

Deep Learning

16 Self-Attention & Transformers

Dr. Konda Reddy Mopuri
Dept. of AI, IIT Hyderabad
Jan-May 2023

- ① Why does one need to think beyond LSTMs?

Motivation

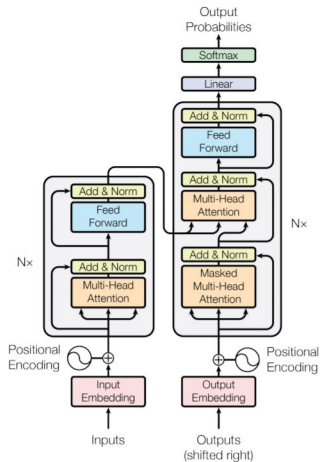
- ① Why does one need to think beyond LSTMs?
- ② Sequential processing doesn't allow parallelization

- ① Why does one need to think beyond LSTMs?
- ② Sequential processing doesn't allow parallelization
- ③ Despite the LSTM/GRU, RNNs need attention to deal with long-range dependencies

- ① Why does one need to think beyond LSTMs?
- ② Sequential processing doesn't allow parallelization
- ③ Despite the LSTM/GRU, RNNs need attention to deal with long-range dependencies
- ④ Since attention enables accesses to any state, do we need RNNs?

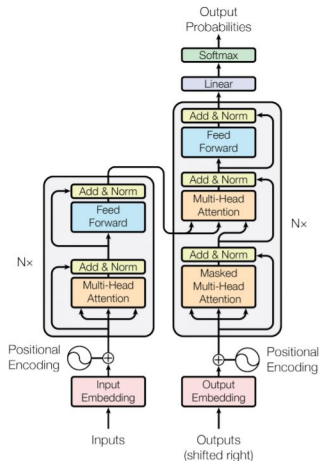
Transformers

- 1 Introduced by Vaswani et al.
NeurIPS 2017



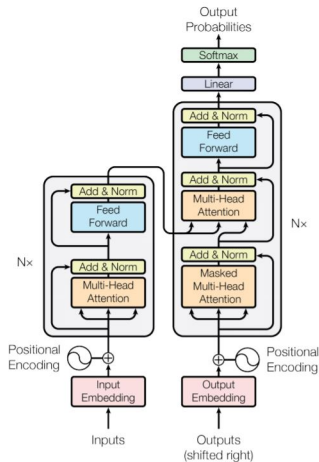
Transformers

- 1 Introduced by Vaswani et al.
NeurIPS 2017
- 2 Sequence to sequence modelling
without RNNs



Transformers

- 1 Introduced by Vaswani et al. NeurIPS 2017
- 2 Sequence to sequence modelling without RNNs
- 3 Transformer model is built on self-attention (no recurrent architectures!)



Transformers

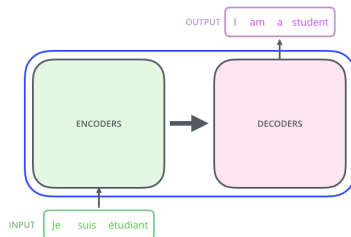


Credits: Jay Alammar

Transformers

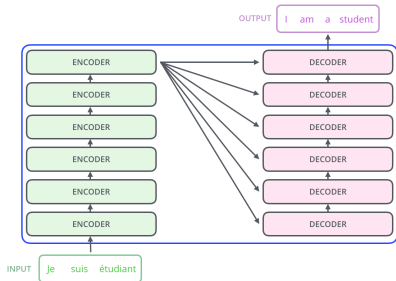


Credits: Jay Alammar



Credits: Jay Alammar

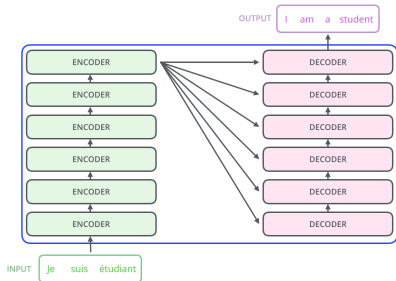
Transformers



Credits: Jay Alammar

- 1 Encoding module has a stack of encoders

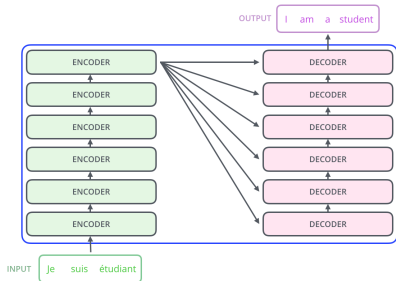
Transformers



Credits: Jay Alammar

- 1 Encoding module has a stack of encoders
- 2 Same structure different parameters

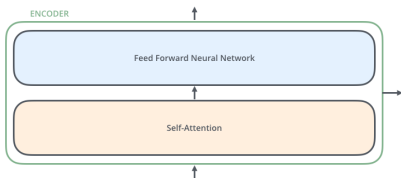
Transformers



Credits: Jay Alammar

- 1 Encoding module has a stack of encoders
- 2 Same structure different parameters
- 3 Similarly the decoding module (same number of components in the stack as encoder)

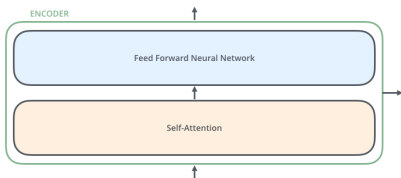
Transformers



- 1 Encoder first has a self-attention layer

Credits: Jay Alammar

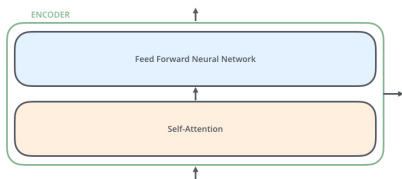
Transformers



Credits: Jay Alammar

- ① Encoder first has a self-attention layer
- ② Looks at the other words while encoding a specific word

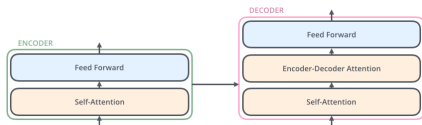
Transformers



Credits: Jay Alammar

- ① Encoder first has a self-attention layer
- ② Looks at the other words while encoding a specific word
- ③ Next a (same) feed-forward NN is applied at all positions

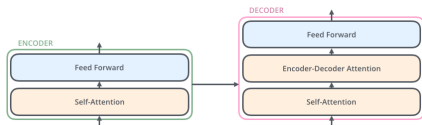
Transformers



Credits: Jay Alammar

- 1 Decoder also has both the layers

Transformers

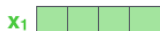


Credits: Jay Alammar

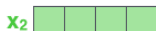
- 1 Decoder also has both the layers
- 2 But, in the middle it has an encoder-decoder attention layer

Transformers-Encoding

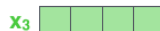
- ① Start with turning each word into a vector at the bottom-most encoder



Je



suis

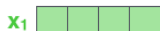


étudiant

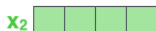
Credits: Jay Alammar

Transformers-Encoding

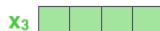
- ① Start with turning each word into a vector at the bottom-most encoder
- ② Others receive a list of vectors from the encoder immediately below



Je



suis

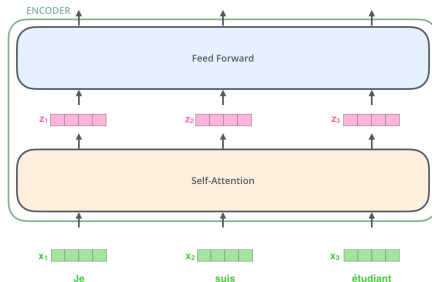


étudiant

Credits: Jay Alammar

Transformers-Encoding

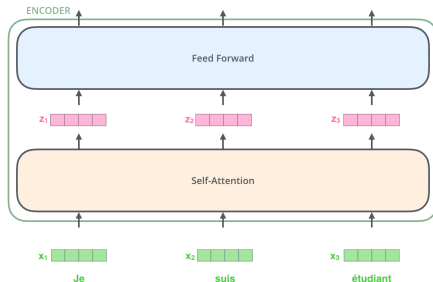
- ① Each word flows through the two layers of the encoder through its own path



Credits: Jay Alammar

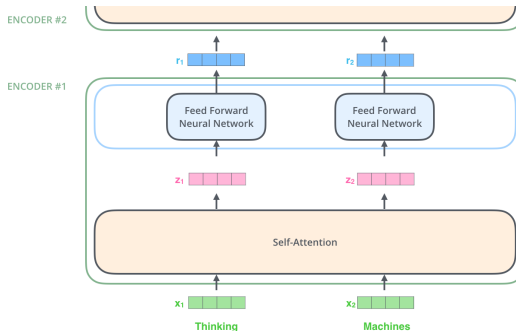
Transformers-Encoding

- ① Each word flows through the two layers of the encoder through its own path
- ② Self-attention layer has dependencies among them, but not the feed-forward layer (which can be parallelized)



Credits: Jay Alammar

Transformers-Encoding



Credits: Jay Alammar

Self-Attention

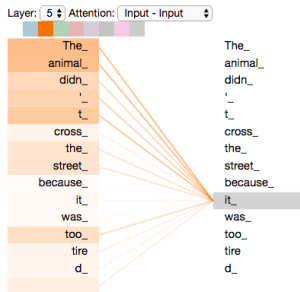
- ① The animal didn't cross the street because it was too tired
- ② The animal didn't cross the street because it was too wide

- ① The animal didn't cross the street because it was too tired
- ② The animal didn't cross the street because it was too wide
- ③ What does 'it' refers to?

- ① The animal didn't cross the street because it was too tired
- ② The animal didn't cross the street because it was too wide
- ③ What does 'it' refers to?
- ④ Easy for humans, but not so much for the traditional Seq2Seq models

Self-Attention

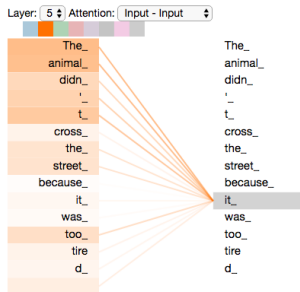
- 1 As the model processes each word, self-attention attends other positions in the i/p sequence to encoder better



Credits: Jay Alammar

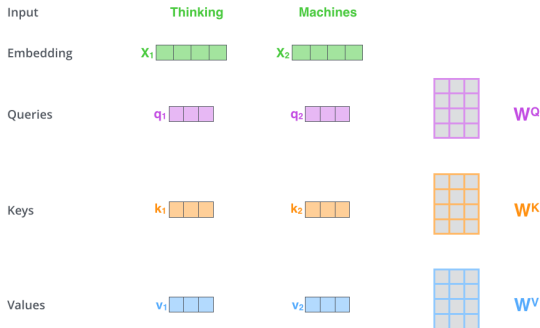
Self-Attention

- ① As the model processes each word, self-attention attends other positions in the i/p sequence to encode better
- ② Unlike RNNs, here we don't keep hidden states from previous positions!



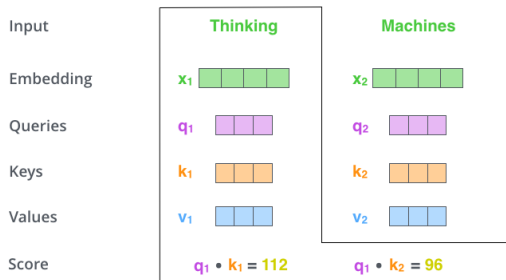
Credits: Jay Alamar

Self-Attention



Credits: Jay Alammar

Self-Attention



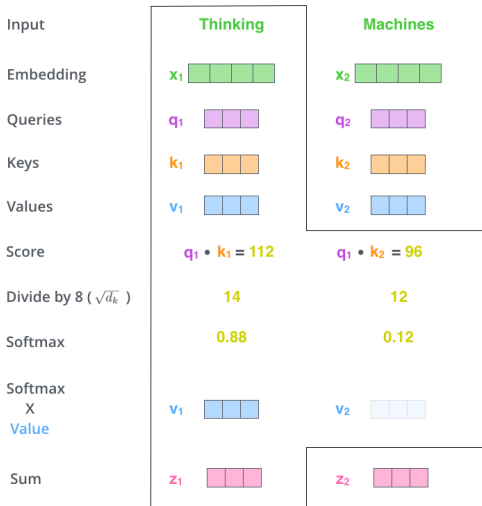
Credits: Jay Alammar

Self-Attention

| 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 |

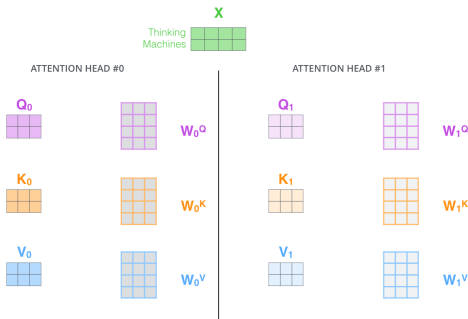
Credits: Jay Alammar

Self-Attention



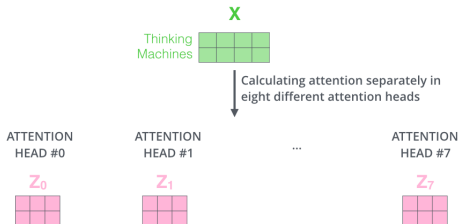
Credits: Jay Alammar

Multi-headed Self-Attention



Credits: Jay Alammar

Multi-headed Self-Attention



Credits: Jay Alammar

Multi-headed Self-Attention

- ① Expands the model's ability to focus on different relevant positions in the i/p

Multi-headed Self-Attention

- ① Expands the model's ability to focus on different relevant positions in the i/p
- ② Enables different 'representational subspace'

Multi-headed Self-Attention

1) Concatenate all the attention heads



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

\times

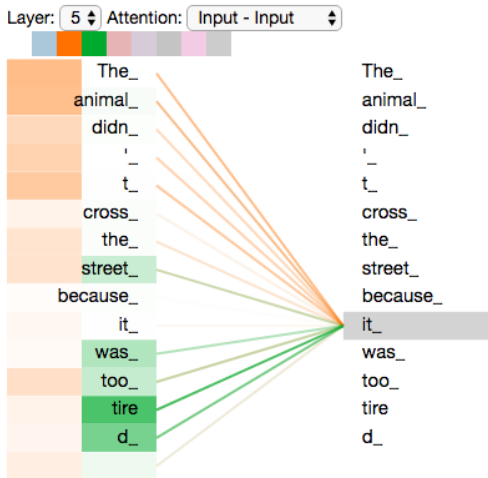


3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



Credits: Jay Alammar

Multi-headed Self-Attention



Credits: Jay Alammur

Positional Encoding

- 1 Unlike RNN and CNN encoders, attention encoder o/p's don't depend on the order

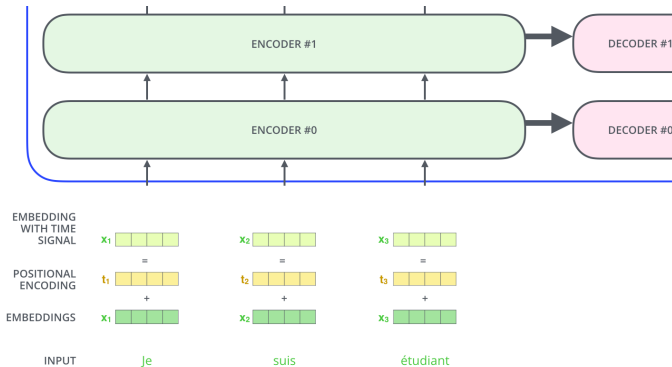
Positional Encoding

- ① Unlike RNN and CNN encoders, attention encoder o/p's don't depend on the order
- ② However, order the sequence conveys vital information in some applications

Positional Encoding

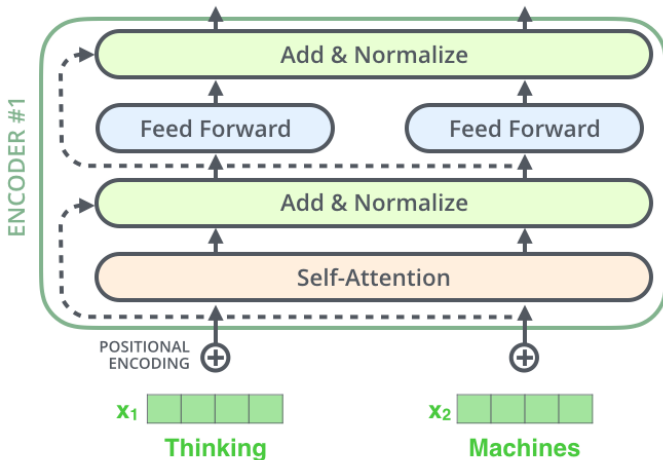
- ① Unlike RNN and CNN encoders, attention encoder o/p/s don't depend on the order
- ② However, order the sequence conveys vital information in some applications
- ③ Solution: Add positional information of the i/p words into their embedding vectors

Positional Encoding



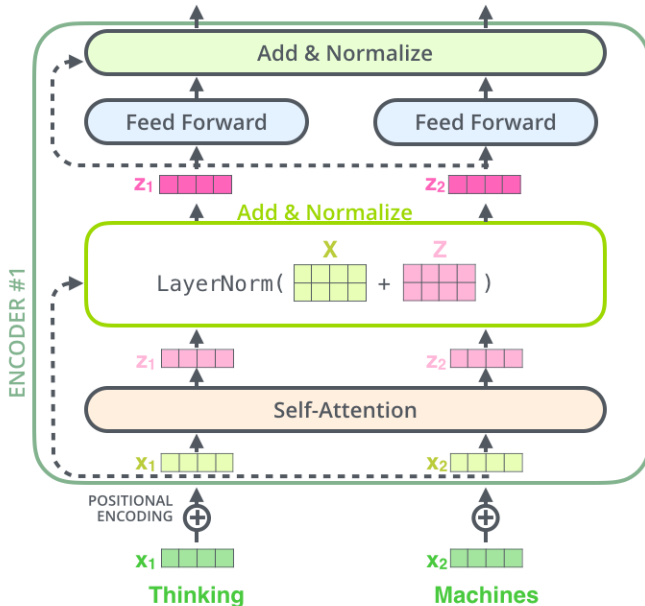
Credits: Jay Alammar

Residuals in the Encoder



Credits: Jay Alammarr

Residuals in the Encoder



The Decoder

- ① Uses the top encoder's K and V vectors for its' encoder-decoder attention

The Decoder

- ① Uses the top encoder's K and V vectors for its' encoder-decoder attention
- ② Self-attention here works in a slightly different way \rightarrow masks the future positions

The Decoder

- ① Uses the top encoder's K and V vectors for its' encoder-decoder attention
- ② Self-attention here works in a slightly different way → masks the future positions
- ③ Encoder-decoder attention layer borrows the queries from the layer below it