# Attention Is All You Need

Vaswani et al. NeurIPS 2017
Presented by Luke Song

# Abstract

- Presents a new neural architecture named the Transformer

- Based solely on the attention mechanism widely used in SEQ2SEQ models

- More parallelizable compared to existing state-of-the-art (SOTA) models

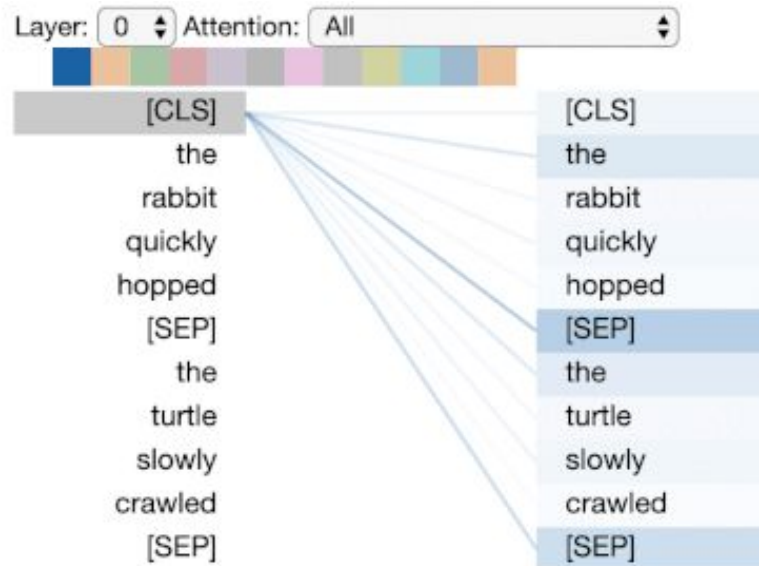- Achieves SOTA in 2 machine translation datasets

# Outline

1. Important Background
2. Model Architecture
3. Experimental Results
4. Model Variation Study
5. Conclusion & Limitation
6. Discussion Time :)

# Important Background
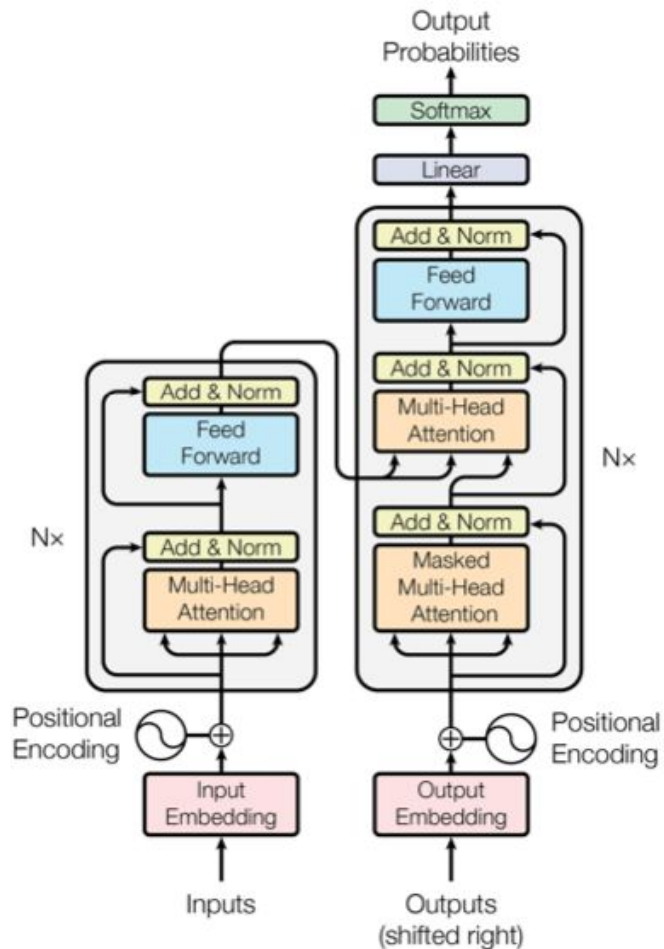
What is Attention Mechanism?

- Mechanism used to let individual tokens "attend" to other tokens regardless of the distance between them

- Transformer uses only self-attention which is attention onto the same sentence

- Think of self-attention as recalculating the representation of each token based on how its meaning is influenced by other tokens in the same sentence



Source: https://github.com/jessevig/bertviz

# Model Architecture

High Level

- Input embedding is first added with Positional Encoding

- 3 components in each encoder/decoder: (Masked) Multi-Head Attention, Addition & Normalization, Feed Forward Network
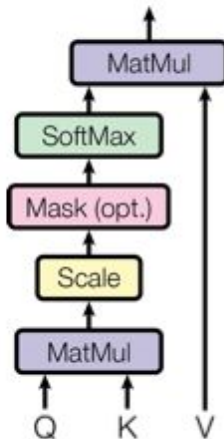


Source: Attention Is All You Need

# Model Architecture

Attention Function

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

- Mapping a query and set of key-value pairs to an output, where the query, keys, values, and output are all vectors

- Q: Queries
  K: Keys
  V: Values
  d_k: dimension of k (64 in the paper)

- Uses a dot-product attention due to its empirical speed/space advantage

- Scale dot product by 1/sqrt(d_k) b/c large values of d_k may push softmax function to region where it has extremely small gradients

Scaled Dot-Product Attention

Input

**Thinking**  **Machines**

Embedding

$X_1$ [ ][ ][ ][ ]     $X_2$ [ ][ ][ ][ ]

Queries

$q_1$ [ ][ ][ ]     $q_2$ [ ][ ][ ]     $W^Q$

Keys

$k_1$ [ ][ ][ ]     $k_2$ [ ][ ][ ]     $W^K$

Values

$v_1$ [ ][ ][ ]     $v_2$ [ ][ ][ ]     $W^V$

| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |

| | Thinking | Machines |
|---|---|---|
| Input | | |
| 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 |

| | Thinking | Machines |
|---|---|---|
| Input | | |
| 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 | $z_1$ | $z_2$ |

Adding it all together...

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \quad V$$

$$= \quad Z$$

# Model Architecture

### Multi-Head Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

- Apply attention to different versions of Q, K, V

- Expands model's ability to focus on different positions

- Generates a multiple "representation subspaces" in order to give the model better representation of the input

- Uses 8 attention heads which are concatenated and fed into a linear layer at the end



Multi-Head Attention

Source: Attention Is All You Need

X

Thinking
Machines

Calculating attention separately in
eight different attention heads

ATTENTION
HEAD #0

ATTENTION
HEAD #1

...

ATTENTION
HEAD #7

$Z_0$

$Z_1$

$Z_7$

1) Concatenate all the attention heads

$Z_0$   $Z_1$   $Z_2$   $Z_3$   $Z_4$   $Z_5$   $Z_6$   $Z_7$

2) Multiply with a weight matrix $W^O$ that was trained jointly with the model

X

$W^O$

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

Z

=

Combining everything attention-wise...

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

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

$Q_0$
$K_0$
$V_0$

$Z_0$

$W^O$

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

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

$Q_1$
$K_1$
$V_1$

$Z_1$

Z

R

...

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

...

$Q_7$
$K_7$
$V_7$

...

$Z_7$

# Before moving on..

- In encoder, all queries, keys, and values come from the same place

- In encoder-decoder attention layer, queries come from the previous decoder layer and keys and values come from the output of the encoder

- This mimics the typical encoder-decoder attention mechanism

- In decoder to ensure auto-regressive property, the model masks everything right to the current token being attended

# Model Architecture

Positional Encoding

- Since attention mechanism in the Transformer does not attend each word auto-regressively (no recurrence nor convolution), model needs something to let it know the relative position of tokens in the sentence

- Positional Encoding is the combination of sine and cosine functions of different frequencies

- Advantages include distance between tokens being symmetrical and being easier to calculate distance between tokens

$$\overrightarrow{p_t} = \begin{bmatrix} \sin(\omega_1.t) \\ \cos(\omega_1.t) \\ \\ \sin(\omega_2.t) \\ \cos(\omega_2.t) \\ \vdots \\ \sin(\omega_{d/2}.t) \\ \cos(\omega_{d/2}.t) \end{bmatrix}_{d\times 1}$$

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



| | | | |
|---|---|---|---|
| dim 4 | | | |
| dim 5 | | | |
| dim 6 | | | |
| dim 7 | | | |

POSITIONAL ENCODING

| 0 | 0 | 1 | 1 |

| 0.84 | 0.0001 | 0.54 | 1 |

| 0.91 | 0.0002 | -0.42 | 1 |

+      +      +

EMBEDDINGS   x₁    x₂    x₃

INPUT    Je    suis    étudiant

# Model Architecture

$$LayerNorm(x + Sublayer(x))$$

Layer Normalization & Residual Connection

- Layer normalization (Ba et al. 2016) is applied to output of sub-layer + input to sub-layer

- Layer normalization normalizes the input across the features

- Empirically shown to reduce training time

- Residual connection means there is a connection that skips few layers (in here 1)

# Model Architecture

Position-wise Feed Forward Networks

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

- Fully connected feed-forward network

- Two linear transformations with a ReLU activation in between

- Inner layer has dimensionality of 2048

- Applied to each position separately and identically

Combining all elements...

Decoding time step: (1) 2 3 4 5 6          OUTPUT



Linear + Softmax

ENCODER

ENCODER

DECODER

DECODER

EMBEDDING
WITH TIME
SIGNAL

EMBEDDINGS

INPUT          Je          suis       étudiant

Decoding time step: 1 (2) 3 4 5 6

OUTPUT    I



**K**encdec    **V**encdec

Linear + Softmax

ENCODERS

DECODERS

EMBEDDING
WITH TIME
SIGNAL

EMBEDDINGS

INPUT        Je      suis    étudiant          PREVIOUS
OUTPUTS      I

# In case you are curious



Which word in our vocabulary is associated with this index?

am

Get the index of the cell with the highest value (argmax)

5

log_probs

0 1 2 3 4 5 ... vocab_size

Softmax

logits

0 1 2 3 4 5 ... vocab_size

Linear

Decoder stack output

# Why Self-Attention?

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

- Less total computational complexity per layer

- More parallelizable than existing fully autoregressive models

- Shorten the path between tokens to enable model to learn long-term dependency better

- Tang et al. (EMNLP 2018) claims that self-attention outperforms RNN/CNN as a semantic feature extractor and empirically show that it excels on word sense disambiguation task (but not subject-verb agreement over long distance!)

# Experimental Results

- Achevies SOTA on 2 machine translation dataset

- Less training cost than existing SOTA models

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 [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 | $\mathbf{3.3 \cdot 10^{18}}$ | |
| Transformer (big) | **28.4** | **41.0** | $2.3 \cdot 10^{19}$ | |

# Model Variation Study

- Attention key size is important

- More heads doesn't necessary mean better performance

- Learned positional embedding is not better than sinusoidal positional encoding

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.

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

# Conclusion & Limitation

- Introduces a groundbreaking new model that is solely based on attention

- Faster and better than existing models

- Still not fully parallelized due to decoder being auto-regressive

- Context is fixed length and cannot attend long-term dependency

- Stacking more encoders/decoders might lead to vanishing gradients

# References

- "*Attention Is All You Need*," Vaswani et al. NeurIPS 2017

- "*Why Self-Attention? A Targeted Evaluation of Neural Machine Translation Architectures*," Tang et al. EMNLP 2018

- "*The Annotated Transformer*," https://nlp.seas.harvard.edu/2018/04/03/attention.html

- "*The Illustrated Transformer*," http://jalammar.github.io/illustrated-transformer/

- "*Positional Embedding*," https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

- "*BertViz*," https://github.com/jessevig/bertviz

# Thank you! & Discussion Time :)