

An Efficient Transformer Decoder with Compressed Sub-layers

Yanyang Li^{1*}, Ye Lin^{1*}, Tong Xiao^{1,2†}, Jingbo Zhu^{1,2}

¹NLP Lab, School of Computer Science and Engineering, Northeastern University, Shenyang, China

²NiuTrans Research, Shenyang, China

{blamedrlee, linye2015}@outlook.com, {xiaotong, zhujingbo}@mail.neu.edu.cn

Abstract

The large attention-based encoder-decoder network (Transformer) has become prevailing recently due to its effectiveness. But the high computation complexity of its decoder raises the inefficiency issue. By examining the mathematic formulation of the decoder, we show that under some mild conditions, the architecture could be simplified by compressing its sub-layers, the basic building block of Transformer, and achieves a higher parallelism. We thereby propose *Compressed Attention Network*, whose decoder layer consists of only one sub-layer instead of three. Extensive experiments on 14 WMT machine translation tasks show that our model is $1.42\times$ faster with performance on par with a strong baseline. This strong baseline is already $2\times$ faster than the widely used standard baseline without loss in performance.

Introduction

Transformer is an attention-based encoder-decoder model (Vaswani et al. 2017). It has shown promising results in machine translation tasks recently (Wang et al. 2019; Aharoni, Johnson, and Firat 2019; Dehghani et al. 2019). Nonetheless, Transformer suffers from the inefficiency issue at inference. This problem is attributed to the Transformer decoder for two reasons: 1) the decoder is deep (Kasai et al. 2020). It consists of multiple layers and each layer contains three sub-layers, including two attentions and a feed-forward network; 2) the attention has a high (quadratic time) complexity (Zhang, Xiong, and Su 2018), as it needs to compute the correlation between any two input words.

Previous work has focused on improving the complexity of the attention in the decoder to accelerate the inference. For example, AAN uses the averaging operation to avoid computing the correlation between input words (Zhang, Xiong, and Su 2018). SAN share the attention results among layers (Xiao et al. 2019). On the other hand, we learn that vanilla attention runs faster in training than in inference thanks to its parallelism. This offers a new direction: a higher degree of parallelism could speed up the inference. The most representative work of this type is the non-autoregressive approach (Gu et al. 2018). Its decoder

predicts all words in parallel, but fails to model the word dependencies. Despite of their successes, all these systems still have a deep decoder.

In this work, we propose to parallelize the sub-layers to obtain a shallow autoregressive decoder. This way does not suffer from the poor result of directly reducing depths and avoids the limitation of non-autoregressive approaches. We prove that the two attention sub-layers in a decoder layer could be parallelized if we assume their inputs are close to each other. This assumption holds and thereby we compress these two attentions into one. Furthermore, we show that the remaining feed-forward network could also be merged into the attention due to their linearity. To the end, we propose *Compressed Attention Network* (CAN for short). The decoder layer of CAN possesses a single attention sub-layer that does the previous three sub-layers' jobs in parallel. As another "bonus", CAN is simple and easy to be implemented.

In addition, Kasai et al. (2020) empirically discover that existing systems are not well balancing the encoder and decoder depths. Based on their work, we build a system with a deep encoder and a shallow decoder, which is $2\times$ faster than the widely used standard baseline without loss in performance. It requires neither the architecture modification nor adding extra parameters. This system serves as a stronger baseline for a more convincing comparison.

We evaluate CAN and the stronger baseline in 14 machine translation tasks, including WMT14 English \leftrightarrow {German, French} (En \leftrightarrow {De, Fr}) and WMT17 English \leftrightarrow {German, Finnish, Latvian, Russian, Czech} (En \leftrightarrow {De, Fi, Lv, Ru, Cs}). The experiments show that CAN is up to $2.82\times$ faster than the standard baseline with almost no loss in performance. Even comparing to our stronger baseline, CAN still has a $1.42\times$ speed-up, while other acceleration techniques such as SAN and AAN are $1.12\sim 1.16\times$ in the same case.

To summarize, our contributions are as follows:

- We propose CAN, a novel architecture that accelerates Transformer by compressing its sub-layers for a higher degree of parallelism. CAN is easy to be implemented.
- Our work is based on a stronger baseline, which is $2\times$ faster than the widely used standard baseline.
- The extensive experiments on 14 WMT machine translation tasks show that CAN is $1.42\times$ faster than the stronger

*Authors contributed equally.

†Corresponding author.

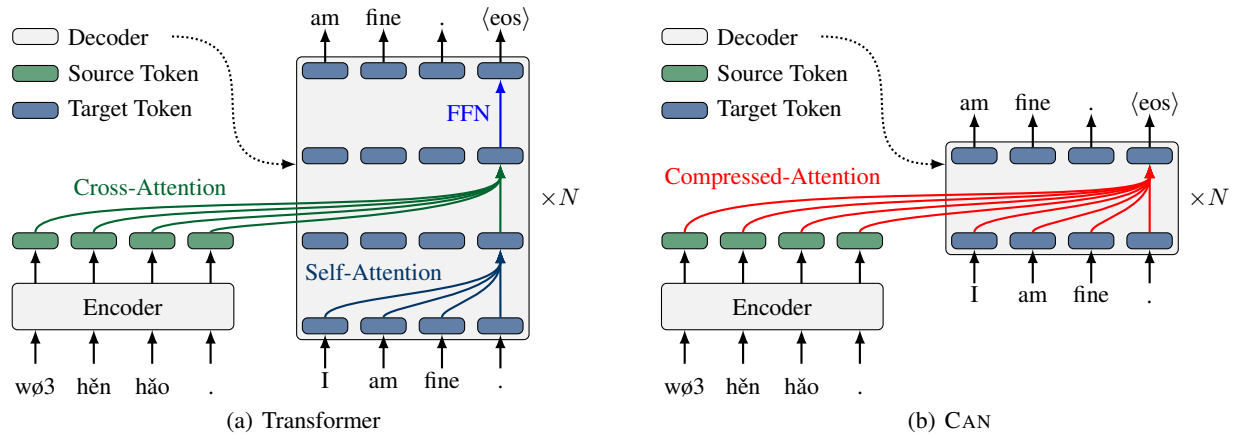


Figure 1: Transformer vs. CAN (Chinese pinyin→English: “w03 h3n h3o .” → “I am fine .”).

baseline and $2.82\times$ for the standard baseline. CAN also outperforms other approaches such as SAN and AAN.

Background: Transformer

Transformer is one of the state-of-the-art neural models in machine translation. It consists of a N -layer encoder and a N -layer decoder, where $N = 6$ in most cases. The encoder maps the source sentence to a sequence of continuous representations and the decoder maps these representations to the target sentence. All layers in the encoder or decoder are identical to each other.

The layer in the decoder consists of three sub-layers, including the self-attention, the cross-attention and the feed-forward network (FFN). The self-attention takes the output X of the previous sub-layer as its input and produces a tensor with the same size as its output. It computes the attention distribution A_x and then averages X by A_x . We denote the self-attention as $Y_x = \text{Self}(X)$, where $X \in \mathbb{R}^{t \times d}$, t is the target sentence length and d is the dimension of the hidden representation:

$$A_x = \text{SoftMax}\left(\frac{XW_{q1}W_{k1}^TX^T}{\sqrt{d}}\right) \quad (1)$$

$$Y_x = A_xXW_{v1} \quad (2)$$

where $W_{q1}, W_{k1}, W_{v1} \in \mathbb{R}^{d \times d}$.

The cross-attention is similar to the self-attention, except that it takes the encoder output H as an additional input. We denote the cross-attention as $Y_h = \text{Cross}(X, H)$, where $H \in \mathbb{R}^{s \times d}$, s is the source sentence length:

$$A_h = \text{SoftMax}\left(\frac{XW_{q2}W_{k2}^TH^T}{\sqrt{d}}\right) \quad (3)$$

$$Y_h = A_hHW_{v2} \quad (4)$$

where $W_{q2}, W_{k2}, W_{v2} \in \mathbb{R}^{d \times d}$.

The FFN applies non-linear transformation to its input X . We denote the FFN as $Y_f = \text{FFN}(X)$:

$$Y_f = \text{ReLU}(XW_1 + b_1)W_2 + b_2 \quad (5)$$

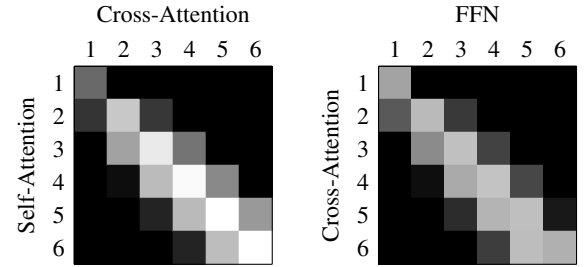


Figure 2: The cosine similarity of inputs for every two adjacent sub-layers on WMT14 En-De translation task (a dark cell means the inputs are dissimilar).

where $W_1 \in \mathbb{R}^{d \times 4d}$, $b_1 \in \mathbb{R}^{4d}$, $W_2 \in \mathbb{R}^{4d \times d}$ and $b_2 \in \mathbb{R}^d$.

All sub-layers are coupled with the residual connection (He et al. 2016a), i.e., $Y = f(X) + X$ where f could be any sub-layer. Their inputs are also preprocessed by the layer normalization first (Ba, Kiros, and Hinton 2016). Fig. 1(a) shows the architecture of Transformer decoder. For more details, we refer the reader to Vaswani et al. (2017).

Compressed Attention Network

Compressing Self-Attention and Cross-Attention

As suggested by Huang et al. (2016), the output of one layer in the residual network can be decomposed into the sum of all outputs from previous layers. For the adjacent self-attention and cross-attention, we can write their final output as $Y = X + \text{Self}(X) + \text{Cross}(X', H)$, where X is the input of self-attention and $X' = X + \text{Self}(X)$ is the input of cross-attention. If X and X' are identical, we are able to accelerate the computation of Y by parallelizing these two attentions, as X' do not need to wait $\text{Self}(X)$ to finish.

Previous work (He et al. 2016b) has shown that inputs of adjacent layers are similar. This implies that X and X' are close and the parallelization is possible. We empirically verify this in the left part of Fig. 2 by examining the cosine similarity between inputs of every self-attention and cross-attention pairs. It shows that X and X' are indeed close to

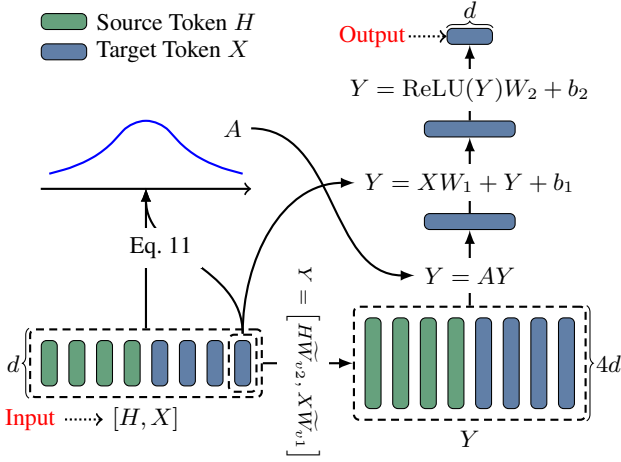


Figure 3: Compressed-Attention.

each other (a high similarity > 0.9 for the diagonal entries). Therefore we could assume X and X' are identical (we omit the layer normalization for simplicity):

$$Y = X + \text{Self}(X) + \text{Cross}(X, H) \quad (6)$$

By observing that Eq. 2 and Eq. 4 are essentially matrix multiplications, we could rewrite $\text{Self}(X) + \text{Cross}(X, H)$ as a single matrix multiplication:

$$A = [A_x^T, A_h^T]^T \quad (7)$$

$$\text{Self}(X) + \text{Cross}(X, H) = A [XW_{v1}, HW_{v2}] \quad (8)$$

$[\cdot]$ is the concatenation operation along the first dimension.

Xiao et al. (2019) shows that some attention distributions A_x and A_h are duplicate. This means that there exists a certain redundancy in $\{W_{q1}, W_{k1}\}$ and $\{W_{q2}, W_{k2}\}$. Thus we could safely share W_{q1} and W_{q2} to parallelize the computation of the attention distribution A :

$$\bar{A} = (XW_q [XW_{k1}, HW_{k2}]^T) / \sqrt{d} \quad (9)$$

$$A = [\text{SoftMax}(\bar{A}_{:,1:t}^T), \text{SoftMax}(\bar{A}_{:,t+1:t+s}^T)]^T \quad (10)$$

However, A consists of two SoftMax distributions and is used in Eq. 8 without normalization. The output variance is then doubled and leads to poor optimization (Glorot and Bengio 2010). It is advised to divide A by $\sqrt{2}$ to preserve the variance. This way resembles a single distribution. So we use one SoftMax instead and this works well:

$$A = \text{SoftMax}\left(\frac{XW_q [XW_{k1}, HW_{k2}]^T}{\sqrt{d}}\right) \quad (11)$$

Now, we can compute Y in Eq. 6 efficiently by using Eq. 11 as well as Eq. 8 to compute $\text{Self}(X) + \text{Cross}(X, H)$.

Compressing Attention and FFN

It is natural to consider to merge the attention and FFN with the same approach for further speed-up. As suggested by the right part of Fig. 2, the similarities between inputs of

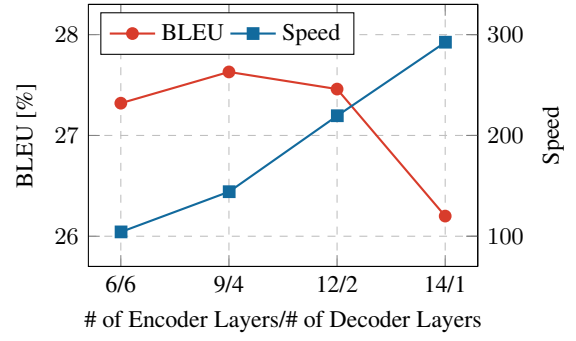


Figure 4: Performance (BLEU) and translation speed (token/sec) vs. the numbers of encoder and decoder layers on WMT14 En-De translation task.

the adjacent cross-attention and FFN are low (dark diagonal entries). This implies that it is not ideal to make the identical input assumption to parallelize the cross-attention and FFN.

Here we provide another solution. Given that attention is merely a weighted sum and FFN performs a linear projection first, we can merge them by exploiting the linearity. This way not only parallelizes the computation of attention and FFN but also removes redundant matrix multiplications.

We substitute X in Eq. 5 by Y in Eq. 6:

$$Y_f = \text{ReLU}(XW_1 + A [XW_{v1}, HW_{v2}] W_1 + b_1) W_2 + b_2 \quad (12)$$

We can combine W_1 with W_{v1} as well as W_{v2} into $\widetilde{W}_{v1}, \widetilde{W}_{v2} \in \mathbb{R}^{d \times 4d}$, as these matrices are learnable and matrix multiplied together:

$$Y_f = \text{ReLU}(XW_1 + A [X\widetilde{W}_{v1}, H\widetilde{W}_{v2}] + b_1) W_2 + b_2 \quad (13)$$

Furthermore, XW_1 can be computed in parallel with other transformations such as XW_q .

This eventually gives us an more efficient decoder layer architecture, named *Compressed-Attention*. The whole computation process is shown in Fig. 3: it first computes the attention distribution A by Eq. 11, then performs the attention operation via Eq. 13, and produces Y_f as the final result. The proposed *Compressed Attention Network* (CAN) stacks compressed-attentions to form its decoder. Fig. 1 shows the difference between Transformer and CAN.

Balancing Encoder and Decoder Depths

Based on the findings of Kasai et al. (2020), we learn that a shallow decoder could offer a great speed gain, while a deep encoder could make up of the loss of a shallow decoder without adding a heavy computation overhead. Since their work is based on knowledge distillation (Hinton, Vinyals, and Dean 2015), here we re-examine this idea under the standard training setting (without knowledge distillation).

Fig. 4 shows the performance and speed if we gradually reduce the decoder depth while adding more encoder layers. We see that although the overall number of parameters remains the same, the baseline can be $2 \times$ faster without losing

| Source | Lang. | Train | | Valid | | Test | |
|--------|-------|-------|------|-------|------|-------|------|
| | | sent. | word | sent. | word | sent. | word |
| WMT14 | En↔De | 4.5M | 220M | 3000 | 110K | 3003 | 114K |
| | En↔Fr | 35M | 2.2B | 26K | 1.7M | 3003 | 155K |
| WMT17 | En↔De | 5.9M | 276M | 8171 | 356K | 3004 | 128K |
| | En↔Fi | 2.6M | 108M | 8870 | 330K | 3002 | 110K |
| | En↔Lv | 4.5M | 115M | 2003 | 90K | 2001 | 88K |
| | En↔Ru | 25M | 1.2B | 8819 | 391K | 3001 | 132K |
| | En↔Cs | 52M | 1.2B | 8658 | 354K | 3005 | 118K |

Table 1: Data statistics (# of sentences and # of words).

any performance (12/2 vs. 6/6). This justifies the previous idea. We thereby choose a stronger baseline with a 12-layer encoder and a 2-layer decoder for a more convincing comparison. This setting is also applied to CAN.

Experiments

Experimental Setup

Datasets We evaluate our methods on 14 machine translation tasks (7 datasets \times 2 translation directions each), including WMT14 En↔{De, Fr} and WMT17 En↔{De, Fi, Lv, Ru, Cs}.

WMT14 En↔{De, Fr} datasets are tokenized by a script from Moses¹. We apply BPE (Sennrich, Haddow, and Birch 2016) with 32K merge operations to segment words into subword units. Sentences with more than 250 subword units are removed. The first two rows of Table 1 are the detailed statistics of these two datasets. For En-De, we share the source and target vocabularies. We choose *newstest-2013* as the validation set and *newstest-2014* as the test set. For En-Fr, we validate the system on the combination of *newstest-2012* and *newstest-2013*, and test it on *newstest-2014*.

All WMT17 datasets are the official preprocessed version from WMT17 website². BPE with 32K merge operations is similarly applied to these datasets. We use the concatenation of all available preprocessed validation sets in WMT17 datasets as our validation set:

- En↔De. We use the concatenation of *newstest2014*, *newstest2015* and *newstest2016* as the validation set.
- En↔Fi. We use the concatenation of *newstest2015*, *newsdev2015*, *newstest2016* and *newstestB2016* as the validation set.
- En↔Lv. We use *newsdev2016* as the validation set.
- En↔Ru. We use the concatenation of *newstest2014*, *newstest2015* and *newstest2016* as the validation set.
- En↔Cs. We use the concatenation of *newstest2014*, *newstest2015* and *newstest2016* as the validation set.

We use *newstest2017* as the test set for all WMT17 datasets. Detailed statistics of these datasets are shown in Table 1. For all 14 translation tasks, we report case-sensitive tokenized BLEU scores³.

¹<https://github.com/moses-smt/mosesdecoder/blob/master/scripts/tokenizer/tokenizer.perl>

²<http://data.statmt.org/wmt17/translation-task/preprocessed/>

³<https://github.com/moses-smt/mosesdecoder/blob/master/scripts/generate/mosesdecoder.perl>

| | System | Test | Δ_{BLEU} | Valid | Speed | Δ_{Speed} |
|-------|----------|-------|-----------------|-------|--------|------------------|
| En-De | Baseline | 27.32 | - | 26.56 | 104.27 | - |
| | Balanced | 27.46 | 0.00 | 26.81 | 219.53 | 0.00% |
| | SAN | 26.91 | -0.55 | 26.04 | 229.89 | +4.72% |
| | AAN | 27.36 | -0.10 | 26.11 | 233.58 | +6.40% |
| De-En | CAN | 27.32 | -0.14 | 26.47 | 290.08 | +32.14% |
| | Baseline | 30.50 | - | 30.34 | 103.97 | - |
| | Balanced | 30.76 | 0.00 | 30.37 | 206.00 | 0.00% |
| | SAN | 30.09 | -0.67 | 30.11 | 240.52 | +16.76% |
| En-Fr | AAN | 30.15 | -0.61 | 30.07 | 232.08 | +12.66% |
| | CAN | 30.37 | -0.39 | 30.17 | 293.16 | +42.31% |
| | Baseline | 40.82 | - | 46.80 | 104.65 | - |
| | Balanced | 40.55 | 0.00 | 46.87 | 206.54 | 0.00% |
| Fr-En | SAN | 40.45 | -0.10 | 46.69 | 208.68 | +1.04% |
| | AAN | 40.50 | -0.05 | 46.57 | 210.29 | +1.82% |
| | CAN | 40.25 | -0.30 | 46.56 | 263.83 | +27.74% |
| | Baseline | 36.33 | - | 47.03 | 105.85 | - |
| Fr-En | Balanced | 36.86 | 0.00 | 46.89 | 201.13 | 0.00% |
| | SAN | 36.73 | -0.13 | 46.82 | 213.30 | +6.05% |
| | AAN | 36.52 | -0.34 | 46.74 | 215.97 | +7.38% |
| | CAN | 36.67 | -0.19 | 46.63 | 266.63 | +32.57% |

Table 2: Comparison of BLEU scores [%] and translation speeds (token/sec) of different attention models on WMT14 En↔{De, Fr} translation tasks.

Model Setup Our baseline system is based on the open-source implementation of the Transformer model presented in Ott et al. (2019). For all machine translation tasks, the standard Transformer baseline (Baseline) consists of a 6-layer encoder and a 6-layer decoder. The embedding size is set to 512. The number of attention heads is 8. The FFN hidden size equals to $4 \times$ embedding size. Dropout with the value of 0.1 is used for regularization. We adopt the inverse square root learning rate schedule with 8,000 warmup steps and 0.0007 learning rate. We stop training until the model stops improving on the validation set. All systems are trained on 8 NVIDIA TITIAN V GPUs with mixed-precision training (Mickevicus et al. 2018) and a batch size of 4,096 tokens per GPU. We average model parameters in the last 5 epochs for better performance. At test time, the model is decoded with a beam of width 4 and half-precision. For an accurate speed comparison, we decode with a batch size of 1 to avoid paddings. The stronger balanced baseline (Balanced) shares the setting with this standard baseline, except that its encoder depth is 12 and decoder depth is 2.

We compare CAN and other model acceleration approaches with our baselines. We choose Sharing Attention Network (SAN) (Xiao et al. 2019) and Average Attention Network (AAN) (Zhang, Xiong, and Su 2018) for comparison, as they have been proven to be effective in various machine translation tasks (Birch et al. 2018). All hyperparameters of CAN, SAN and AAN are identical to the balanced baseline system. Results are the average of 3 runs.

Results

Table 2 shows the results of various systems on WMT14 En↔{De, Fr}. Our balanced baseline has nearly the same

| | System | Test | Δ_{BLEU} | Valid | Speed | Δ_{Speed} |
|-------|----------|-------|-----------------|-------|--------|------------------|
| En-De | Baseline | 28.40 | - | 31.30 | 106.58 | - |
| | Balanced | 28.65 | 0.00 | 31.39 | 218.35 | 0.00% |
| | CAN | 28.30 | -0.35 | 30.94 | 280.57 | +28.50% |
| De-En | Baseline | 34.48 | - | 35.36 | 103.04 | - |
| | Balanced | 34.38 | 0.00 | 35.16 | 220.05 | 0.00% |
| | CAN | 33.99 | -0.39 | 34.82 | 286.23 | +30.07% |
| En-Fi | Baseline | 21.28 | - | 18.31 | 103.84 | - |
| | Balanced | 21.38 | 0.00 | 18.67 | 207.73 | 0.00% |
| | CAN | 21.14 | -0.24 | 18.19 | 286.36 | +37.85% |
| Fi-En | Baseline | 25.54 | - | 21.32 | 106.59 | - |
| | Balanced | 25.63 | 0.00 | 21.29 | 209.88 | 0.00% |
| | CAN | 25.25 | -0.38 | 21.31 | 287.57 | +37.02% |
| En-Lv | Baseline | 16.14 | - | 21.33 | 107.20 | - |
| | Balanced | 15.98 | 0.00 | 21.21 | 219.02 | 0.00% |
| | CAN | 15.90 | -0.08 | 20.75 | 287.33 | +31.19% |
| Lv-En | Baseline | 18.74 | - | 24.79 | 106.25 | - |
| | Balanced | 18.69 | 0.00 | 24.54 | 216.06 | 0.00% |
| | CAN | 18.21 | -0.48 | 24.16 | 275.89 | +27.69% |
| En-Ru | Baseline | 30.44 | - | 30.67 | 106.46 | - |
| | Balanced | 30.28 | 0.00 | 30.59 | 214.52 | 0.00% |
| | CAN | 29.89 | -0.39 | 30.28 | 287.13 | +33.85% |
| Ru-En | Baseline | 34.44 | - | 32.39 | 107.24 | - |
| | Balanced | 34.24 | 0.00 | 32.22 | 213.78 | 0.00% |
| | CAN | 33.95 | -0.29 | 31.92 | 287.86 | +34.65% |
| En-Cs | Baseline | 24.00 | - | 28.09 | 106.18 | - |
| | Balanced | 23.69 | 0.00 | 28.03 | 212.65 | 0.00% |
| | CAN | 23.59 | -0.10 | 27.71 | 272.37 | +28.08% |
| Cs-En | Baseline | 30.00 | - | 33.01 | 104.00 | - |
| | Balanced | 30.06 | 0.00 | 32.86 | 202.96 | 0.00% |
| | CAN | 29.87 | -0.19 | 32.99 | 269.70 | +32.88% |

Table 3: BLEU scores [%] and translation speeds (token/sec) on WMT17 En \leftrightarrow {De, Fi, Lv, Ru, Cs} translation tasks.

performance as the standard baseline, but its speed is $2\times$ faster on average. A similar phenomenon is also observed from WMT17 experiments in Table 3. This observation indicates that existing systems do not well balance the encoder and decoder depths. We also report the performance of AAN, SAN and the proposed CAN. All three approaches have similar BLEU scores and slightly underperform the balanced baseline. CAN is more stable than the others, as its maximum Δ_{BLEU} is -0.39, while SAN is -0.67 and AAN is -0.61. For speeds of these systems, SAN and AAN have a similar level of acceleration (1~16%) over the balanced baseline. CAN, on the other hand, provides a higher level of acceleration (27~42%). Interestingly, we find that the acceleration is more obvious in De-En than in others, e.g., 42% in De-En and 27% in En-Fi for CAN. We find that the length ratio between the translation and the source sentence in De-En is higher than others, e.g., 1.0 for De-En and 0.981 for En-Fi. In this case the decoder tends to predict more words and consumes more time in De-En, and thus acceleration approaches that work on the decoder are more effective.

More experimental results to justify the effectiveness of CAN are presented in Table 3. We evaluate the balanced

| System | Before KD | | After KD | |
|----------|-----------|-----------------|----------|-----------------|
| | Test | Δ_{BLEU} | Test | Δ_{BLEU} |
| Balanced | 27.46 | 0.00 | 27.82 | 0.00 |
| SAN | 26.91 | -0.55 | 27.76 | -0.06 |
| AAN | 27.36 | -0.10 | 27.85 | +0.03 |
| CAN | 27.32 | -0.14 | 28.08 | +0.26 |

Table 4: BLEU scores [%] of applying knowledge distillation (KD) on WMT14 En-De translation task.

| System | Test | Δ_{BLEU} | Speed | Δ_{Speed} |
|----------------------|-------|-----------------|--------|------------------|
| Balanced | 27.46 | 0.00 | 219.53 | 0.00% |
| + Compress Attention | 27.09 | -0.37 | 263.64 | +20.09% |
| + Compress FFN | 27.69 | +0.23 | 233.17 | +6.21% |
| + Compress All | 27.32 | -0.14 | 290.08 | +32.14% |

Table 5: Ablation study on WMT14 En-De translation task (Compress Attention: compress the self-attention and cross-attention only; Compress FFN: compress the cross-attention and FFN only; Compress All: compress the self-attention, cross-attention and FFN).

baseline as well as CAN on five WMT17 language pairs. The results again show that the balanced baseline is indeed a strong baseline with BLEU scores close to the standard baseline and is consistently $2\times$ faster. CAN also shows a similar trend that it slightly underperforms the balanced baseline (< 0.5 BLEU scores) but is $> 27\%$ faster.

Analysis

Knowledge Distillation

Although SAN, AAN and CAN offer considerable speed gain over the balanced baseline, they all suffer from the performance degradation as shown in Table 2 and Table 3. The popular solution to this is knowledge distillation (KD). Here we choose sequence-level knowledge distillation (Kim and Rush 2016) for better performance in machine translation tasks. The balanced baseline is used to generate the pseudo data for KD.

Table 4 shows that KD closes the performance gap between the fast attention models (SAN, AAN and CAN) and the balanced baseline. This fact suggests that all three systems have enough capacity for a good performance, but training from scratch is not able to reach a good convergence state. It suggests that these systems might require a more careful hyper-parameters tuning or a better optimization method.

Ablation Study

To investigate in which part CAN contributes the most to the acceleration as well as the performance loss, we only compress the self-attention and cross-attention or compress the cross-attention and FFN for study. Table 5 shows the results of this ablation study. We can see that compressing the two attentions provides a 20.09% speed-up, while only 6.21% for compressing attention and FFN. This is because FFN is already highly parallelized and accelerating itself does not

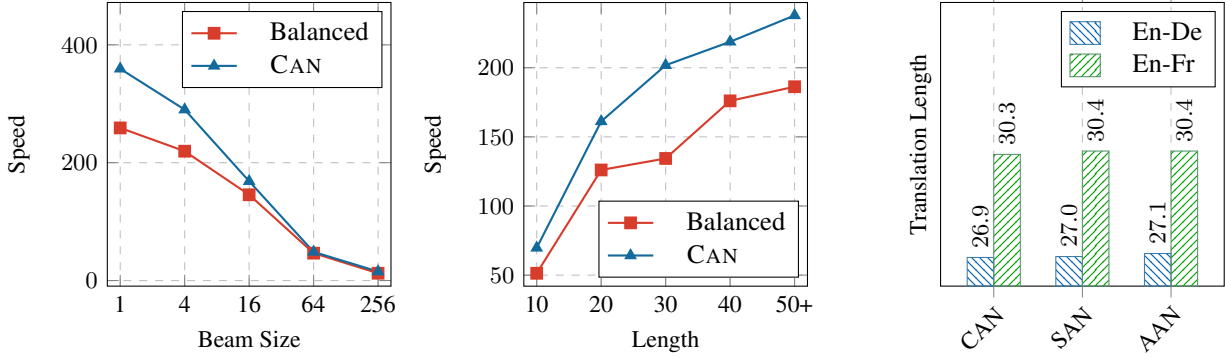


Figure 5: Translation speed (token/sec) vs. beam size and translation length on WMT14 En-De translation task.

bring much gain. On the other hand, compressing attentions brings the most performance loss, which shows that the identical input assumption is strong. Fig. 2 shows that inputs of the adjacent layers are not very similar in lower layers. Therefore using CAN in low layers might bring a great loss. We also find that compressing attention and FFN has an even better result. This might be that we remove the redundant parameters in the model.

Sensitivity Analysis

We study how the speed could be affected by other factors in Fig. 5, e.g., the beam size and the translation length. The left of Fig. 5 shows that CAN is consistently faster than the balanced baseline with different beam size. As the acceleration provided by CAN is constantly proportional to the speed of the baseline, it becomes less obvious when the baseline is slow, i.e., translating with a large beam. An opposite trend happens in the middle of Fig. 5 for the translation length. This is because overheads such as data preparation dominate the translation time of short sentences. This way results in a slow speed even when the translation time is short. As both the baseline and CAN get faster when generating longer translations, one might suspect that the superior acceleration of CAN over other approaches comes from the fact that CAN generates longer translations. Further analysis is conducted and shown in the right of Fig. 5. We see that CAN, SAN and AAN generate translations with similar lengths in two WMT14 translation tasks. This observation justifies that the superior acceleration brought by CAN did come from its design rather than translation lengths.

Error Analysis

As shown in Table 2 and Table 3, the acceleration of CAN comes at the cost of performance. Here we conduct experiments to better understand in which aspect CAN sacrifices for speed-up. We first evaluate the sentence-level BLEU score for each translation, then cluster these translations according to their averaged word frequencies or lengths.

Fig. 6 shows the results. The left of Fig. 6 indicates that CAN did well on sentences with low frequencies, but not on those with high frequencies. The right of Fig. 6 shows that CAN does not translate short sentences well but is quite

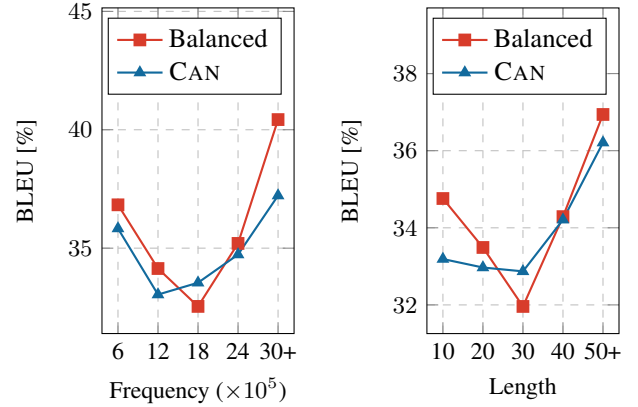


Figure 6: BLEU score [%] vs. word frequency ($\times 10^5$) and translation length on WMT14 En-De translation task.

good at translating long sentences. These facts are counter-intuitive as one might expect a poor model could do well on easy samples (high frequency and short sentences) but not on hard ones (low frequency and long sentences). This might due to the identical input assumptions we used to derive CAN are critical to easy samples. We left this for the future exploration.

Parallelism Study

A simple approach to obtain a higher parallelism without modifying the architecture is to increase the batch size at inference. Fig. 7 compares the inference time of the balanced baseline and CAN by varying the batch size. We can see that both systems run faster with a larger batch size and CAN is consistently faster than the balanced baseline. But the acceleration of CAN over the baseline Δ_{Speed} diminishes when the batch size gets larger. In this case we observe that CAN reaches the highest parallelism (a nearly 100% GPU utility) in a smaller batch size (≥ 32) than the baseline (> 64). This means that enlarging the batch size no longer provides acceleration for CAN, while the baseline can still be further speeded up. We expect CAN could be faster if more tensor cores are available in the future.

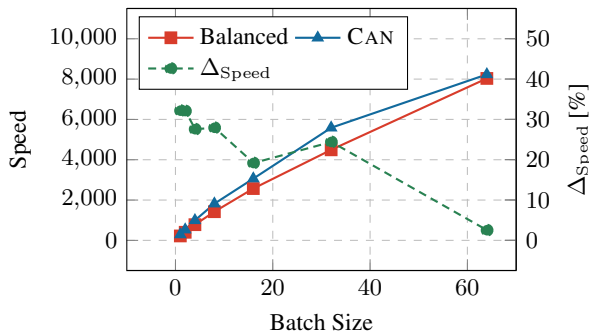


Figure 7: Speed (token/sec) and Δ_{Speed} [%] vs. batch size on WMT14 En-De translation task.

Training Study

We plot the training and validation loss curve of the standard baseline, the balanced baseline and CAN in Fig. 8 for studying their convergence. We can see that all systems converge stably. The balanced baseline has a higher loss than the standard baseline in both the training and validation sets, but their BLEU scores are close as shown in Table 2. This is due to the shallow decoder in the balanced baseline. Since the loss is determined by the decoder, a shallow decoder with less capacity would have a higher loss. Wang et al. (2019) indicates that the encoder depth has a greater impact than the decoder on BLEU scores, therefore the deep encoder makes up the performance loss of the shallow decoder. We also see that CAN has a higher loss than the balanced baseline because we compress the decoder. Since we do not enhance the encoder, the BLEU score drops accordingly.

Related Work

Model Acceleration

Large Transformer has demonstrated its effectiveness on various natural language processing tasks, including machine translation (Vaswani et al. 2017), language modelling (Baevski and Auli 2019) and etc. The by-product brought by this huge network is the slow inference speed. Previous work focuses on improving model efficiency from different perspectives. For example, knowledge distillation approaches treat the large network output as the ground truth to train a small network (Kim and Rush 2016). Low-bit quantization approaches represent and run the model with 8-bit integer (Lin et al. 2020). Our work follows another line of researches, which pursues a more efficient architecture.

Chen et al. (2018) show that the attention of Transformer benefits the encoder the most and the decoder could be safely replaced by a recurrent network. This way reduces the complexity of the decoder to linear time but incurs a high cost in training. Zhang, Xiong, and Su (2018) show that the self-attention is not necessary and a simple averaging is enough. Xiao et al. (2019) indicate that most attention distributions are redundant and thus share these distributions among layers. Kitaev, Kaiser, and Levskaya (2020) use locality-sensitive hashing to select a constant number of words and perform attention on them.

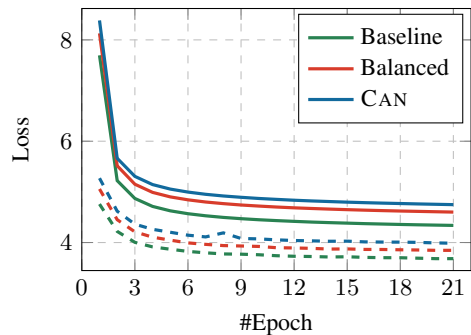


Figure 8: Loss vs. # of epochs on WMT14 En-De translation task (solid lines are the training losses, dashed lines are the validation losses).

Fan, Grave, and Joulin (2020) train a large Transformer and drop some layers at testing for fast inference. Gu et al. (2018) use a non-autoregressive decoder to predict the whole sentence at one time instead of generating it word by word. This approach makes a linear time translation process to constant time via the parallel computation. Perhaps the most related works are He et al. (2018); Zhang, Titov, and Sennrich (2019). They merge the self-attention and cross-attention and share their parameters. We, on the other hand, use different sets of parameters for each attention and mathematically prove that this way is equivalent to the standard Transformer under some mild conditions. We further show that the attention and FFN could also be merged together due to their linearity.

Deep Transformer

Recent studies have shown that deepening the Transformer encoder is more beneficial than widening the encoder or deepening the decoder (Bapna et al. 2018). Wang et al. (2019) show that placing the layer normalization before (Pre-Norm) rather than behind (Post-Norm) the sub-layer allows us to train deep Transformer. Xiong et al. (2020) prove that the success of the Pre-Norm network relies on its well-behaved gradient. Zhang, Titov, and Sennrich (2019) suggest that a proper initialization is enough to train a deep Post-Norm network. Kasai et al. (2020) similarly exploit this observation but to build a faster instead of a better model. They show that using knowledge distillation, a deep encoder and shallow decoder model could run much faster without losing any performance. Based on their work, we use this model as our baseline system and evaluate it on extensive machine translation tasks without knowledge distillation.

Conclusion

In this work, we propose CAN, whose decoder layer consists of only one attention. CAN offers consistent acceleration by providing a high degree of parallelism. Experiments on 14 WMT machine translation tasks show that CAN is $2.82\times$ faster than the baseline. We also use a stronger baseline for comparison. It employs a deep encoder and a shallow decoder, and is $2\times$ faster than the standard Transformer base-

line without loss in performance.

Acknowledgments

This work was supported in part by the National Science Foundation of China (Nos. 61876035 and 61732005), the National Key R&D Program of China (No. 2019QY1801). The authors would like to thank anonymous reviewers for their comments.

References

- Aharoni, R.; Johnson, M.; and Firat, O. 2019. Massively Multilingual Neural Machine Translation. In Burstein, J.; Doran, C.; and Solorio, T., eds., *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, 3874–3884. Association for Computational Linguistics. doi:10.18653/v1/n19-1388. URL <https://doi.org/10.18653/v1/n19-1388>.
- Ba, L. J.; Kiros, J. R.; and Hinton, G. E. 2016. Layer Normalization. *CoRR* abs/1607.06450. URL <http://arxiv.org/abs/1607.06450>.
- Baevski, A.; and Auli, M. 2019. Adaptive Input Representations for Neural Language Modeling. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net. URL <https://openreview.net/forum?id=ByxZX20qFQ>.
- Bapna, A.; Chen, M. X.; Firat, O.; Cao, Y.; and Wu, Y. 2018. Training Deeper Neural Machine Translation Models with Transparent Attention. In Riloff, E.; Chiang, D.; Hockenmaier, J.; and Tsujii, J., eds., *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, 3028–3033. Association for Computational Linguistics. doi:10.18653/v1/d18-1338. URL <https://doi.org/10.18653/v1/d18-1338>.
- Birch, A.; Finch, A. M.; Luong, M.; Neubig, G.; and Oda, Y. 2018. Findings of the Second Workshop on Neural Machine Translation and Generation. In Birch, A.; Finch, A. M.; Luong, M.; Neubig, G.; and Oda, Y., eds., *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation, NMT@ACL 2018, Melbourne, Australia, July 20, 2018*, 1–10. Association for Computational Linguistics. doi:10.18653/v1/w18-2701. URL <https://doi.org/10.18653/v1/w18-2701>.
- Chen, M. X.; Firat, O.; Bapna, A.; Johnson, M.; Macherey, W.; Foster, G. F.; Jones, L.; Schuster, M.; Shazeer, N.; Parmar, N.; Vaswani, A.; Uszkoreit, J.; Kaiser, L.; Chen, Z.; Wu, Y.; and Hughes, M. 2018. The Best of Both Worlds: Combining Recent Advances in Neural Machine Translation. In Gurevych, I.; and Miyao, Y., eds., *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, 76–86. Association for Computational Linguistics. doi:10.18653/v1/P18-1008. URL <https://www.aclweb.org/anthology/P18-1008>.
- Dehghani, M.; Gouws, S.; Vinyals, O.; Uszkoreit, J.; and Kaiser, L. 2019. Universal Transformers. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net. URL <https://openreview.net/forum?id=HyzdRiR9Y7>.
- Fan, A.; Grave, E.; and Joulin, A. 2020. Reducing Transformer Depth on Demand with Structured Dropout. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. URL <https://openreview.net/forum?id=SylO2yStDr>.
- Glorot, X.; and Bengio, Y. 2010. Understanding the difficulty of training deep feedforward neural networks. In Teh, Y. W.; and Titterton, D. M., eds., *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*, volume 9 of *JMLR Proceedings*, 249–256. JMLR.org. URL <http://proceedings.mlr.press/v9/glorot10a.html>.
- Gu, J.; Bradbury, J.; Xiong, C.; Li, V. O. K.; and Socher, R. 2018. Non-Autoregressive Neural Machine Translation. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net. URL <https://openreview.net/forum?id=B118BtlCb>.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016a. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, 770–778. IEEE Computer Society. doi:10.1109/CVPR.2016.90. URL <https://doi.org/10.1109/CVPR.2016.90>.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016b. Identity Mappings in Deep Residual Networks. In Leibe, B.; Matas, J.; Sebe, N.; and Welling, M., eds., *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, volume 9908 of *Lecture Notes in Computer Science*, 630–645. Springer. doi:10.1007/978-3-319-46493-0_38. URL https://doi.org/10.1007/978-3-319-46493-0_38.
- He, T.; Tan, X.; Xia, Y.; He, D.; Qin, T.; Chen, Z.; and Liu, T. 2018. Layer-Wise Coordination between Encoder and Decoder for Neural Machine Translation. In Bengio, S.; Wallach, H. M.; Larochelle, H.; Grauman, K.; Cesa-Bianchi, N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada*, 7955–7965. URL <http://papers.nips.cc/paper/8019-layer-wise-coordination-between-encod>.
- Hinton, G. E.; Vinyals, O.; and Dean, J. 2015. Distilling the Knowledge in a Neural Network. *CoRR* abs/1503.02531. URL <http://arxiv.org/abs/1503.02531>.
- Huang, G.; Sun, Y.; Liu, Z.; Sedra, D.; and Weinberger, K. Q. 2016. Deep Networks with Stochastic Depth. In Leibe, B.; Matas, J.; Sebe, N.; and Welling, M., eds., *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings*,

- Part IV, volume 9908 of *Lecture Notes in Computer Science*, 646–661. Springer. doi:10.1007/978-3-319-46493-0_39. URL https://doi.org/10.1007/978-3-319-46493-0_39.
- Kasai, J.; Pappas, N.; Peng, H.; Cross, J.; and Smith, N. A. 2020. Deep Encoder, Shallow Decoder: Reevaluating the Speed-Quality Tradeoff in Machine Translation. *CoRR* abs/2006.10369. URL <https://arxiv.org/abs/2006.10369>.
- Kim, Y.; and Rush, A. M. 2016. Sequence-Level Knowledge Distillation. In Su, J.; Carreras, X.; and Duh, K., eds., *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, 1317–1327. The Association for Computational Linguistics. doi:10.18653/v1/d16-1139. URL <https://doi.org/10.18653/v1/d16-1139>.
- Kitaev, N.; Kaiser, L.; and Levskaya, A. 2020. Reformer: The Efficient Transformer. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. URL <https://openreview.net/forum?id=rkgNKkHtvB>.
- Lin, Y.; Li, Y.; Liu, T.; Xiao, T.; Liu, T.; and Zhu, J. 2020. Towards Fully 8-bit Integer Inference for the Transformer Model. In Bessiere, C., ed., *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, 3759–3765. ijcai.org. doi:10.24963/ijcai.2020/520. URL <https://doi.org/10.24963/ijcai.2020/520>.
- Micikevicius, P.; Narang, S.; Alben, J.; Diamos, G. F.; Elsen, E.; García, D.; Ginsburg, B.; Houston, M.; Kuchaiev, O.; Venkatesh, G.; and Wu, H. 2018. Mixed Precision Training. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net. URL <https://openreview.net/forum?id=r1gs9JgRZ>.
- Ott, M.; Edunov, S.; Baevski, A.; Fan, A.; Gross, S.; Ng, N.; Grangier, D.; and Auli, M. 2019. fairseq: A Fast, Extensible Toolkit for Sequence Modeling. In Ammar, W.; Louis, A.; and Mostafazadeh, N., eds., *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Demonstrations*, 48–53. Association for Computational Linguistics. doi:10.18653/v1/n19-4009. URL <https://doi.org/10.18653/v1/n19-4009>.
- Sennrich, R.; Haddow, B.; and Birch, A. 2016. Neural Machine Translation of Rare Words with Subword Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics. doi:10.18653/v1/p16-1162. URL <https://doi.org/10.18653/v1/p16-1162>.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention is All you Need. In Guyon, I.; von Luxburg, U.; Bengio, S.; Wallach, H. M.; Fergus, R.; Vishwanathan, S. V. N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, 5998–6008. URL <http://papers.nips.cc/paper/7181-attention-is-all-you-need>.
- Wang, Q.; Li, B.; Xiao, T.; Zhu, J.; Li, C.; Wong, D. F.; and Chao, L. S. 2019. Learning Deep Transformer Models for Machine Translation. In Korhonen, A.; Traum, D. R.; and Màrquez, L., eds., *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, 1810–1822. Association for Computational Linguistics. doi:10.18653/v1/p19-1176. URL <https://doi.org/10.18653/v1/p19-1176>.
- Xiao, T.; Li, Y.; Zhu, J.; Yu, Z.; and Liu, T. 2019. Sharing Attention Weights for Fast Transformer. In Kraus, S., ed., *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, 5292–5298. ijcai.org. doi:10.24963/ijcai.2019/735. URL <https://doi.org/10.24963/ijcai.2019/735>.
- Xiong, R.; Yang, Y.; He, D.; Zheng, K.; Zheng, S.; Xing, C.; Zhang, H.; Lan, Y.; Wang, L.; and Liu, T. 2020. On Layer Normalization in the Transformer Architecture. *CoRR* abs/2002.04745. URL <https://arxiv.org/abs/2002.04745>.
- Zhang, B.; Titov, I.; and Sennrich, R. 2019. Improving Deep Transformer with Depth-Scaled Initialization and Merged Attention. In Inui, K.; Jiang, J.; Ng, V.; and Wan, X., eds., *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 898–909. Association for Computational Linguistics. doi:10.18653/v1/D19-1083. URL <https://doi.org/10.18653/v1/D19-1083>.
- Zhang, B.; Xiong, D.; and Su, J. 2018. Accelerating Neural Transformer via an Average Attention Network. In Gurevych, I.; and Miyao, Y., eds., *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, 1789–1798. Association for Computational Linguistics. doi:10.18653/v1/P18-1166. URL <https://www.aclweb.org/anthology/P18-1166/>.