

Attention is All You Need

Vaswani et al.

2017

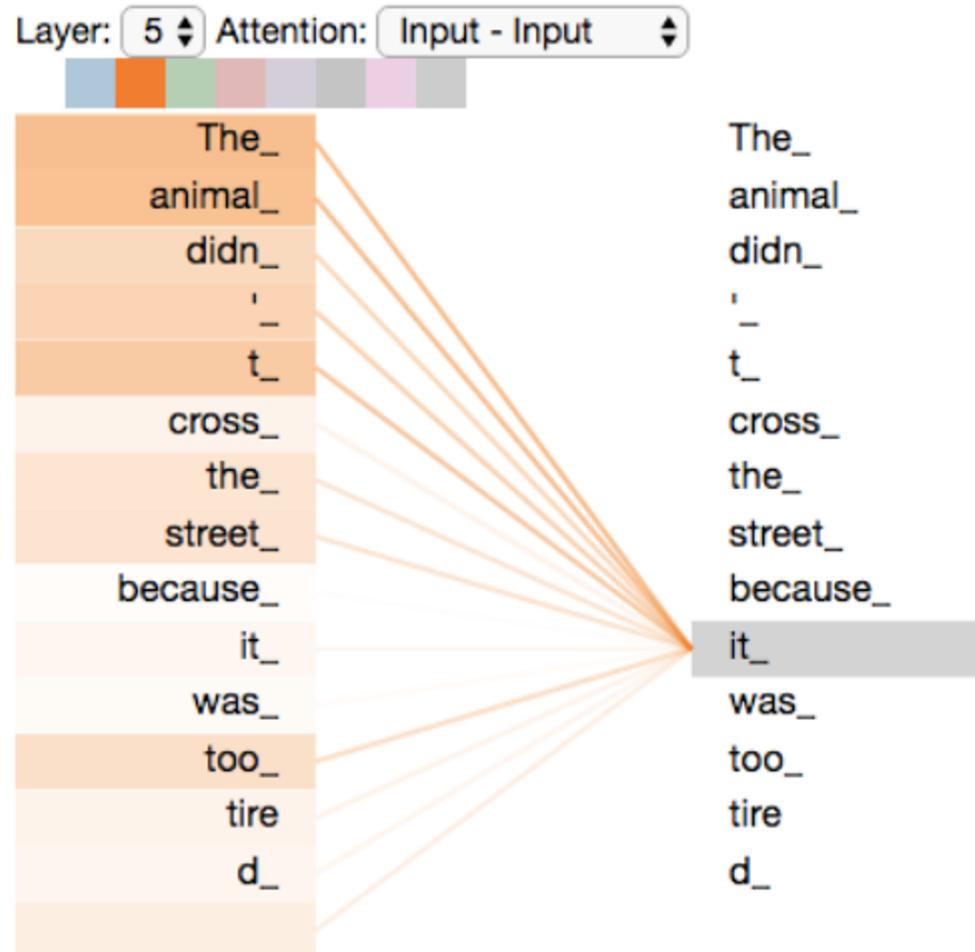


Presented by Eli Lifland, 2/23/2020

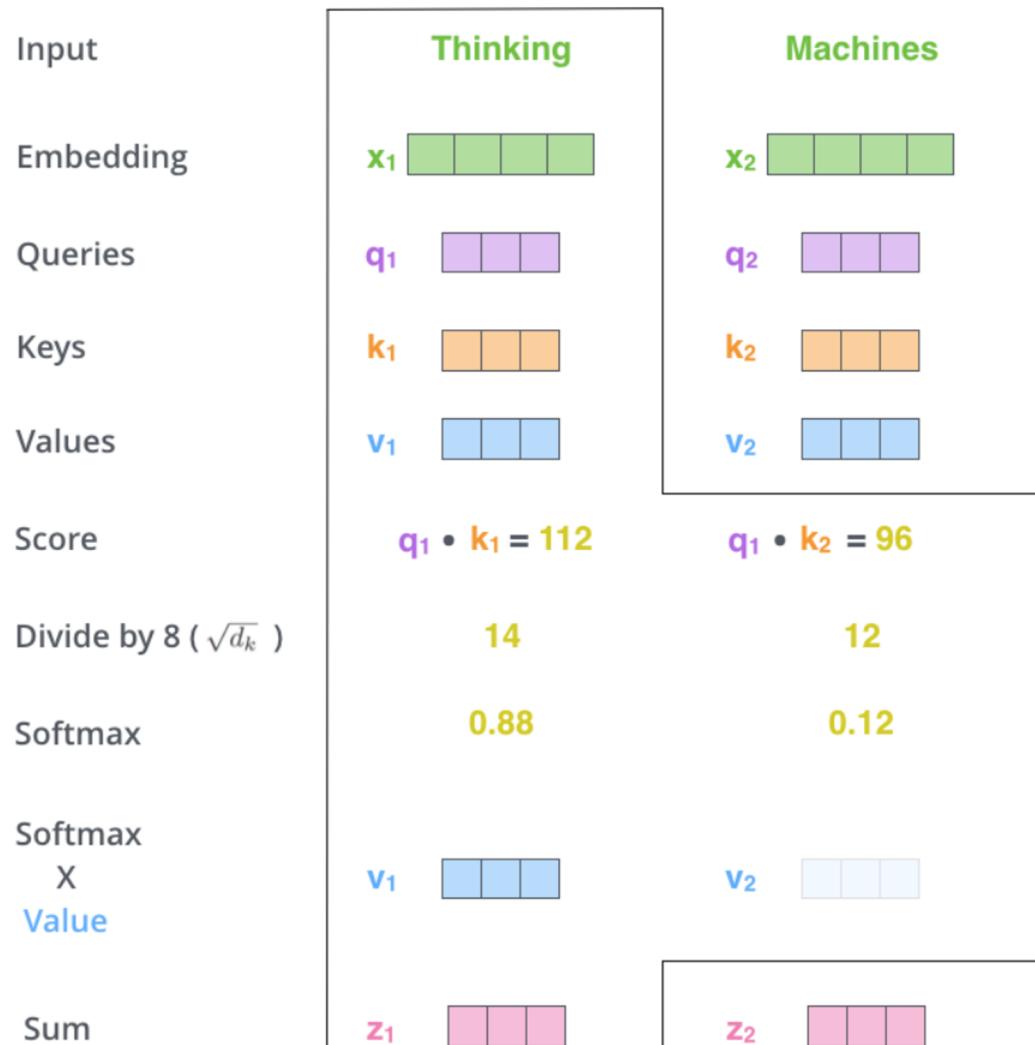
Background

- CNN-based methods popular
 - Difficult to learn dependencies between distant positions
- RNNs SOTA: LSTMs, Gated RNNs
 - Generate sequence of hidden states h_t as function of h_{t-1} , input at position t
 - Sequential nature doesn't allow for parallelization within training examples
- Encoder-decoder architectures
- Attention allows modeling of dependencies without regard to distance between positions
 - Before: used with RNNs
 - But attention is all you need! Transformers use no RNNs or convolution

Attention: General Idea



Attention: Scaled Dot-Product



Attention: Matrix Operations

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$

Diagram illustrating the matrix multiplication of input \mathbf{X} (green 4x4 grid) by weight matrix \mathbf{W}^Q (purple 4x4 grid) to produce query matrix \mathbf{Q} (purple 4x4 grid).

$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$

Diagram illustrating the matrix multiplication of input \mathbf{X} (green 4x4 grid) by weight matrix \mathbf{W}^K (orange 4x4 grid) to produce key matrix \mathbf{K} (orange 4x4 grid).

$$\text{softmax}\left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}}\right) = \mathbf{Z}$$

Diagram illustrating the computation of attention weights \mathbf{Z} (pink 4x4 grid) using the query matrix \mathbf{Q} (purple 4x4 grid) and the transpose of the key matrix \mathbf{K}^T (orange 4x4 grid), normalized by the square root of the dimension d_k .

$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$

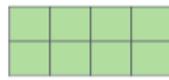
Diagram illustrating the matrix multiplication of input \mathbf{X} (green 4x4 grid) by weight matrix \mathbf{W}^V (blue 4x4 grid) to produce value matrix \mathbf{V} (blue 4x4 grid).

Attention: Multi-Head

- 1) This is our input sentence* X
- 2) We embed each word* R
- 3) Split into 8 heads. We multiply X or R with weight matrices W_0^Q, W_0^K, W_0^V , W_1^Q, W_1^K, W_1^V , ..., W_7^Q, W_7^K, W_7^V
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

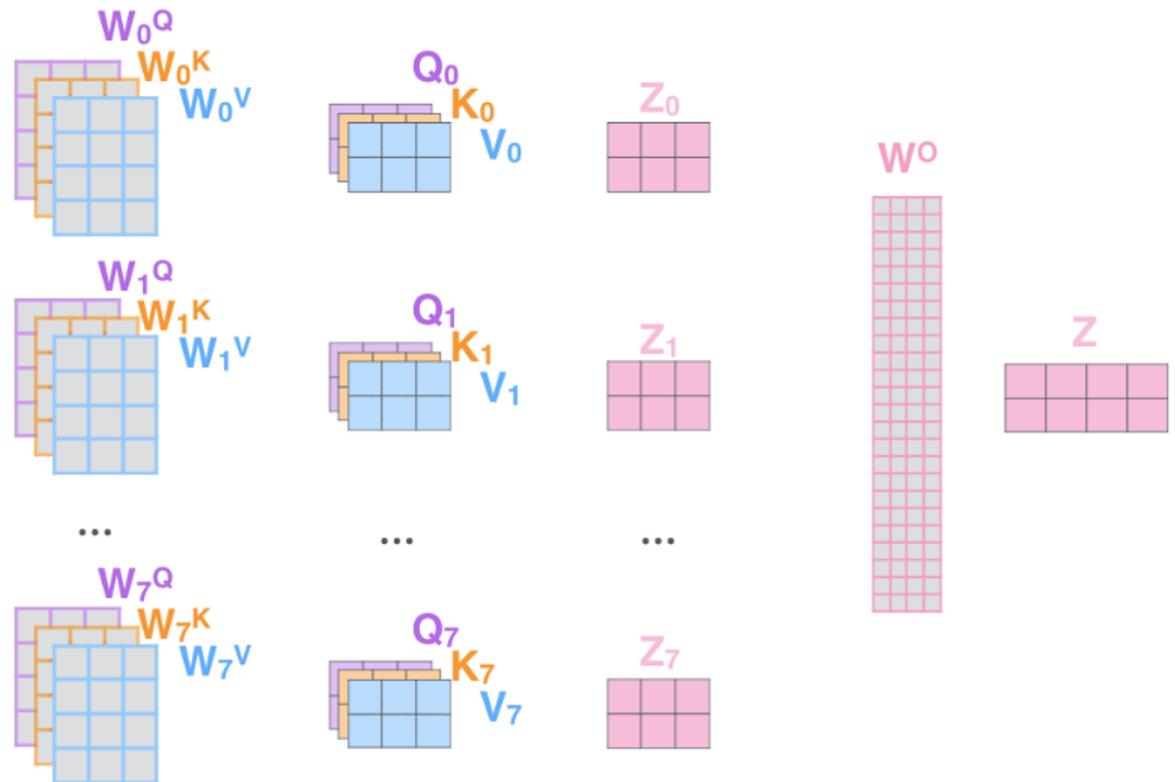
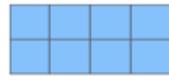
Thinking
Machines

X



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

R



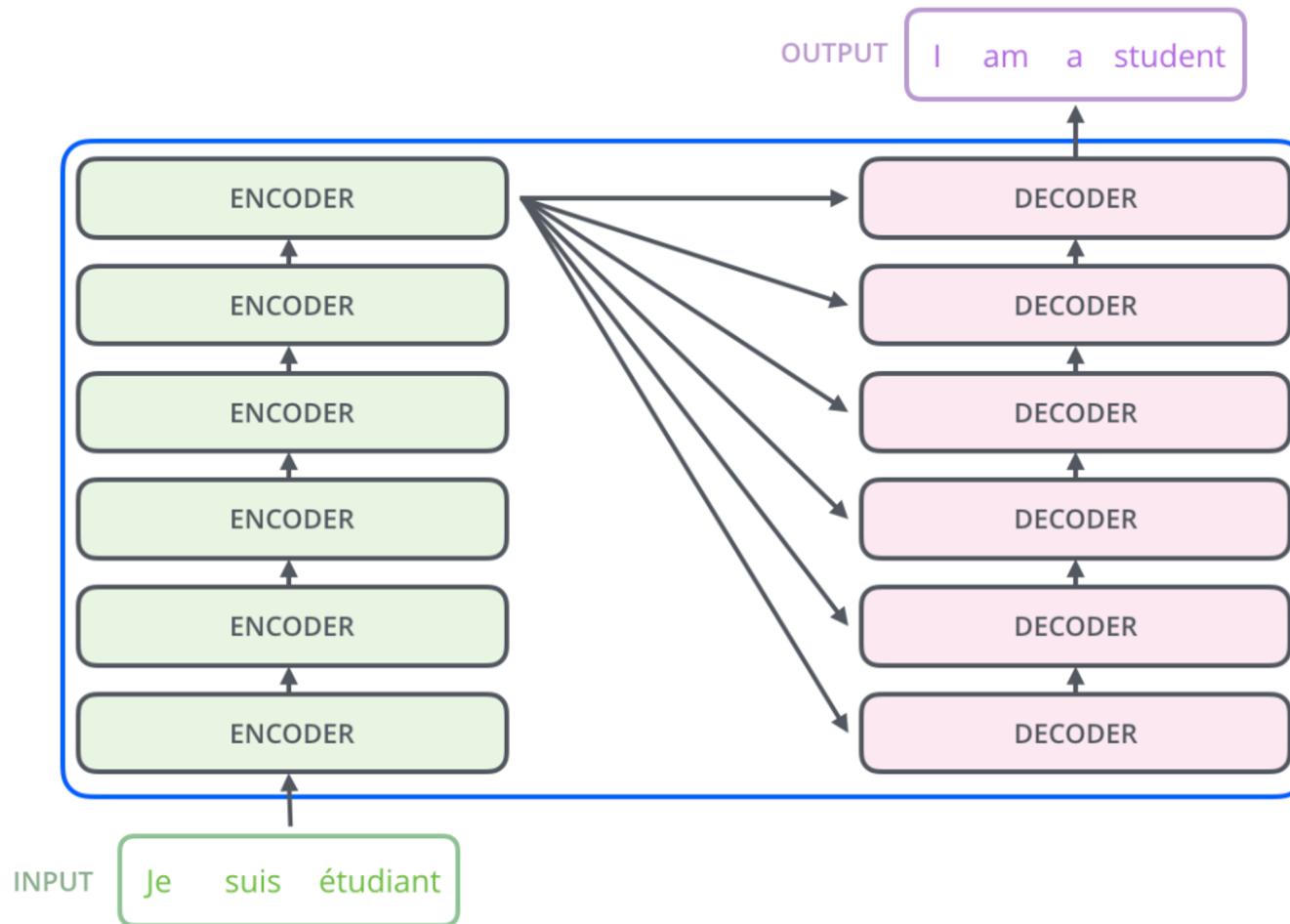
Why Attention?

- Low computational complexity per layer
- High level of parallelizability (low number of sequential operations required)
- Low path length between long-range dependencies
 - Easier to learn long-range dependencies even though in principle possible with RNNs, CNNs

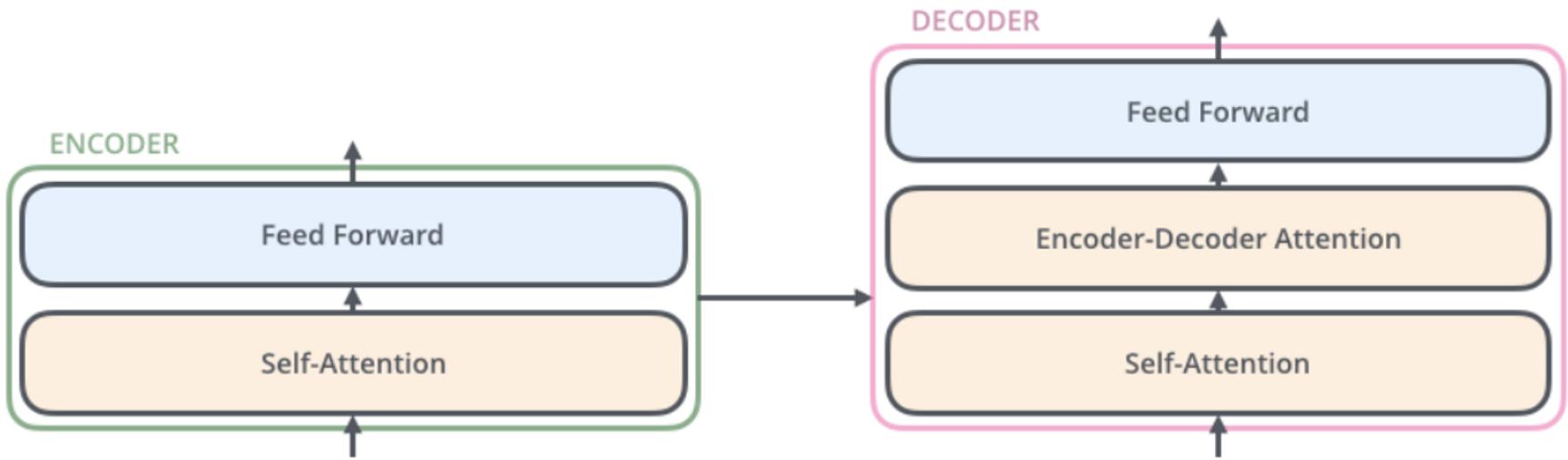
Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

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)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

Transformer Architecture: Encoder-Decoder Structure

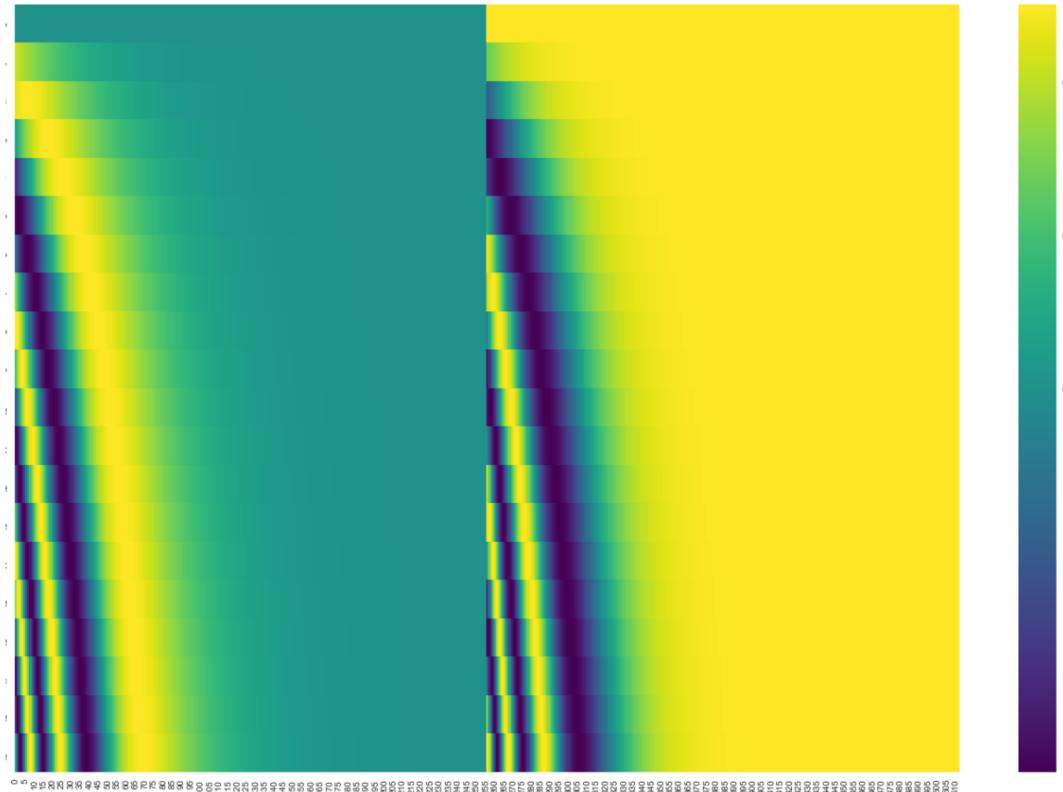


Transformer Architecture: Architecture of Encoders, Decoders



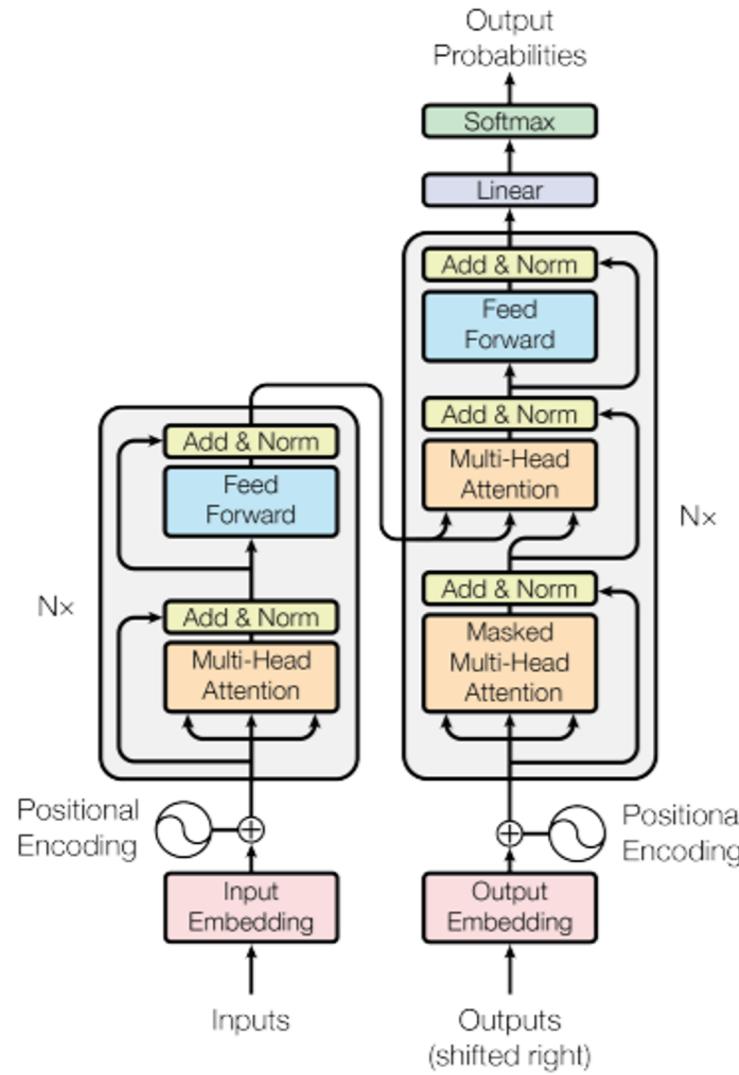
- Modify self-attention layers in decoder to prevent positions from attending to subsequent positions (masking)
- This unidirectionality later removed for BERT

Transformer Architecture: Positional Encoding



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

Transformer Architecture: Full



Results

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	$2.3 \cdot 10^{19}$	

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

Parser	Training	WSJ 23 F1
Vinyals & Kaiser el 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 el 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]	Adaptive	93.3

<https://qdata.github.io/deep2Read/>

Results: Ablation

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)										5.29	24.9	
										5.00	25.5	
										4.91	25.8	
										5.01	25.4	
(B)									16	5.16	25.1	58
									32	5.01	25.4	60
(C)									2	6.11	23.7	36
									4	5.19	25.3	50
									8	4.88	25.5	80
									256	5.75	24.5	28
									1024	4.66	26.0	168
									1024	5.12	25.4	53
									4096	4.75	26.2	90
									0.0	5.77	24.6	
(D)									0.2	4.95	25.5	
									0.0	4.67	25.3	
									0.2	5.47	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