# **Deep Learning**

## 16 Self-Attention & Transformers

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

# Motivation

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
   - Path length $= \mathbb{O}(n)$
   - RNNs need at most $\mathbb{O}(n)$ sequential computations to access each element
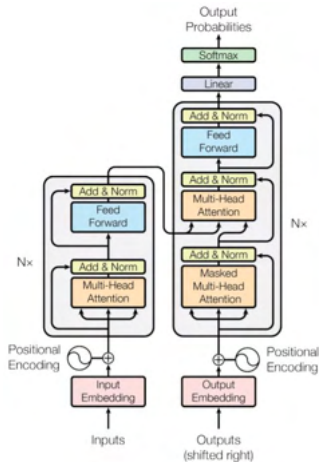
# Motivation

1. Despite the LSTM/GRU, RNNs need attention to deal with long-range dependencies

# Motivation

1. Despite the LSTM/GRU, RNNs need attention to deal with long-range dependencies

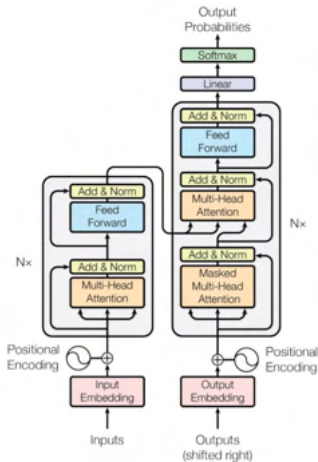2. Since attention enables accesses to any state, do we need RNNs?

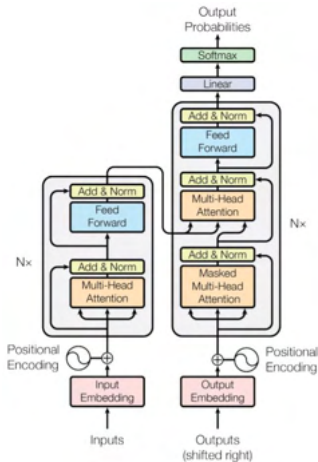# Transformers

1. Introduced by Vaswani et al. NeurIPS 2017

# Transformers

1. Introduced by Vaswani et al. NeurIPS 2017

2. Sequnce to sequence modelling without RNNs

# Transformers

1. Introduced by Vaswani et al. NeurIPS 2017

2. Sequnce to sequence modelling without RNNs

3. Transformer model is built on self-attention (no recurrent architectures!)

# Transformers
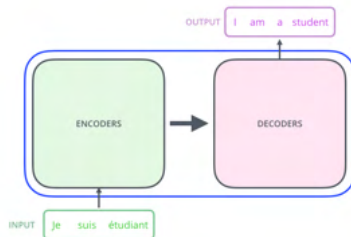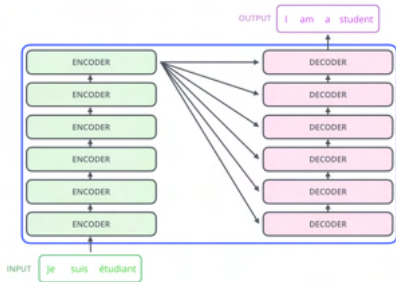


Credits: Jay Alammar
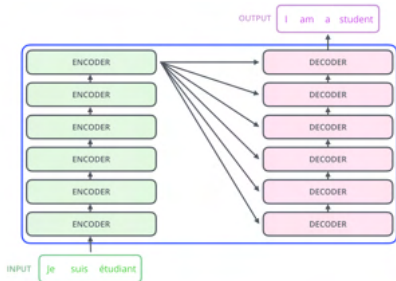
# Transformers

Credits: Jay Alammar



Credits: Jay Alammar

# Transformers
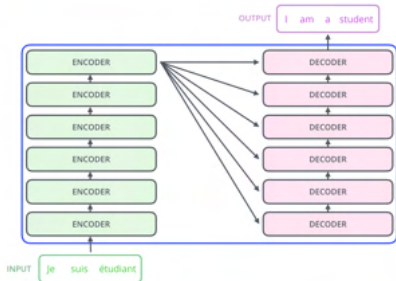
Credits: Jay Alammar

1. Encoding module has a stack of encoders

# Transformers



Credits: Jay Alammar

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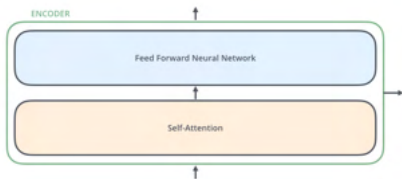2. Same structure different parameters

# Transformers

Credits: Jay Alammar

1. Encoding module has a stack of encoders

2. Same structure different parameters

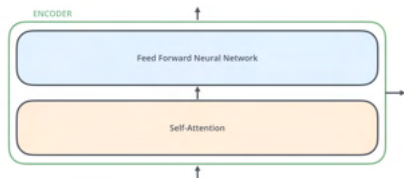3. Similarly the decoding module (same number of components in the stack as encoder)

# Transformers
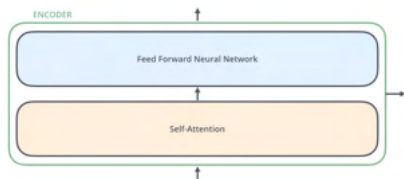
1. Encoder first has a self-attention layer



Credits: Jay Alammar

# Transformers
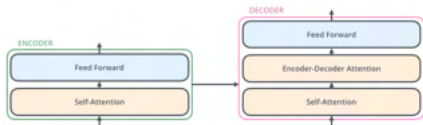
Credits: Jay Alammar

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

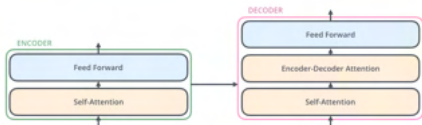3. Next a (same) feed-forward NN is applied at all positions

# 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

$x_1$ ☐☐☐☐

**Je**

$x_2$ ☐☐☐☐

**suis**

$x_3$ ☐☐☐☐

**étudiant**

Credits: Jay Alammar

# Transformers-Encoding

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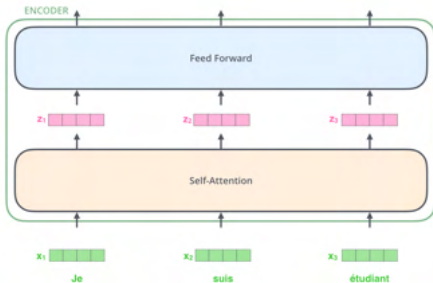2. Others receive a list of vectors from the encoder immediately below



x₁ **Je**

x₂ **suis**

x₃ **étudiant**
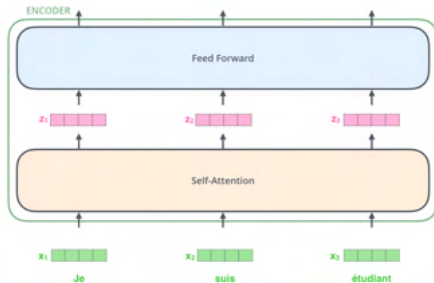
Credits: Jay Alammar

# Transformers-Encoding

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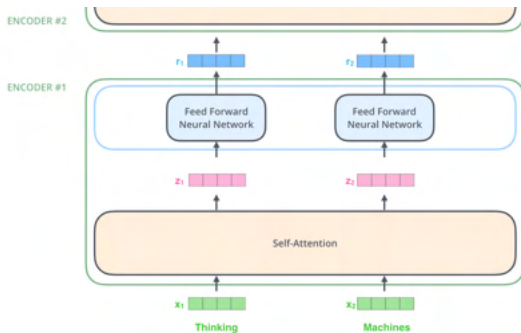


Credits: Jay Alammar

# 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, however, the path length is $\mathbb{O}(1)$



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
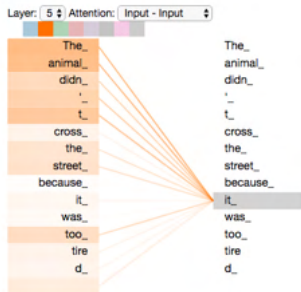
3. What does 'it' refer to?

# 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
3. What does 'it' refer to?
4. Easy for humans, but not so much for the traditional Seq2Seq models
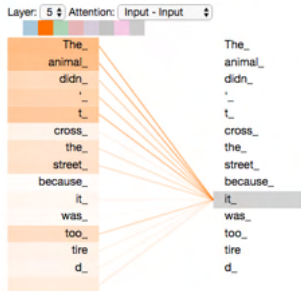
# Self-Attention

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



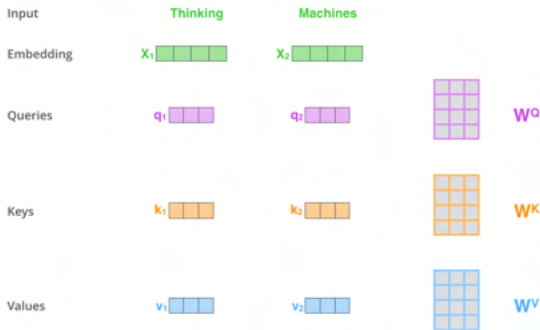Credits: Jay Alammar

# Self-Attention

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

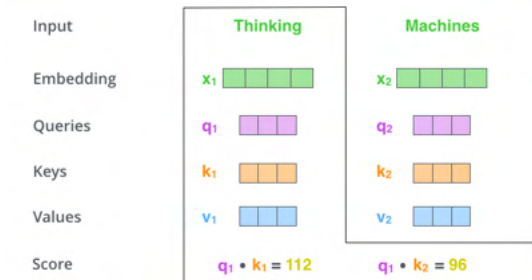2. Unlike RNNs, here we don't keep hidden states from previous positions!



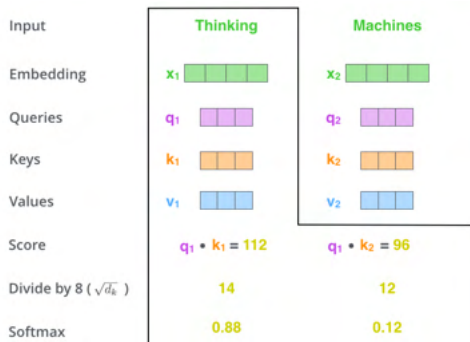Credits: Jay Alammar

# Self-Attention



Credits: Jay Alammar

# Self-Attention

Credits: Jay Alammar

# Self-Attention

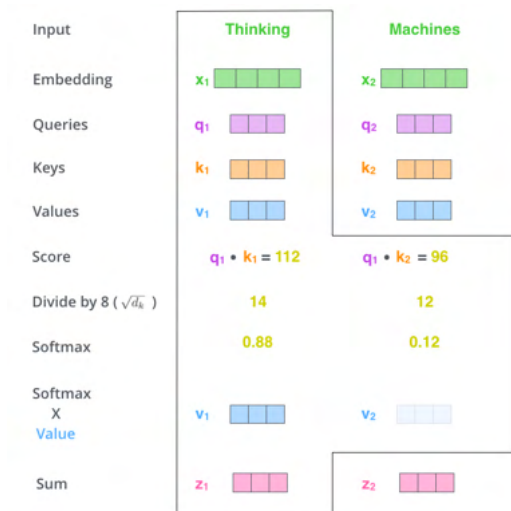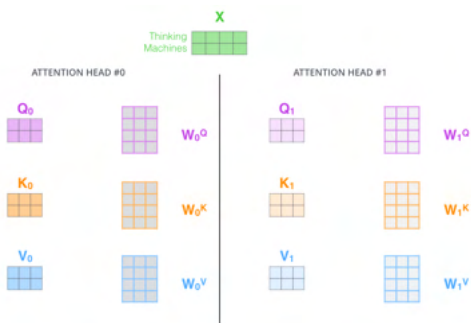| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

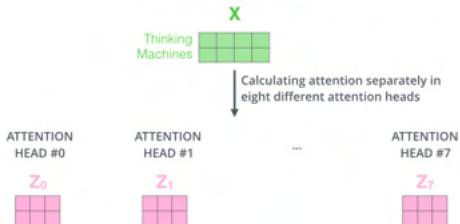Credits: Jay Alammar

# 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

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

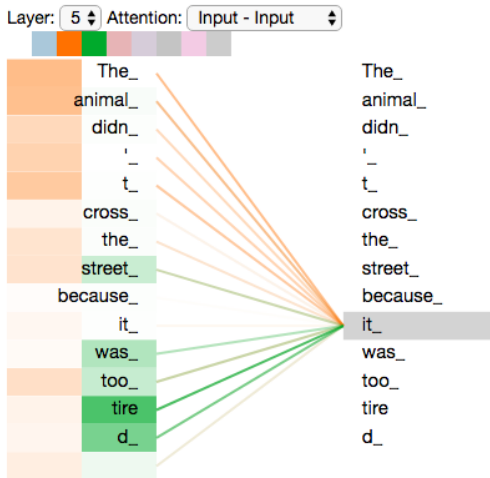2) Multiply with a weight matrix $W^O$ that was trained jointly with the model

x

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

$W^O$

Credits: Jay Alammar

# Multi-headed Self-Attention



Credits: Jay Alammar

# Positional Encoding

1. Unlike RNN and CNN encoders, attention encoder outputs don't depend on the order
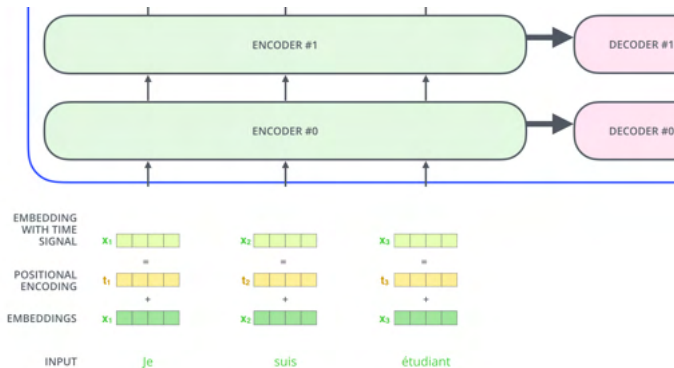
# Positional Encoding

1. Unlike RNN and CNN encoders, attention encoder outputs don't depend on the order

2. However, order the sequence conveys vital information in some applications
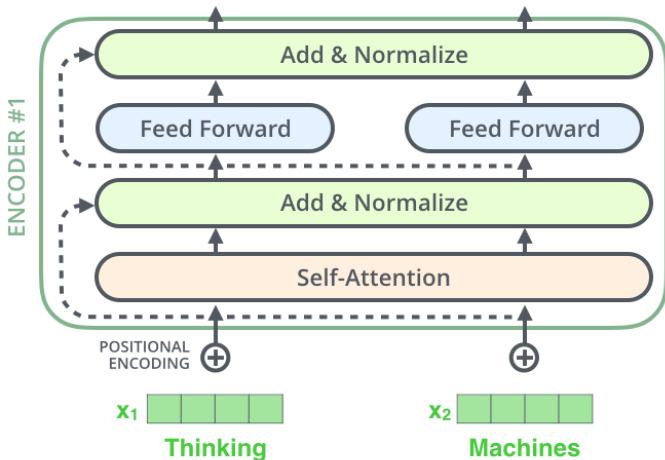
# Positional Encoding

1. Unlike RNN and CNN encoders, attention encoder outputs 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
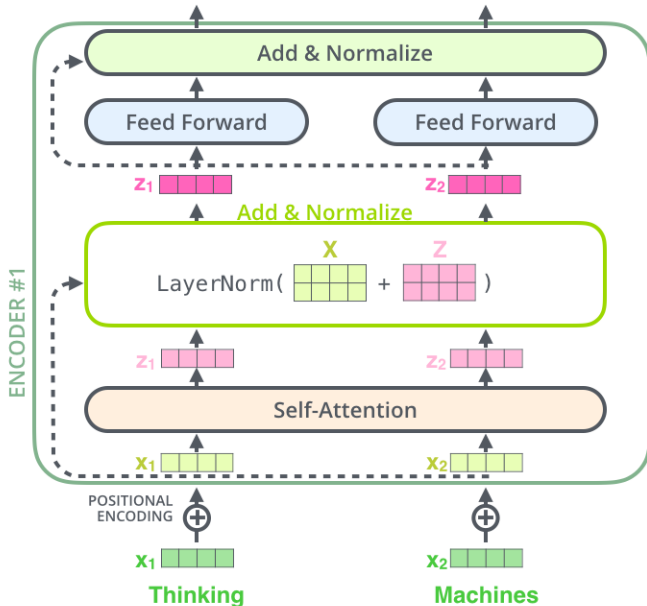
# Positional Encoding

Credits: Jay Alammar
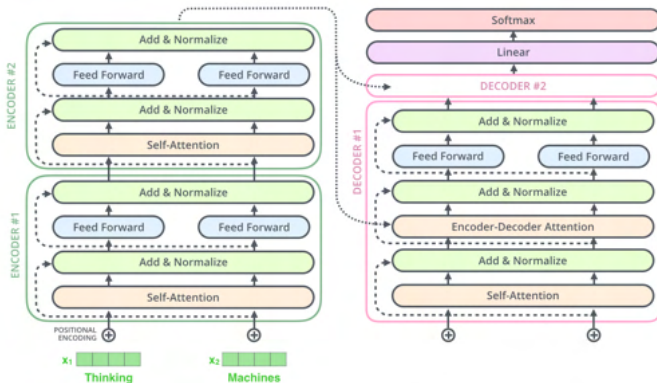
# Residuals in the Encoder



Credits: Jay Alammar

# Residuals in the Encoder

# Tranformer-Decoder

Credits: Jay Alammar

# Transformer-Decoder

1. Self-attention here works in a slightly different way $\rightarrow$ masks the future positions

# Transformer-Decoder

1. Self-attention here works in a slightly different way $\rightarrow$ masks the future positions

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

# Transformer-Decoder

1. Self-attention here works in a slightly different way $\rightarrow$ masks the future positions

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

3. Encoder-decoder attention layer borrows the queries from the layer below it

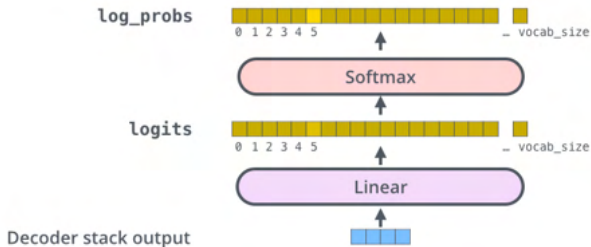# Transformer-Decoder

# Transformer-Decoder

# Final o/p

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

Credits: Jay Alammar