# Attention Is All You Need

Vaswani et al (2017)
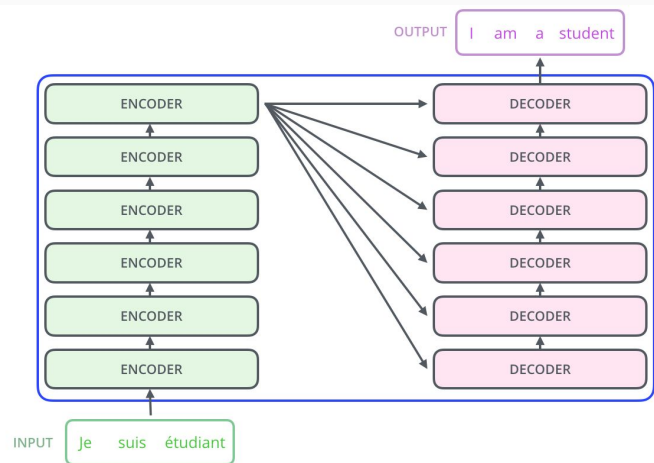
NLP reading club 2021/06/24

# Key Points

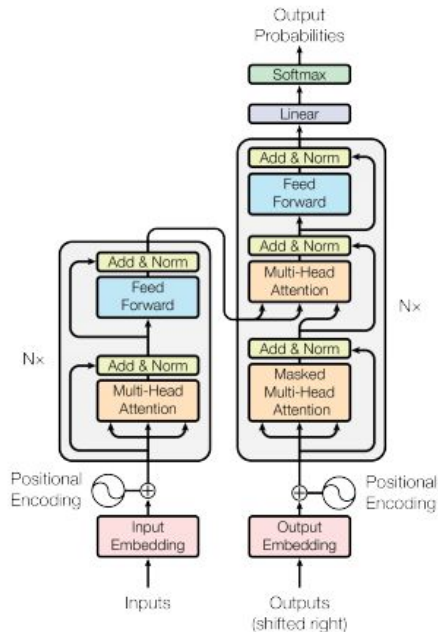The paper introduces the Transformer architecture, which:

- Rely completely on attention to learn the sequence of input instead of using recurrent/convolutional units
- Uses multiheaded attention over the same input to get different projections to improve model performance
- Positional encodings to encode word ordering
- Faster training times as compared to RNNs/CNNs trained to do the same (translation/seq2seq) task

# Model architecture



OUTPUT: I am a student

INPUT: Je suis étudiant
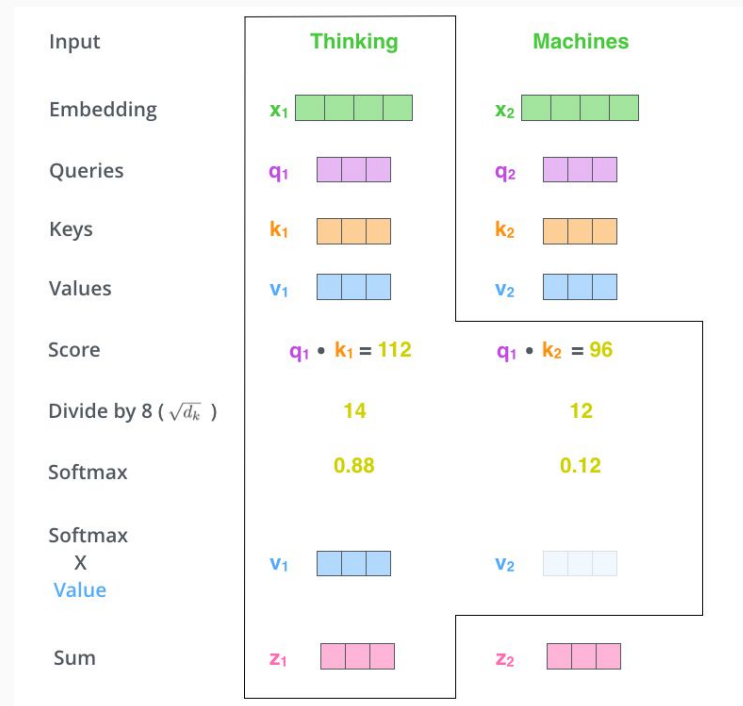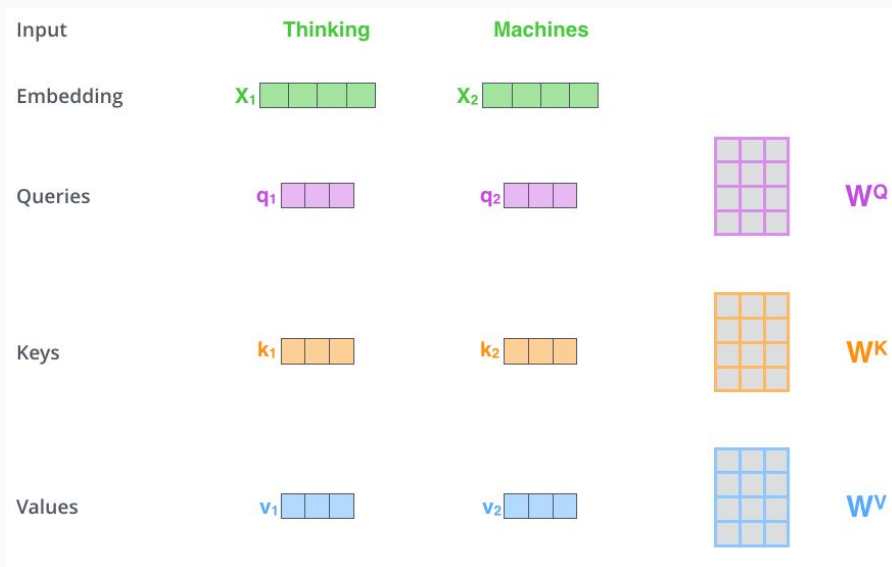
from http://jalammar.github.io/illustrated-transformer/

Similar to previous seq2seq models, the transformer has a encoder-decoder setup. However, multiple words can now pass through the model at once as opposed to RNN architectures where words have to be fed sequentially

# Why self-attention

- Consider constant number of positions regardless of length of input ( as opposed to RNNs, which have to consider $O(n)$ positions
- Easily parallelizable
- Path length from distance part of the sentence is the same as those that are closer (unless you restrict the attention to a certain distance)

# Attention mechanism



| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

from http://jalammar.github.io/illustrated-transformer/

# Multihead attention

**1)** This is our input sentence*

**2)** We embed each word*

**3)** Split into 8 heads. We multiply $X$ or $R$ with weight matrices

**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
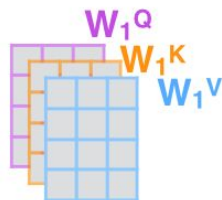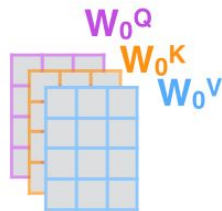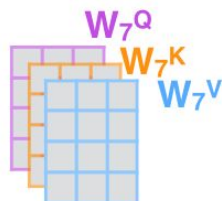
Thinking Machines

$X$

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

$R$

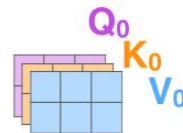$W_0^Q$
$W_0^K$
$W_0^V$

$W_1^Q$
$W_1^K$
$W_1^V$

$\cdots$

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_0$
$K_0$
$V_0$

$Q_1$
$K_1$
$V_1$

$\cdots$
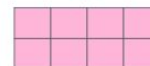
$Q_7$
$K_7$
$V_7$

$Z_0$

$Z_1$

$\cdots$

$Z_7$

$W^O$

$Z$

# Experiment 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 | $\mathbf{3.3 \cdot 10^{18}}$ | |
| Transformer (big) | **28.4** | **41.8** | $2.3 \cdot 10^{19}$ | |

## Parameter variations

| | $N$ | $d_{\text{model}}$ | $d_{\text{ff}}$ | $h$ | $d_k$ | $d_v$ | $P_{drop}$ | $\epsilon_{ls}$ | 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 | 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 | | | | | | | | | 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 |

# Further resources

- [The illustrated transformer](#)
- [The Annotated Transformer](#)
- [https://github.com/xmu-xiaoma666/External-Attention-pytorch](https://github.com/xmu-xiaoma666/External-Attention-pytorch)