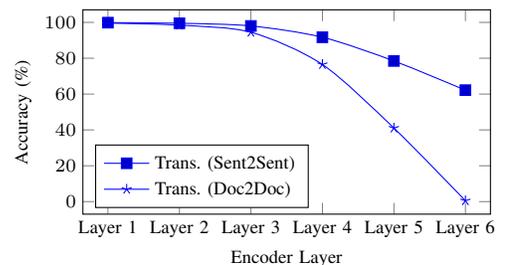


Fig. 1. Accuracy of absolute position probing task for Transformer-based Sent2Sent and Doc2Doc NMT. We take probing on different encoder layers to access the position information encoded. The higher the probing accuracy, the more position information is encoded in the hidden state of this layer.

Figure 1 compares the classification accuracy for hidden states in different encoder layers of both Sent2Sent and Doc2Doc NMT. From the figure we have the following observations:

⁵<https://github.com/facebookresearch/SentEval>

⁶When we apply the three probing tasks on middle-scale English-to-German Europarl dataset, experimental results show that Doc2Doc Transformer can preserve position information.

- For both Sent2Sent and Doc2Doc NMT, the accuracy is very high for the hidden states of low layers (layer 1 ~ 3), suggesting that position information is strongly preserved. From layer 4 to layer 6, there is an obvious trend that the accuracy starts to decrease, suggesting that the position information gets weakened or even vanished once it reaches the upper layers.
- Let us focus on the accuracy of the top layer, i.e., layer 6 where the hidden states will be used as the final output of the encoder. The accuracy for Sent2Sent NMT is 62.2% while the accuracy for Doc2Doc NMT is as extremely low as 0.5%. It indicates that the final output of the encoder for Doc2Doc NMT almost does not contain absolute position information. Therefore, this will confuse the cross-attention in the decoder and result in flat attention weights.

Even when we loose the metric to allow approximate matching, i.e., a prediction is correct if the predicted position is in the ± 3 window-size of the correct position, the accuracy of the top layer for Sent2Sent NMT increases to 95.1% while it is still as low as 1.4% for Doc2Doc NMT. This further suggests that position information has almost vanished in the final encoder output of Doc2Doc NMT.

B. Relative Position Probing Task

Given a pair of hidden states $\mathbf{h}_i, \mathbf{h}_j \in \mathbb{R}^D$ ($i \neq j$) of the i -th and the j -th source-side words, the relative position probing task is to predict their relative position, i.e., $i - j$, between the word pair. Similar to the absolute position probing, we use a neural network to map the two hidden states into a $2K$ -class probability distribution via:

$$P_R(\mathbf{k}|\mathbf{h}_i, \mathbf{h}_j) = \text{Softmax}(\mathbf{W}_R(\mathbf{h}_i - \mathbf{h}_j)), \quad (2)$$

where $\mathbf{W}_R \in \mathbb{R}^{2K \times D}$ is a trainable matrix, K is the maximum relative distance, and D is the hidden state size. $\mathbf{k} \in \{-K, -K+1, \dots, -1, 1, \dots, K-1, K\}$ is relative pairwise distance.

Given a sequence (s_1, s_2, \dots, s_I) with I words, the total number of word pairs from it is $I(I-1)$. Instead of using all word pairs, we sample I word pairs by setting h_i as h_1, h_2, \dots , and h_I , respectively while randomly sampling h_j from $\pm K$ window-size of h_i . Figure 2 shows the accuracy curves of the relative position probing on the hidden states of the top encoder layer, where the maximum relative distance $K = 20$. From Figure 2, we observe that:

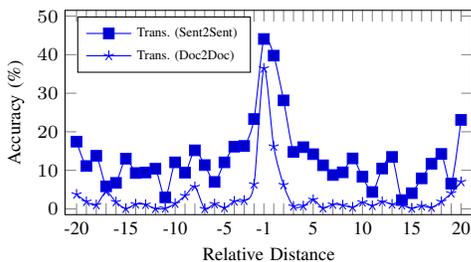


Fig. 2. Accuracy of relative position probing task for Transformer-based Sent2Sent and Doc2Doc NMT.

- For both Sent2Sent and Doc2Doc NMT, the accuracy is high among hidden states of the local area with a relative distance of ± 1 , suggesting that relative position information is well preserved.
- For relative distance from ± 2 to longer, there is a visible trend that the accuracy starts to decrease, which is more serious in the long sequence input of Doc2Doc NMT. This suggests that the relative position information gets weakened once the word pairs in a long distance.
- The probing results indicate that the relative position information is also weak or even vanishing in long sequence Doc2Doc NMT encoder output.

C. Word Order Probing Task

Given a pair of hidden states $\mathbf{h}_i, \mathbf{h}_j \in \mathbb{R}^D$ ($i \neq j$) of the i -th and the j -th source-side words, the word order probing task is a two-class classification problem by predicting the word order between the word pair, i.e., 0 for $i > j$ and 1 for $i < j$. Similar to the relative position probing task, we use a neural network to map the two hidden states into a two-class probability distribution via:

$$P_O(\mathbf{k}|\mathbf{h}_i, \mathbf{h}_j) = \text{Softmax}(\mathbf{W}_O(\mathbf{h}_i - \mathbf{h}_j)), \quad (3)$$

where $\mathbf{W}_O \in \mathbb{R}^{2 \times D}$ is a training matrix, $\mathbf{k} \in \{0, 1\}$, and D is the hidden state size.

Same as the relative position probing task, we sample I word pairs from I -length input sequence by constraining h_j from $\pm K$ window-size of h_i . Figure 3 shows the accuracy curves of the word order probing task on the hidden states of the top encoder layer, where the maximum relative distance $K = 20$. From it, we observe that:

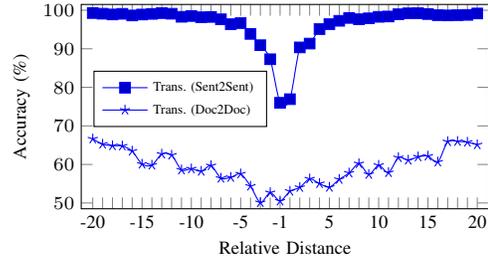


Fig. 3. Accuracy of word order probing task for Transformer-based Sent2Sent and Doc2Doc NMT.

- For Sent2Sent NMT, the probing accuracy is about 76% on the hidden states of local distance (-1 or 1), suggesting that word order information is well preserved in the local area. In the relative distance from 2 to the longer sequence, there is a visible trend that the accuracy starts to increase, reaching about 96% ~ 100%. This suggests that the word order information is strongly preserved at a long distance.
- For Doc2Doc NMT, the probing accuracy is low (about 50% ~ 55%) on the hidden states of local distance (-5 ~ 5), suggesting that word order information is very weak in the local area. The relative distance over ± 5 to the longer sequence, there is a trend that the accuracy starts to increase.

- The word order information of Sent2Sent NMT encoder output is strongly preserved, while it is very weak in Doc2Doc NMT encoder output.

III. P-TRANSFORMER: POSITION-AWARE TRANSFORMER

In Transformer [10], the attention function is viewed as a mapping between a query and a set of key-value pairs, to an output. The Transformer applies the attention function in two self-attention modules and one cross-attention module. To enhance position information in the attention function, we propose position-aware attention in which the query-key pairs are explicitly equipped with their corresponding position information. Moreover, we also propose to integrate relative position into the self-attention module to further improve the ability to perceive position information. In the next, we present the position-aware self-attention, cross-attention, and the relative position encoding for self-attention in detail.

A. Position-Aware Self-Attention

For the self-attention module, we define its input as $\mathbf{H} \in \mathbb{R}^{I \times D}$, where I is the input length and D is the size of hidden states. The original self-attention [10] computes the input as:

$$\text{Softmax} \left(\frac{(\mathbf{H}\mathbf{W}_Q)(\mathbf{H}\mathbf{W}_K)^T}{\sqrt{D}} \right) (\mathbf{H}\mathbf{W}_V), \quad (4)$$

where $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{D \times D}$ denote the query, key, and value projection matrix, respectively. In Eq. 4, position information is implicitly encoded in the hidden states \mathbf{H} . To enable the self-attention to be aware of the absolute position of the input, we simply add the position embeddings to the hidden states of query-key pairs, which we refer to as *position-aware self-attention*. The updated Eq. 4 is computed as:

$$\text{Softmax} \left(\frac{((\mathbf{H}+\mathbf{P})\mathbf{W}_Q)((\mathbf{H}+\mathbf{P})\mathbf{W}_K)^T}{\sqrt{D}} \right) (\mathbf{H}\mathbf{W}_V), \quad (5)$$

where $\mathbf{P} \in \mathbb{R}^{I \times D}$ is the absolute position embeddings, for which we use the sinusoidal position encoding proposed in Vaswani et al. [10]. Therefore, the proposed position-aware self-attention does not introduce new parameters.

Adding position information to self-attention is not novel. For example, related studies in [13–16] propose relative position or more sophisticated position encoding to the self-attention module. Compared to them, our approach is simple enough and does not need any new position encoding mechanism. Moreover, it does not change other parts of the Transformer.

B. Position-Aware Cross-Attention

In Transformer, cross-attention, also called encoder-decoder attention, passes information from the encoder to the decoder. The original cross-attention [10] computes the input as:

$$\text{Softmax} \left(\frac{(\mathbf{H}_s \mathbf{W}_Q)(\mathbf{H}_t \mathbf{W}_K)^T}{\sqrt{D}} \right) (\mathbf{H}_s \mathbf{W}_V), \quad (6)$$

where $\mathbf{H}_s \in \mathbb{R}^{I_s \times D}$, $\mathbf{H}_t \in \mathbb{R}^{I_t \times D}$ denote the output of the encoder and the output of self-attention of the decoder side, respectively, while I_s and I_t denote the lengths of source-side and target-side sequences, respectively. Similar to position-aware self-attention, we enhance the cross-attention to be aware of the absolute position of the hidden states of the encoder and the decoder, which we refer to as *position-aware cross-attention*. The updated Eq. 6 is computed as:

$$\text{Softmax} \left(\frac{((\mathbf{H}_t+\mathbf{P}_t)\mathbf{W}_Q)((\mathbf{H}_s+\mathbf{P}_s)\mathbf{W}_K)^T}{\sqrt{D}} \right) (\mathbf{H}_s \mathbf{W}_V), \quad (7)$$

where $\mathbf{P}_s \in \mathbb{R}^{I_s \times D}$, $\mathbf{P}_t \in \mathbb{R}^{I_t \times D}$ denote the corresponding sinusoidal position embeddings of \mathbf{H}_s and \mathbf{H}_t , respectively.

After adding absolute position information to the query-key pairs, the cross-attention will learn to effectively extract useful context information over the source document by comparing the absolute position of the query against the position of the keys. Compared with the methods of integrating local attention and global attention [6, 17], our position-aware cross-attention is much simpler without introducing new parameters.

C. Relative Position for Self-attention

From the relative position probing result in Figure 2, we find that the relative position information will decrease rapidly for pairs of long distance. To enhance relative position information in self-attention, we follow Huang et al. [18] to integrate the relative position embeddings into the query-key component in the self-attention network.

Given an input $\mathbf{H} \in \mathbb{R}^{I \times D}$, the relative position-aware self-attention is computed as:

$$\text{Softmax} \left(\frac{(\mathbf{H}\mathbf{W}_Q)(\mathbf{H}\mathbf{W}_K)^T + \mathbf{S}_{rel}}{\sqrt{D}} \right) (\mathbf{H}\mathbf{W}_V), \quad (8)$$

$$\mathbf{S}_{rel} = (\mathbf{H}\mathbf{W}_Q)\mathbf{R}^T, \quad (9)$$

where $\mathbf{R} \in \mathbb{R}^{I \times I \times D}$ is a tensor of relative positional embeddings.⁷ Note that unlike Shaw et al. [13], here we do not clip the minimum/maximum relative position (i.e., ± 512 , 512 is the maximum sequence length) to a certain value. When we apply relative position for position-aware self-attention, Eq. 8 and Eq. 9 are updated as:

$$\text{Softmax} \left(\frac{((\mathbf{H}+\mathbf{P})\mathbf{W}_Q)((\mathbf{H}+\mathbf{P})\mathbf{W}_K)^T + \mathbf{S}_{rel}}{\sqrt{D}} \right) (\mathbf{H}\mathbf{W}_V), \quad (10)$$

$$\mathbf{S}_{rel} = ((\mathbf{H} + \mathbf{P}) \mathbf{W}_Q) \mathbf{R}^T. \quad (11)$$

The relative position [13] provides the relative pairwise distance of the input sequence, which has been successfully applied to the sentence-level NMT model [13] and other pre-training models [19]. In this paper, we enhance the document-level seq2seq model to be aware of relative position information in a long input sequence. Note that we integrate relative position into the self-attention module both in the encoder and decoder side.

⁷In implementation, the result of $\mathbf{H}\mathbf{W}_Q$ will be reshaped as $\mathbb{R}^{I \times I \times D}$ before multiplying with \mathbf{R}^T .

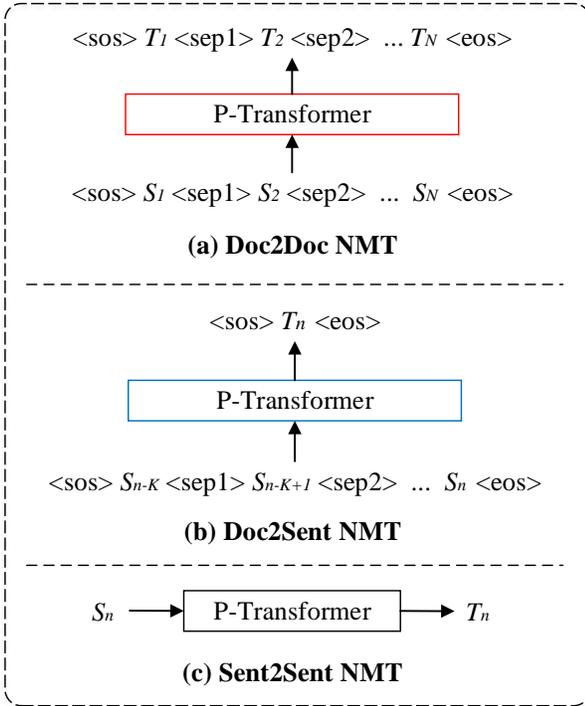


Fig. 4. Overview of our proposed P-Transformer-based Doc2Doc, Doc2Sent, and Sent2Sent NMT models.

D. Applying P-Transformer to Multiple NMT Models

We use a parallel document pair $(S, T) = \{S_n, T_n\}_{n=1}^N$ with N sentence pairs to illustrate how we apply the proposed P-Transformer to the following three types of seq2seq-based translation:

Doc2Doc Translation. For Doc2Doc translation, we concatenate the sentences within a document, both source-side and target-side, together into a long sequence, as shown in Figure 4 (a). We insert $\langle \text{sos} \rangle$, $\langle \text{eos} \rangle$, $\langle \text{sep}I \rangle$ ($I \in \{1, 2, 3, \dots\}$) into the sequence to denote the start of a document, the end of a document, and the separator between the I -th and the $(I + 1)$ -th sentences, respectively. With the sentence separators $\langle \text{sep}I \rangle$, it is easy to get sentence-to-sentence alignment between a source document and its translation.

Doc2Sent Translation. We follow previous work [1, 2] to use K previous sentences as context for Doc2Sent translation. As shown in Figure 4 (b), the source side input is a sequence consisting of the K previous sentences, plus the current sentence while the target side output is the translation of the current sentence. Similar to Doc2Doc translation, we insert $\langle \text{sos} \rangle$, $\langle \text{eos} \rangle$, $\langle \text{sep}I \rangle$ ($I \in \{1, 2, \dots, K\}$) in the source-side sequence. For the target side sequence, we only insert $\langle \text{sos} \rangle$ and $\langle \text{eos} \rangle$.

Sent2Sent Translation. As shown in Figure 4 (c) we can also directly apply P-Transformer to Sent2Sent translation.

IV. EXPERIMENTATION

A. Experimental Setup

Datasets. We carry out extensive experiments on various document-level translation datasets covering 7 language pairs, including Chinese-to-English (Zh-En), English-to-German (En-De), French-to-English (Fr-En), Spanish-to-English (Es-En), Russian-to-English (Ru-En), English-to-Russian (En-Ru) and English-to-French (En-Fr). The details of the datasets are listed as follows:

- **Zh-En:** we use the parallel document corpus (PDC) released by Sun et al. [1] for Zh-En translation.⁸ The dataset contains middle-scale data with documents from different domains, including politics, finance, health, culture, etc. The training, validation, and test sets are including 1.39M / 2.0M / 4.9K sentences, respectively.
- **En-De:** we use three datasets provided by Maruf et al. [4], including TED, News and Europarl.⁹ The TED dataset is from IWSLT 2017¹⁰ [20] and we use tst2016-2017 as the test set, and other as the validation set. The News dataset is from the News Commentary v11 corpus.¹¹ We use WMT newstest2015 as the validation set and newstest2016 as the test set, respectively. The Europarl dataset is extracted from Europarl v7, and the training, validation, and test sets are obtained by randomly splitting the corpus. The training, validation, and test sets of TED are including 0.21M / 9.0K / 2.3K sentences, respectively. The training, validation, and test sets of News are including 0.24M / 2.0K / 3.0K sentences, respectively. The training, validation, and test sets of Europarl are including 1.67M / 3.6K / 5.1K sentences, respectively.
- **Fr-En, Es-En and Ru-En:** we use News Commentary v14 (News14) from WMT 2019¹² as their training data. The validation sets / test sets are newstest2012 / newstest2013 for Es-En translation, newstest2013 / newstest2014 for FR-EN translation, newstest2018 / newstest2019 for Ru-En translation, respectively. The training, validation, and test sets of Fr-En are including 0.33M / 3.0K / 3.0K sentences, respectively. The training, validation, and test sets of Es-En are including 0.38M / 3.0K / 3.0K sentences, respectively. The training, validation, and test sets of Ru-En are including 0.28M / 3.0K / 2.0K sentences, respectively.
- **En-Ru and En-Fr:** we choose the widely used large-scale OpenSubtitles 2018 (Subtitles) datasets.¹³ The original En-Ru and En-Fr datasets contain 25.9M / 41.8M sentences. we split long documents into sub-documents with up to 512 tokens, resulting in 0.48M / 0.81M sub-documents. For En-Ru and En-Fr translations, we randomly selected 0.36M / 2K / 1K and 0.60M / 2K / 1K documents as the training set, validation set, and test set

⁸https://github.com/sunzewei2715/Doc2Doc_NMT

⁹<https://github.com/sameenmaruf/selective-attn/tree/master/data>

¹⁰<https://wit3.fbk.eu/>

¹¹<http://www.casmacat.eu/corpus/news-commentary.html>

¹²<https://www.statmt.org/wmt19/translation-task.html>

¹³<https://opus.nlpl.eu/OpenSubtitles2018.php>

TABLE I
STATISTICS ON THE PRE-PROCESSED DOCUMENT-LEVEL TRANSLATION DATA SETS. “S”, “M” AND “L” DENOTE THAT THE CORRESPONDING TRAINING DATA IS IN SMALL-, MIDDLE- OR LARGE-SCALED, RESPECTIVELY.

Dataset	Type	#Doc			Avg # Token/Doc		
		(train / valid / test)					
Zh-En	PDC	M	119.8K / 255 / 320	383 / 282 / 378			
En-De	TED	S	10.9K / 468 / 115	451 / 444 / 445			
	News	S	17.9K / 165 / 249	405 / 365 / 338			
	Europarl	M	150.8K / 322 / 458	334 / 343 / 342			
Fr-En	News14	S	24.7K / 201 / 274	430 / 436 / 352			
Es-En	News14	S	28.0K / 229 / 192	424 / 414 / 443			
Ru-En	News14	S	21.4K / 266 / 213	411 / 345 / 258			
En-Ru	Subtitles	L	360K / 2000 / 1000	486 / 487 / 488			
En-Fr	Subtitles	L	600K / 2000 / 1000	478 / 479 / 478			

from the processed document-level instances, and there is none-overlapping between them.

In pre-processing, we tokenize all sentences by Moses toolkit¹⁴ while the Chinese sentences are segmented by the widely used Jieba toolkit.¹⁵ Then, the source and target sentences are segmented into sub-words by a joint BPE model [24] with 32K merged operations. Note that the pre-processing is identical to the related work of G-Transformer [6] and MR Doc2Doc [1]. Finally, following G-Transformer [6], we split long documents into sub-documents with up to 512 tokens in Doc2Doc NMT experiments. Table I shows the detailed statistics of datasets after pre-processing.

Model Setting. We use *Fairseq* [25] as the implementation of Transformer models. We follow the standard Transformer base model setting [10], in which we use 6 layers, 8 heads, 512 dimension outputs, and 2048 dimension hidden states. Note that the position information we add to attention is always sinusoidal PEs. In all experiments, we use the learning rate decay policy proposed by Vaswani et al. [10] (warm-up step 4K) with label smoothing of 0.1, and the dropout is 0.3. We share bilingual vocabulary to reduce the computation cost. In inference, we choose the best checkpoint on the validation set to evaluate the translation performance. We set the beam size as 5, and the three parameters to control the generation length $\text{lenpen} / \text{max-len-a} / \text{max-len-b}$ as 1.0 / 1.1 / 7, respectively. We run all experiments 3 times with 3 different random seeds on the small- and middle-scaled Zh-En, En-De, En-Fr, En-Es, and En-Ru tasks and reported averaged BLEU scores. Since the big load of computation, the experiments on large-size datasets only run once.

Training. All the models are trained on a single Tesla V100 GPU. We set the gradient update frequency / token size as 4 / 8192 for the Transformer model on middle- and large-scaled PDC, Europarl, and Subtitles datasets, 4 / 4096 for other small datasets. Finally, we stop training by using an early stopping strategy on the validation set for small- and middle-scaled datasets. For large Subtitles datasets, we train NMT models with 20 epochs on both En-Ru and En-Fr translations. As fine-tuning the sentence-level NMT model benefits

Doc2Doc model [6, 11], we also combine both Doc2Sent (or Doc2Doc) training instances and Sent2Sent training instances when training P-Transformer-based Doc2Sent (or Doc2Doc) models.

Evaluation. Following previous related work, we report case-sensitive document-level BLEU [11] (d-BLEU) for the document-level NMT systems, which is computed by matching n-grams in the whole document after removing the special tokens of $\langle \text{sos} \rangle$, $\langle \text{eos} \rangle$, and $\langle \text{sep}I \rangle$. Thanks to the proposed sentence separators, we can obtain the sentence-to-sentence alignments between a source document and its translation. So we also report the conventional case-sensitive sentence-level BLEU (s-BLEU).¹⁶

B. Main Results of Sent2Sent Translation and Doc2Sent Translation

Table II (Sent2Sent) shows the BLEU scores (s-BLEU) of the Sent2Sent NMT models. The performance of our Transformer baseline (#3) is comparable with that of related studies (#1 and #2). Moreover, our proposed P-Transformer (#4) improves the translation quality of Sent2Sent NMT on all datasets, which suggests that even though the position information is well preserved at the top layer of the Transformer encoder, it also benefits from our proposed position-aware attention.

Table II (Doc2Sent) presents the experimental results of Doc2Sent models. Note that in the Doc2Sent instances, the source side input is a sequence consisting of the 4 previous sentences, plus one current sentence. From the table, we observe that with the help of contextual sentences, the translation performance of P-Transformer (ctx) (#12) is improved on En-De translations. This shows that a larger source context benefits translation performance. By including sentence-level training data, P-Transformer (ctx + sent) (#13) further improves the translation quality on both Zh-En and En-De translations. Compared with related studies in Doc2Sent, our model achieves the best performance.

C. Main Results of Doc2Doc NMT

Results of Zh-En and En-De translations. As shown in Table II (Doc2Doc), it fails to train a Doc2Doc model with vanilla Transformer (#18) on small-scaled datasets (e.g., TED and News), while it can be successfully trained on the middle-scaled PDC and Europarl datasets. This is consistent with the findings in G-Transformer[6]. Different from the vanilla Transformer, the proposed P-Transformer (#19) can be successfully trained on all datasets. Compared with Sent2Sent (#3) baseline systems, P-Transformer (#19) trained directly on document-level data achieves comparable performance on small-scaled TED and News datasets, or significantly better performance on middle-size PDC and Europarl datasets, which shows good capability of P-Transformer in document-level context modeling. Moreover, P-Transformer (#19) outperforms both G-Transformer and SR Doc2Doc (#14 and #16) when direct

¹⁴<https://github.com/marian-nmt/moses-scripts>

¹⁵<https://github.com/fxsjy/jieba>

¹⁶As shown in TableVII, very few source sentences have no alignment in the translation. For these source sentences, we set their translation as empty when calculating s-BLEU.

TABLE II
BLEU SCORES ON ZH-EN AND EN-DE TRANSLATIONS. †/‡: SIGNIFICANT OVER TRANSFORMER (SENT) BASELINES AT 0.05/0.01.

Type	#	Model	Zh-En PDC		TED		En-De News		Europarl		
			s-BLEU	d-BLEU	s-BLEU	d-BLEU	s-BLEU	d-BLEU	s-BLEU	d-BLEU	
Existing Models											
Sent2Sent	1	Transformer [6]	-	-	24.82	-	25.19	-	31.37	-	
	2	Transformer [1]	-	-	25.19	-	24.98	-	31.70	-	
	Our Models										
	3	Transformer (sent)	25.80	-	24.83	-	25.13	-	31.55	-	
	4	P-Transformer (sent)	26.13 †	-	25.26 †	-	25.57 †	-	31.85 †	-	
Existing Models											
Doc2Sent	5	DocT [2]	-	-	24.00	-	23.08	-	29.32	-	
	6	HAN [21]	-	-	24.58	-	25.03	-	28.60	-	
	7	SAN [4]	-	-	24.42	-	24.48	-	29.75	-	
	8	QCN [22]	-	-	25.19	-	22.37	-	29.82	-	
	9	Flat-Tran. [9]	-	-	24.87	-	23.55	-	30.09	-	
	10	MCN [23]	-	-	25.10	-	24.91	-	30.40	-	
	11	MR Doc2Sent [1]	-	-	25.24	29.20	25.00	26.70	32.11	34.18	
	Our Models										
		12	P-Transformer (ctx)	25.50	-	25.47‡	-	25.15	-	31.80†	-
		13	P-Transformer (ctx + sent)	26.43 ‡	-	25.73 ‡	-	25.71 ‡	-	32.15 ‡	-
	Existing Models										
Doc2Doc	14	G-Transformer [6]	-	-	23.53	25.84	23.55	25.23	32.18	33.87	
	15	G-Transformer (fine-tuned) [6]	-	-	25.12	27.17	25.52	27.11	32.39	34.08	
	16	SR Doc2Doc [1]	-	24.33	-	4.70	-	21.18	-	34.16	
	17	MR Doc2Doc [1]	-	27.80	-	29.27	-	26.71	-	34.48	
	Our Models										
		18	Transformer (doc)	25.38	27.39	-	0.76	-	0.60	31.54	33.21
	19	P-Transformer (doc)	26.53‡	28.53	24.91	27.09	24.92	27.03	32.37‡	34.14	
	20	P-Transformer (doc + sent)	27.14 ‡	29.26	25.67 ‡	27.94	25.93 ‡	27.67	32.62 ‡	34.49	

Note: (sent), (ctx), (doc) indicate the model is trained on Sent2Sent, Doc2Sent and Doc2Doc training instances, respectively. (ctx/doc + sent) indicate the model is trained on Doc2Sent/Doc2Doc and Sent2Sent training instances.

TABLE III
BLEU SCORES OF DOCUMENT-LEVEL NMT MODELS ON FR-EN, ES-EN, AND RU-EN TRANSLATION TASKS. †/‡: SIGNIFICANT OVER TRANSFORMER (SENT) BASELINES AT 0.05/0.01.

Type	Model	Fr-En		Es-En		Ru-En	
		s-BLEU	d-BLEU	s-BLEU	d-BLEU	s-BLEU	d-BLEU
Sent2Sent	Transformer (sent)	28.11	-	28.21	-	22.03	-
Doc2Doc	SR Doc2Doc [1]	-	23.86	-	26.79	-	16.47
	MR Doc2Doc [1]	-	28.85	-	29.37	-	23.98
	Transformer (doc)	-	0.79	-	0.73	-	0.06
	P-Transformer (doc)	28.29	29.94	28.10	30.65	21.56	23.96
	P-Transformer (doc + sent)	29.16 ‡	30.79	28.98 ‡	31.26	22.72 ‡	24.39

training on document-level data, especially on small-scale datasets. Finally, compared with G-Transformer (#15) and MR Doc2Doc (#17), P-Transformer with additional sentence-level training data (#20) achieves a new state-of-the-art on PDC, TED, News, and Europarl. Note that our model (#20) is lower than both MR Doc2sent (#11) and MR Doc2Doc (#17) models on the TED data set measured by d-BLEU, but it outperforms MR Doc2Sent (#11) in s-BLEU, while the two systems (#11 and #17) get similar performance in d-BLEU.

Results of Fr-En, Es-En and Ru-En translations. Table III shows the translation performance on Fr-En, Es-En, and Ru-En translations. It is unsurprising to observe that vanilla Transformer fails to be successfully trained on the three small-scale datasets. When being trained on document-level data, P-Transformer (doc) achieves similar performance on Fr-En and Es-En translation compared to sentence-level Transformer while it achieves better or comparable performance than MR Doc2Doc [1]. By including sentence-level training data, P-

TABLE IV
BLEU SCORES OF DOCUMENT-LEVEL NMT MODELS ON LARGE-SCALE EN-RU AND EN-FR TRANSLATIONS, RESPECTIVELY. †/‡: SIGNIFICANT OVER TRANSFORMER (SENT) BASELINES AT 0.05/0.01.

Models	En-Ru		En-Fr	
	s-BLEU	d-BLEU	s-BLEU	d-BLEU
Transformer (sent)	25.80	-	33.89	-
Transformer (doc)	24.32	28.16	32.82	36.70
P-Transformer (doc)	24.94	28.81	34.24‡	38.07
P-Transformer (doc + sent)	25.98 †	29.83	34.36 ‡	38.19

Transformer (doc + sent) achieves the best performance over the three translation tasks.

Results of En-Ru and En-Fr translations. Table IV shows the translation performance on large-scale En-Ru and En-Fr translation tasks. From it, we observe that our proposed P-Transformer (doc) significantly outperforms Transformer (doc) in terms of both s-BLEU and d-BLEU. For example, P-

Transformer (doc) achieves gains of 1.42 / 1.37 on s-BLEU / d-BLEU over Transformer (doc) in En-Fr translation. Similarly, by taking advantage of sentence-level training data, P-Transformer (doc + sent) achieves the best performance, with further improved translation performance. From the above experimental results, we see that the proposed P-Transformer has a strong generalization ability on context modeling over different language pairs and dataset scales.

D. Parameters and Training Speed

TABLE V
PARAMETER (IN MILLIONS) AND TRAINING SPEED (SECOND / EPOCH).

#	Model	#Param.	Speed
1	Transformer (sent)	59.38M	460s
2	P-Transformer (sent)	59.45M	473s
3	-rel-self	59.38M	462s
4	P-Transformer (doc)	59.45M	647s
5	-rel-self	59.38M	562s

Note: “-rel-self” denotes removing the relative position encoding in self-attention.

We use models on the En-De TED translation task as an example to analyze model parameters and training speed, as shown in Table V. Comparing #1 and #3, we observe that our proposed position-aware self-attention and cross-attention bring no additional parameters and negligible computation cost. Moreover, the relative position for self-attention only introduces a matrix of relative position embeddings with $1025 \times 64 = 65,600$ additional parameters and requires slightly more computation cost (as shown in #2 VS. #3, and #4 VS. #5). It is well known that the computational complexity of self-attention is $I^2 \times D$, where I is the sequence length, and D is the hidden state size. Compared with the sentence-level baseline model, the training time of an epoch for P-Transformer (doc) increases by about 40% (as shown in #1 VS. #4).

V. ANALYSIS

In this section, we take Doc2Doc NMT models as representatives to discuss the effectiveness of P-Transformer, which unless otherwise specified, is trained on document-level instances, i.e., P-Transformer (doc).

A. Analysis on Different Components

We explore the contributions of position-aware self-attention, cross-attention, and relative position encoding for Doc2Doc P-Transformer. Table VI shows the performance in d-BLEU on the three En-De translation tasks.

On the one hand, P-Transformer fails on TED and News when we disable position-aware attention in the cross-attention module (-cross-attn). On the other hand, P-Transformer still can be properly trained when we disable position-aware attention in the three self-attention modules (-rel-self, -self-src, -self-tgt). The performance trend suggests that the cross-attention module is the key for Doc2Doc NMT while the two self-attention modules can also benefit from position information.

TABLE VI
PERFORMANCE (D-BLEU) ON THE THREE EN-DE TRANSLATION TASKS.

Model	TED	News	Europarl	Drop
P-Trans. (doc)	27.09	27.03	34.14	-
-cross-attn	0.70	0.83	33.67	-17.53
-self-src, -self-tgt	24.05	26.05	33.78	-1.34
-self-src	25.67	25.12	33.90	-1.11
-self-tgt	26.62	26.12	33.93	-0.46
-rel-self	26.75	26.61	34.01	-0.25

Note: “-cross-attn”, “-self-src”, and “-self-tgt” denote removing the absolute position embeddings in the cross-attention, source-side self-attention, and target-side self-attention, respectively while “-rel-self” denotes removing the relative position embeddings in self-attention.

TABLE VII
STATISTICS ON THE TEST SETS, REGARDING DOCUMENTS WITHOUT SENTENCE-LEVEL COVERAGE ISSUE IN Doc2Doc TRANSLATION.

Dataset	Correct	Total	Percentage
Zh-En (PDC)	316	317	99.7%
En-De (TED)	114	115	99.1%
En-De (News)	245	249	98.4%
En-De (Europarl)	458	458	100%
Fr-En (News14)	273	274	99.6%
Es-En (News14)	188	192	97.9%
Ru-En (News14)	209	213	98.1%
En-Ru (Subtitles)	991	1000	99.1%
En-Fr (Subtitles)	985	1000	98.5%
Overall	3779	3818	99.0%

B. Analysis on Sentence Separator in Doc2Doc Translation

In this section, we first analyze the sentence-level coverage issue in Doc2Doc translation. Then we compare the performance of using different sentence separator methods.

In Doc2Doc translation, we insert the separator $\langle \text{sep}I \rangle$ between the I -th and the $I+1$ -th sentences in both the source and target sides. Therefore, we recover sentence-level translation from these separators. We say a document having no sentence-level coverage issues if we can perfectly recover sentence-level translation for all its source sentences. From Table VII, we find that 97.9% to 100% documents are properly translated without sentence-level coverage issues. For those few documents having sentence-level coverage issues, further analysis reveals that most of them miss one or two sentences at the end of translation.

TABLE VIII
PERFORMANCE (S-BLEU / D-BLEU) COMPARISON OF USING DIFFERENT SENTENCE SEPARATOR METHODS.

Model	Sep.	TED	News	Europarl
P-Trans.	$\langle \text{sep}I \rangle$	24.91 / 27.09	24.92 / 27.03	32.37 / 34.14
P-Trans.	$\langle \text{sep} \rangle$	24.37 / 27.25	23.23 / 26.71	32.06 / 33.96

Note: Here P-Trans. indicates P-Transformer (doc) model.

Instead of using index-aware sentence separators, i.e., $\langle \text{sep}I \rangle$, we compare its performance with another method that uses a uniform sentence separator $\langle \text{sep} \rangle$. As shown in Table VIII, we find that overall the uniform sentence separator method under-performs the index-aware sentence separator method. Moreover, further analysis of the test sets reveals that

the uniform sentence separator suffers more from the sentence-level coverage issue.

C. One Model Learns both Sent2Sent and Doc2Doc Translation

As mentioned, our P-Transformer (doc + sent) in Table II (Doc2Doc) is trained on both sentence-level and document-level sequences. Therefore, the trained model could translate both sentences and documents. Table IX compares the performance of s-BLEU when the input unit in inference is a sentence or document.

TABLE IX
PERFORMANCE COMPARISON ON S-BLEU SCORE OF SENT2SENT AND DOC2DOC NMT MODELS.

Model	Input	PDC	TED	News	Europarl
Trans. (sent)	sent	25.80	24.83	25.13	31.55
P-Trans. (sent)	sent	26.13	25.26	25.57	31.85
P-Trans. (doc + sent)	sent	26.40	25.21	25.37	31.99
P-Trans. (doc + sent)	doc	27.14	25.67	25.93	32.62

Table IX shows that no matter whether the input unit is a sentence or document, P-Transformer (doc + sent) achieves better performance than sentence-level Transformer baseline. This shows that our proposed document-level seq2seq model is beneficial for both document- and sentence-level NMT. This result demonstrates that a single P-Transformer (doc + sent) model can be used for both Sent2Sent and Doc2Doc translation tasks. Moreover, it shows Doc2Doc translation benefits from using document-level context.

D. Analysis on Discourse Phenomena

To verify whether Doc2Doc NMT truly learns the useful contextual information to improve discourse coherence, we use the linguistic feature test set provided by Sun et al. [1] to evaluate different discourse phenomena, including tense consistency (TC), conjunction presence (CP), and pronoun translation (PT). TCP is an overall score calculated as the geometric mean of TC, CP, and PT, which has a strong correlation with human evaluation [1].

TABLE X
DISCOURSE PHENOMENA EVALUATION ON THE ZH-EN TEST SET.

Model	s-BLEU	TC	CP	PT	TCP
Transformer (sent)	25.80	57.3	34.7	61.6	49.6
Transformer (doc)	25.38	55.7	31.9	60.8	47.6
P-Transformer (doc)	26.53	57.7	37.4	62.2	51.2
P-Transformer (doc + sent)	27.14	58.2	33.6	64.9	50.3

As shown in Table X, compared with Transformer (doc) and Transformer (sent) systems, our proposed P-Transformer (doc) improves not only the translation performance in BLEU but also the performance of discourse phenomena. By including sentence-level data, our P-Transformer (doc + sent) further improves TC and PT, but decreases both CP and overall TCP. This suggests that though sentence-level data improves the document-level NMT quality in BLEU, it has an uncertain

effect on various discourse phenomena, e.g., improving performance in TC and PT while decreasing performance in CP and overall TCP.

E. Analysis on Transformer Encoder Depth

As our probing task result suggests that position information gradually vanished from the bottom encoder layers to the up layers in vanilla Transformer. Therefore, decreasing its encoder depth might make the training successful, especially on small datasets.

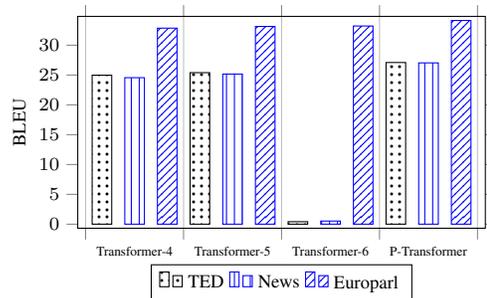


Fig. 5. BLUE scores when using different encoder depths for Doc2Doc NMT on the three En-De translation tasks. “Transformer-X” represents vanilla Transformer with X encoder layers.

As shown in Figure 5, vanilla Transformer models with encoder depth of either 4 or 5 achieve good performance on all datasets. However, decreasing encoder depth would weaken the capability of capturing useful information from input sequences, thus lowering the translation performance.

F. Analysis on Input Length

Next, we investigate the effect of the input length. To this end, we split long documents of TED into sub-documents with up to I tokens, where $I \in \{64, 128, 256, \dots\}$.

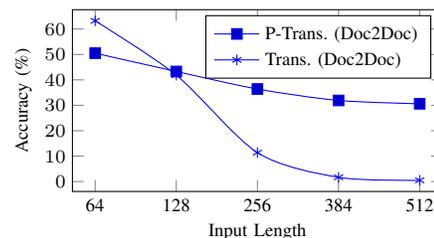


Fig. 6. The accuracy of position information probing task with respect to different input lengths.

Figure 6 compares the accuracy of the absolute position probing task of the top encoder layer. It shows that 1) for Transformer, the accuracy significantly decreases when the input length exceeds 128 and the position information has almost vanished when the length is over 256; 2) for P-Transformer, the encoder still captures certain position information even when the input length increases to 512.

Figure 7 compares the translation performance on the TED dataset. It shows that the vanilla Transformer benefits from increasing the input length from 64 to 128 and its performance

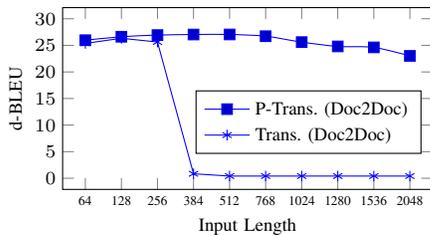


Fig. 7. d-BLEU scores with respective to the different input lengths.

starts to decrease when increasing the length to 256. Moreover, translation is completely failed when the input length is up to 384 or longer. On the other hand, P-Transformer is less sensitive to the input length, as its translation performance is getting higher when increasing the input length from 64 to 768. Even when the input length is as long as 2048, it still achieves a 23 d-BLEU score.

TABLE XI
PERFORMANCE OF P-TRANSFORMER ON ZH-EN PDC DATASET.

Input Length	s-BLEU	d-BLEU
512	26.53	28.53
1024	26.74	29.17
2018	25.91	28.58

To further verify the performance of P-Transformer on long document translation, we split documents of Zh-En PDC dataset into sub-documents with up to J tokens, where $J \in \{512, 1024, 2048\}$. From Table XI, we find that in the middle-scale PDC dataset, the translation performance is stable and robust over different lengths, in which the document of 1024 length achieves the best performance.

G. Analysis on Training Data Scale

To test the performance of the P-Transformer on small training data, we train models on training datasets with different data sizes by randomly selecting 1K ~ 10K sub-documents from En-De TED translation.

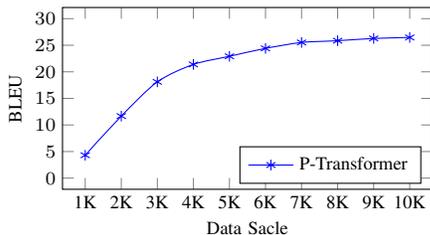


Fig. 8. BLEU scores of P-Transformer on various data scales.

As shown in Figure 8, P-Transformer has a smooth curve of performance when increasing the size of training data from 1K to 10K. It suggests that P-Transformer can be successfully trained for document-level NMT even on extremely small datasets.

VI. RELATED WORK

In this section, we discuss related work from three perspectives: Doc2Sent NMT models, Doc2Doc NMT models, and position encoding for Transformer.

A. Doc2Sent NMT Models

Previous studies have proposed various Doc2Sent NMT models and have achieved great success by generating target-side sentences independently. These studies could be further roughly categorized into two groups. Studies of the first group [2, 7, 9, 22, 26–34] consider partial source-side context, such as a few previous sentences. These models are strictly sentence-level as document-level context could be viewed as extra input. Some of those studies [2, 7, 9, 22, 26, 27, 31–34] aim to extract general useful information from context to help for translating current sentences while some [28–30] aim to resolve discourse phenomena. Studies of the second group [4, 35–39] extend the source-side context from the local context of a few sentences into larger context within the whole document. These models usually take documents as input units and extract useful information from the global context to help translate current sentences. However, one obvious disadvantage of these Doc2Sent models is that it is difficult to properly use the target-side context, which is potentially useful for translation.

B. Doc2Doc NMT Models

Much effort has also been devoted to Doc2Doc translation. Concatenating multiple sentences into a unit is the preliminary exploration of Doc2Doc NMT [7–9, 40]. However, these approaches are limited within a short context of up to 4 ~ 6 sentences. Recent studies successfully train vanilla Transformer models for Doc2Doc translation by taking advantage of either large augmented datasets [1, 3, 41] or pre-trained models [11]. Alternatively, Bao et al. [6] find that the training failure is not caused by over-fitting, but by sticking around local minima. Consequently, they propose G-Transformer by introducing local bias to attention to constrain a target sentence to attend to its corresponding source sentence. G-Transformer is the first Doc2Doc NMT model which can be trained from scratch even on a small dataset. However, this approach requires sentence-to-sentence alignment between the source and target sides, which limits its generality for other seq2seq document-level NLP tasks, like text summarization. Following this research line, we find that the failure of training Doc2Doc NMT models from scratch is due to the vanishing of position information at the encoder output. To this end, we propose a position-aware Transformer to enhance both the absolute and relative position information in attention modules. Our position-aware self-attention and cross-attention do not introduce any parameters and more importantly, can be applied to other document-level seq2seq tasks.

C. Position Encoding for Transformer

Position information plays a crucial role in Transformer to model the word order of the input sequence. The absolute

position information [10, 42] is properly propagated to higher layers through residual connections. The relative positional encoding [13, 43–45] extends the self-attention that can be used to incorporate relative position information among tokens within a sequence. Later, Transformer-based pre-trained models, including BERT [46] and BART [47] use fully learnable PEs. Chen et al. [16] show the gains of the relative positional encoding coming from moving positional information from the input to attention. So Chen et al. [16] propose the decoupled positional attention to encode position and segment information into the Transformer models. Different from the above studies, in this paper we find that position information has almost vanished in the output of the encoder when the input becomes long on small datasets. And this could be easily resolved by letting the attention modules be explicitly aware of the position information for query-key pairs.

VII. CONCLUSION

In this paper, we have investigated the main reasons for training failure on the document-to-document Transformer, especially on small datasets. Through the position information probing tasks, we find that position information in the source representation has almost vanished at the top layer. Therefore, we propose a simple position-aware self-attention and cross-attention by explicitly adding the position embeddings to the query-key pairs in the attention function. The proposed P-Transformer is a general, truly Doc2Doc model that directly translates a source document to its corresponding target document. Experimental results on several datasets show that P-Transformer significantly outperforms the strong baseline and achieves state-of-the-art performance.

Though the P-Transformer achieves good performance on Doc2Doc translation. However, memory usage will be a bottleneck for documents with thousands of tokens. So the computation efficiency on long-range sequence processing in Doc2Doc translation is the limitation that remains to be further explored. Moreover, in the future, we will evaluate the proposed P-Transformer against other document-level seq2seq tasks, e.g., text summarization [48], text simplification [49] or chat translation [50].

REFERENCES

- [1] Z. Sun, M. Wang, H. Zhou, C. Zhao, S. Huang, J. Chen, and L. Li, “Rethinking document-level neural machine translation,” in *Findings of the Association for Computational Linguistics: ACL 2022*. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 3537–3548. [Online]. Available: <https://aclanthology.org/2022.findings-acl.279>
- [2] J. Zhang, H. Luan, M. Sun, and et al., “Improving the transformer translation model with document-level context,” in *Proceedings of EMNLP*, 2018, pp. 533–542. [Online]. Available: <https://aclanthology.org/D18-1049>
- [3] M. Junczys-Dowmunt, “Microsoft translator at WMT 2019: Towards large-scale document-level neural machine translation,” in *Proceedings of WMT*, 2019, pp. 225–233. [Online]. Available: <https://aclanthology.org/W19-5321>
- [4] S. Maruf, A. F. T. Martins, and G. Haffari, “Selective attention for context-aware neural machine translation,” in *Proceedings of NAACL*, 2019, pp. 3092–3102. [Online]. Available: <https://aclanthology.org/N19-1313>
- [5] E. Voita, R. Sennrich, and I. Titov, “When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion,” in *Proceedings of ACL*, 2019, pp. 1198–1212. [Online]. Available: <https://aclanthology.org/P19-1116>
- [6] G. Bao, Y. Zhang, Z. Teng, B. Chen, and W. Luo, “G-transformer for document-level machine translation,” in *Proceedings of ACL*, 2021, pp. 3442–3455. [Online]. Available: <https://aclanthology.org/2021.acl-long.267>
- [7] J. Tiedemann and Y. Scherrer, “Neural machine translation with extended context,” in *Proceedings of WMT*, 2017, pp. 82–92. [Online]. Available: <https://aclanthology.org/W17-4811>
- [8] R. Agrawal, M. Turchi, and M. Negri, “Contextual handling in neural machine translation: Look behind, ahead and on both sides,” in *In Proceedings of EACL*, 2018, pp. 11–20.
- [9] S. Ma, D. Zhang, and M. Zhou, “A simple and effective unified encoder for document-level machine translation,” in *Proceedings of ACL*, 2020, pp. 3505–3511. [Online]. Available: <https://aclanthology.org/2020.acl-main.321>
- [10] 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 *Proceedings of NIPS*, 2017, pp. 5998–6008. [Online]. Available: <https://arxiv.org/pdf/1706.03762.pdf>
- [11] Y. Liu, J. Gu, N. Goyal, and et al., “Multilingual denoising pre-training for neural machine translation,” *Transactions of the Association for Computational Linguistics*, pp. 726–742, 2020. [Online]. Available: <https://aclanthology.org/2020.tacl-1.47>
- [12] A. Conneau, G. Kruszewski, G. Lample, L. Barrault, and M. Baroni, “What you can cram into a single \mathbb{R}^d vector: Probing sentence embeddings for linguistic properties,” in *Proceedings of ACL*, 2018. [Online]. Available: <https://www.aclweb.org/anthology/P18-1198>
- [13] P. Shaw, J. Uszkoreit, and A. Vaswani, “Self-attention with relative position representations,” in *Proceedings of NAACL*, 2018, pp. 464–468. [Online]. Available: <https://aclanthology.org/N18-2074>
- [14] X. Wang, Z. Tu, L. Wang, and S. Shi, “Self-attention with structural position representations,” in *Proceedings of EMNLP-IJCNLP*, 2019, pp. 1403–1409. [Online]. Available: <https://aclanthology.org/D19-1145>
- [15] Z. Yang, Z. Dai, Y. Yang, and et al., “Xlnet: Generalized autoregressive pretraining for language understanding,” in *Proceedings of NeurIPS*, vol. 32, 2019, pp. 1–11.
- [16] P.-C. Chen, H. Tsai, S. Bhojanapalli, and et al., “A simple and effective positional encoding for transformers,” in *Proceedings of EMNLP*, 2021, pp. 2974–2988. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.236>
- [17] I. Beltagy, M. E. Peters, and A. Cohan, “Longformer: The long-document transformer,” *arXiv:2004.05150*, 2020.
- [18] C.-Z. A. Huang, A. Vaswani, J. Uszkoreit, and et al, “Music transformer: Generating music with long-term structure,” in *ICLR 2019*, 2019, pp. 1–15.
- [19] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. Le, and R. Salakhutdinov, “Transformer-XL: Attentive language models beyond a fixed-length context,” in *Proceedings of the ACL 2019*. Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 2978–2988. [Online]. Available: <https://aclanthology.org/P19-1285>
- [20] M. Cettolo, C. Girardi, and M. Fed-erico, “Wit3: Web inventory of transcribed and translated talks,” in *Proceedings of EAMT*, 2012, pp. 261–268.
- [21] L. Miculicich, D. Ram, N. Pappas, and J. Henderson, “Document-level neural machine translation with hierarchical attention networks,” in *Proceedings of EMNLP*, pp. 2947–2954. [Online]. Available: <https://aclanthology.org/D18-1325>
- [22] Z. Yang, J. Zhang, F. Meng, S. Gu, Y. Feng, and J. Zhou, “Enhancing context modeling with a query-guided capsule network for document-level translation,” in *Proceedings of EMNLP-IJCNLP*, 2019, pp. 1527–1537.
- [23] Z. Zheng, X. Yue, S. Huang, J. Chen, and A. Birch, “Towards making the most of context in neural machine translation,” in *Proceedings of IJCAI*, 2020, pp. 3983–3989. [Online]. Available: <https://doi.org/10.24963/ijcai.2020/551>
- [24] R. Sennrich, B. Haddow, and A. Birch, “Neural machine translation of rare words with subword units,” in *Proceedings of ACL*, 2016, pp. 1715–1725. [Online]. Available: <https://www.aclweb.org/anthology/P16-1162>
- [25] M. Ott, S. Edunov, D. Grangier, and M. Auli, “Scaling neural machine translation,” in *Proceedings of the 3th WMT*, 2018, pp. 1–9. [Online]. Available: <https://www.aclweb.org/anthology/W18-6301>
- [26] S. Jean, S. Lauly, O. Firat, and K. Cho, “Does neural machine translation benefit from larger context?” *Computing Research Repository*, vol. arXiv:1704.05135, 2017.
- [27] L. Wang, Z. Tu, A. Way, and Q. Liu, “Exploiting cross-sentence context for neural machine translation,” in *Proceedings of EMNLP*, 2017, pp.

- 2826–2831.
- [28] R. Bawden, R. Sennrich, A. Birch, and B. Haddow, “Evaluating discourse phenomena in neural machine translation,” in *Proceedings of NAACL-HLT*, 2018, pp. 1304–1313.
- [29] E. Voita, P. Serdyukov, R. Sennrich, and I. Titov, “Context-aware neural machine translation learns anaphora resolution,” in *Proceedings of ACL*, 2018, pp. 1264–1274.
- [30] K. Wong, S. Maruf, and G. Haffari, “Contextual neural machine translation improves translation of cataphoric pronouns,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, 2020, pp. 5971–5978.
- [31] H. Yun, Y. Hwang, and K. Jung, “Improving context-aware neural machine translation using self-attentive sentence embedding,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020, pp. 9498–9506.
- [32] B. Li, H. Liu, Z. Wang, Y. Jiang, T. Xiao, J. Zhu, T. Liu, and C. Li, “Does multi-encoder help? a case study on context-aware neural machine translation,” in *Proceedings of ACL*, 2020, p. 3512–3518.
- [33] X. Kang, Y. Zhao, J. Zhang, and C. Zong, “Dynamic context selection for document-level neural machine translation via reinforcement learning,” in *Proceedings of EMNLP*, 2020, pp. 2242–2254.
- [34] L. Zhang, T. Zhang, H. Zhang, B. Yang, W. Ye, and S. Zhang, “Multi-hop transformer for document-level machine translation,” in *Proceedings of NAACL*, Online, 2021, pp. 3953–3963.
- [35] S. Maruf and G. Haffari, “Document context neural machine translation with memory networks,” in *Proceedings of ACL*, 2018, pp. 1275–1284. [Online]. Available: <https://aclanthology.org/P18-1118>
- [36] X. Tan, L. Zhang, D. Xiong, and G. Zhou, “Hierarchical modeling of global context for document-level neural machine translation,” in *Proceedings of the EMNLP-IJCNLP*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 1576–1585. [Online]. Available: <https://aclanthology.org/D19-1168>
- [37] X. Tan, L. Zhang, and G. Zhou, “Coupling context modeling with zero pronoun recovering for document-level natural language generation,” in *Proceedings of the ACL 2021*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 2530–2540. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.197>
- [38] X. Lyu, J. Li, Z. Gong, and M. Zhang, “Encouraging lexical translation consistency for document-level neural machine translation,” in *Proceedings of EMNLP*, 2021, p. 3265–3277.
- [39] M. Xu, L. Li, D. F. Wong, Q. Liu, and L. S. Chao, “Document graph for neural machine translation,” in *Proceedings of the EMNLP*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, 2021, pp. 8435–8448.
- [40] P. Zhang, B. Chen, N. Ge, and et al., “Long-short term masking transformer: A simple but effective baseline for document-level neural machine translation,” in *Proceedings of EMNLP*, 2020, pp. 1081–1087. [Online]. Available: <https://aclanthology.org/2020.emnlp-main.81>
- [41] L. Lupo, M. Dinarelli, and L. Besacier, “Focused concatenation for context-aware neural machine translation,” in *expriarXiv: 2210.13388v1*, 2022, pp. 1–14.
- [42] G. Ke, D. He, and T.-Y. Liu, “Rethinking positional encoding in language pre-training,” in *Proceedings of ICLR*, May 2021.
- [43] A. Qu, J. Niu, and S. Mo, “Explore better relative position embeddings from encoding perspective for transformer models,” in *Proceedings of the EMNLP*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 2989–2997. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.237>
- [44] J. Ainslie, S. Ontanon, C. Alberti, and et al., “ETC: Encoding long and structured inputs in transformers,” in *Proceedings of EMNLP*, 2020, pp. 268–284. [Online]. Available: <https://aclanthology.org/2020.emnlp-main.19>
- [45] Z. Huang, D. Liang, P. Xu, and B. Xiang, “Improve transformer models with better relative position embeddings,” in *Findings of EMNLP 2020*. Online: Association for Computational Linguistics, Nov. 2020, pp. 3327–3335. [Online]. Available: <https://aclanthology.org/2020.findings-emnlp.298>
- [46] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *Proceedings of NAACL*, 2019, pp. 4171–4186. [Online]. Available: <https://aclanthology.org/N19-1423>
- [47] M. Lewis, Y. Liu, N. Goyal, and et al., “Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension,” in *Proceedings of ACL*, 2020, p. 7871–7880. [Online]. Available: <https://aclanthology.org/2020.acl-main.703/>
- [48] A. Bajaj, P. Dangati, K. Krishna, and et al., “Long document summarization in a low resource setting using pretrained language models,” in *Proceedings of ACL*, 2021, pp. 71–80.
- [49] R. Sun, H. Jin, and X. Wan, “Document-level text simplification: Dataset, criteria and baseline,” in *Proceedings of EMNLP 2021*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 7997–8013. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.630>
- [50] Y. Liang, F. Meng, Y. Chen, J. Xu, and J. Zhou, “Modeling bilingual conversational characteristics for neural chat translation,” in *Proceedings of the ACL 2021 (Volume 1: Long Papers)*. Online: Association for Computational Linguistics, Aug. 2021, pp. 5711–5724. [Online]. Available: <https://aclanthology.org/2021.acl-long.444>