

Attention Is All You Need

Muhammad Umar Salman
umar.salman@mbzuai.ac.ae

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

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez*[†]
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin*[‡]
illia.polosukhin@gmail.com

NIPS 2017

Agenda

Motivation

Summary

Strengths

Weaknesses

Methodology

Motivation

- In S2S problems e.g. Neural Machine Translation, RNN with Encoder + Decoder architecture was the go-to solution.
- However, working with long sequences, the architecture loses its ability to retain information from first elements.
- Decoder has access to last hidden state of Encoder stack making it lose or forget early information, for this reason **Attention Mechanism** was introduced.
- The decoder now looks at all the states of the encoder, being able to access information about all the elements of the input sequence.

Summary

- **Transformers**, first sequence transduction model replacing recurrent layers used in Encoder-Decoder architecture to Multi-head self-attention.
- Its architecture makes use of stacked self-attention and point-wise, fully connected layers
- Experiments on WMT-2014 datasets (En-Fr & En-Ger) show that this approach achieves higher BLEU score compared to previous models or reported ensembles
- Approach has had significant impact on introducing new state-of-the-art language models and vision models such as BERT, GPT-2 and ViT.
- Observed Transformer generalized well when applied to other tasks such as English constituency parsing.

Strengths

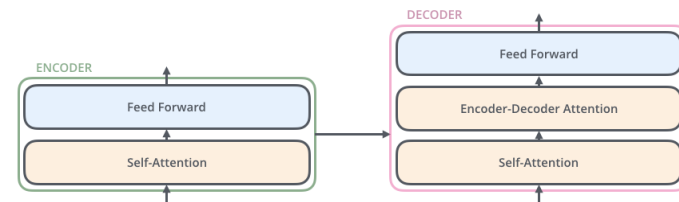
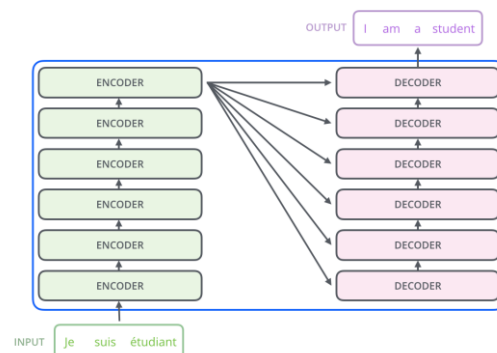
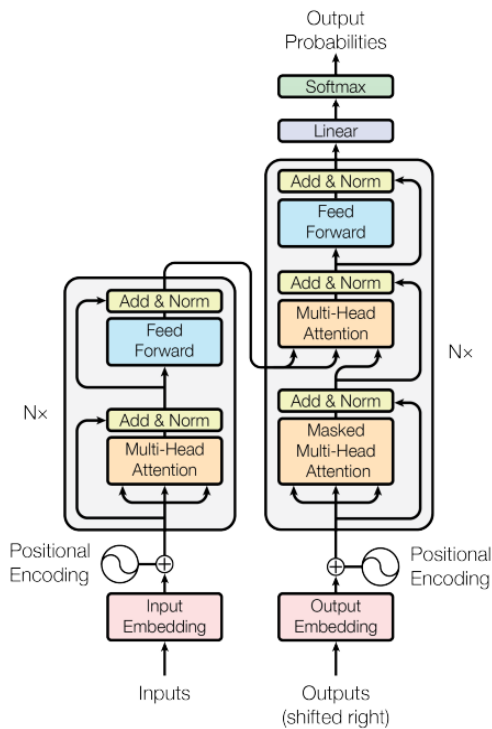
- Propose a **novel architecture** based on Attention Mechanism
- Established a **new state-of-the-art BLEU score** on machine translation tasks
- The architecture is **highly efficient** as architecture reduces model complexity which **reduces the training cost**
- Achieves **parallelizability** within the architecture to speed up training.
- Well-written, well-structured and experiments are clearly set-up for replication purposes.

Weaknesses

- $O(n^{2.d})$, larger the input the higher the time it will take to make inferences/train.
- Paper lacks in-depth mathematical and methodological details by giving a high overview. Reader has to go over code to fully understand the nitty-gritty details.
- Paper argues that retains information for long input sequences. It should have maybe specifically evaluated on longer sequences and its compare it with other models
- Although Tables are very well formatted table column names should be properly described in the footnote of the image for ease of the reader.

Methodology

General Architecture

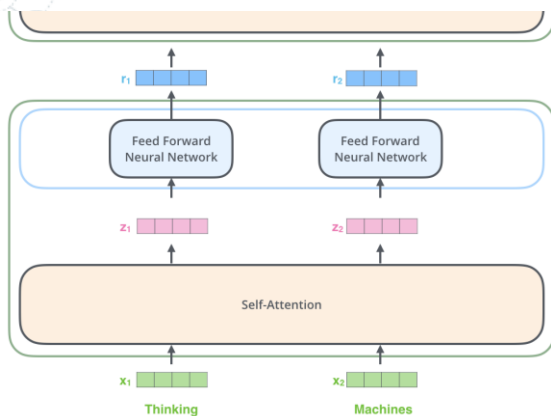


Methodology

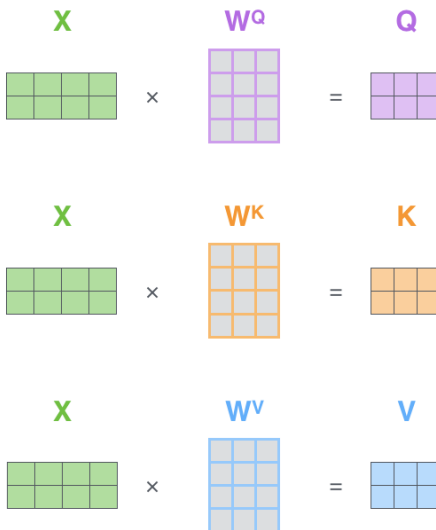
Embedding + Positional Encoding

ENCODER #2

ENCODER #1

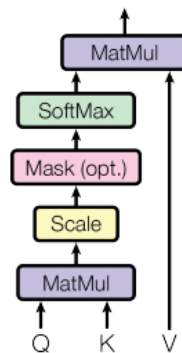


Key, Value & Query



Self Attention / Attention Score / Attention Head

Scaled Dot-Product Attention



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

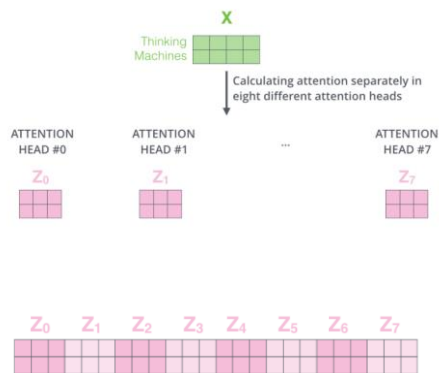
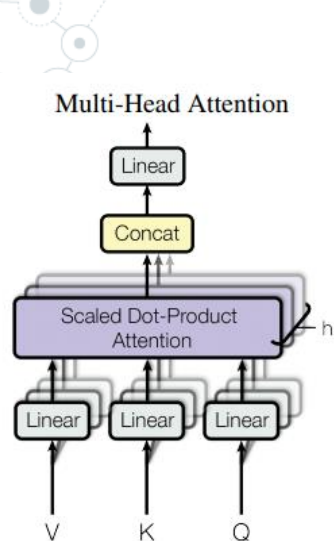
$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \times V = Z$$

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

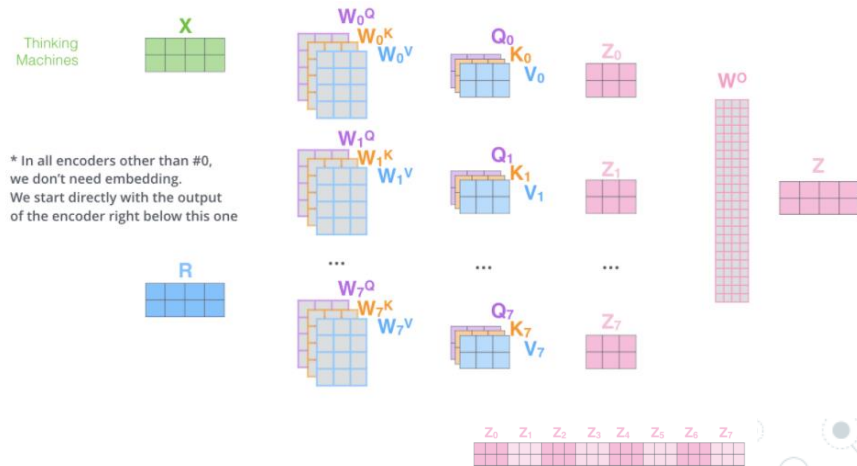
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Methodology

Multi-Head Attention



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

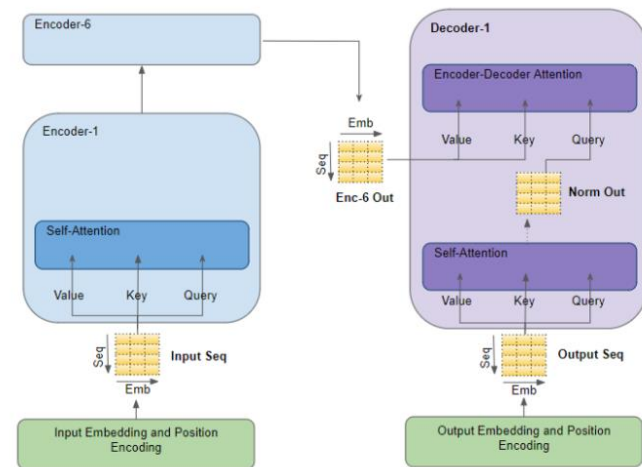
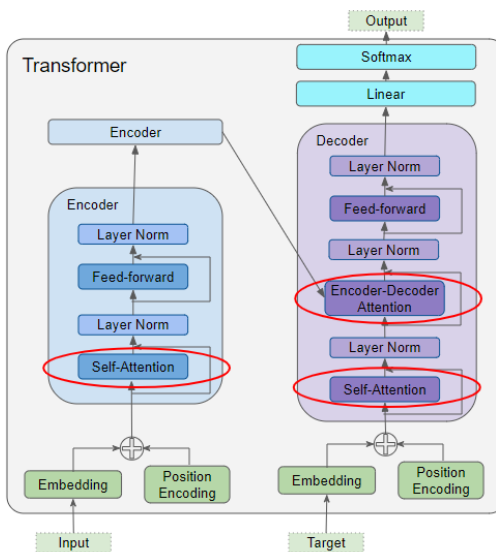
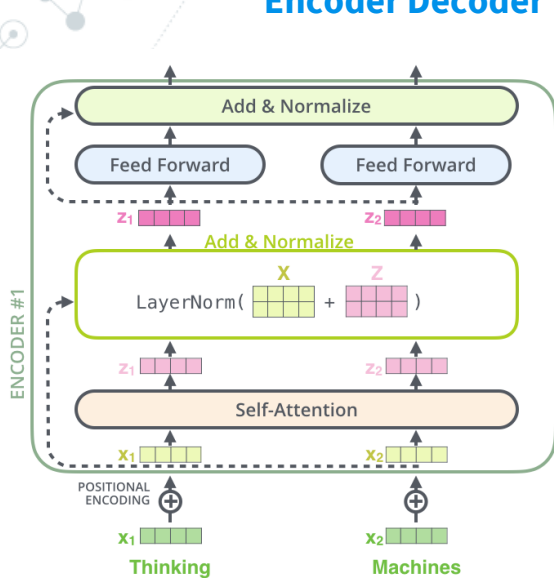


$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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

Methodology

Encoder Decoder



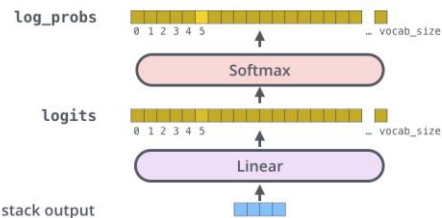
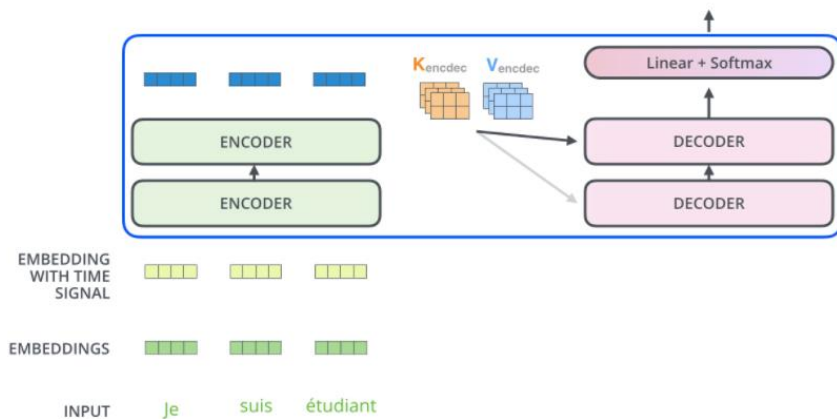
$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

Methodology

Inference

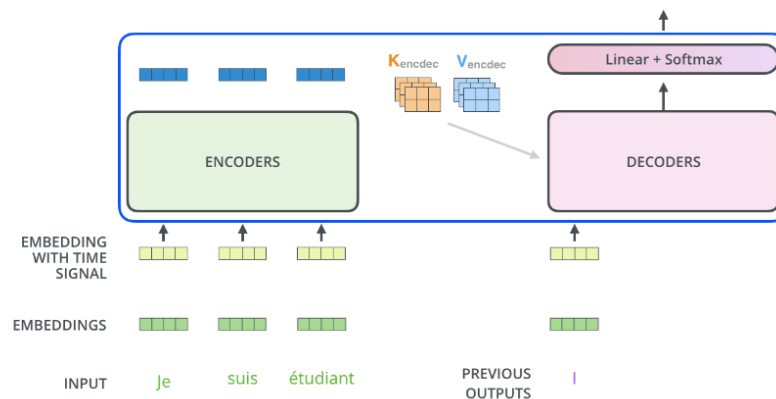
Decoding time step: 1 2 3 4 5 6

OUTPUT |



Decoding time step: 1 2 3 4 5 6

OUTPUT | am



The background of the slide is a light gray network pattern. It consists of numerous small circles, some of which are solid gray and others are hollow with a gray outline. These circles are interconnected by a web of thin, light gray lines, creating a complex, organic structure that resembles a molecular or digital network.

Thank you

Muhammad Umar Salman
umar.salman@mbzau.ac.ae