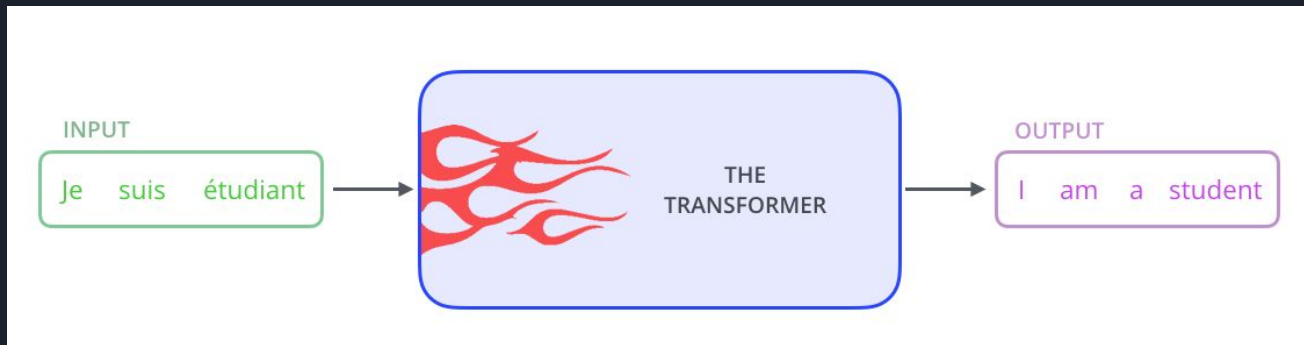# Deep dive into Transformers

Amirhassan Amirmahani - Gholamreza Dar
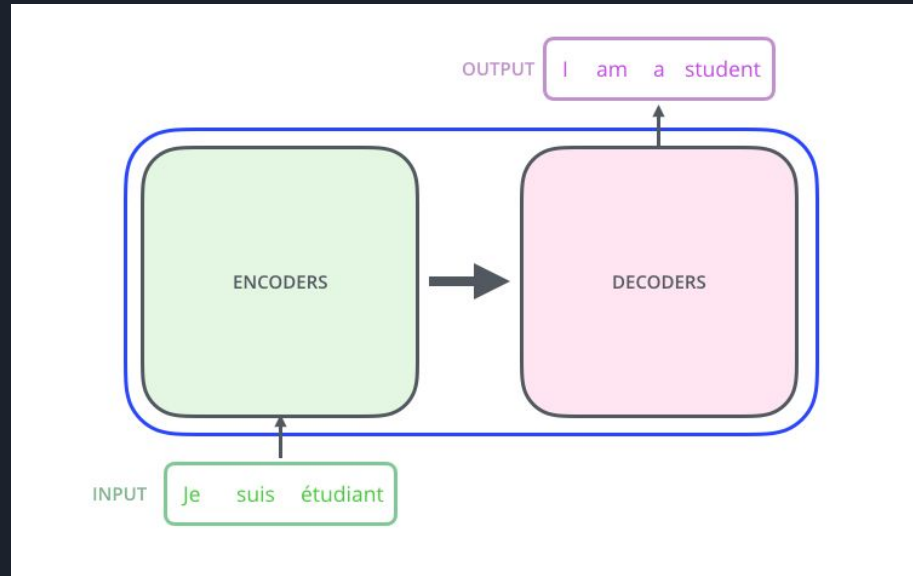
# Transformer Architecture

- Transformer architecture
- Encoder and Decoder stack
- Self-attention
- Multihead attention
- Positional encoding
- Residual connections
- Encoder - Decoder Attention
- Encoder Attention Vs Decoder Attention (Masking)
- The final Linear and Softmax layers
- Training
  - One Hot Encoding
  - Loss Function
  - Another Training Method
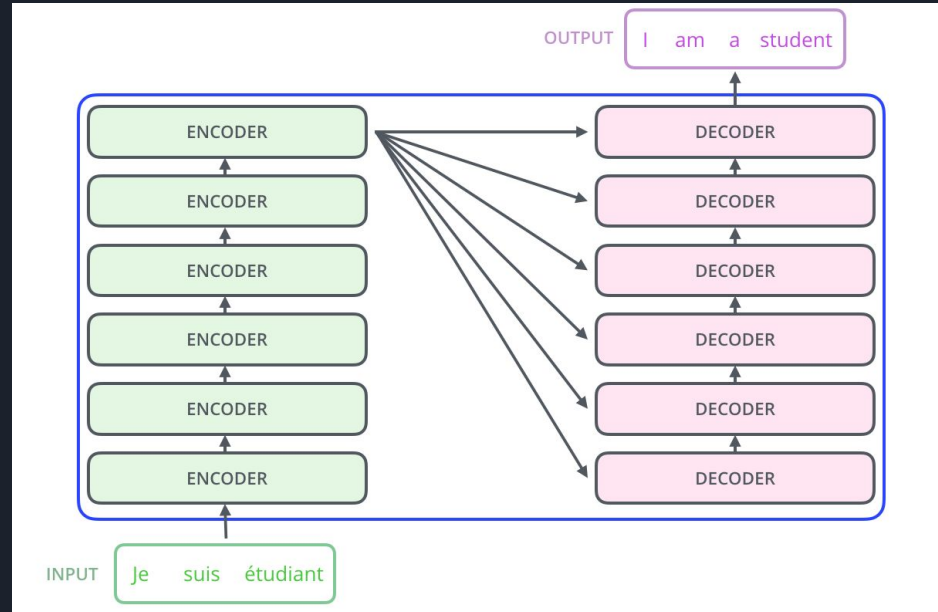
# Transformer as a black box
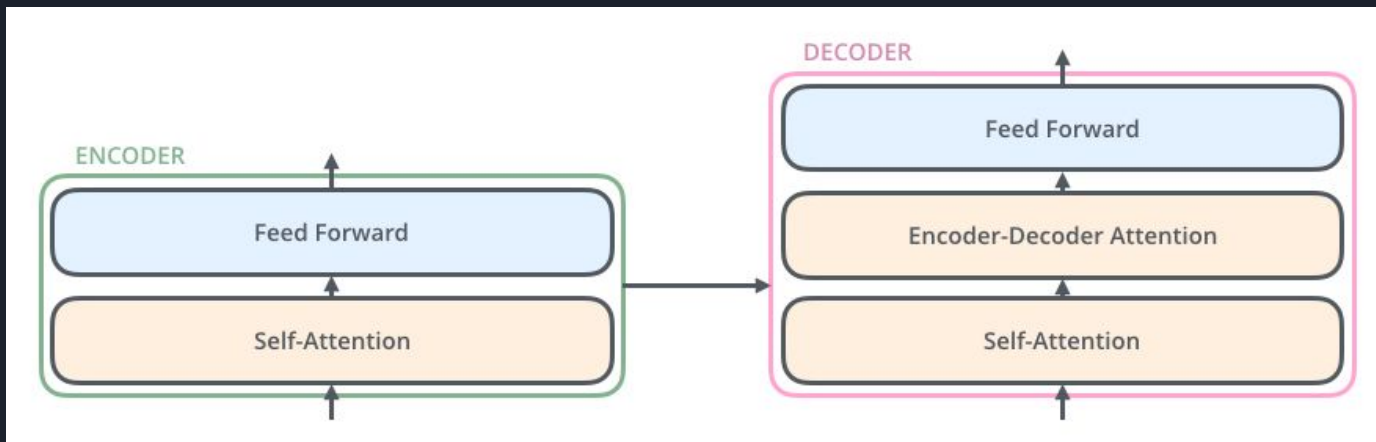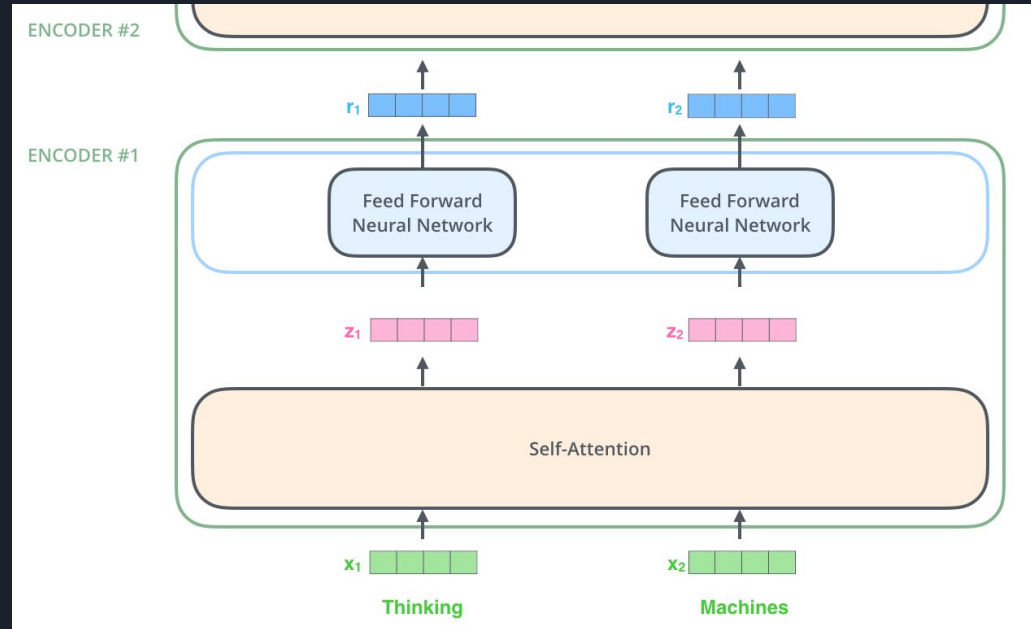
# Inside a Transformer



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

# Encoder and Decoder Stacks



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

# Inside Encoder and Decoder Units



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

# Encoding



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

# Self-attention

"The animal didn't cross the street because it was too tired"

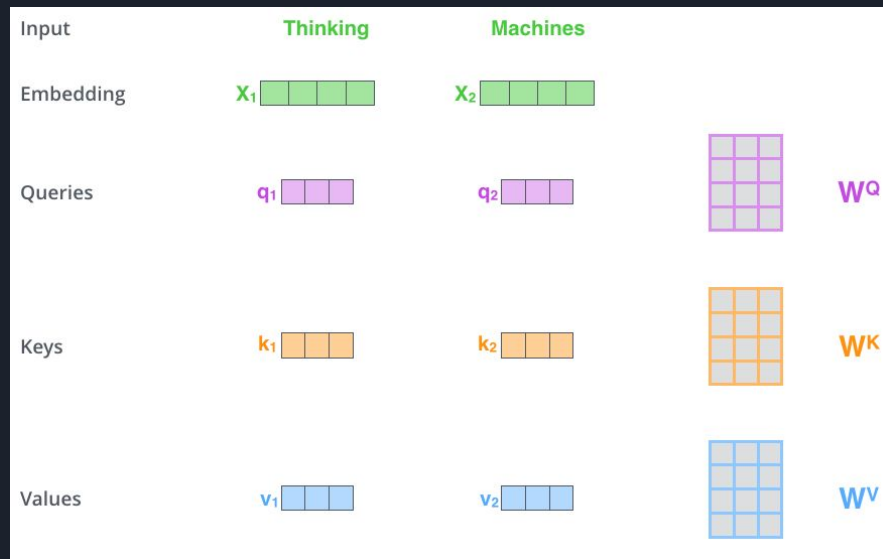

http://jalammar.github.io/illustrated-transformer/

# Tensor2Tensor by Google Brain team

Tensor2Tensor, or T2T for short, is a library of deep learning models and datasets designed to make deep learning more accessible and accelerate ML research.
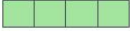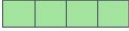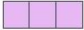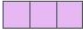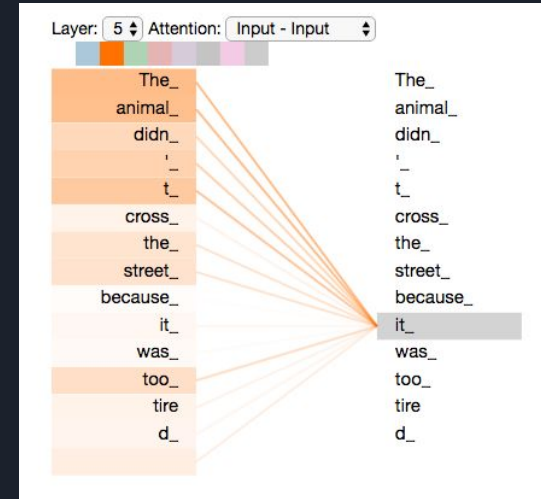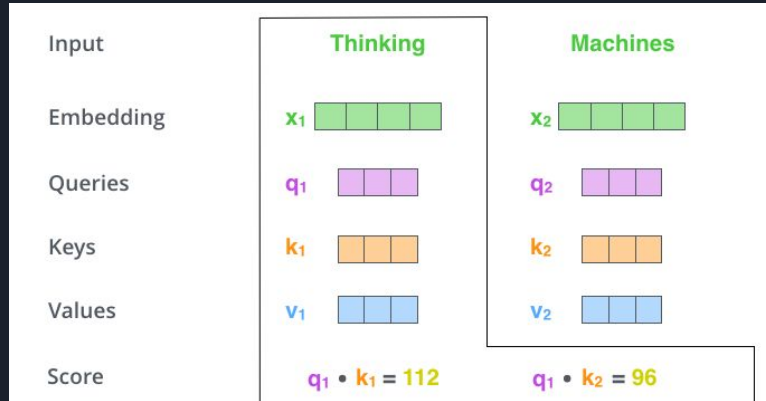
# Self-attention in more detail

# Second step : Score

# Second step : Score





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

# Third step : Divide by 8



| 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 |

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

# Fourth step : Softmax



google.com

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

# Matrix Representations



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

# Self-attention Formula

# Multi-headed attention

# Multi-headed attention

# We only need one Z

# Concatenation and Wo weight matrix



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

=

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

# Multi-headed attention visualization

1 Head

2 Heads

8 Heads

# Positional Encoding



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

https://arxiv.org/abs/1706.03762

# Residual Connections



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

# Encoder Decoder Attention



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

# Encoder Decoder Attention

# Encoder Decoder Attention

Creates Queries matrix from the layer below it, and takes the Keys and Values matrix from the encoder stack.

Encoder Attention
Vs
Decoder Attention

# Masking

# The Final Linear and Softmax Layer



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

https://arxiv.org/abs/1706.03762

# Training

- One Hot Encoding
- Loss Function
- Another Training Method

# One Hot Encoding



Output Vocabulary
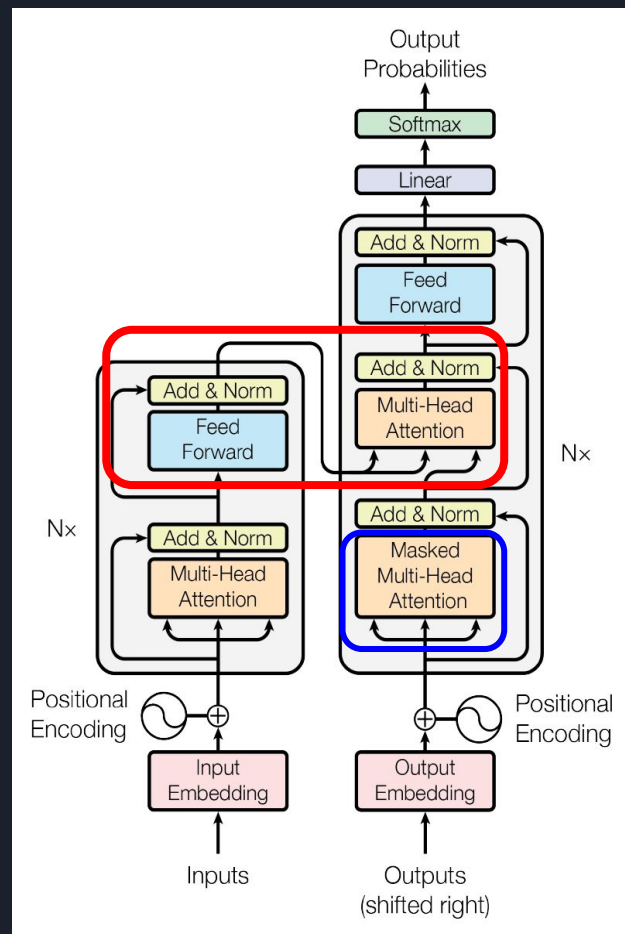
| WORD | a | am | I | thanks | student | <eos> |
|------|---|----|----|--------|---------|-------|
| INDEX | 0 | 1 | 2 | 3 | 4 | 5 |

One-hot encoding of the word "am"

| 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
|-----|-----|-----|-----|-----|-----|

# Loss Function



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

# Loss Function



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

# Another Training Method

Another way to do it would be to hold on to, say, the top two words (say, 'I' and 'a' for example), then in the next step, run the model twice

# Resources

[Illustrated-transformer](#)

[Illustrated Guide to Transformers Neural Network: A step by step explanation](#)

[A Deep Dive Into the Transformer Architecture – The Development of Transformer Models](#)

[Attention Is All You Need](#)

[Transformers and Language Models - YouTube Playlist](#)

# Thanks

hasanmahani08@gmail.com

rezadar1378@gmail.com