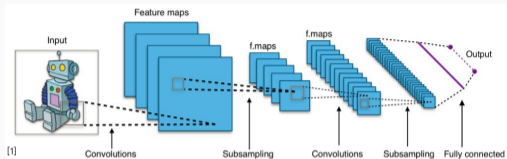


Transformer

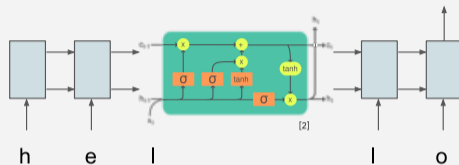
Computer Vision

Convolutional NNs (+ResNets)



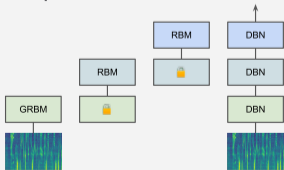
Natural Lang. Proc.

Recurrent NNs (+LSTMs)



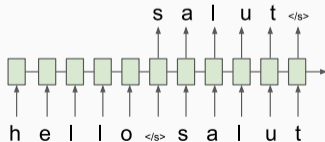
Speech

Deep Belief Nets (+non-DL)



Translation

Seq2Seq



RL

BC/GAIL

Algorithm 1 Generative adversarial imitation learning

- Input:** Expert trajectories $\tau_E \sim \pi_E$, initial policy and discriminator parameters θ_0, ω_0
- for** $i = 0, 1, 2, \dots$ **do**
- Sample trajectories $\tau_i \sim \pi_{\theta_i}$
- Update the discriminator parameters from ω_i to ω_{i+1} with the gradient

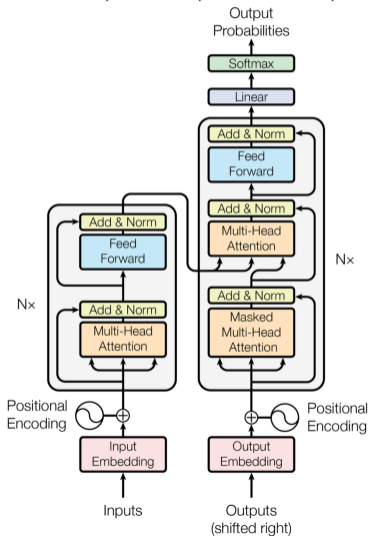
$$\hat{\mathbb{E}}_{\tau_i}[\nabla_{\omega} \log(D_{\omega}(s, a))] + \hat{\mathbb{E}}_{\tau_E}[\nabla_{\omega} \log(1 - D_{\omega}(s, a))] \quad (17)$$
- Take a policy step from θ_i to θ_{i+1} , using the TRPO rule with cost function $\log(D_{\omega_{i+1}}(s, a))$. Specifically, take a KL-constrained natural gradient step with

$$\hat{\mathbb{E}}_{\tau_i}[\nabla_{\theta} \log \pi_{\theta}(a|s)Q(s, a)] - \lambda \nabla_{\theta} H(\pi_{\theta}), \quad (18)$$
 where $Q(\bar{s}, \bar{a}) = \hat{\mathbb{E}}_{\tau_i}[\log(D_{\omega_{i+1}}(s, a)) | s_0 = \bar{s}, a_0 = \bar{a}]$
- end for**

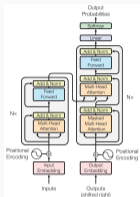
[1] CNN image CC-BY-SA by Apex34 for Wikipedia https://commons.wikimedia.org/wiki/File:Typical_cnn.png
 [2] RNN image CC-BY-SA by GChE for Wikipedia https://commons.wikimedia.org/wiki/File:The_LSTM_Cell.svg

Attention is All you Need

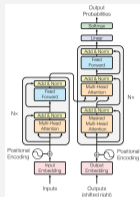
A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin (2017)



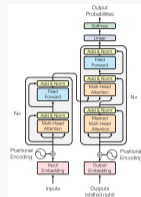
Computer Vision



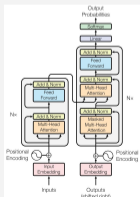
Natural Lang. Proc.



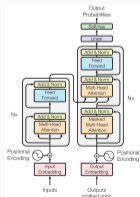
Reinf. Learning



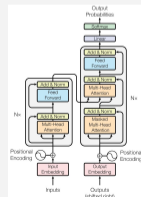
Speech



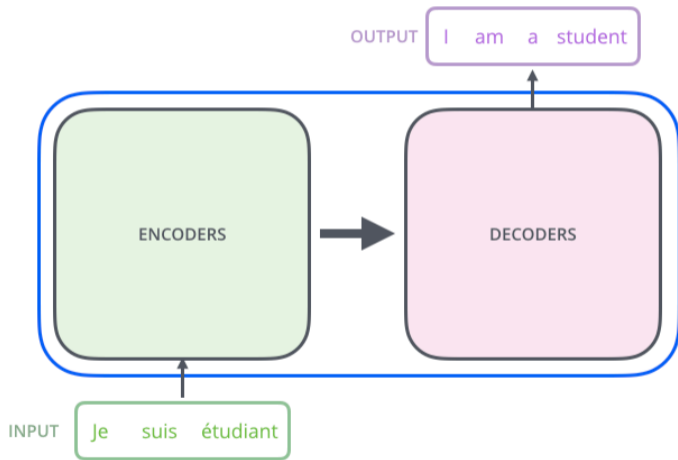
Translation



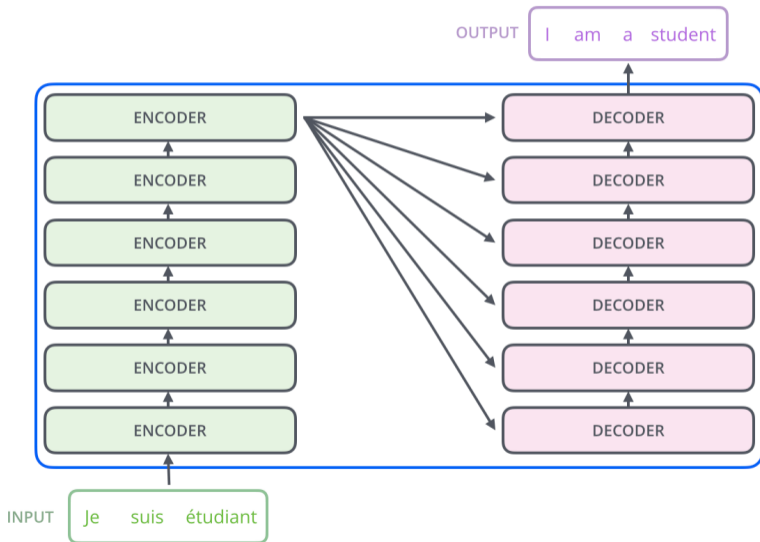
Graphs/Science



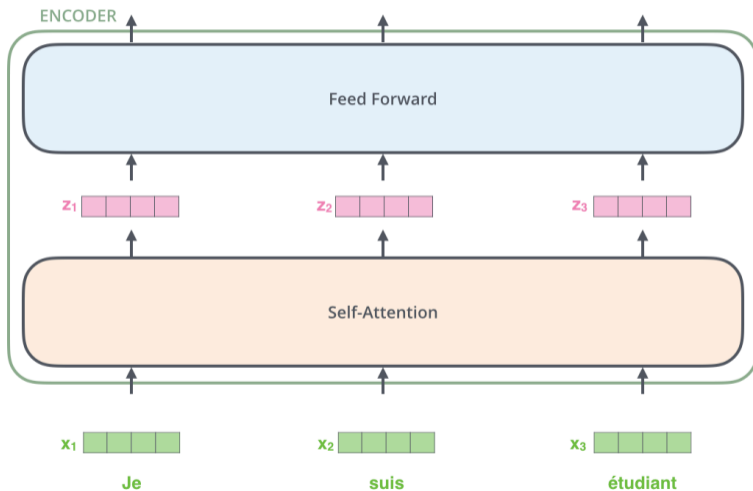
A big picture



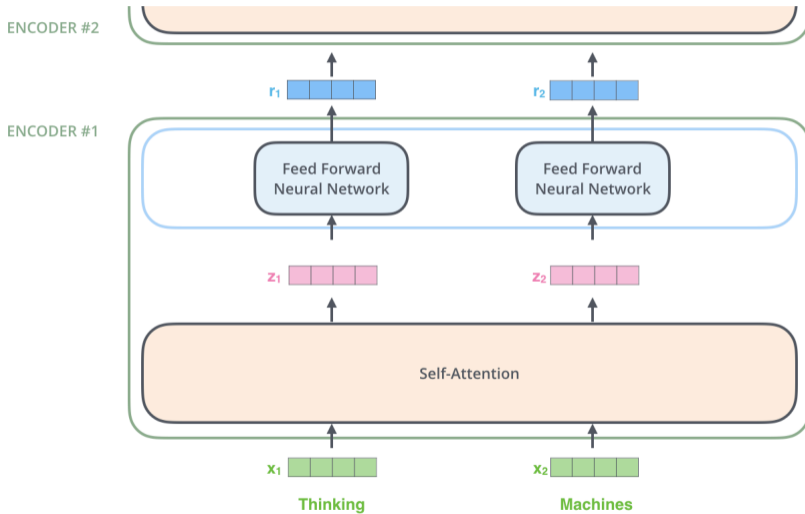
A smaller picture



Inside an Encoder



Zooming in on an Encoder



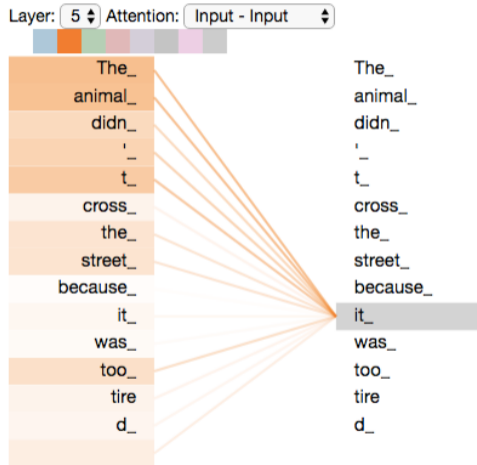
Idea behind Self-Attention

Consider an input sentence:

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

What does “it” in this sentence refer to? The animal, the cross or the street?

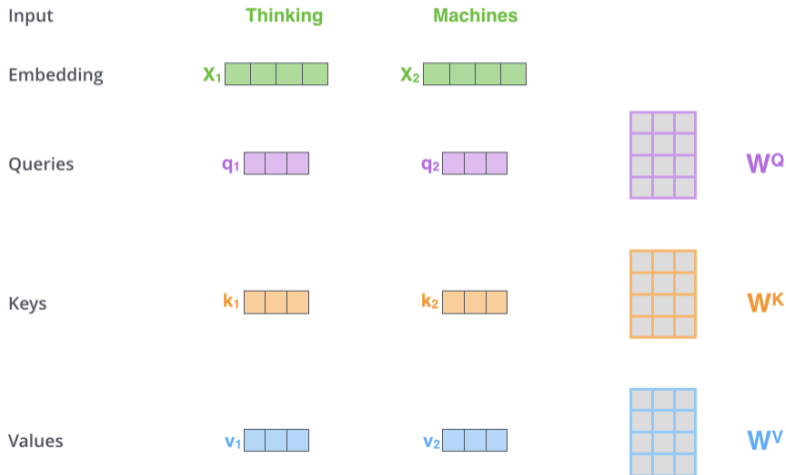
Self-Attention as finding relevant words



Try it yourself

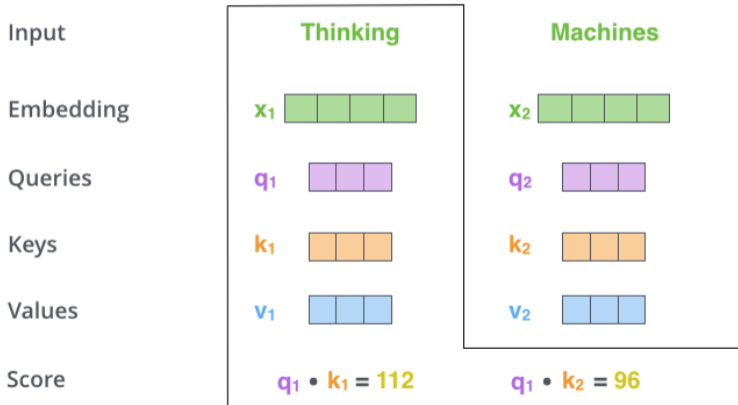
Calculating Self-Attention

Step 1: From an input vectors, create **Q**uery vector, **K**ey vector and **V**alue vector



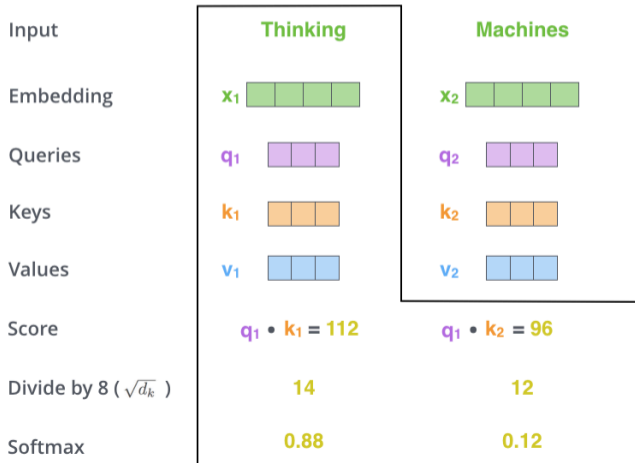
Calculating Self-Attention

Step 2: Compute the **scores** of each word by taking dot-product of its **query** and the **keys** of the all words



Calculating Self-Attention

Step 3: Divide by 8 (or the square root of the dimension of the key vectors) and apply Softmax



Calculating Self-Attention

Step 4: Multiply each value vector by the softmax score, then sum the vectors. Irrelevant words with low softmax scores will not contribute much to the sum

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 

Self-Attention with Multiple inputs

We can write multiple dot-products as a matrix multiplication

$$X \times W^Q = Q$$


A diagram illustrating the calculation of the query matrix Q. It shows a 2x4 green matrix X multiplied by a 4x3 purple matrix W^Q, resulting in a 2x3 purple matrix Q.

$$X \times W^K = K$$


A diagram illustrating the calculation of the key matrix K. It shows a 2x4 green matrix X multiplied by a 4x3 orange matrix W^K, resulting in a 2x3 orange matrix K.

$$X \times W^V = V$$


A diagram illustrating the calculation of the value matrix V. It shows a 2x4 green matrix X multiplied by a 4x3 blue matrix W^V, resulting in a 2x3 blue matrix V.

Self-Attention in one equation

$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \mathbf{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \end{matrix} \right) \begin{matrix} \mathbf{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

=

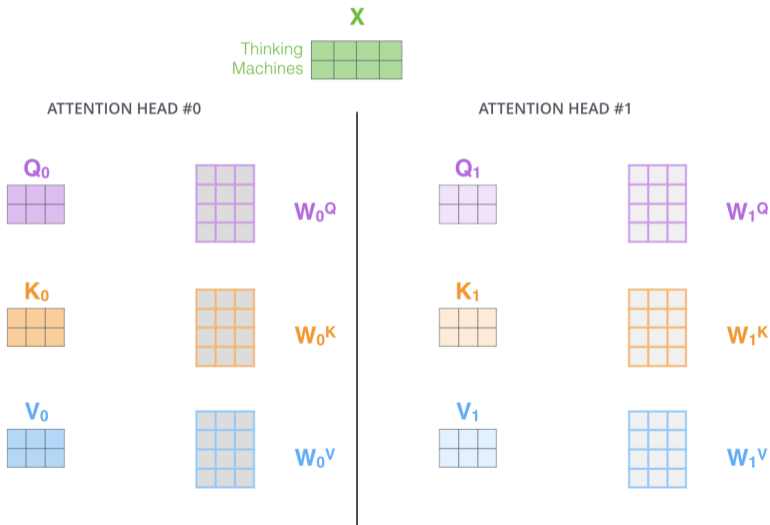
$$\begin{matrix} \mathbf{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

Multi-head Attention

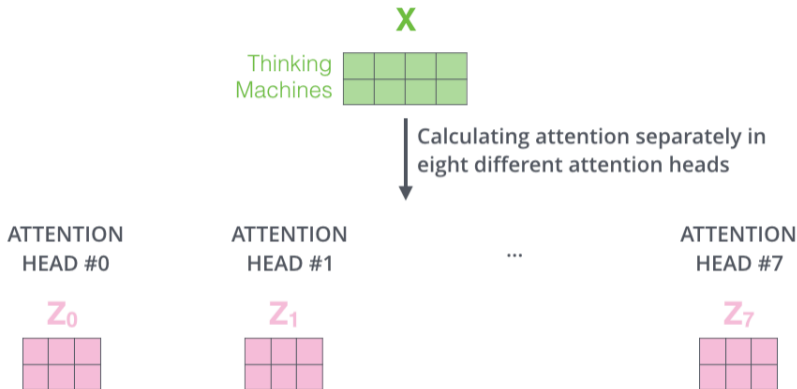
Motivation: want multiple Query/Key/Value combinations

Multi-head Attention

Motivation: want multiple Query/Key/Value combinations



Multi-head Attention



but the feed-forward layer only takes a single matrix. How do we combine these into a single matrix?

Combining matrices

1) Concatenate all the attention heads



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



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

X



Summary of Multi-head Self-Attention

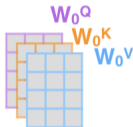
1) This is our input sentence*

Thinking
Machines

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



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



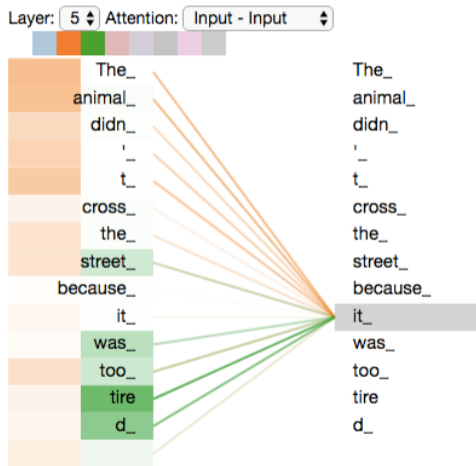
W^O



Z



Revisit the Visualization



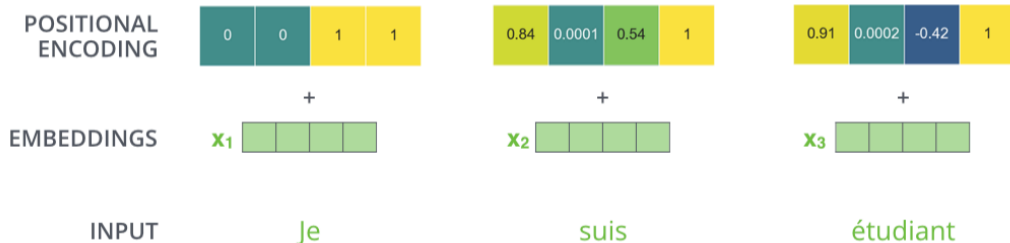
Positional Encoding

Issue: The model does not know the order of words

Positional Encoding

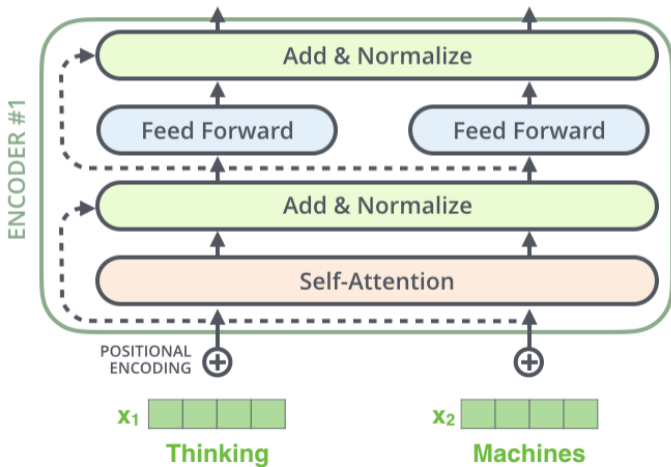
Issue: The model does not know the order of words

Solution: Add different vectors to the sequence of vectors (**positional encoding**)

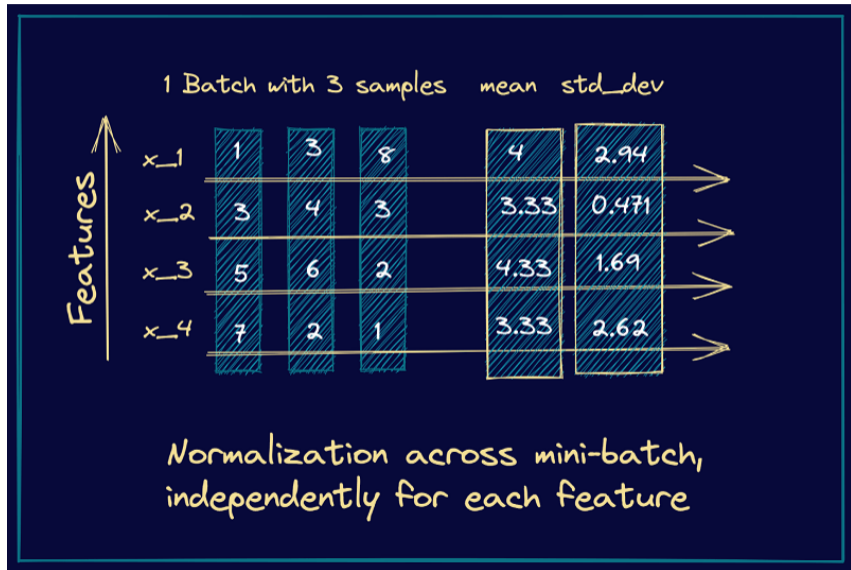


Missing details in the encoder

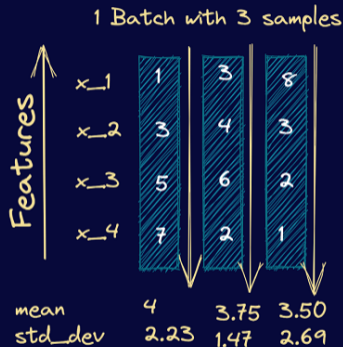
After Self-Attention and Feed-forward, apply **Residual Connection** and **Layer normalization**



Review: Batch normalization

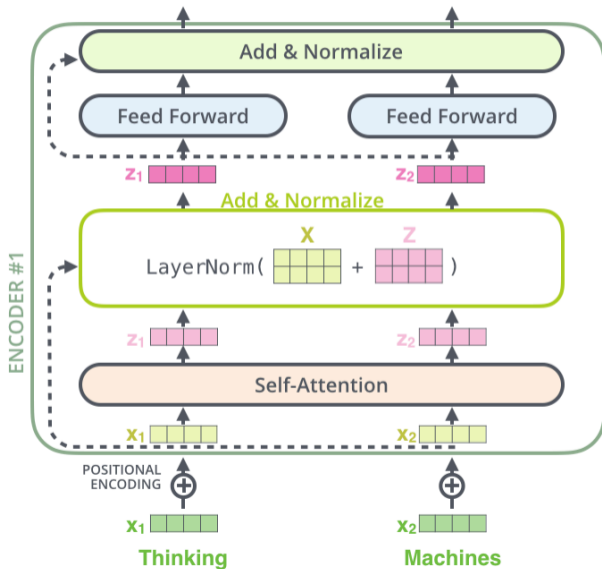


Layer normalization

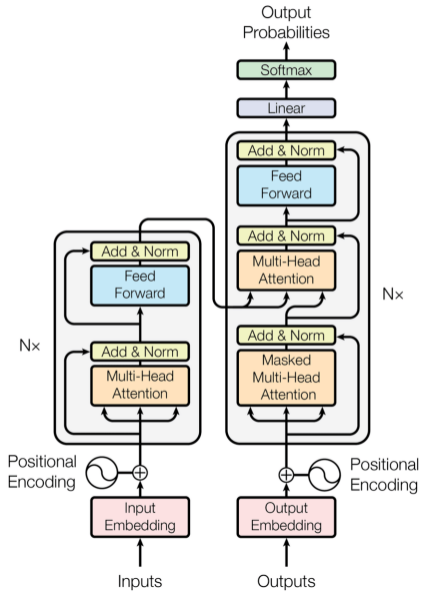


Normalization across features,
independently for each sample

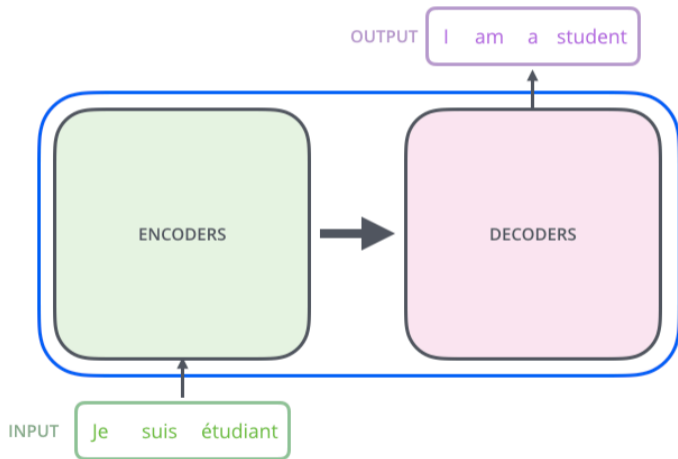
Inside the Encoder



The Decoder: Cross-Attention



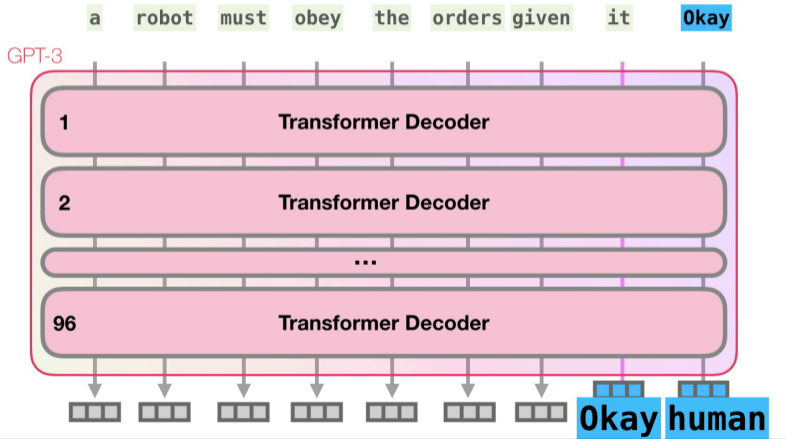
Decoder's outputs



Transformer-based Language Models

GPT3

96 stacks of Transformer Decoder



AI DUNGEON

You arrive at a small wooden store with sheet music in the window. The door is locked so you knock. The music suddenly stops and you hear footsteps coming to the door.

"Hello?" A voice says behind the door.

"I'm looking for the book of essence," You say.

"What do you want it for?" The voice answers, a bit of an edge to it.

> You tell them you are a humble seeker of knowledge.

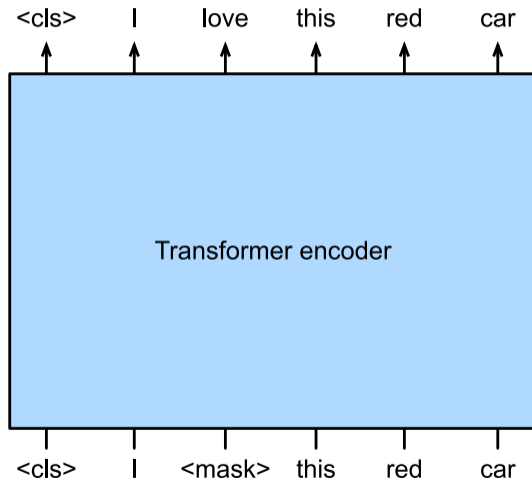
"I am a humble seeker of knowledge," You answer.

"Ah, I see. You seek to learn about the essence of reality, and the universe." The voice says. "I know of this book. I am the guard of this library. The only way you can get to it, is if you get past me. Are you up for the challenge?"

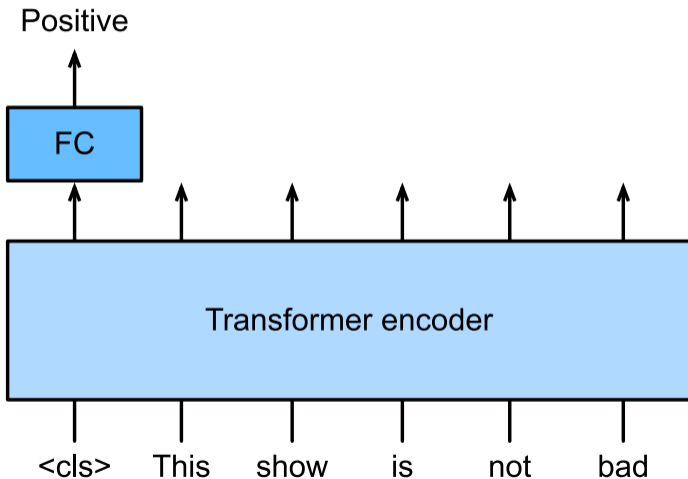
> You ask them what the challenge is.

"The challenge is, you have to win a game of chess against me. If you win, then you may try to take the book. I will get the board." The door unlocks and opens.

BERT



Fine-tuning BERT



References

- Jay Alammar. The Illustrated Transformer.
<https://jalammar.github.io/illustrated-transformer>
- Bala Priya. Build Better Deep Learning Models with Batch and Layer Normalization.
<https://www.pinecone.io/learn/batch-layer-normalization>