

# Attention Is All You Need

Ashish Vaswani  
google brain  
avaswani@google.com

Noam Shazeer  
google brain  
noam@google.com

Vik Parmar  
google research  
vikip@google.com

Jakob Uszkoreit  
google research  
juszko@google.com

Lionel Jones  
google research  
ljones@google.com

Aidan Gomez  
University of Toronto  
aidan@cs.toronto.edu

Lukas Kaiser  
google brain  
lukasz.kaiser@google.com

Ilya Polosukhin  
ilya.polosukhin@gmail.com

## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show that models of this type are superior in quality while being more parallelizable and requiring significantly less time & train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results including ensembles by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 5.5 days on eight GPUs. A small fraction of the training costs of the best models from the literature.

## 1 Introduction

Recurrent neural networks (long short-term memory [12] and gated recurrent [13] neural networks) in particular have been firmly established as state of the art approaches to sequence modeling and transduction problems such as language modeling and machine translation [29, 9, 31]. Numerous efforts have since continued to push the boundaries of recurrent language models and encoder-decoder architectures [33, 44, 3].

Equal contribution. Ashish, Noam, and Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish and Ilya designed and implemented the first transformer model and has been crucially involved in every aspect of this work. Noam proposed scaled dot-product attention, multi-head attention and the parameter-free position representation and became the other person involved in nearly every detail. Nik designed, implemented, tuned and evaluated countless model variants; it was his original codebase and tensor2tensor [10] that experimented with novel model variants, was responsible for the initial codebase and efficient inference and visualizations. Lukasz and Aidan spent countless long days designing various parts of the implementing tensor2tensor replacing our earlier codebase, greatly improving results and massively accelerating our research.

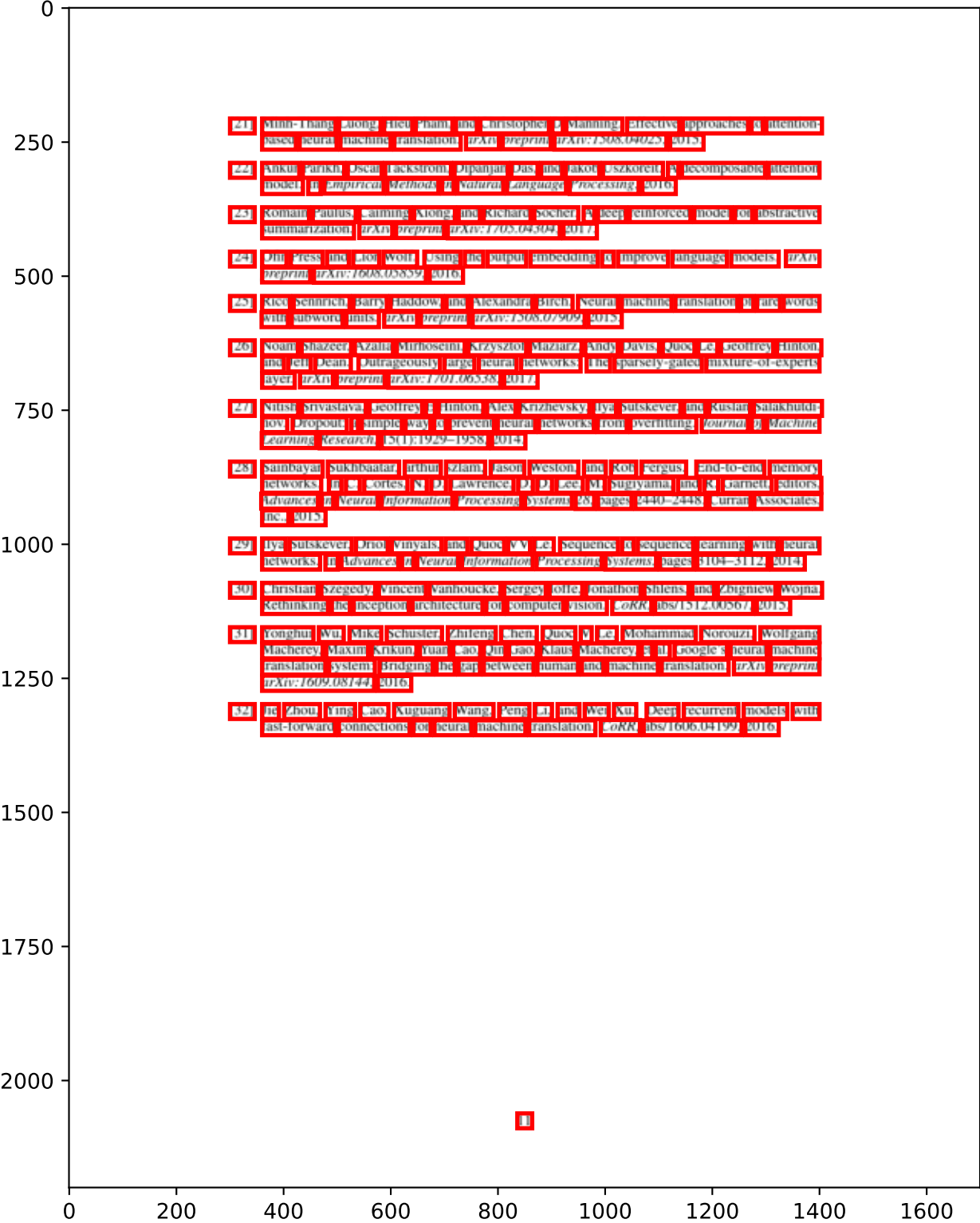
Work performed while at Google Brain.

Work performed while at Google Research.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

## References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.00420*, 2016.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2014.
- [3] Penny Britz, Xiang Sordle, Minh-Thang Luong, and Quoc V. Le. Massive exploration in neural machine translation architectures. *CoRR*, abs/1703.03906, 2017.
- [4] Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory networks for machine reading. *arXiv preprint arXiv:1601.06733*, 2016.
- [5] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fionn Flauger, Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *CoRR*, abs/1406.1078, 2014.
- [6] François Chollet. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint arXiv:1610.02357*, 2016.
- [7] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks in sequence modeling. *CoRR*, abs/1412.3555, 2014.
- [8] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. Convolutional sequence to sequence learning. *arXiv preprint arXiv:1705.03122v2*, 2017.
- [9] Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.
- [10] Kaiyang He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [11] Sepp Hochreiter, Yoshua Bengio, Paolo Frasconi, and Jürgen Schmidhuber. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies. 2001.
- [12] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1755–1780, 1997.
- [13] Karel Jezekowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. Exploring the limits of language modeling. *arXiv preprint arXiv:1602.02410*, 2016.
- [14] Łukasz Kaiser and Ilya Sutskever. Neural GPUs learn algorithms. In *International Conference on Learning Representations (ICLR)*, 2016.
- [15] Na Kalchbrenner, Jesse Espeholt, Karen Simonyan, Aaron van der Oord, Alex Graves, and Koray Kavukcuoglu. Neural machine translation in linear time. *arXiv preprint arXiv:1610.10099v2*, 2017.
- [16] Soohyun Kim, Carsten Donlon, Luong Hoang, and Alexander M. Rush. Structured attention networks. In *International Conference on Learning Representations*, 2017.
- [17] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [18] Diederik Kingma and Boris Ginsburg. Factorization tricks for LSTM networks. *arXiv preprint arXiv:1703.10722*, 2017.
- [19] Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Md Yu Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. *arXiv preprint arXiv:1703.03130*, 2017.
- [20] Samy Bengio, Łukasz Kaiser, and Carsten Donlon. Replace attention. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.



Recurrent models typically do the computation along the symbol positions in the input and output sequences, assigning the positions to steps in computation time they generate a sequence of hidden states  $h$  as a function of the previous hidden state  $h$  and the input or position  $i$ . This inherently sequential nature precludes parallelization within training examples which becomes critical in longer sequence lengths as memory constraints limit batching across examples. Recent work has achieved significant improvements in computational efficiency through factorization tricks [16] and conditional computation [40] while also improving model performance by use of the later the fundamental constraint of sequential computation, however, remains.

Attention mechanisms have become an integral part of competitive sequence modeling and transduction models in various tasks allowing modeling of dependencies without regard to their distance in the input or output sequences [4, 16] in a few cases [22], however, such attention mechanisms are used in conjunction with a recurrent network.

In this work we propose the transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The transformer allows for significantly more parallelization and so far, sets a new state of the art in translation quality after being trained for a little less than a month on 160K GPUs.

## 1 Background

The goal of reducing sequential computation also forms the foundation of the extended Neural GPU [40], ByteNet [15] and convS2S [8], all of which use convolutional neural networks as basic building blocks, computing hidden representations at parallel or in input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows as the distance between positions, linearly for convS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [11] as the transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions. In effect, we counteract with Multi-Head Attention as described in section 3.4.

Self-attention, sometimes called intra-attention, is an attention mechanism relating different positions of a single sequence to order a compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment, and learning task-independent sentence representations [41, 44, 45, 49].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [28].

To the best of our knowledge, however, the transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence aligned CNNs or convolution. In the following sections we will describe the transformer motivate self-attention and discuss its advantages over models such as [14, 15] and [8].

## 2 Model Architecture

Most competitive neural sequence transduction models have an encoder-decoder structure [5, 9, 42]. Here, the encoder maps an input sequence of symbol representations  $x_1, \dots, x_n$  to a sequence of continuous representations  $h_1, \dots, h_n$ . The decoder then generates an output sequence  $y_1, \dots, y_n$  of symbols one element at a time. At each step the model is auto-regressive [47] consuming the previously generated symbols as additional input when generating the next.

The transformer follows this overall architecture using stacked self-attention and point-wise fully connected layers for both the encoder and decoder shown in the left and right halves of figure 1 respectively.

### 2.1 Encoder and Decoder Stacks

**Encoder:** The encoder is composed of  $N$  identical layers, each layer has two sub-layers. The first is a multi-head self-attention mechanism and the second is a simple position-

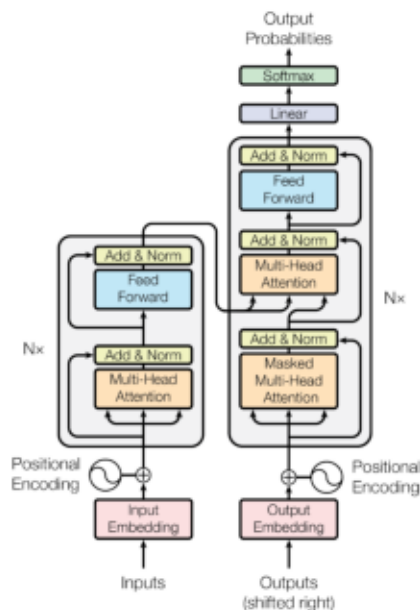


Figure 1 The Transformer model architecture

wise fully connected feed-forward network. We employ residual connection (10) around each of the two sub-layers, followed by layer normalization (11) that is the output of each sub-layer:  $\text{LayerNorm}(x + \text{Sublayer}(x))$  where  $\text{Sublayer}(x)$  is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers of the model as well as the embedding layers produce outputs of dimension  $d_{\text{model}} = d_1/2$ .

**Decoder:** The decoder is also composed of stack of  $N_D$  identical layers. In addition to the two sub-layers, in each encoder layer the decoder inserts a third sub-layer which performs multi-head attention over the output of the encoder stack, similar to the encoder we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer of the decoder stack to prevent positions from attending to subsequent positions. This masking combined with fact that the output embeddings are offset by one position ensure that the predictions for position  $i$  can depend only on the known outputs at positions less than  $i$ .

## 3.2 Attention

An attention function can be describes as mapping a query and a set of key-value pairs to an output, where the query keys, values and output are all vectors. The output is computed as a weighted sum of the values, where the weight assigned to each value is computed by a compatibility function of the query with the corresponding key.

### 3.2.1 Scaled Dot-Product Attention

We call our particular attention 'Scaled Dot-Product Attention' (Figure 9). The input consists of queries and keys of dimension  $d_k$  and values of dimension  $d_v$ . We compute the dot products of the

Figure 9. Left: Scaled Dot-Product Attention. Right: Multi-Head Attention consists of several attention layers running in parallel.

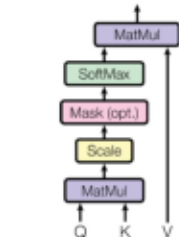


Figure 9. Left: Scaled Dot-Product Attention. Right: Multi-Head Attention consists of several attention layers running in parallel.

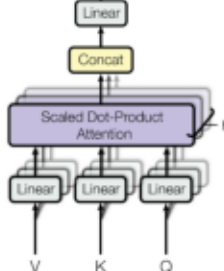


Figure 9. Left: Scaled Dot-Product Attention. Right: Multi-Head Attention consists of several attention layers running in parallel.

query with all keys, divide each by  $\sqrt{d_k}$ , and apply the softmax function to obtain the weights on the values.

In practice, we compute the attention function on a set of queries simultaneously packed together into a matrix  $Q$ ; the keys and values are also packed together into matrices  $K$  and  $V$ . We compute the matrix of outputs as:

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

The two most commonly used attention functions are additive attention [4] and dot-product multiplicative attention. Dot-product attention is identical to our algorithm except for the scaling factor in (9). Additive attention computes the compatibility function using a feed-forward network with a single hidden layer. While the two are similar in theoretical complexity, dot-product attention is much faster and more space-efficient in practice, since it can be implemented using highly optimized matrix multiplication code.

While for small values of  $d_k$ , the two mechanisms perform similarly, additive attention outperforms dot-product attention without scaling for larger values of  $d_k$  [3]. We suspect that for large values of  $d_k$ , the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients. To counteract this effect, we scale the dot products by  $\frac{1}{\sqrt{d_k}}$ .

### 3.2.2 Multi-Head Attention

Instead of performing a single attention function with  $d_{\text{model}}$  dimensional keys, values and queries, we found it beneficial to linearly project the queries, keys and values  $h$  times with different learned linear projections to  $d_k$ ,  $d_v$  and  $d_k$  dimensions respectively. In fact, in these projected versions of queries, keys and values we then perform the attention function in parallel, yielding  $h$  dimensional output values. These are concatenated and once again projected, resulting in the final values, as depicted in Figure 9.

Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. With a single attention head, averaging inhibits this.

It illustrates why the dot products are large, assuming that the components of  $q$  and  $k$  are independent random variables with mean 0 and variance 1. Then their dot product  $q \cdot k = \sum_{i=1}^d q_i k_i$  has mean 0 and variance  $d$ .

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

where the projections are parameter matrices  $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$  and  $W^O \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$ .

in this work we employ  $h = 8$  parallel attention layers in heads. For each  $d$ , these are  $d_k = d_v = d_{\text{model}}/h = 64$ . Due to the reduced dimension of each head the total computational cost is similar to that of single-head attention with full dimensionality.

### 3.4. Applications of Attention in our Model

The transformer uses multi-head attention in three different ways:

- In encoder-decoder attention layers, the queries come from the previous decoder layer and the memory keys and values come from the output of the encoder. This allows every position in the decoder to attend over all positions in the input sequence. This mimics the typical encoder-decoder attention mechanisms in sequence-to-sequence models such as [20, 21].
- The encoder contains self-attention layers. In a self-attention layer all of the keys, values and queries come from the same place. In this case, the output of the previous layer in the encoder. Each position in the encoder can attend to all positions in the previous layer of the encoder.
- Similarly self-attention layers in the decoder allow each position in the decoder to attend to all positions in the decoder up to and including that position. We need to prevent selfward information flow in the decoder to preserve the auto-regressive property. We implement this inside the scaled dot-product attention by masking out settings to  $-\infty$  all values in the input of the softmax which correspond to illegal connections. See figure 4.

### 3.5. Position-wise Feed-Forward Networks

In addition to attention sub-layers, each of the layers in our encoder and decoder contains a fully connected feed-forward network, which is applied to each position separately and identically. This consists of two linear transformations with a ReLU activation in between:

$$\text{FFN}(x) = \max(0, W_1 x + b_1) W_2 + b_2 \quad (2)$$

While the linear transformations are the same across different positions, they use different parameters from layer to layer. Another way of describing this is as two convolutions with kernel size 1. The dimensionality of input and output is  $d_{\text{model}} = 512$  and the inner-layer has dimensionality  $d_{\text{ffn}} = 2048$ .

### 3.6. Embeddings and Softmax

Similarly to other sequence transduction models, we use learned embeddings to convert the input tokens and output tokens in vectors of dimension  $d_{\text{model}}$ . We also use the usual learned linear transformation and softmax function to convert the decoder output to predicted next-token probabilities. In our model we share the same weight matrix between the two embedding layers and the pre-softmax linear transformation, similar to [24]. In the embedding layers, we multiply those weights by  $\sqrt{d_{\text{model}}}$ .

### 3.7. Positional Encoding

Since our model contains no recurrence and no convolution, in order for the model to make use of the order of the sequence, we must inject some information about the relative or absolute position of the tokens in the sequence. To this end, we add 'positional encodings' to the input embeddings in the



Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations on different layer types, as a function of sequence length  $L$ , the representation dimension  $d$ , the kernel size  $k$ , convolutional and  $l$  its size  $n$ , the neighborhood  $n$  restricted self-attention

Layer type	Complexity per layer	Sequential operations	Maximum path length
Self-Attention	$\mathcal{O}(n^2)$	$\mathcal{O}(1)$	$\mathcal{O}(L)$
Recurrent	$\mathcal{O}(L)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$
Convolutional	$\mathcal{O}(L)$	$\mathcal{O}(1)$	$\mathcal{O}(L)$
Self-Attention restricted	$\mathcal{O}(L)$	$\mathcal{O}(1)$	$\mathcal{O}(L)$

bottoms of the encoder and decoder stacks. The positional encodings have the same dimension  $d_{\text{model}}$  as the embeddings, so that the two can be summed. There are many choices of positional encodings learned and fixed [5].

In this work, we use sine and cosine functions of different frequencies:

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

where  $pos$  is the position and  $i$  is the dimension. That is each dimension of the positional encoding corresponds to  $\sin$  or  $\cos$ . The wavelengths form a geometric progression from  $2\pi$  to  $10000 \cdot 2\pi$ . We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset  $k$ ,  $PE_{(pos+k)}$  can be represented as a linear function of  $PE_{pos}$ .

We also experimented with using learned positional embeddings [5]. Instead, we found that the two versions produced nearly identical results (see Table 1 row 11). We chose the sinusoidal version because it may allow the model to extrapolate to sequence lengths longer than the ones encountered during training.

## 1 Why Self-Attention

In this section we compare various aspects of self-attention layers to the recurrent and convolutional layers commonly used for mapping one variable-length sequence of symbol representations to another of equal length [1, 2, 3, 4, 10, 11] with  $L \leq 10^5$  such as in [1] hidden layers in a typical sequence transduction encoder or decoder. Motivating our use of self-attention we consider three desiderata.

One of the main computational complexity per layer. Another is the amount of computation that can be parallelized, as measured by the minimum number of sequential operations required.

The third is the path length between long-range dependencies in the network. Learning long-range dependencies is a key challenge in many sequence transduction tasks. One key factor affecting the ability to learn such dependencies is the length of the paths forward and backward signals have to traverse in the network. The shorter these paths between any combination of positions in the input and output sequences, the easier it is to learn long-range dependencies [11]. Hence we also compare the maximum path length between any two input and output positions in networks composed of the different layer types.

As noted in Table 1, self-attention layers connect all positions with a constant number of sequentially executed operations, whereas a recurrent layer requires  $\mathcal{O}(n)$  sequential operations. In terms of computational complexity, self-attention layers are faster than recurrent layers when the sequence length  $L$  is smaller than the representation dimensionality  $d$ , which is most often the case with sentence representations used by state-of-the-art models of machine translations such as word-piece [11] and byte-pair [25] representations. To improve computational performance on tasks involving very long sequences, self-attention could be restricted to considering only a neighborhood of size  $n$  in



the input sequence centered around its respective output position. This would increase the maximum path length to  $2(n/2)$ . We plan to investigate this approach during a future work.

A single convolutional layer with kernel width  $k$  does not connect all pairs of input and output positions, doing so requires a stack of  $2(n/k)$  convolutional layers in the case of contiguous kernels of  $2(\log_2 n)$ . In the case of dilated convolutions [12], increasing the length of the longest paths between any two positions in the network, convolutional layers are generally more expensive than recurrent layers by a factor of  $n$ . Separable convolutions [9], however, decrease the complexity considerably to  $2(n/k)$  even with  $k = 1$ . However, the complexity of 1D separable convolution is equal to the combination of a self-attention layer and a point-wise feed-forward layer. The approach we take in our model.

As side benefit, self-attention could yield more interpretable models. We inspect attention distributions from our models and present and discuss examples in the appendix. Not only do individual attention heads clearly learn to perform different tasks, many appear to exhibit behavior related to the syntactic and semantic structure of the sentences.

## 5 Training

This section describes the training regime of our models.

### 5.1 Training Data and Batching

We trained on the standard WMT 2014 English-German dataset consisting of about 46 million sentence pairs. Sentences were encoded using byte-pair encoding [3], which has a shared source-target vocabulary of about 37000 tokens. For English-French, we used the significantly larger WMT 2014 English-French dataset consisting of 56M sentences and split tokens into a 52000 word-piece vocabulary [31]. Sentence pairs were batched together by approximate sequence length. Each training batch contained 1 set of sentence pairs containing approximately 25000 source tokens and 25000 target tokens.

### 5.2 Hardware and Schedule

We trained our models on one machine with 4 NVIDIA P100 GPUs. For our base models using the hyperparameters described throughout the paper, each training step took about 0.4 seconds. We trained the base models on a total of 100,000 steps or 2 hours. For our big models, described at the bottom line of table 3, the step time was 10 seconds. The big models were trained on 800,000 steps (32 days).

### 5.3 Optimizer

We used the Adam optimizer [17] with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$  and  $\epsilon = 10^{-9}$ . We varied the learning rate over the course of training according to the formula:

$$\text{rate} = \frac{1}{\sqrt{\text{step\_num}}} \min(\text{step\_num}^{-0.5}, \text{step\_num}^{-\text{warmup\_steps}^{-1.5}}) \quad (3)$$

This corresponds to increasing the learning rate linearly for the first *warmup\_steps* training steps and decreasing it thereafter proportionally to the inverse square root of the step number. We used *warmup\_steps* = 1000.

### 5.4 Regularization

We employ three types of regularization during training.

**Residual Dropout** We apply dropout [20] to the output of each sub-layer before it is added to the sub-layer input and normalized. In addition, we apply dropout to the sum of the embeddings and the positional encodings in both the encoder and decoder stacks. For the base model we use a rate of  $p_{\text{resid}} = 0.1$ .

Table 3: The Transformer achieves better BLEU scores than previous state-of-the-art models in the English-to-German and English-to-French newstest2013 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	25.7			
DeepAt + PosUnk [34]		39.2	31	10%
ENMT1 + RL [31]	25.1	39.92	31	10%
CONVS2S [8]	25.13	40.46	31	10%
MoE [20]	30.05	40.54	31	10%
DeepAt + PosUnk ensemble [34]		40.5	31	10%
ENMT1 + RL ensemble [31]	30.34	41.14	31	10%
CONVS2S ensemble [8]	30.34	41.25	31	10%
Transformer (base model)	27.3	38.1	28	10%
Transformer (big)	28.4	41.3	28	10%

**Label Smoothing.** During training we employed label smoothing of value  $\epsilon = 0.1$  [30]. This hurts perplexity as the model learns to be more unsure, but improves accuracy and BLEU score.

## 2 Results

### 2.1 Machine Translation

On the WMT 2014 English-to-German translation task the big transformer model (Transformer big in table 3) outperforms the best previously reported models (including ensembles of more than 20 BLEU) establishing a new state-of-the-art BLEU score of 28.4. The configuration of this model is listed in the bottom line of table 3. Training cost for large  $n = 100$  GPUs. Even our base model surpasses all previously published models and ensembles at a fraction of the training cost of any of the competitive models.

On the WMT 2014 English-to-French translation task our big model achieves a BLEU score of 41.3, outperforming all of the previously published single models as well as 77% the training cost of the previous state-of-the-art model. The transformer big model trained on English-to-French uses dropout rate  $\epsilon_{\text{norm}} = 0.1$  instead of 0.3.

For the base models we used a single model obtained by averaging the last 5 checkpoints which were written at 10-minute intervals. For the big models we averaged the last 50 checkpoints. We used beam search with beam size of 1 and length penalty  $\beta = 0.0$  [31]. These hyperparameters were chosen after experimentation on the development set. We set the maximum output length during inference to input length + 50 but terminate early when possible [31].

Table 3 summarizes our results and compares our translation quality and training costs to other model architectures from the literature. We estimate the number of floating point operations used to train a model by multiplying the training time, the number of GPUs used and an estimate of the sustained single-precision floating-point capacity of each GPU [1].

### 2.2 Model Variations

To evaluate the importance of different components of the Transformer we varied our base model in different ways, measuring the change in performance on English-to-German translation on the development set newstest2013. We used beam search as described in the previous section but in checkpoint averaging we present these results in Table 3.

In table 3, rows A) we vary the number of attention heads and the attention key and value dimensions keeping the amount of computation constant as described in section 3.2.2. While single-head attention is 0.3 BLEU worse than the best setting, quality also drops off with too many heads.

We need values of 9.8, 10.7, 33 and 95 TFLOPs or 380, 540, 540 and 1100 respectively.

