Attention is all you need (2017)

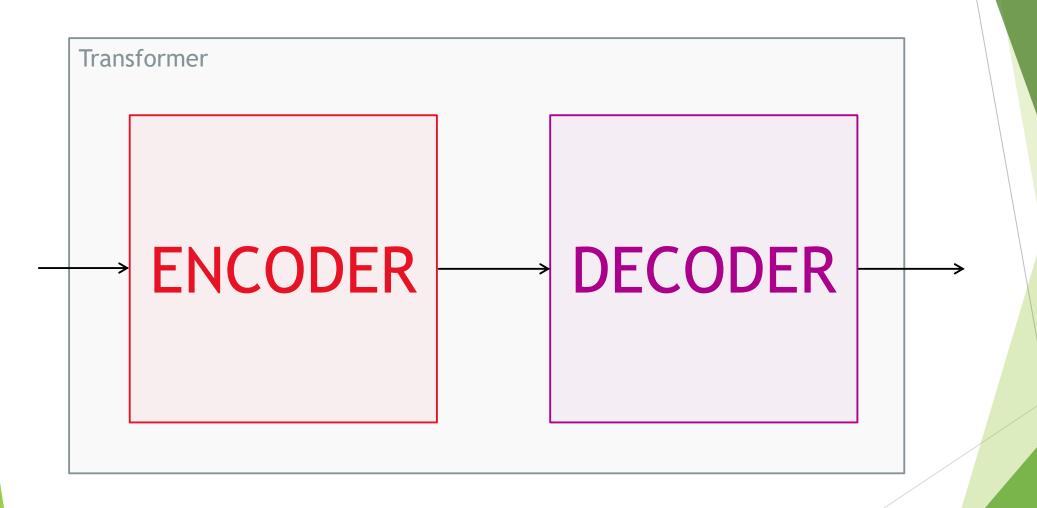
How do transformers work?

Alejandro C. Parra Garcia

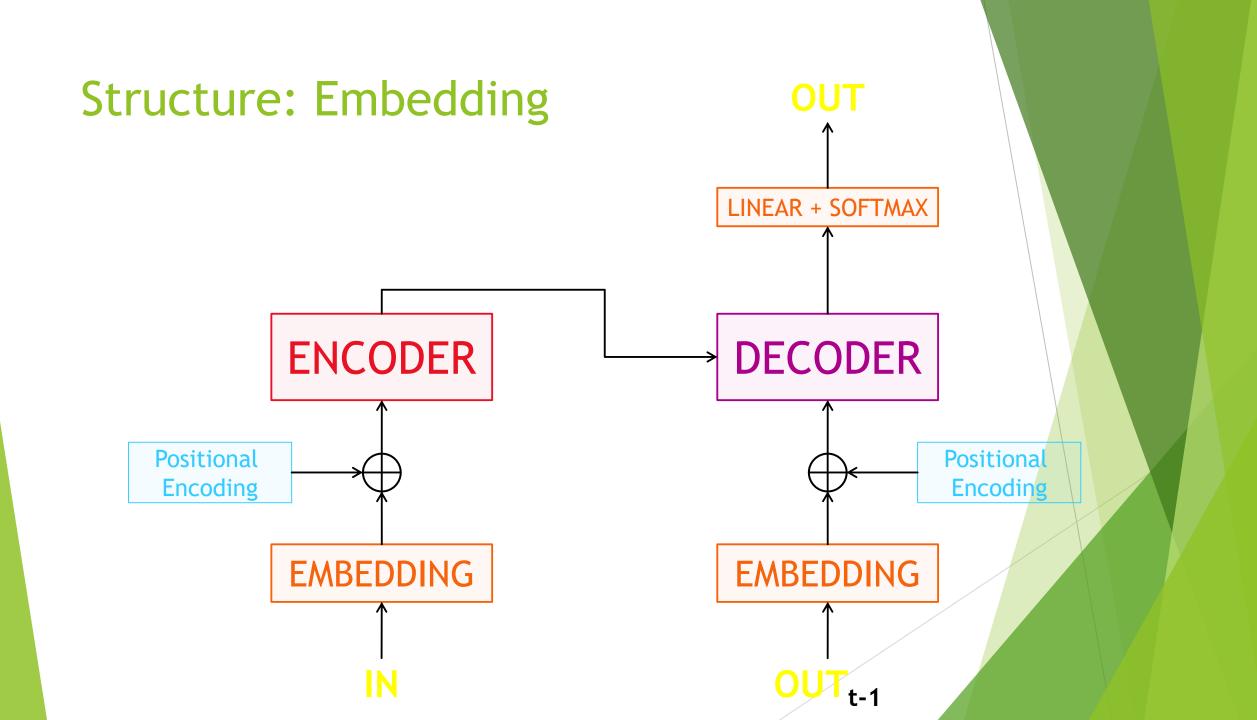
Translation model



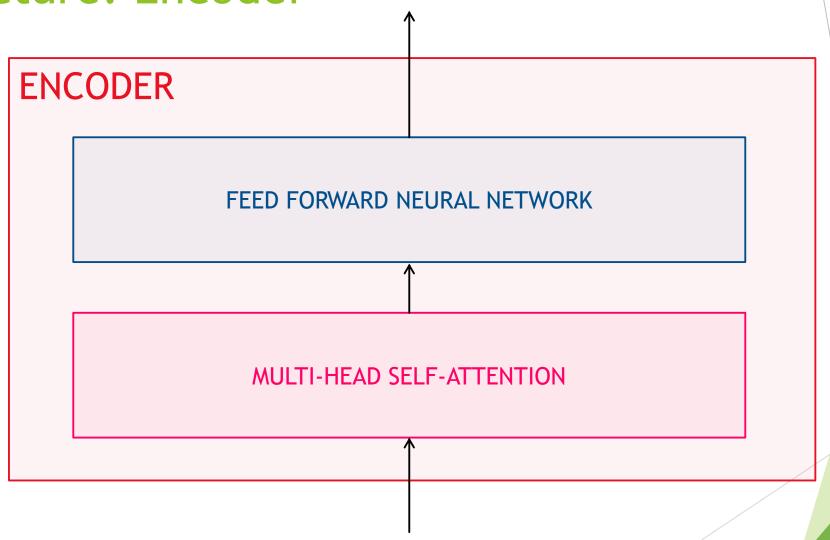
Structure



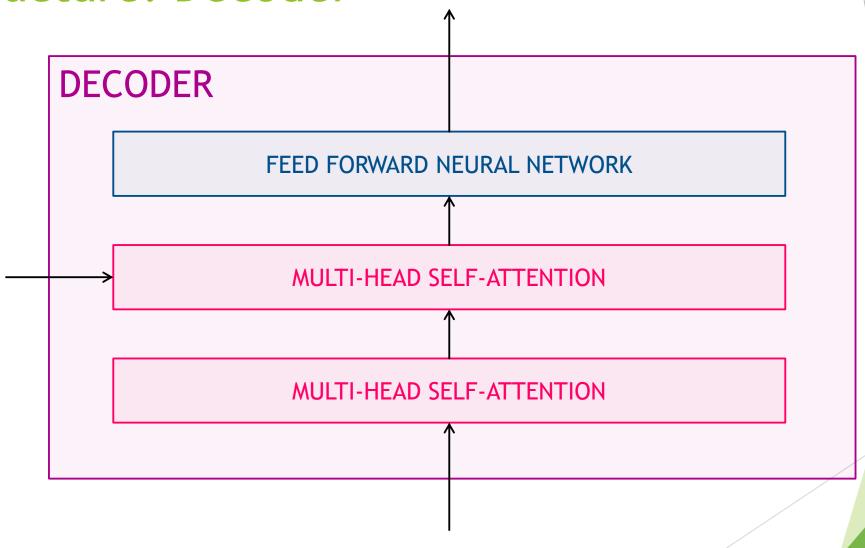
Structure **OUT ENCODER DECODER ENCODER DECODER ENCODER DECODER ENCODER DECODER ENCODER DECODER ENCODER DECODER**



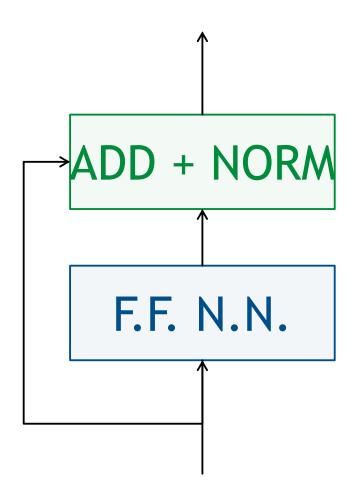
Structure: Encoder

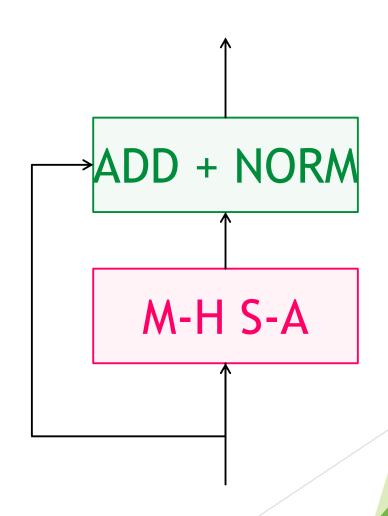


Structure: Decoder



Structure: Residuals





Output Probabilities **Complete Structure** Softmax Linear Add & Norm Feed Forward Add & Norm Add & Norm Multi-Head Feed Attention Forward N× Add & Norm N× Add & Norm Masked Multi-Head Multi-Head Attention Attention Positional Positional Encoding Encoding Output Input Embedding Embedding Inputs Outputs

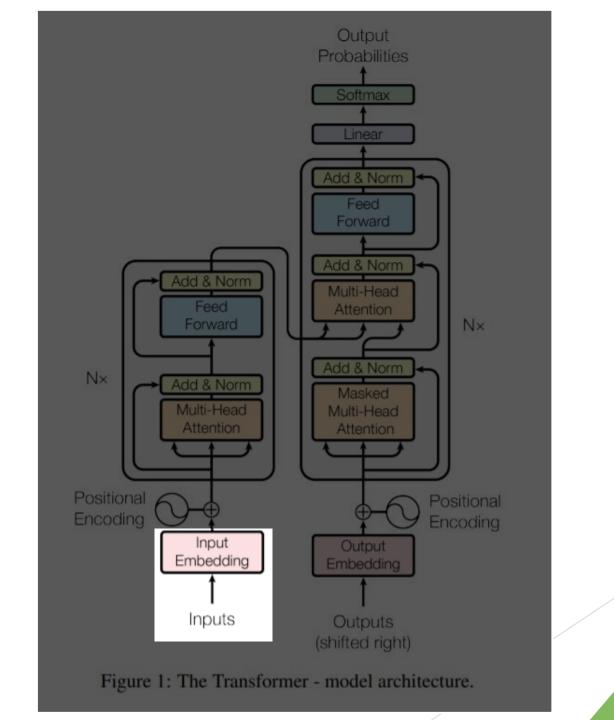
Figure 1: The Transformer - model architecture.

(shifted right)

How it works

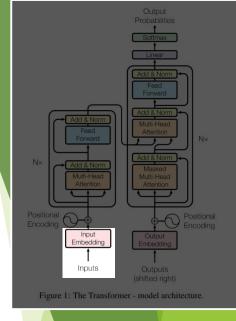


Embedding

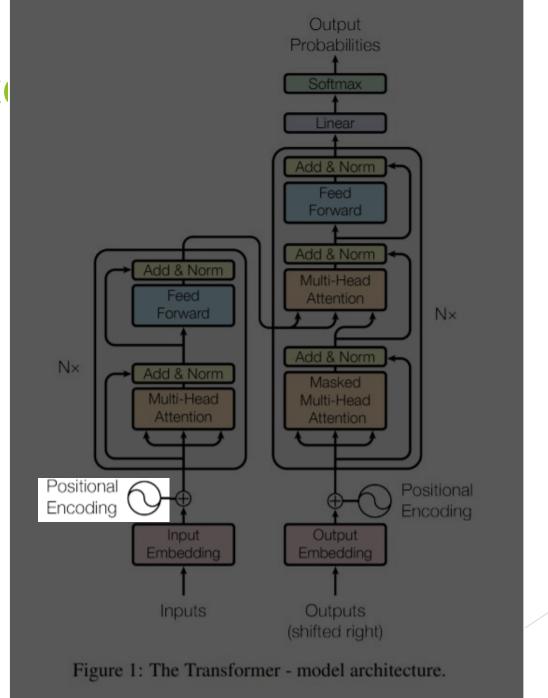


Embedding





Positional Ence



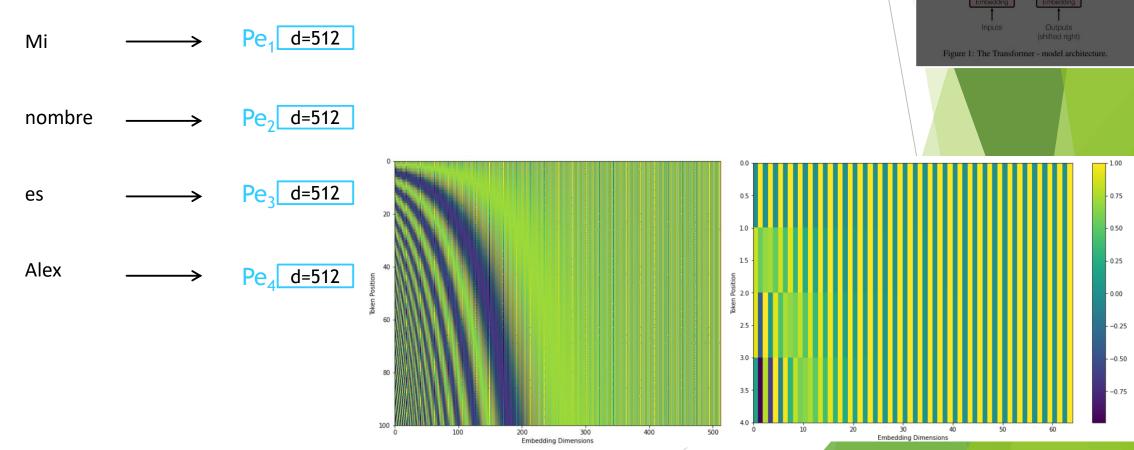
Positional Encoding

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

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

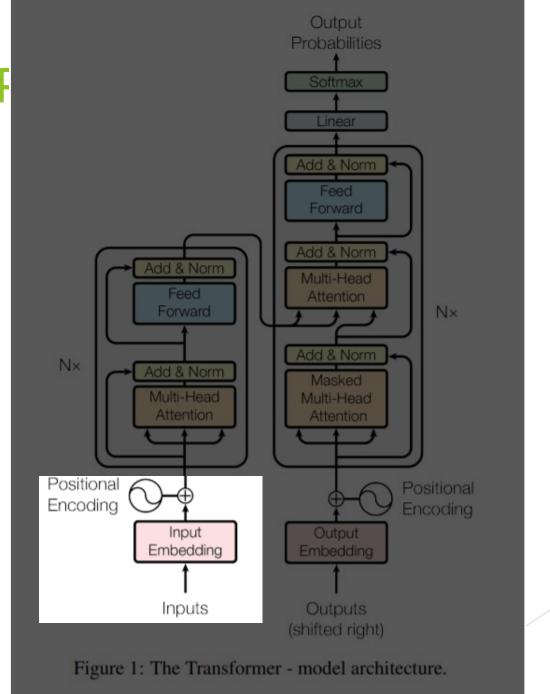
where pos is the position and i is the dimension. That is, each dimension of the positional encoding corresponds to a sinusoid. The wavelengths form a geometric progression from 2π to $10000 \cdot 2\pi$. We chose this function because we hypothesized it would allow the model to easily learn to attend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

Positional Encoding

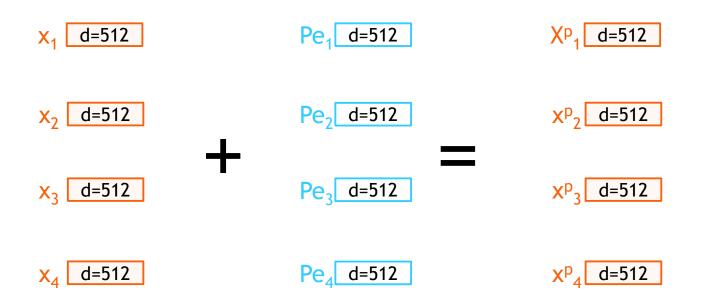


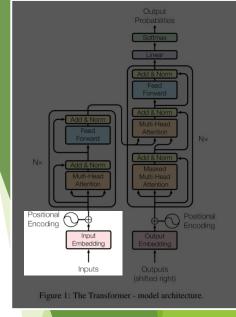
https://github.com/jalammar/jalammar.github.io/blob/master/notebookes/transformer_positional_encoding_graph.ipynb

Embedding + F

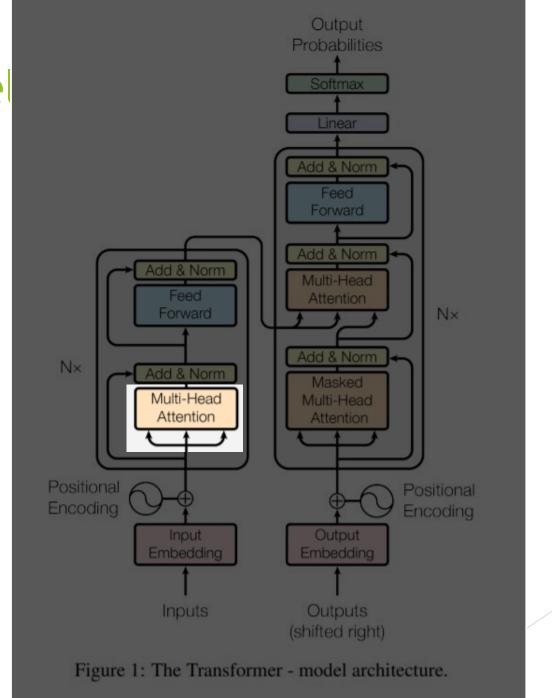


Embedding + Pos. Encoding



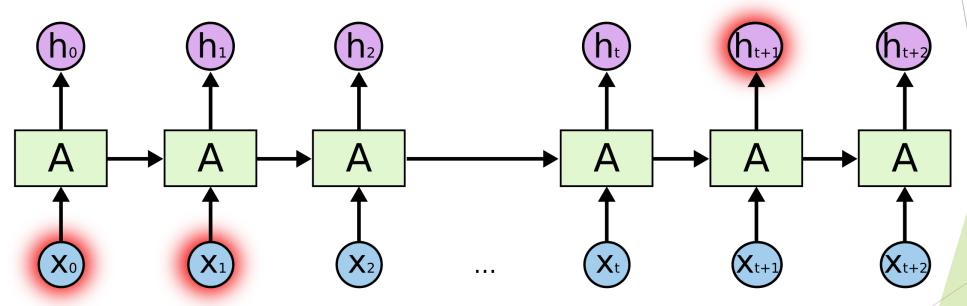


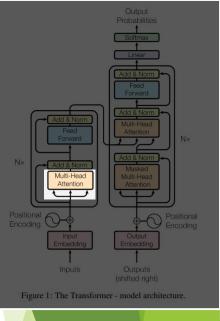
Multi-Head Sel



Substitute for RNN

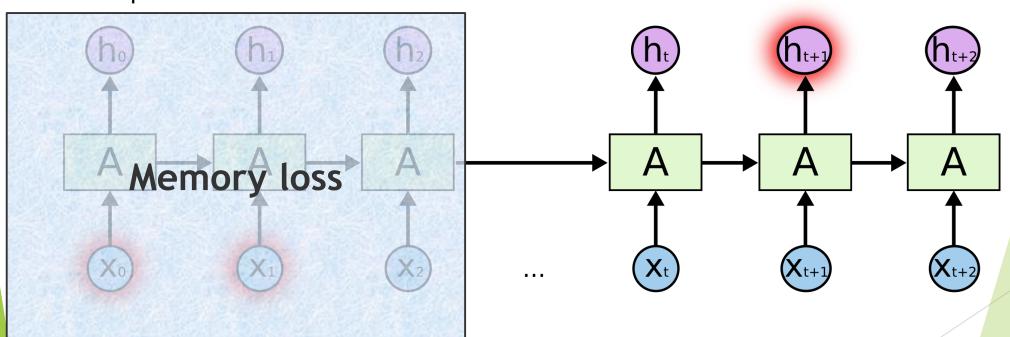


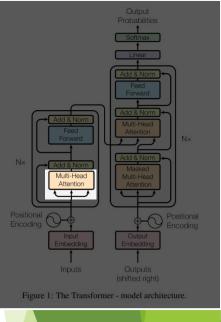




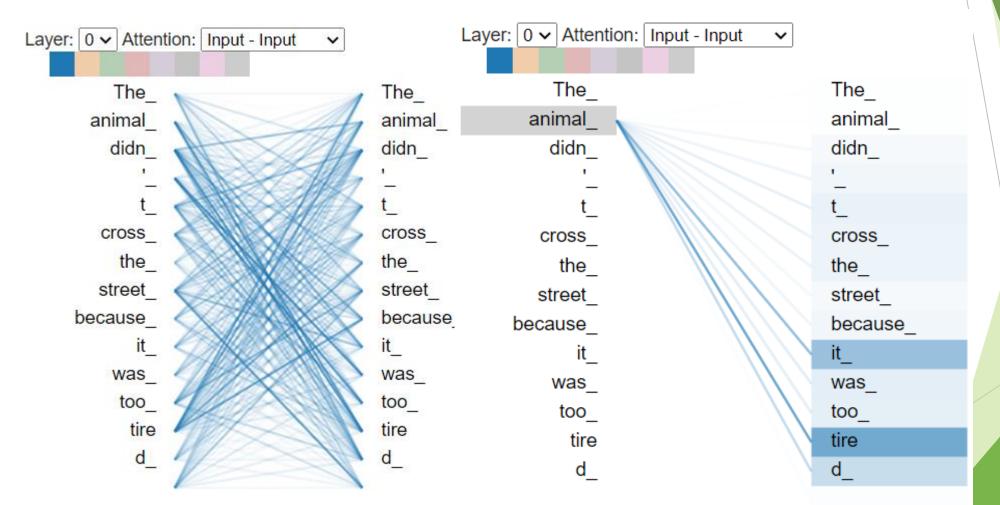
Substitute for RNN

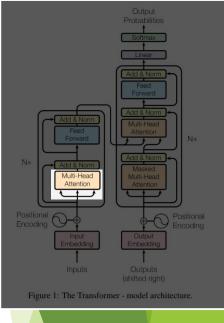
RNN problem





Substitute for RNN





Self-Attention

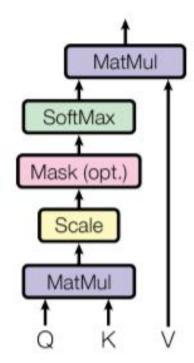
The sublayers needs 3 vectors for each word (d=64)

Queries

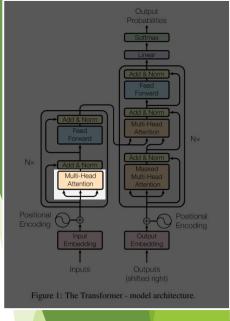
Keys

Values

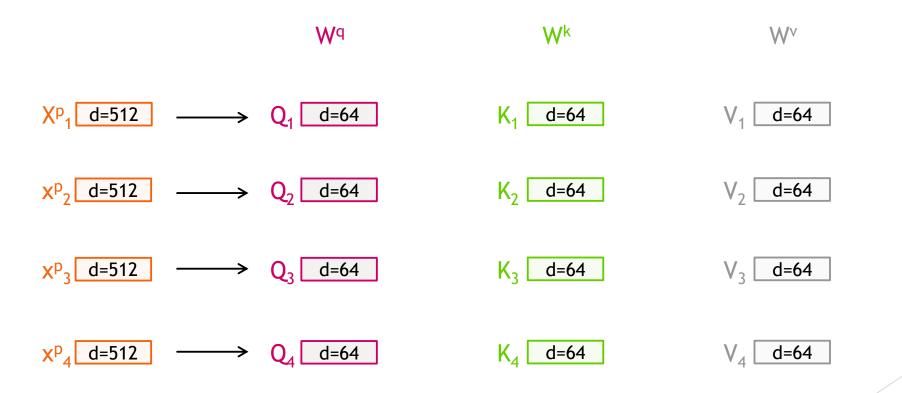
Scaled Dot-Product Attention

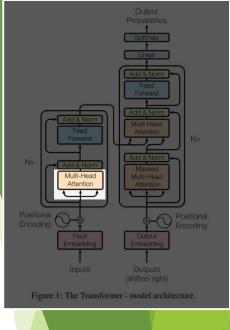


Attention
$$(Q, K, V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

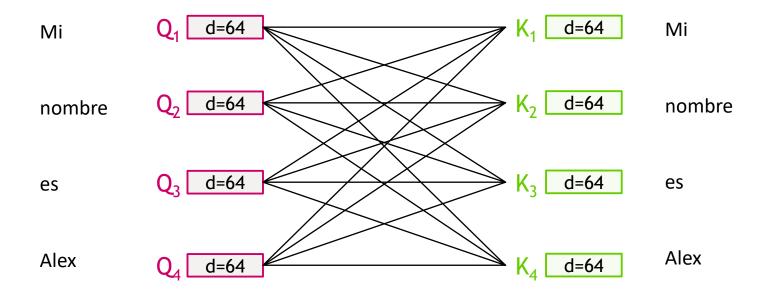


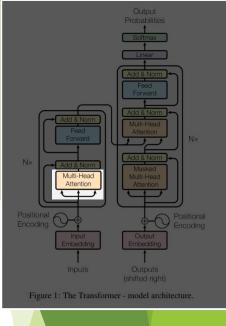
Self-Attention





Self-Attention

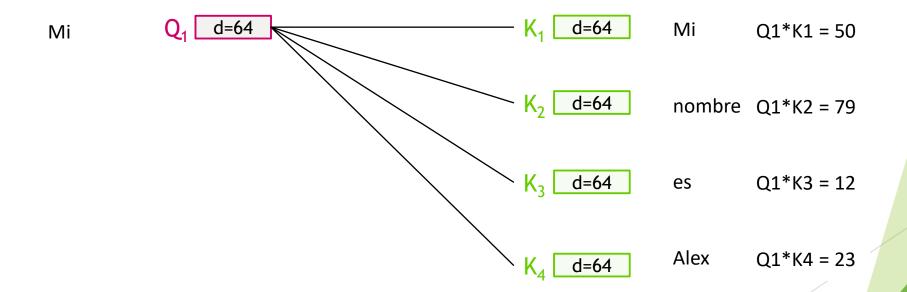


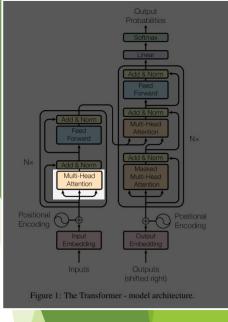


Self-Attention

Example

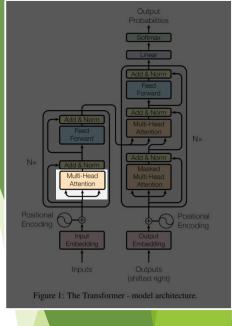
Attention
$$(Q, K, V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$



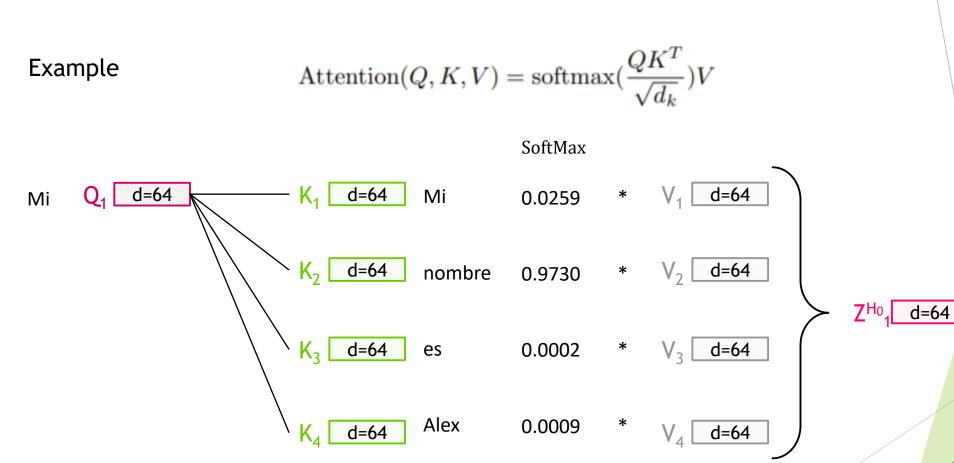


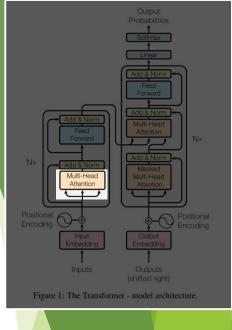
Self-Attention

 $\operatorname{Attention}(Q,K,V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$ Example SoftMax d=64 d=64 Mi Q1*K1 = 5050/8 = 6.25 Mi 0.0259 d=64 nombre Q1*K2 = 7979/8 = 9.875 0.9730 d=64 12/8 = 1.5Q1*K3 = 120.0002 Alex Q1*K4 = 2323/8 = 2.8750.0009



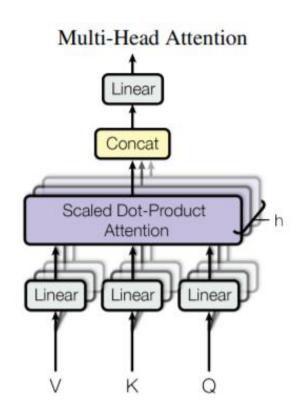
Self-Attention

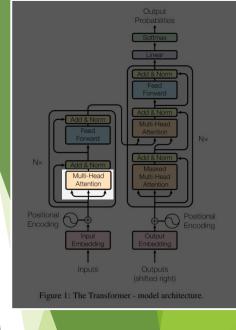




Multi-Head

8 Heads





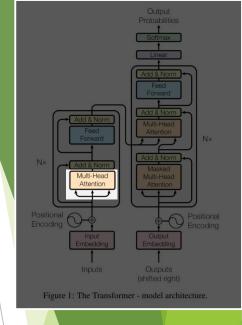
 $\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, ..., \text{head}_{\text{h}}) W^O \\ \text{where head}_{\text{i}} &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned}$

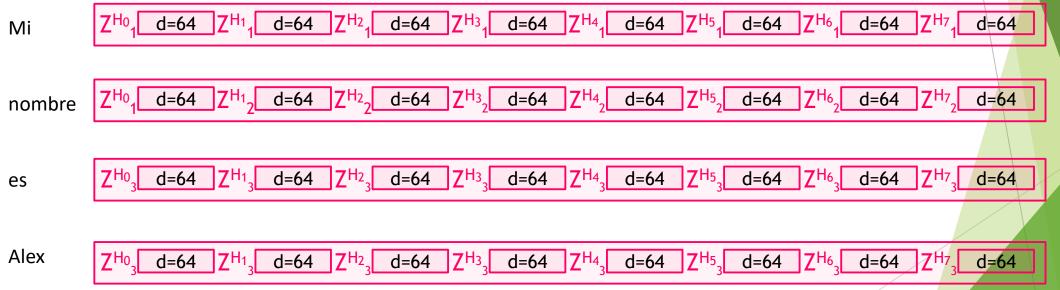
Multi-Head

Example

$$MultiHead(Q, K, V) = Concat(head_1, ..., head_h)W^O$$

$$where head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$





Multi-Head

Example

$$MultiHead(Q, K, V) = Concat(head_1, ..., head_h)W^O$$

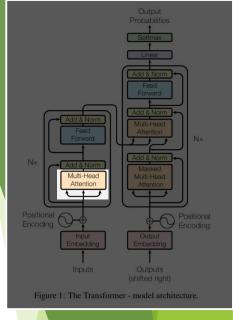
$$where head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

Mi Z_1 d=512

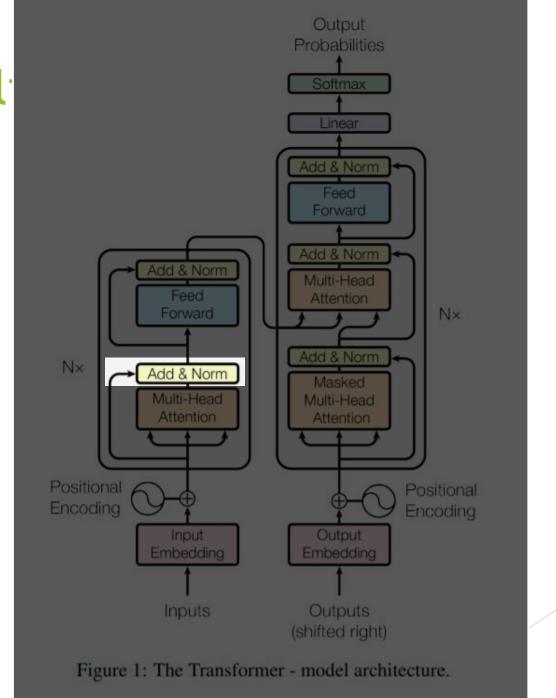
nombre Z_2 d=512

es Z_3 d=512

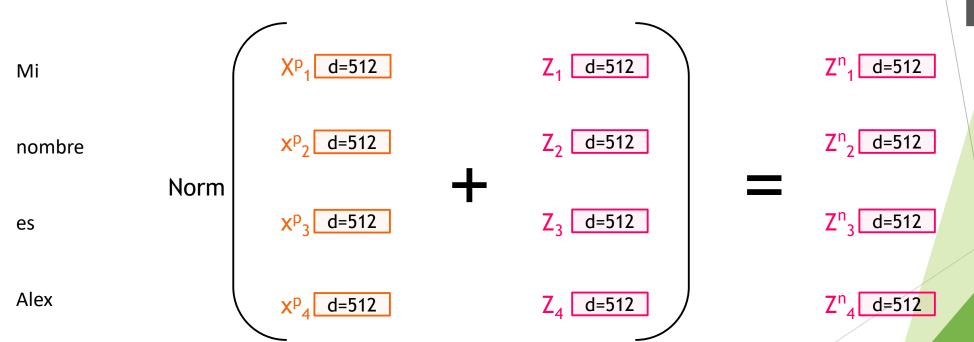
Alex Z_4 d=512

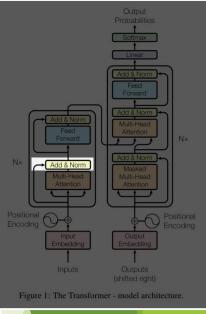


Add & Normal

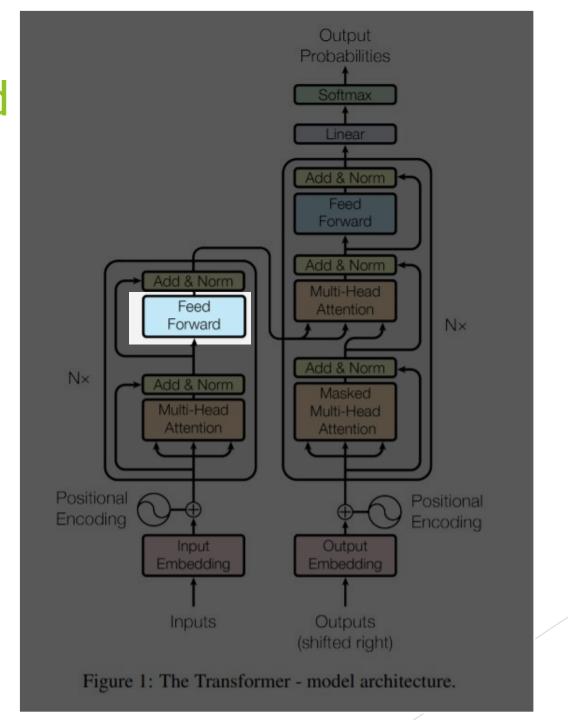


Add & Normalize

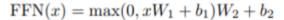


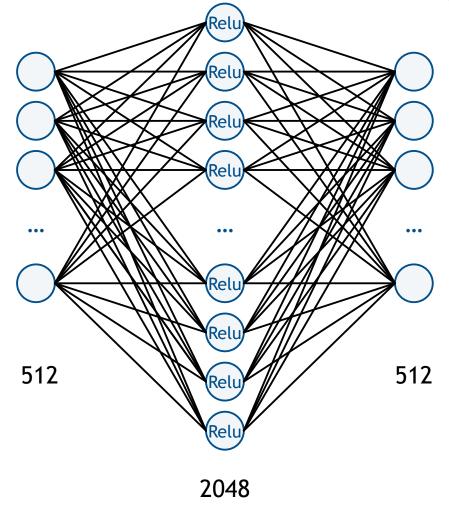


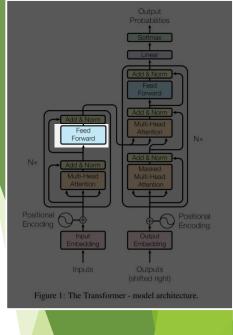
Feed Forward



Feed Forward Neural Network

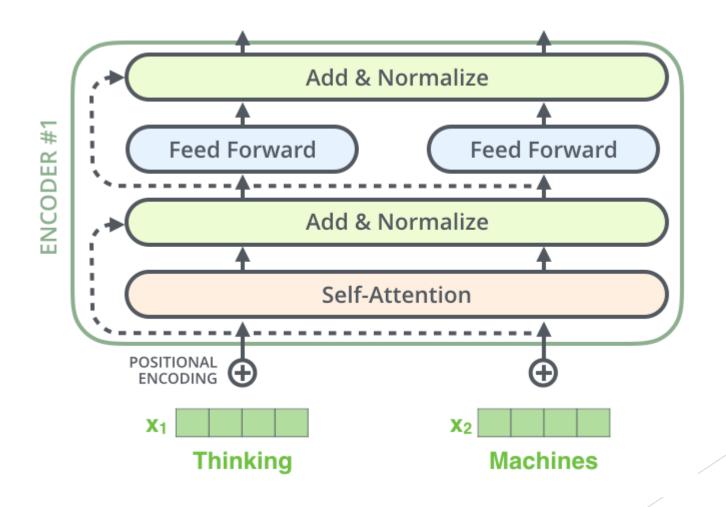






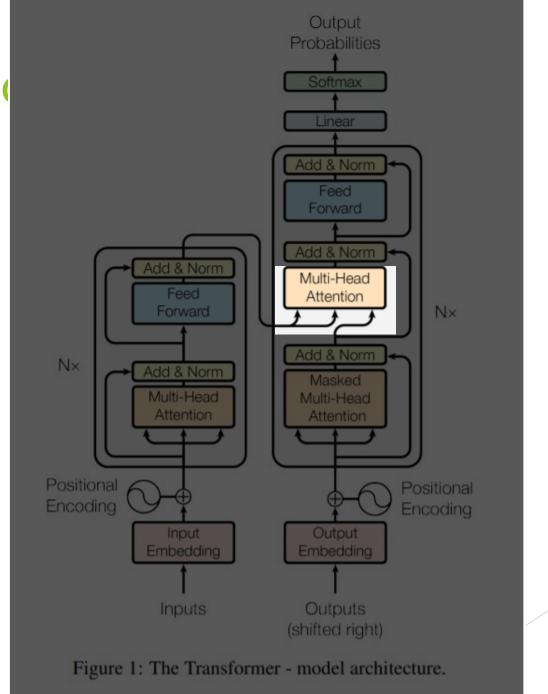
The FFNN is applied to each word separately

Recap



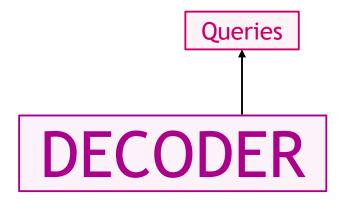
Source: <u>Jay Alammar</u>

Encoder-Deco

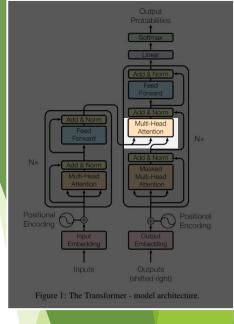


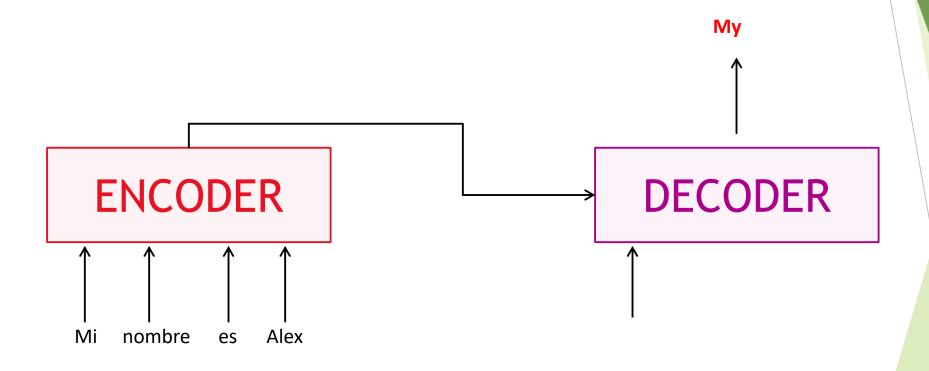
Encoder-Decoder M-H S-A

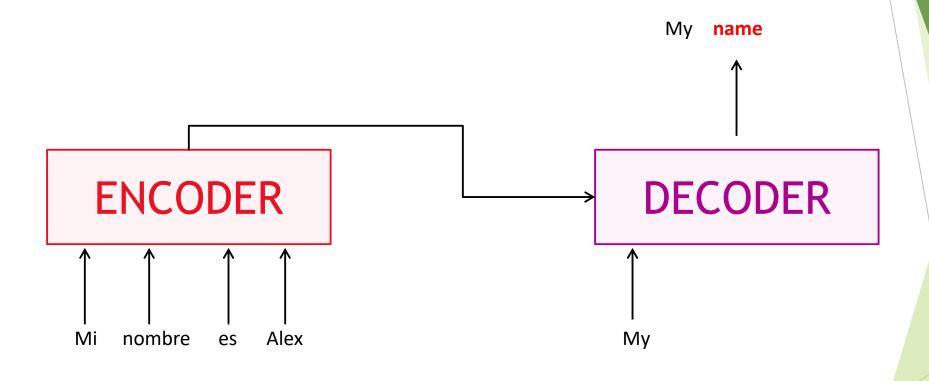
The sublayers needs 3 vectors for each word (d=64)

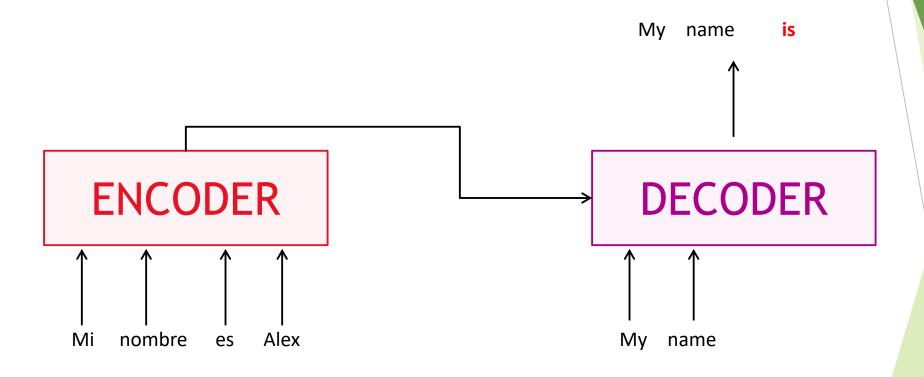


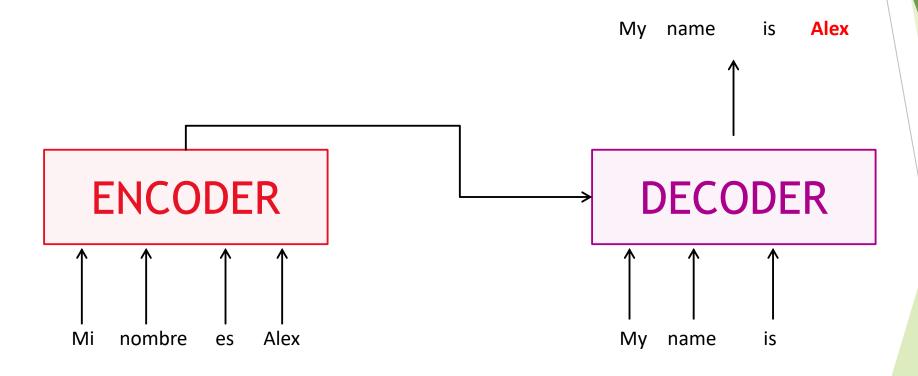


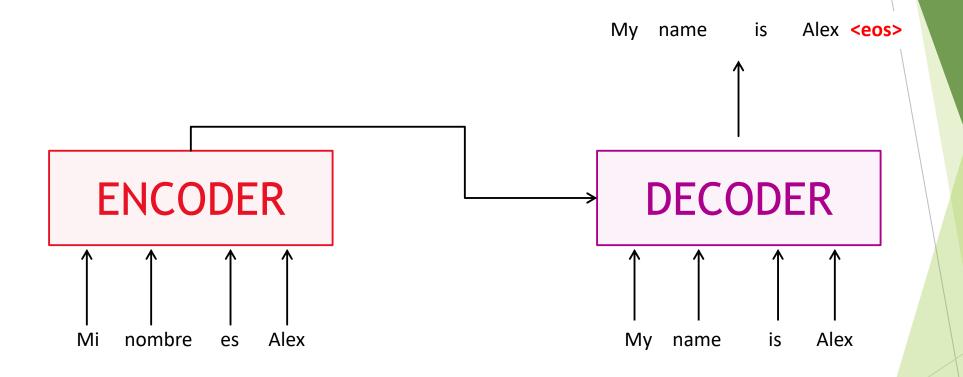












Optimizer

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:

$$lrate = d_{\text{model}}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup_steps^{-1.5})$$
 (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 $warmup_steps = 4000$.

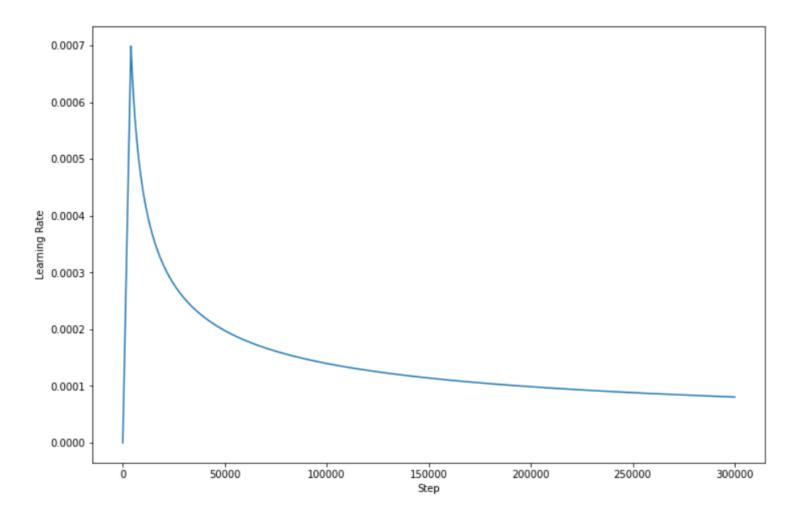


Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BL	EU	Training Cost (FLOPs)		
Model	EN-DE EN-FR		EN-DE	EN-FR	
ByteNet [15]	23.75				
Deep-Att + PosUnk [32]		39.2		$1.0 \cdot 10^{20}$	
GNMT + RL [31]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$	
ConvS2S [8]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$	
MoE [26]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$	
Deep-Att + PosUnk Ensemble [32]		40.4		$8.0 \cdot 10^{20}$	
GNMT + RL Ensemble [31]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$	
ConvS2S Ensemble [8]	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	41.0	2.3 ·	10^{19}	

	N	$d_{ m model}$	$d_{ m ff}$	h	d_k	d_v	P_{drop}	ϵ_{ls}	train steps	PPL (dev)	BLEU (dev)	params ×10 ⁶
base	6	512	2048	8	64	64	0.1	0.1	100K	4.92	25.8	65
(A)				1	512	512				5.29	24.9	
				4	128	128				5.00	25.5	
				16	32	32				4.91	25.8	
				32	16	16				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 8									5.19	25.3	50
	8									4.88	25.5	80
		256			32	32				5.75	24.5	28
		1024			128	128				4.66	26.0	168
			1024							5.12	25.4	53
			4096							4.75	26.2	90
(D)							0.0			5.77	24.6	
							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

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 [8], and observe nearly identical results to the base model.

END

- https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1 c4a845aa-Paper.pdf
- https://jalammar.github.io/illustrated-transformer/
- https://colab.research.google.com/github/jalammar/jalammar.github.io/blo b/master/notebookes/transformer/transformer_positional_encoding_graph.ip ynb#scrollTo=SB9-56JWKv6T
- https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb#scrollTo=OJKU36QAfqOC