

Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

在提供适当署名的情况下, 谷歌特此授权复制本文中的表格和图表, 仅限于用于新闻或学术作品。

Attention Is All You Need

Attention Is All You Need

Ashish Vaswani* Google Brain avaswani@google.com	Noam Shazeer* Google Brain noam@google.com	Niki Parmar* Google Research nikip@google.com	Jakob Uszkoreit* Google Research usz@google.com
---	---	--	--

Ashish Vaswani* Google Brain avaswani@google.com	Noam Shazeer* Google Brain noam@google.com	Niki Parmar* Google Research nikip@google.com	Jakob Uszkoreit* Google Research huz@google.com
---	---	--	--

Llion Jones* Google Research llion@google.com	Aidan N. Gomez* † University of Toronto aidan@cs.toronto.edu	Lukasz Kaiser* Google Brain lukaszkaizer@google.com
--	---	--

Llion Jones* Google Research llion@google.com	Aidan N. Gomez* †多 伦多大学 aidan@cs.toronto.edu	Lukasz Kaiser*Google Brain lukaszkaizer@google.com
--	---	---

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

伊利亚·波洛苏金* ‡
illia.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 these models to be superior in quality while being more parallelizable and requiring significantly less time to 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 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

Abstract

主导的序列转导模型基于复杂的循环或卷积神经网络, 这些网络包含编码器和解码器。表现最佳的模型还通过注意力机制连接编码器和解码器。我们提出了一种新的简单网络架构——Transformer, 它完全基于注意力机制, 摒弃了循环和卷积。在两个机器翻译任务上的实验表明, 这些模型在质量上更优越, 同时更易于并行化, 且训练时间显著减少。我们的模型在WMT 2014英语到德语翻译任务上取得了28.4的BLEU分数, 比现有的最佳结果(包括集成模型)提高了超过2个BLEU。在WMT 2014英语到法语翻译任务上, 我们的模型在八块GPU上训练3.5天后, 建立了一个新的单模型SOTA BLEU分数41.8, 其训练成本只是文献中最佳模型的极小一部分。我们通过将其成功应用于英语依存句法分析(无论使用大量还是有限训练数据), 证明了Transformer在其他任务上的泛化能力。

*Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, designed and implemented the first Transformer models 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. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor2tensor. Llion also experimented with novel model variants, was responsible for our initial codebase, and efficient inference and visualizations. Lukasz and Aidan spent countless long days designing various parts of and 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.

*贡献均等。排序顺序随机。Jakob提议用自注意力机制替代循环神经网络, 并开始评估这一想法。

Ashish与Illia设计并实现了首批Transformer模型, 并在这项工作的各个方面都发挥了关键作用。

Noam提出了缩放点积注意力、多头注意力和无参数位置表示, 并成为参与几乎所有细节的另一位关键人物。Niki设计、实现、调优和评估了我们原始代码库和tensor2tensor中的无数模型变体。

Llion也实验了新型模型变体, 负责我们的初始代码库、高效推理和可视化。Lukasz和Aidan花费了无数漫长的日子设计和实现tensor2tensor的各个部分, 替换了我们早期的代码库, 极大地提升了结果, 并极大地加速了我们的研究。

†在Google Brain期间完成的工作。

‡在Google Research期间完成的工作。

1 Introduction

Recurrent neural networks, long short-term memory [13] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and transduction problems such as language modeling and machine translation [35, 2, 5]. Numerous efforts have since continued to push the boundaries of recurrent language models and encoder-decoder architectures [38, 24, 15].

Recurrent models typically factor computation along the symbol positions of the input and output sequences. Aligning the positions to steps in computation time, they generate a sequence of hidden states h_t , as a function of the previous hidden state h_{t-1} and the input for position t . This inherently sequential nature precludes parallelization within training examples, which becomes critical at longer sequence lengths, as memory constraints limit batching across examples. Recent work has achieved significant improvements in computational efficiency through factorization tricks [21] and conditional computation [32], while also improving model performance in case of the latter. The fundamental constraint of sequential computation, however, remains.

Attention mechanisms have become an integral part of compelling sequence modeling and transduction models in various tasks, allowing modeling of dependencies without regard to their distance in the input or output sequences [2, 19]. In all but a few cases [27], 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 can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In 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, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to 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 [4, 27, 28, 22].

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 [34].

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 RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

3 Model Architecture

Most competitive neural sequence transduction models have an encoder-decoder structure [5, 2, 35]. Here, the encoder maps an input sequence of symbol representations (x_1, \dots, x_n) to a sequence of continuous representations $\mathbf{z} = (z_1, \dots, z_n)$. Given \mathbf{z} , the decoder then generates an output sequence (y_1, \dots, y_m) of symbols one element at a time. At each step the model is auto-regressive [10], consuming the previously generated symbols as additional input when generating the next.

1 简介

循环神经网络，特别是长短期记忆 [13] 和门控循环 [7] 神经网络，已被牢固确立为序列建模和转导问题（如语言建模和机器翻译 [35, 2, 5]）的SOTA方法。此后，人们持续投入大量努力，不断拓展循环语言模型和编码器-解码器架构 [38, 24, 15] 的边界。

循环模型通常沿着输入和输出序列的符号位置分解计算。将位置与计算时间中的步骤对齐，它们生成一个隐藏状态序列 h_t ，作为先前隐藏状态 h_{t-1} 和位置 t 输入的函数。这种固有的顺序性质阻止了在训练示例内部的并行化，这在较长的序列长度时变得至关重要，因为内存限制限制了跨示例的批处理。最近的工作通过分解技巧 [21] 和条件计算 [32]，实现了计算效率的显著改进，同时也在后者的情况下提高了模型性能。然而，顺序计算的 **fundamental constraint** 仍然存在。

注意力机制已成为各种任务中引人入胜的序列建模和转导模型的一个组成部分，允许在不考虑它们在输入或输出序列 [2, 19] 中的距离的情况下建模依赖关系。但在少数情况下 [27]，这种注意力机制是与循环网络一起使用的。

在这项工作中，我们提出了 **Transformer**，一种模型架构，它避免了递归，而是完全依赖于注意力机制来在输入和输出之间建立全局依赖关系。**Transformer** 允许显著更多的并行化，并且在经过在八个 P100 GPU 上训练仅十二小时后，可以实现翻译质量的新的 **state of the art**。

2 背景

减少顺序计算的目标也构成了扩展神经GPU[16]、ByteNet [18] 和ConvS2S [9] 的基础，所有这些模型都使用卷积神经网络作为基本构建块，并行计算所有输入和输出位置的隐藏表示。在这些模型中，将两个任意输入或输出位置的信号关联起来所需的操作数量随着位置之间距离的增长而增加，对于ConvS2S是线性增长，对于ByteNet是对数增长。这使得学习远处位置之间的依赖关系变得更加困难 [12]。在Transformer中，这减少到了一个常数数量的操作，尽管这会由于对注意力加权位置进行平均而导致有效分辨率降低，而我们通过第3.2节中描述的多头注意力来抵消这种效应。

自注意力，有时称为序列内注意力，是一种关联单个序列不同位置的注意力机制，目的是计算序列的表示。自注意力已成功应用于各种任务，包括阅读理解、抽象摘要、文本蕴涵和学习任务无关的句子表示 [4, 27, 28, 22]。

端到端记忆网络基于循环注意力机制，而非序列对齐的递归，并且已被证明在简单语言问答和语言建模任务 [34] 上表现良好。

据我们所知，然而，**Transformer** 是首个完全依赖自注意力机制来计算其输入和输出表示的转换模型，而不使用序列对齐的循环神经网络或卷积。在接下来的部分中，我们将描述Transformer，阐述自注意力的动机，并讨论其相对于 [17, 18] 和 [9] 等模型的优点。

3 模型架构

大多数具有竞争力的神经序列转换模型采用编码器-解码器结构 [5, 2, 35]。在此，编码器将符号表示的输入序列 (x_1, \dots, x_n) 映射为连续表示的序列 $\mathbf{z} = (z_1, \dots, z_n)$ 。给定 \mathbf{z} ，解码器随后逐个元素地生成输出序列 (y_1, \dots, y_m) 。在每一步，模型都是自回归 [10] 的，在生成下一个元素时将先前生成的符号作为附加输入进行消耗。

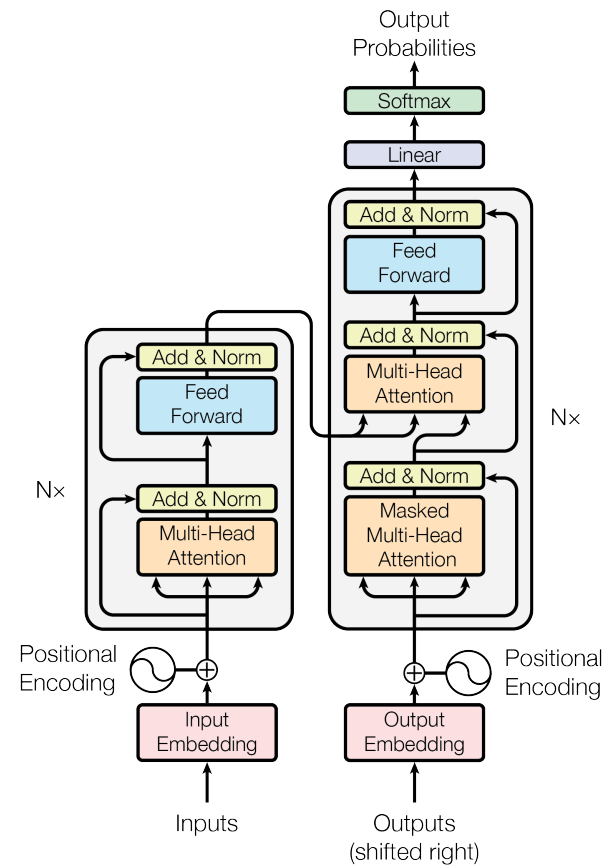


Figure 1: The Transformer - model architecture.

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.

3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ 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 in 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, ensures 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 described 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

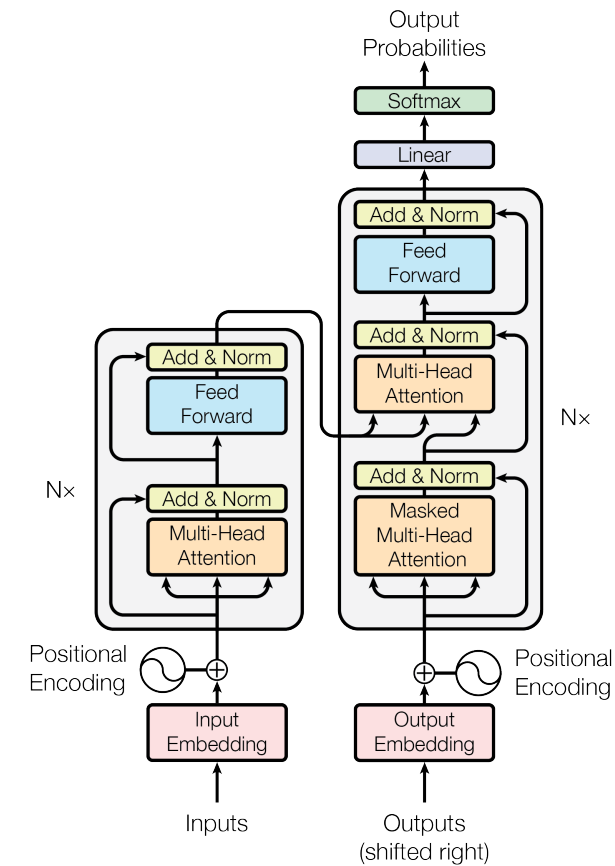


图1: Transformer - 模型架构。

Transformer 采用堆叠的自注意力机制和逐点、全连接层来构建编码器和解码器，分别如图1的左半部和右半部所示。

3.1 Encoder and Decoder Stacks

编码器: 编码器由一个堆叠的 $N = 6$ 个相同层组成。每一层有两个子层。第一个是多头自注意力机制，第二个是一个简单、位置感知的全连接前馈网络。我们在每个子层周围使用残差连接 [11]，然后进行层归一化 [1]。也就是说，每个子层的输出是 $\text{LayerNorm}(x + \text{Sublayer}(x))$ ，其中 $\text{Sublayer}(x)$ 是子层本身实现的函数。为了便于这些残差连接，模型中的所有子层以及嵌入层都产生维度为 $d_{\text{model}} = 512$ 的输出。

解码器: 解码器也由一个堆叠的 $N = 6$ 个相同层组成。除了编码器层中的两个子层外，解码器还插入一个第三个子层，该子层对编码器堆栈的输出执行多头注意力。与编码器类似，我们在每个子层周围使用残差连接，然后进行层归一化。我们还修改了解码器堆栈中的自注意力子层，以防止位置关注后续位置。这种掩码，结合输出嵌入偏移一个位置的事实，确保位置 i 的预测只能依赖于位置小于 i 的已知输出。

3.2 注意力

注意力函数可以被描述为将一个查询和一组键值对映射到一个输出，其中查询、键、值和输出都是向量。输出计算为值的加权求和



Figure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several attention layers running in parallel.

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 2). The input consists of queries and keys of dimension d_k , and values of dimension d_v . We compute the dot products of the query with all keys, divide each by $\sqrt{d_k}$, and apply a 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:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

The two most commonly used attention functions are additive attention [2], and dot-product (multiplicative) attention. Dot-product attention is identical to our algorithm, except for the scaling factor of $\frac{1}{\sqrt{d_k}}$. 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 gradients⁴. 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_{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_k and d_v dimensions, respectively. On each of these projected versions of queries, keys and values we then perform the attention function in parallel, yielding d_v -dimensional

⁴To illustrate why the dot products get large, assume that the components of q and k are independent random variables with mean 0 and variance 1. Then their dot product, $q \cdot k = \sum_{i=1}^{d_k} q_i k_i$, has mean 0 and variance d_k .

图2：左侧为缩放点积注意力。右侧为多头注意力由多个并行运行的注意力层组成。

其中分配给每个值的权重由查询与相应键的兼容性函数计算得出。

3.2.1 缩放点积注意力

我们将这种注意力称为“缩放点积注意力”（图2）。输入由维度为 d_k 的查询和键，以及维度为 d_v 的值组成。我们将查询与所有键的点积计算出来，除以 $\sqrt{d_k}$ ，并应用softmax函数来获得值上的权重。

在实践中，我们同时在一组查询上计算注意力函数，将它们打包成一个矩阵 Q 。键和值也被打包成矩阵 K 和 V 。输出的矩阵计算如下：

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

两种最常用的注意力函数是加性注意力和点积（乘积）注意力。点积注意力与我们的算法相同，除了缩放因子 $\frac{1}{\sqrt{d_k}}$ 。加性注意力使用单隐层的前馈网络来计算兼容性函数。虽然两者在理论复杂度上相似，但在实践中，点积注意力要快得多，也更节省空间，因为它可以使用高度优化的矩阵乘法代码来实现。

虽然对于较小的 d_k 值，这两种机制表现相似，但对于较大的 d_k [3]值，加性注意力在无缩放的情况下优于点积注意力。我们推测，对于较大的 d_k 值，点积的绝对值会变得很大，将softmax函数推入梯度极小的区域⁴。为了抵消这种效应，我们通过 $\frac{1}{\sqrt{d_k}}$ 对点积进行缩放。

3.2.2 多头注意力

与其使用具有 d_{model} 维键、值和查询的单个注意力函数，我们发现将查询、键和值分别通过不同的、已学习的线性投影进行 h 次线性投影到 d_k 、 d_k 和 d_v 维是有益的。然后，我们对这些投影版本的查询、键和值并行执行注意力函数，产生 d_v 维

⁴为了说明为什么点积会变得很大，假设 q 和 k 的分量是独立的随机变量，均值为0，方差为1。那么它们的点积 $q \cdot k = \sum_{i=1}^{d_k} q_i k_i$ 的均值为0，方差为 d_k 。

output values. These are concatenated and once again projected, resulting in the final values, as depicted in Figure 2.

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.

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

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

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 $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$.

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

3.2.3 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 [38, 2, 9].
- 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 leftward information flow in the decoder to preserve the auto-regressive property. We implement this inside of scaled dot-product attention by masking out (setting to $-\infty$) all values in the input of the softmax which correspond to illegal connections. See Figure 2.

3.3 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.

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (2)$$

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 $d_{ff} = 2048$.

3.4 Embeddings and Softmax

Similarly to other sequence transduction models, we use learned embeddings to convert the input tokens and output tokens to vectors of dimension d_{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 [30]. In the embedding layers, we multiply those weights by $\sqrt{d_{\text{model}}}$.

输出值。这些值被连接起来，再次投影，最终得到图2中所示的值。

多头注意力允许模型在不同位置同时关注来自不同表示子空间的信息。使用单个注意力头时，平均操作会抑制这种能力。

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

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

其中投影是参数矩阵 $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$ 、 $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$ 、 $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$ 和 $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$ 。

在这项工作中，我们采用 $h = 8$ 并行注意力层，或称为头。对于每一个，我们使用 $d_k = d_v = d_{\text{model}}/h = 64$ 。由于每个头的维度降低，总计算成本与全维度单头注意力相似。

3.2.3 注意力机制在我们模型中的应用

Transformer 以三种不同的方式使用多头注意力：

- 在“编码器-解码器注意力”层中，查询来自前一个解码器层，而内存键和值来自编码器的输出。这允许解码器中的每个位置关注输入序列中的所有位置。这模拟了序列到序列模型（如 [38, 2, 9]）中典型的编码器-解码器注意力机制。
- 编码器包含自注意力层。在自注意力层中，所有的键、值和查询都来自同一位置，在本例中，来自编码器前一层的输出。编码器中的每个位置可以关注编码器前一层中的所有位置。
- 类似地，解码器中的自注意力层允许解码器中的每个位置关注解码器中所有位置，包括该位置本身。我们需要防止解码器中的左向信息流动，以保持自回归特性。我们在缩放点积注意力内部通过将softmax输入中对应非法连接的所有值置为 $-\infty$ 来实现这一点。见图2。

3.3 位置前馈网络

除了注意力子层之外，我们编码器和解码器中的每一层都包含一个全连接前馈网络，该网络分别且相同地应用于每个位置。这由两个线性变换组成，中间有一个ReLU激活函数。

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (2)$$

虽然线性变换在不同位置上是相同的，但它们在层与层之间使用不同的参数。另一种描述方式是将其视为两个内核大小为1的卷积。输入和输出的维度是 $d_{\text{model}} = 512$ ，内部层的维度是 $d_{ff} = 2048$ 。

3.4 嵌入和Softmax

同样地，与其他序列转导模型一样，我们使用学习到的嵌入将输入标记和输出标记转换为维度为 d_{model} 的向量。我们还使用通常的学习到的线性变换和softmax函数将解码器输出转换为预测的下一个标记概率。在我们的模型中，我们在两个嵌入层和预softmax线性变换之间共享相同的权重矩阵，类似于 [30]。在嵌入层中，我们将这些权重乘以 $\sqrt{d_{\text{model}}}$ 。

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

3.5 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 at the bottoms of the encoder and decoder stacks. The positional encodings have the same dimension d_{model} as the embeddings, so that the two can be summed. There are many choices of positional encodings, learned and fixed [9].

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 a sinusoid. The wavelengths form a geometric progression from 2π 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 [9] instead, and found that the two versions produced nearly identical results (see Table 3 row (E)). We chose the sinusoidal version because it may allow the model to extrapolate to sequence lengths longer than the ones encountered during training.

4 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 (x_1, \dots, x_n) to another sequence of equal length (z_1, \dots, z_n) , with $x_i, z_i \in \mathbb{R}^d$, such as a hidden layer in a typical sequence transduction encoder or decoder. Motivating our use of self-attention we consider three desiderata.

One is the total 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 [12]. 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, a self-attention layer connects all positions with a constant number of sequentially executed operations, whereas a recurrent layer requires $O(n)$ sequential operations. In terms of computational complexity, self-attention layers are faster than recurrent layers when the sequence

表1: 不同层类型的最大路径长度、每层复杂度和最小顺序操作数。 n 是序列长度, d 是表示维度, k 是卷积的核大小, r 是受限自注意力中的邻域大小。

层类型	每层复杂度	顺序操作	最大路径长度
自注意力机制	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
循环	$O(n \cdot d^2)$	$O(n)$	$O(n)$
卷积	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
自注意力机制 (受限)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

3.5 位置编码

由于我们的模型不包含递归和卷积, 为了让模型能够利用序列的顺序, 我们必须注入一些关于序列中标记的相对位置或绝对位置的信息。为此, 我们在编码器和解码器堆栈底部的输入嵌入中添加了“位置编码”。位置编码的维度与嵌入相同 d_{model} , 以便两者可以相加。位置编码有多种选择, 包括学习型和固定型 [9]。

在本工作中, 我们使用不同频率的正弦和余弦函数:

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

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

pos 是位置, i 是维度。也就是说, 位置编码的每个维度都对应一个正弦波。波长从 2π 到 $10000 \cdot 2\pi$ 形成几何级数。我们选择这个函数, 因为我们假设它能让模型容易地通过相对位置来学习注意力, 因为对于任何固定的偏移量 k , PE_{pos+k} 可以表示为 PE_{pos} 的线性函数。

我们还尝试使用学习位置嵌入 [9], 发现这两个版本产生的结果几乎相同 (见表3行(E))。我们选择正弦版本, 因为它可能允许模型推广到训练中遇到的序列长度更长的序列。

4 为什么使用自注意力机制

在本节中, 我们将自注意力层与通常用于将一个变量长度的符号表示序列 (x_1, \dots, x_n) 映射到另一个等长序列 (z_1, \dots, z_n) 的循环层和卷积层进行比较, 其中 $x_i, z_i \in \mathbb{R}^d$, 例如典型序列转换编码器或解码器中的一个隐藏层。为了说明我们使用自注意力机制的原因, 我们考虑了三个期望。

一个是每层的总计算复杂度。另一个是可以并行化的计算量, 以所需的最小顺序操作数衡量。

三是网络中长期距离依赖之间的路径长度。学习长期距离依赖是许多序列转换任务中的一个关键挑战。影响学习这种依赖能力的一个关键因素是前向和后向信号在网络中必须穿越的路径长度。输入和输出序列中任何位置组合之间的路径越短, 学习长期距离依赖就越容易 [12]。因此, 我们还比较由不同层类型组成的网络中任何两个输入和输出位置之间的最大路径长度。

如表1所示, 自注意力层通过固定数量的顺序执行操作连接所有位置, 而循环层需要 $O(n)$ 个顺序操作。在计算复杂度方面, 当序列长度 O 小于表示维度 n 时, 自注意力层比循环层更快, 这在机器翻译中使用的当前最佳模型的句子表示中最为常见, 例如词片和字节对表示。为了提高涉及非常长序列的任务的计算性能, 自注意力可以限制为仅考虑输入序列中围绕相应输出位置的大小为 r 的邻域。这将使最大路径长度增加到 $O(n/r)$ 。我们计划在未来的工作中进一步研究这种方法。

length n is smaller than the representation dimensionality d , which is most often the case with sentence representations used by state-of-the-art models in machine translations, such as word-piece [38] and byte-pair [31] representations. To improve computational performance for tasks involving very long sequences, self-attention could be restricted to considering only a neighborhood of size r in the input sequence centered around the respective output position. This would increase the maximum path length to $O(n/r)$. We plan to investigate this approach further in future work.

A single convolutional layer with kernel width $k < n$ does not connect all pairs of input and output positions. Doing so requires a stack of $O(n/k)$ convolutional layers in the case of contiguous kernels, or $O(\log_k(n))$ in the case of dilated convolutions [18], 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 k . Separable convolutions [6], however, decrease the complexity considerably, to $O(k \cdot n \cdot d + n \cdot d^2)$. Even with $k = n$, however, the complexity of a 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 for our models.

5.1 Training Data and Batching

We trained on the standard WMT 2014 English-German dataset consisting of about 4.5 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 36M sentences and split tokens into a 32000 word-piece vocabulary [38]. Sentence pairs were batched together by approximate sequence length. Each training batch contained a 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 8 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 for a total of 100,000 steps or 12 hours. For our big models,(described on the bottom line of table 3), step time was 1.0 seconds. The big models were trained for 300,000 steps (3.5 days).

5.3 Optimizer

We used the Adam optimizer [20] 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:

$$lrate = d_{\text{model}}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup_steps^{-1.5}) \quad (3)$$

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 $warmup_steps = 4000$.

5.4 Regularization

We employ three types of regularization during training:

长度 n 小于表示维度 d , 这在使用当前最佳模型的机器翻译中最为常见, 例如词片[38] 和字节对 [31] 表示。为了提高涉及非常长序列的任务的计算性能, 自注意力可以限制为仅考虑输入序列中围绕相应输出位置的大小为 r 的邻域。这将使最大路径长度增加到 $O(n/r)$ 。我们计划在未来的工作中进一步研究这种方法。

一个卷积层, 其核宽度为 $k < n$, 并不能连接所有输入和输出位置。要实现这一点, 在连续核的情况下需要堆叠 $O(n/k)$ 个卷积层, 而在扩张卷积的情况下则需要 $O(\log_k(n))$ 个卷积层, 这会增加网络中任意两个位置之间最长路径的长度。卷积层通常比循环层更昂贵, 成本高出 k 倍。然而, 可分离卷积 [6], 可以显著降低复杂度, 降至 $O(k \cdot n \cdot d + n \cdot d^2)$ 。即使如此, 可分离卷积的复杂度仍等于自注意力层和逐点前馈层的组合, 这也是我们模型采用的方法。

作为附加好处, 自注意力机制可以产生更易于解释的模型。我们检查了模型中的注意力分布, 并在附录中呈现和讨论了相关示例。不仅单个注意力头能明显学会执行不同任务, 许多注意力头似乎还表现出与句子句法和语义结构相关的行为。

5 训练

本节描述了我们模型的训练流程。

5.1 训练数据与批处理

我们在标准 WMT 2014 英语-德语数据集上训练, 该数据集包含约 450 万个句子对。句子使用字节对编码 [3], 进行编码, 具有约 37000 个共享的源-目标词汇表。对于英语-法语, 我们使用了显著更大的 WMT 2014 英语-法语数据集, 包含 3600 万个句子, 并将标记拆分为 32000 个词块词汇 [38]。句子对按近似序列长度进行批处理。每个训练批次包含一组包含约 25000 个源标记和 25000 个目标标记的句子对。

5.2 硬件与时间表

我们在一台机器上的 8 个 NVIDIA P100 GPU 上训练我们的模型。对于整个论文中描述的超参数所使用的基模型, 每一步训练大约需要 0.4 秒。我们训练基模型总共 100,000 步或 12 小时。对于我们的大型模型 (如表 3 最下面一行所述), 步长时间为 1.0 秒。大型模型训练了 300,000 步 (3.5 天)。

5.3 优化器

我们使用了 Adam 优化器 [20], 结合 $\beta_1 = 0.9$ 、 $\beta_2 = 0.98$ 和 $\epsilon = 10^{-9}$ 。我们在训练过程中调整学习率, 具体公式如下:

$$lrate = d_{\text{model}}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup_steps^{-1.5}) \quad (3)$$

这相当于在最初的 $warmup_steps$ 个训练步骤中线性增加学习率, 之后按步骤编号的平方根倒数比例递减。我们使用了 $warmup_steps = 4000$ 。

5.4 正则化

我们在训练过程中采用三种正则化方法:

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

Residual Dropout We apply dropout [33] 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 sums of the embeddings and the positional encodings in both the encoder and decoder stacks. For the base model, we use a rate of $P_{drop} = 0.1$.

Label Smoothing During training, we employed label smoothing of value $\epsilon_{ls} = 0.1$ [36]. This hurts perplexity, as the model learns to be more unsure, but improves accuracy and BLEU score.

6 Results

6.1 Machine Translation

On the WMT 2014 English-to-German translation task, the big transformer model (Transformer (big) in Table 2) outperforms the best previously reported models (including ensembles) by more than 2.0 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 took 3.5 days on 8 P100 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 41.0, outperforming all of the previously published single models, at less than 1/4 the training cost of the previous state-of-the-art model. The Transformer (big) model trained for English-to-French used dropout rate $P_{drop} = 0.1$, 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 20 checkpoints. We used beam search with a beam size of 4 and length penalty $\alpha = 0.6$ [38]. 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 [38].

Table 2 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 GPU⁵.

6.2 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

⁵We used values of 2.8, 3.7, 6.0 and 9.5 TFLOPS for K80, K40, M40 and P100, respectively.

表2: Transformer在英语到德语和英语到法语的newstest2014测试中, 以训练成本的一小部分, 取得了比先前当前最佳模型更好的BLEU分数。

模型	BLEU		训练成本 (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk 集成 [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + 集成 [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S集成 [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (基础模型)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (大模型)	28.4	41.8	$2.3 \cdot 10^{19}$	

残差Dropout 我们对每个子层的输出应用Dropout [33], 在它被加到子层输入并归一化之前。此外, 我们还在编码器和解码器堆栈中的嵌入和位置编码的和中应用Dropout。对于基础模型, 我们使用 $P_{drop} = 0.1$ 的比率。

标签平滑 在训练过程中, 我们采用了值为 $\epsilon_{ls} = 0.1$ [36]的标签平滑。这会降低困惑度, 因为模型会变得不那么确定, 但能提高准确率和BLEU分数。

6 结果

6.1 机器翻译

在 WMT 2014 英语到德语翻译任务上, 表2中的大型 Transformer 模型 (Transformer (big)) 比之前报道的所有最佳模型 (包括集成模型) 多出超过 2.0 BLEU, 创下了新的 28.4 BLEU 最佳成绩。该模型的配置列于表3的最底行。在 8 个 P100 GPU 上训练耗时 3.5 天。即使我们的基础模型也超越了所有之前发表的模型和集成模型, 其训练成本远低于任何竞争模型。

在 WMT 2014 英语到法语翻译任务上, 我们的大型模型取得了 41.0 的 BLEU 分数, 超越了所有之前发表的单一模型, 且训练成本不到 1/4 之前最佳模型的成本。用于英语到法语翻译的 Transformer (big) 模型使用了 dropout 率 $P_{drop} = 0.1$, 而不是 0.3。

对于基础模型, 我们使用了通过平均最后5个检查点 (这些检查点以10分钟间隔写入) 得到的单个模型。对于大型模型, 我们平均了最后20个检查点。我们使用了beam search, beam大小为4, 并应用了长度惩罚 $\alpha = 0.6$ [38]。这些超参数是在开发集上进行实验后选择的。我们在推理过程中将最大输出长度设置为输入长度 + 50, 但在可能的情况下提前终止 [38]。

表 2 总结了我们的结果, 并将我们的翻译质量和训练成本与其他文献中的模型架构进行了比较。我们通过将训练时间、使用的 GPU 数量以及每个 GPU 的持续单精度浮点性能估计值相乘, 来估计训练模型所使用的浮点运算次数⁵。

6.2 模型变体

为了评估 Transformer 不同组件的重要性, 我们以不同方式调整了基础模型, 测量了在英语到德语翻译任务上开发集 newstest2013 的性能变化

⁵我们使用了2.8和3.7, 6.0和9.5 TFLOPS用于K80、K40、M40和P100, 分别。

Table 3: Variations on the Transformer architecture. Unlisted values are identical to those of the base model. All metrics are on the English-to-German translation development set, newstest2013. Listed perplexities are per-wordpiece, according to our byte-pair encoding, and should not be compared to per-word perplexities.

	N	d_{model}	d_{ff}	h	d_k	d_v	P_{drop}	ϵ_{ls}	train steps	PPL (dev)	BLEU (dev)	params $\times 10^6$
base	6	512	2048	8	64	64	0.1	0.1	100K	4.92	25.8	65
(A)				1 4 16 32	512 128 32 16	512 128 32 16				5.29 5.00 4.91 5.01	24.9 25.5 25.8 25.4	
(B)					16 32					5.16 5.01	25.1 25.4	58 60
(C)	2 4 8	256 1024			32 128	32 128				6.11 5.19 4.88 5.75 4.66 5.12 4.75	23.7 25.3 25.5 24.5 26.0 25.4 26.2	36 50 80 28 168 53 90
(D)							0.0 0.2	0.0 0.2		5.77 4.95 4.67 5.47	24.6 25.5 25.3 25.7	
(E)		positional embedding instead of sinusoids								4.92	25.7	
big	6	1024	4096	16			0.3		300K	4.33	26.4	213

development set, newstest2013. We used beam search as described in the previous section, but no 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 0.9 BLEU worse than the best setting, quality also drops off with too many heads.

In Table 3 rows (B), we observe that reducing the attention key size d_k hurts model quality. This suggests that determining compatibility is not easy and that a more sophisticated compatibility function than dot product may be beneficial. We further observe in rows (C) and (D) that, as expected, bigger models are better, and dropout is very helpful in avoiding over-fitting. In row (E) we replace our sinusoidal positional encoding with learned positional embeddings [9], and observe nearly identical results to the base model.

6.3 English Constituency Parsing

To evaluate if the Transformer can generalize to other tasks we performed experiments on English constituency parsing. This task presents specific challenges: the output is subject to strong structural constraints and is significantly longer than the input. Furthermore, RNN sequence-to-sequence models have not been able to attain state-of-the-art results in small-data regimes [37].

We trained a 4-layer transformer with $d_{\text{model}} = 1024$ on the Wall Street Journal (WSJ) portion of the Penn Treebank [25], about 40K training sentences. We also trained it in a semi-supervised setting, using the larger high-confidence and BerkleyParser corpora from with approximately 17M sentences [37]. We used a vocabulary of 16K tokens for the WSJ only setting and a vocabulary of 32K tokens for the semi-supervised setting.

We performed only a small number of experiments to select the dropout, both attention and residual (section 5.4), learning rates and beam size on the Section 22 development set, all other parameters remained unchanged from the English-to-German base translation model. During inference, we

表3: Transformer架构的变体。未列出的值与基础模型相同。所有指标均在英语到德语翻译开发集newstest2013上。列出的困惑度是按我们的字节对编码的每个词片段计算的, 不应与每个词的困惑度进行比较。

	N	$d_{\text{模型}}$	d_{ff}	h	d_k	d_v	P_{drop}	ϵ_{ls}	训练步数	PPL (开发)	BLEU (dev)	参数 $\times 10^6$
base	6	512	2048	8	64	64	0.1	0.1	100K	4.92	25.8	65
(A)				1 4 16 32	512 128 32 16	512 128 32 16				5.29 5.00 4.91 5.01	24.9 25.5 25.8 25.4	
(B)					16 32					5.16 5.01	25.1 25.4	58 60
(C)	2 4 8	256 1024			32 128	32 128				6.11 5.19 4.88 5.75 4.66 5.12 4.75	23.7 25.3 25.5 24.5 26.0 25.4 26.2	36 50 80 28 168 53 90
(D)							0.0 0.2	0.0 0.2		5.77 4.95 4.67 5.47	24.6 25.5 25.3 25.7	
(E)		位置嵌入而非正弦函数								4.92	25.7	
big	6	1024	4096	16			0.3		300K	4.33	26.4	213

我们使用了前文描述的 beam search 方法, 但没有使用模型检查点平均。我们将这些结果展示在表 3 中

在表3的行(A)中, 我们变化了注意力头的数量以及注意力键和值的维度, 同时保持计算量恒定, 如第3.2.2节所述。虽然单头注意力比最佳设置差0.9 BLEU, 但过多的头也会导致质量下降。

在表3的行(B)中, 我们观察到减少注意力键大小 d_k 会损害模型质量。这表明确定兼容性并不容易, 并且可能需要比点积更复杂的兼容性函数。我们进一步观察到在行(C)和(D)中, 正如预期的那样, 更大的模型效果更好, dropout在避免过拟合方面非常有帮助。在行(E)中, 我们用学习位置嵌入 [9], 替换了我们的正弦位置编码, 并观察到与基础模型几乎相同的结果。

6.3 英语依存句法分析

为了评估Transformer是否可以泛化到其他任务, 我们在英语依存句法分析上进行了实验。这项任务提出了特定的挑战: 输出受强结构约束, 并且显著长于输入。此外, RNN序列到序列模型在小数据环境下尚未达到最先进的结果 [37]。

我们使用 $d_{\text{model}} = 1024$ 在华尔街日报 (WSJ) 部分的宾夕法尼亚树库 [25], 上训练了一个 4层Transformer, 大约有40K个训练句子。我们还以半监督的方式训练它, 使用来自的更大规模高置信度和伯克利解析器语料库, 大约有17M个句子[37]。对于WSJ仅设置, 我们使用了16K个token的词汇表; 对于半监督设置, 我们使用了32K个token的词汇表。

我们仅进行少量实验来选择 dropout、注意力机制和残差 (第 5.4 节)、学习率以及束大小, 这些实验是在第 22 个开发集上进行的, 所有其他参数与英语到德语的基准翻译模型保持不变。在推理过程中, 我们将最大输出长度调整为输入长度 {v1}。我们使用了 21 的束大小和 {v2}

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

Parser	Training	WSJ 23 F1
Vinyals & Kaiser et al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser et al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

increased the maximum output length to input length + 300. We used a beam size of 21 and $\alpha = 0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8].

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

7 Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goal of ours.

The code we used to train and evaluate our models is available at <https://github.com/tensorflow/tensor2tensor>.

Acknowledgements We are grateful to Nal Kalchbrenner and Stephan Gouws for their fruitful comments, corrections and inspiration.

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 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] Denny Britz, Anna Goldie, Minh-Thang Luong, and Quoc V. Le. Massive exploration of 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.

表4: Transformer对英语依存句法分析泛化效果良好 (结果来自WSJ第23节)

Parser	Training	WSJ 23 F1
Vinyals & Kaiser等 (2014) [37]	WSJ仅, 判别式	88.3
Petrov等人 (2006) [29]	WSJ仅, 判别性	90.4
Zhu等人 (2013) [40]	WSJ仅, 判别性	90.4
Dyer等人 (2016) [8]	WSJ仅, 判别性	91.7
Transformer (4层)	WSJ仅, 判别性	91.3
Zhu等人 (2013) [40]	半监督	91.3
黄 & 哈珀 (2009) [14]	半监督	91.3
McClosky等人 (2006) [26]	半监督	92.1
Vinyals & Kaiser等人 (2014) [37]	半监督	92.1
Transformer (4层)	半监督	92.7
Luong等人 (2015) [23]	多任务	93.0
Dyer等人 (2016) [8]	生成式	93.3

incr将最大输出长度放宽至输入长度 + 300。我们使用了 21 的束大小和 $\alpha = 0.3$ for WSJ仅和半监督设置。 3

表4中的结果显示, 尽管缺乏任务特定调优, 我们的模型表现惊人地好, 其结果优于所有先前报告的模型, 例外是循环神经网络语法 [8]。

与RNN序列到序列模型 [37]不同, Transformer甚至在使用仅包含40K句子的WSJ训练集进行训练时, 也优于Berkeley-Parser [29]。

7 结论

在这项工作中, 我们提出了Transformer, 这是第一个完全基于注意力的序列转换模型, 它用多头自注意力机制替换了编码器-解码器架构中最常用的循环层。

对于翻译任务, Transformer的训练速度比基于循环层或卷积层的架构快得多。在WMT 2014英语到德语和WMT 2014英语到法语的翻译任务中, 我们实现了新的技术突破。在后者任务中, 我们的最佳模型甚至优于所有先前报告的集成模型。

我们对基于注意力的模型的未来充满期待, 并计划将它们应用于其他任务。我们计划将Transformer扩展到涉及文本以外的输入输出模态的问题, 并研究局部受限注意力机制, 以高效处理图像、音频和视频等大型输入输出。使生成过程减少序列性也是我们另一个研究目标。

我们用于训练和评估模型的代码可在 <https://github.com/tensorflow/tensor2tensor> 获取。

致谢 我们感谢 Nal Kalchbrenner 和 Stephan Gouws 的富有成效的评论、修正和启发! comments、修正和启发。

参考文献

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, 以及Geoffrey E Hinton. 层归一化. *arXiv预印本arXiv:1607.06450*, 2016年.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, 以及Yoshua Bengio. 通过联合学习对齐和翻译进行神经机器翻译. *CoRR*, abs/1409.0473, 2014年.
- [3] Denny Britz, Anna Goldie, Minh-Thang Luong, 以及Quoc V. Le. 大规模探索神经机器翻译架构. *CoRR*, abs/1703.03906, 2017年.
- [4] Jianpeng Cheng, Li Dong, 以及Mirella Lapata. 用于机器阅读的长短期记忆网络. *arXiv预印本arXiv:1601.06733*, 2016年.

- [5] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *CoRR*, abs/1406.1078, 2014.
- [6] Francois Chollet. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint arXiv:1610.02357*, 2016.
- [7] Junyoung Chung, Çağlar Gülçehre, Kyunghyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014.
- [8] Chris Dyer, Adhiguna Kuncoro, Miguel Ballesteros, and Noah A. Smith. Recurrent neural network grammars. In *Proc. of NAACL*, 2016.
- [9] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. Convolutional sequence to sequence learning. *arXiv preprint arXiv:1705.03122v2*, 2017.
- [10] Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.
- [11] Kaiming 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.
- [12] Sepp Hochreiter, Yoshua Bengio, Paolo Frasconi, and Jürgen Schmidhuber. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001.
- [13] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [14] Zhongqiang Huang and Mary Harper. Self-training PCFG grammars with latent annotations across languages. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 832–841. ACL, August 2009.
- [15] Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. Exploring the limits of language modeling. *arXiv preprint arXiv:1602.02410*, 2016.
- [16] Łukasz Kaiser and Samy Bengio. Can active memory replace attention? In *Advances in Neural Information Processing Systems*, (NIPS), 2016.
- [17] Łukasz Kaiser and Ilya Sutskever. Neural GPUs learn algorithms. In *International Conference on Learning Representations (ICLR)*, 2016.
- [18] Nal Kalchbrenner, Lasse Espeholt, Karen Simonyan, Aaron van den Oord, Alex Graves, and Koray Kavukcuoglu. Neural machine translation in linear time. *arXiv preprint arXiv:1610.10099v2*, 2017.
- [19] Yoon Kim, Carl Denton, Luong Hoang, and Alexander M. Rush. Structured attention networks. In *International Conference on Learning Representations*, 2017.
- [20] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [21] Oleksii Kuchaiev and Boris Ginsburg. Factorization tricks for LSTM networks. *arXiv preprint arXiv:1703.10722*, 2017.
- [22] Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. *arXiv preprint arXiv:1703.03130*, 2017.
- [23] Minh-Thang Luong, Quoc V. Le, Ilya Sutskever, Oriol Vinyals, and Łukasz Kaiser. Multi-task sequence to sequence learning. *arXiv preprint arXiv:1511.06114*, 2015.
- [24] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015.

- [5] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Fethi Bougares, Holger Schwenk, 以及 Yoshua Bengio. 使用RNN编码器-解码器学习短语表示以进行统计机器翻译. *CoRR*, abs/1406.1078, 2014年. [6] Francois Chollet. Xception: 使用深度可分离卷积的深度学习. *arXiv预印本 arXiv:1610.02357*, 2016年. [7] Junyoung Chung, Çağlar Gülçehre, Kyunghyun Cho, 以及 Yoshua Bengio. 门控循环神经网络在序列建模上的经验评估. *CoRR*, abs/1412.3555, 2014年. [8] Chris Dyer, Adhiguna Kuncoro, Miguel Ballesteros, 以及 Noah A. Smith. 循环神经网络语法. *Proc. of NAACL*, 2016年. [9] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, 以及 Yann N. Dauphin. 卷积序列到序列学习. *arXiv预印本 arXiv:1705.03122v2*, 2017年. [10] Alex Graves. 使用循环神经网络生成序列. *arXiv预印本 arXiv:1308.0850*, 2013年. [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, 以及 Jian Sun. 用于图像识别的深度残差学习. *IEEE 计算机视觉与模式识别会议论文集*, 第770–778页, 2016年. [12] Sepp Hochreiter, Yoshua Bengio, Paolo Frasconi, 以及 Jürgen Schmidhuber. 循环网络中的梯度流: 学习长期依赖的困难, 2001年. [13] Sepp Hochreiter 和 Jürgen Schmidhuber. 长短期记忆. *神经计算*, 9(8):1735–1780, 1997年. [14] Zhongqiang Huang 和 Mary Harper. 跨语言使用潜在标注自训练PCFG语法. *2009年自然语言处理经验方法会议论文集*, 第832–841页. ACL, 2009年8月. [15] Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, 以及 Yonghui Wu. 探索语言建模的极限. *arXiv预印本 arXiv:1602.02410*, 2016年. [16] Łukasz Kaiser 和 Samy Bengio. 主动记忆能替代注意力吗? *神经信息处理系统进展 (NIPS)*, 2016年. [17] Łukasz Kaiser 和 Ilya Sutskever. 神经GPU学习算法. *国际学习表示会议 (ICLR)*, 2016年. [18] Nal Kalchbrenner, Lasse Espeholt, Karen Simonyan, Aaron van den Oord, Alex Graves, 以及 Koray Kavukcuoglu. 线性时间的神经机器翻译. *arXiv预印本 arXiv:1610.10099v2*, 2017年. [19] Yoon Kim, Carl Denton, Luong Hoang, 以及 Alexander M. Rush. 结构化注意力网络. *国际学习表示会议*, 2017年. [20] Diederik Kingma 和 Jimmy Ba. Adam: 一种随机优化方法. *ICLR*, 2015年. [21] Oleksii Kuchaiev 和 Boris Ginsburg. LSTM网络的分解技巧. *arXiv预印本 arXiv:1703.10722*, 2017年. [22] Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, 以及 Yoshua Bengio. 结构化自注意力句子嵌入. *arXiv预印本 arXiv:1703.03130*, 2017年. [23] Minh-Thang Luong, Quoc V. Le, Ilya Sutskever, Oriol Vinyals, 以及 Łukasz Kaiser. 多任务序列到序列学习. *arXiv预印本 arXiv:1511.06114*, 2015年. [24] Minh-Thang Luong, Hieu Pham, 以及 Christopher D Manning. 基于注意力的神经机器翻译的有效方法. *arXiv预印本 arXiv:1508.04025*, 2015年.

- [25] Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 19(2):313–330, 1993.
- [26] David McClosky, Eugene Charniak, and Mark Johnson. Effective self-training for parsing. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 152–159. ACL, June 2006.
- [27] Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model. In *Empirical Methods in Natural Language Processing*, 2016.
- [28] Romain Paulus, Caiming Xiong, and Richard Socher. A deep reinforced model for abstractive summarization. *arXiv preprint arXiv:1705.04304*, 2017.
- [29] Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. Learning accurate, compact, and interpretable tree annotation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the ACL*, pages 433–440. ACL, July 2006.
- [30] Ofir Press and Lior Wolf. Using the output embedding to improve language models. *arXiv preprint arXiv:1608.05859*, 2016.
- [31] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*, 2015.
- [32] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.
- [33] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [34] Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, and Rob Fergus. End-to-end memory networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2440–2448. Curran Associates, Inc., 2015.
- [35] Ilya Sutskever, Oriol Vinyals, and Quoc VV Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112, 2014.
- [36] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. *CoRR*, abs/1512.00567, 2015.
- [37] Vinyals & Kaiser, Koo, Petrov, Sutskever, and Hinton. Grammar as a foreign language. In *Advances in Neural Information Processing Systems*, 2015.
- [38] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- [39] Jie Zhou, Ying Cao, Xuguang Wang, Peng Li, and Wei Xu. Deep recurrent models with fast-forward connections for neural machine translation. *CoRR*, abs/1606.04199, 2016.
- [40] Muhua Zhu, Yue Zhang, Wenliang Chen, Min Zhang, and Jingbo Zhu. Fast and accurate shift-reduce constituent parsing. In *Proceedings of the 51st Annual Meeting of the ACL (Volume 1: Long Papers)*, pages 434–443. ACL, August 2013.

- [25] Mitchell P Marcus, Mary Ann Marcinkiewicz, 和 Beatrice Santorini. 构建大型英语标注语料库: 宾夕法尼亚树库。 *计算语言学*, 19(2):313–330, 1993.[26] David McClosky, Eugene Charniak, 和 Mark Johnson. 高效的自训练句法分析。在 *NAACL 人类语言技术会议主会*, 页面 152–159. ACL, 2006 年 6 月.[27] Ankur Parikh, Oscar Täckström, Dipanjan Das, 和 Jakob Uszkoreit. 可分解注意力模型。在 *自然语言处理经验方法*, 2016.[28] Romain Paulus, Caiming Xiong, 和 Richard Socher. 用于抽象摘要的深度强化模型。 *arXiv 预印本 arXiv:1705.04304*, 2017.[29] Slav Petrov, Leon Barrett, Romain Thibaux, 和 Dan Klein. 学习准确、紧凑且可解释的树标注。在 *第 21 届计算语言学国际会议暨 ACL 第 44 届年会*, 页面 433–440. ACL, 2006 年 7 月.[30] Ofir Press 和 Lior Wolf. 使用输出嵌入改进语言模型。 *arXiv 预印本 arXiv:1608.05859*, 2016.[31] Rico Sennrich, Barry Haddow, 和 Alexandra Birch. 使用子词单元进行罕见词神经机器翻译。 *arXiv 预印本 arXiv:1508.07909*, 2015.[32] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, 和 Jeff Dean. 惊人的大神经网络: 稀疏门控专家混合层。 *arXiv 预印本 arXiv:1701.06538*, 2017.[33] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, 和 Ruslan Salakhutdinov. Dropout: 一种简单的方法来防止神经网络过拟合。 *机器学习研究杂志*, 15(1):1929–1958, 2014.[34] Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, 和 Rob Fergus. 端到端记忆网络。在 C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, 和 R. Garnett, 编辑的 *第 28 届神经信息处理系统进展*, 页面 2440–2448. Curran Associates, Inc., 2015.[35] Ilya Sutskever, Oriol Vinyals, 和 Quoc VV Le. 基于神经网络的序列到序列学习。在 *神经信息处理系统进展*, 页面 3104–3112, 2014.[36] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, 和 Zbigniew Wojna. 重新思考用于计算机视觉的 Inception 架构。 *CoRR*, abs/1512.00567, 2015.[37] Vinyals & Kaiser, Koo, Petrov, Sutskever, 和 Hinton. 语法作为一种外语。在 *神经信息处理系统进展*, 2015.[38] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, 等. Google 的神经机器翻译系统: 弥合人与机器翻译之间的差距。 *arXiv 预印本 arXiv:1609.08144*, 2016.[39] Jie Zhou, Ying Cao, Xuguang Wang, Peng Li, 和 Wei Xu. 用于神经机器翻译的具有快速前馈连接深度循环模型。 *CoRR*, abs/1606.04199, 2016.[40] Muhua Zhu, Yue Zhang, Wenliang Chen, Min Zhang, 和 Jingbo Zhu. 快速且准确的移位-规约成分句法分析。在 *第 51 届 ACL 年会 (第一卷: 长论文)*, 页面 434–443. ACL, 2013 年 8 月.

Attention Visualizations

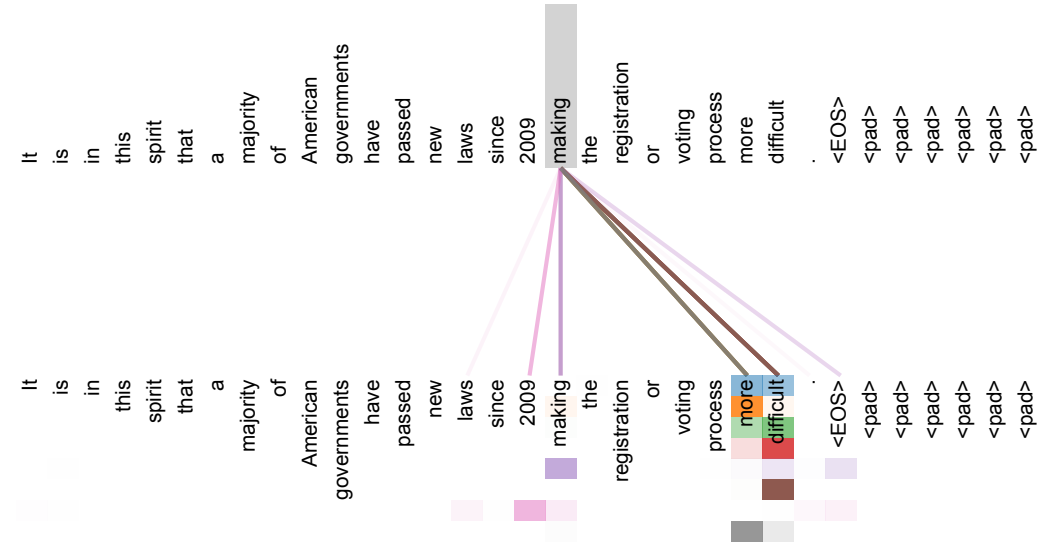


Figure 3: An example of the attention mechanism following long-distance dependencies in the encoder self-attention in layer 5 of 6. Many of the attention heads attend to a distant dependency of the verb 'making', completing the phrase 'making...more difficult'. Attentions here shown only for the word 'making'. Different colors represent different heads. Best viewed in color.

注意力可视化

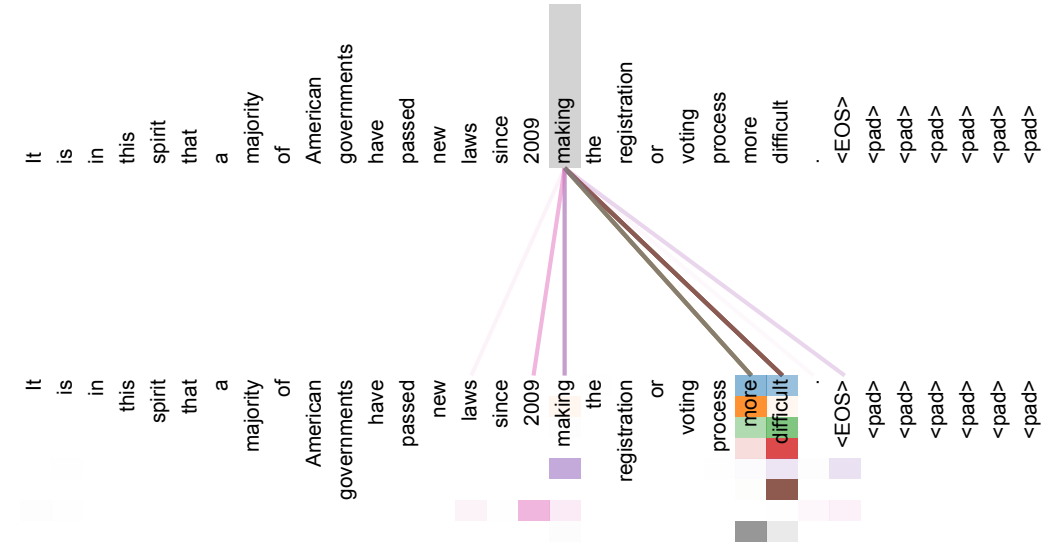


图3: 展示了第6层中的第5层编码器自注意力机制下, 注意力机制如何处理长距离依赖的示例。许多注意力头关注动词'making'的远处依赖, 完成了短语'making...更困难'。此处仅显示单词'making'的注意力。不同颜色代表不同的头。建议彩色查看。

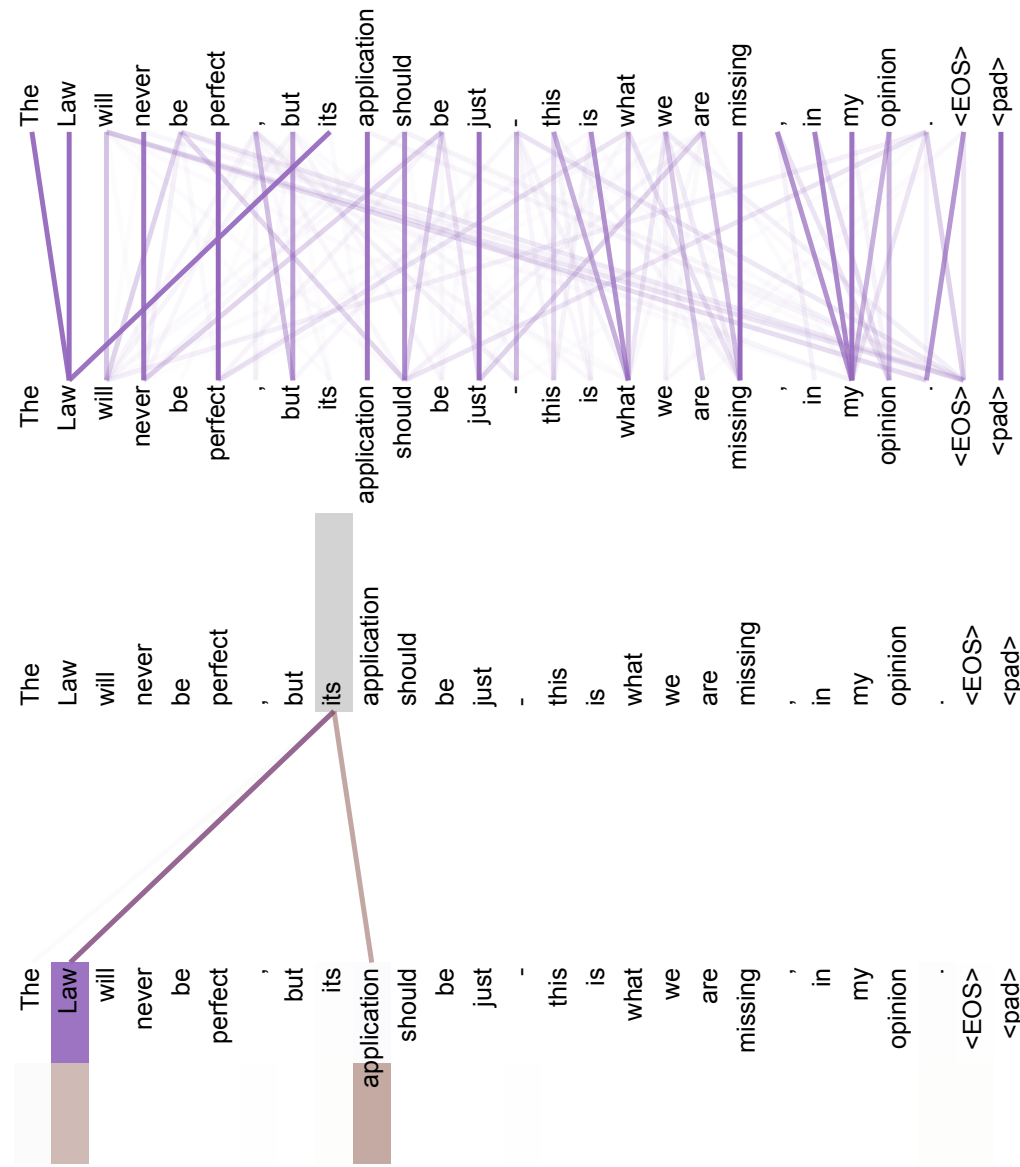


Figure 4: Two attention heads, also in layer 5 of 6, apparently involved in anaphora resolution. Top: Full attentions for head 5. Bottom: Isolated attentions from just the word 'its' for attention heads 5 and 6. Note that the attentions are very sharp for this word.

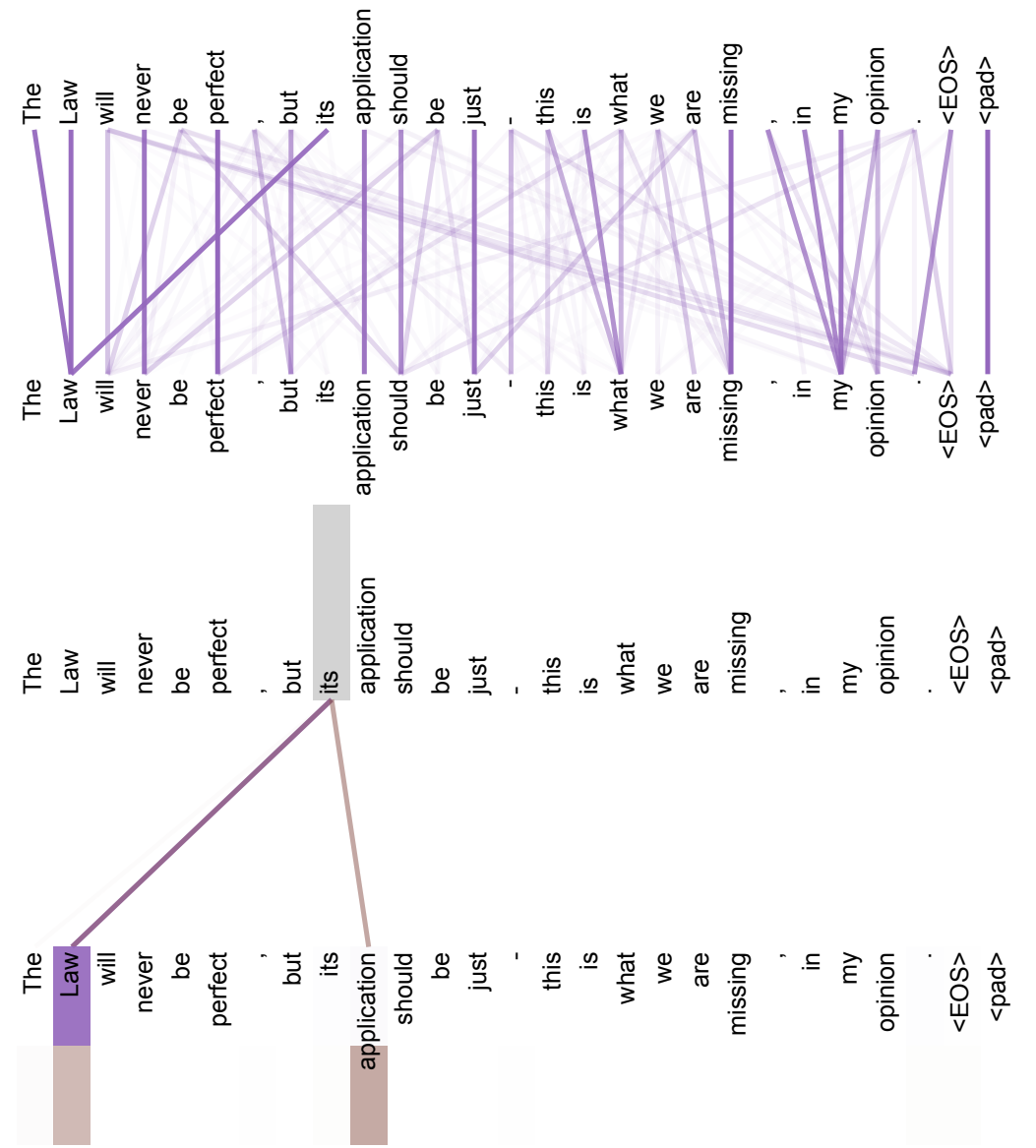


图4: 两个注意力头, 也位于第6层的第5层, 显然参与指代消解。顶部: 头5的完整注意力。底部: 仅从单词‘其’对头5和6的注意力分离。注意, 这个单词的注意力非常集中。

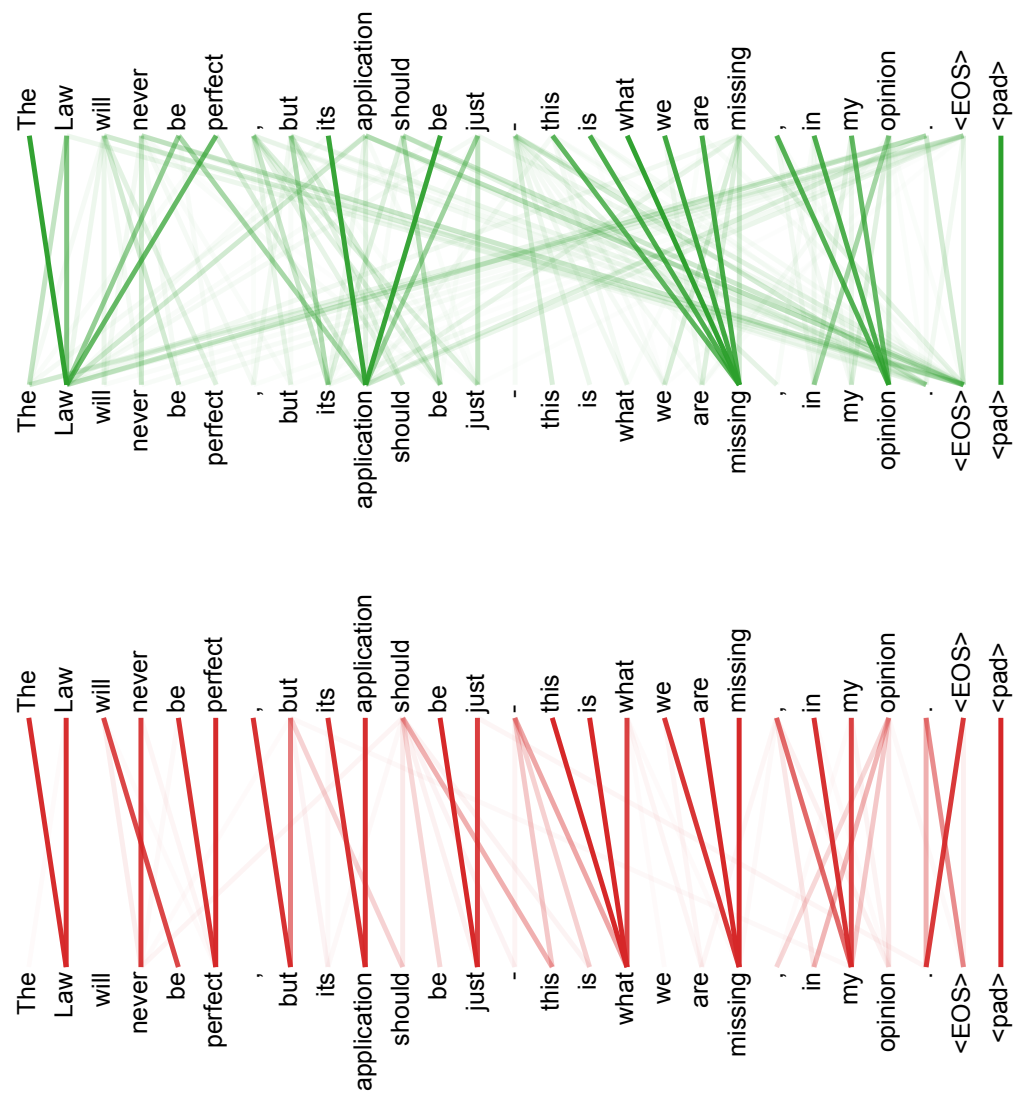


Figure 5: Many of the attention heads exhibit behaviour that seems related to the structure of the sentence. We give two such examples above, from two different heads from the encoder self-attention at layer 5 of 6. The heads clearly learned to perform different tasks.

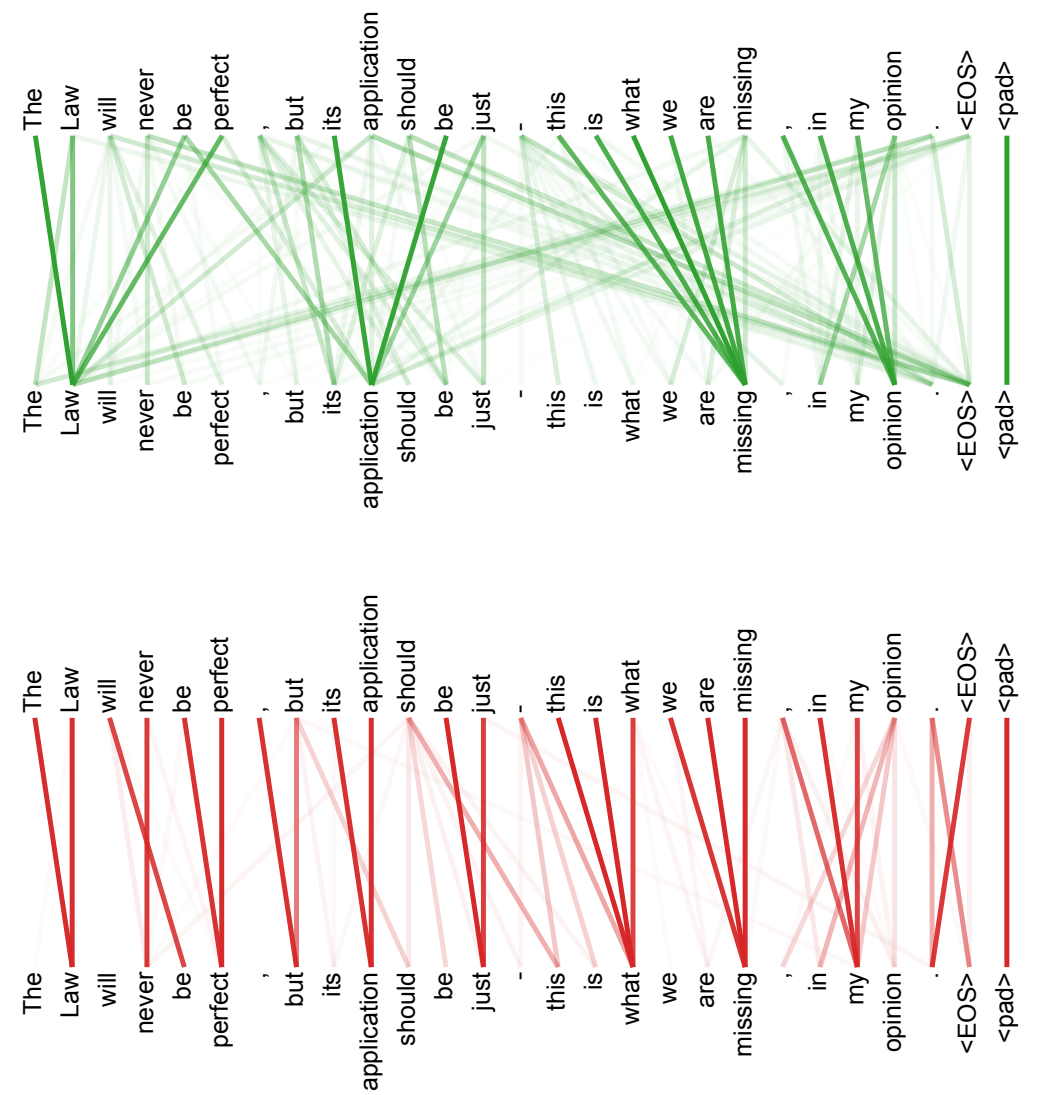


图5：许多注意力头表现出与句子结构相关的行为。我们在上方给出了两个这样的例子，来自编码器自注意力机制在第5层（共6层）的两个不同头。这些头显然学会了执行不同的任务。