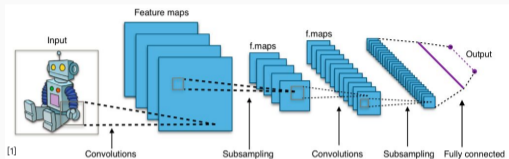


**Transformer**

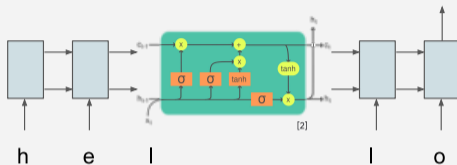
# Computer Vision

Convolutional NNs (+ResNets)



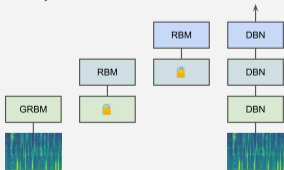
# Natural Lang. Proc.

Recurrent NNs (+LSTMs)



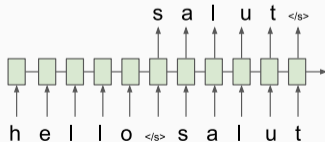
# Speech

Deep Belief Nets (+non-DL)



# Translation

Seq2Seq



# RL

BC/GAIL

**Algorithm 1** Generative adversarial imitation learning

- 1: **Input:** Expert trajectories  $\tau_E \sim \pi_E$ , initial policy and discriminator parameters  $\theta_0, \omega_0$
- 2: **for**  $i = 0, 1, 2, \dots$  **do**
- 3:   Sample trajectories  $\tau_i \sim \pi_{\theta_i}$
- 4:   Update the discriminator parameters from  $\omega_i$  to  $\omega_{i+1}$  with the gradient

$$\hat{\mathbb{E}}_{\tau_i}[\nabla_{\omega} \log(D_{\omega}(s, a))] + \hat{\mathbb{E}}_{\tau_E}[\nabla_{\omega} \log(1 - D_{\omega}(s, a))] \quad (17)$$

- 5:   Take a policy step from  $\theta_i$  to  $\theta_{i+1}$ , 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  $Q(\bar{s}, \bar{a}) = \hat{\mathbb{E}}_{\tau_i}[\log(D_{\omega_{i+1}}(s, a)) | s_0 = \bar{s}, a_0 = \bar{a}]$

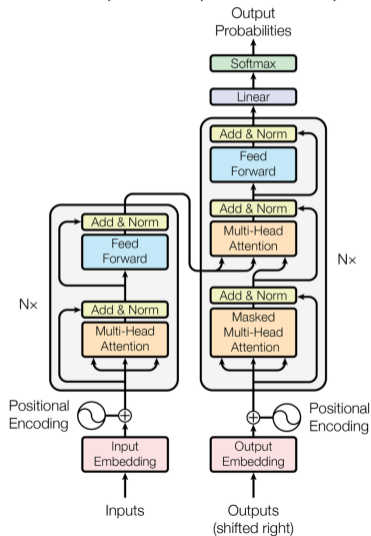
6: **end for**

[1] CNN image CC-BY-SA by Apex34 for Wikipedia [https://commons.wikimedia.org/wiki/File:Typical\\_cnn.png](https://commons.wikimedia.org/wiki/File:Typical_cnn.png)

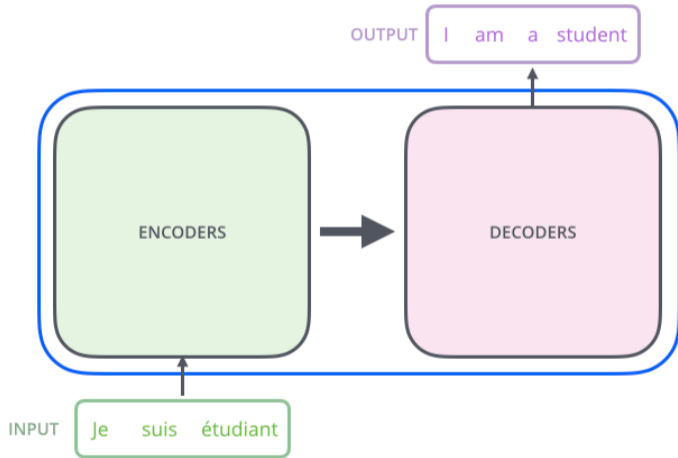
[2] RNN image CC-BY-SA by GChе for Wikipedia [https://commons.wikimedia.org/wiki/File:The\\_LSTM\\_Cell.svg](https://commons.wikimedia.org/wiki/File:The_LSTM_Cell.svg)

# Attention is All you Need

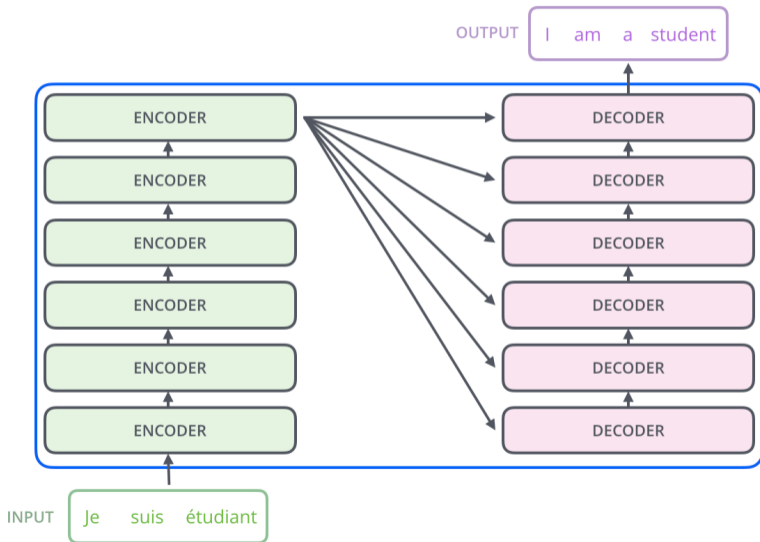
A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin (2017)



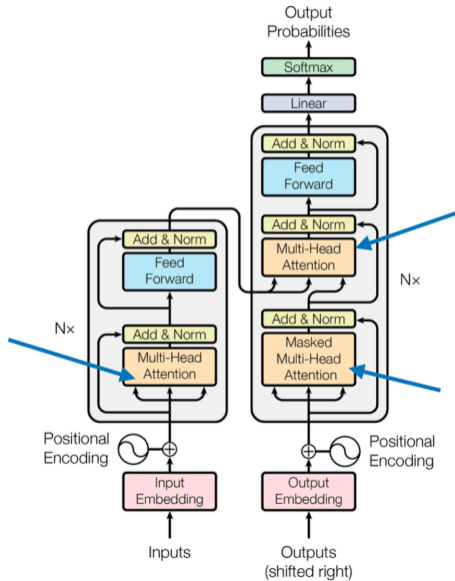
# A big picture



# A smaller picture



# Attention

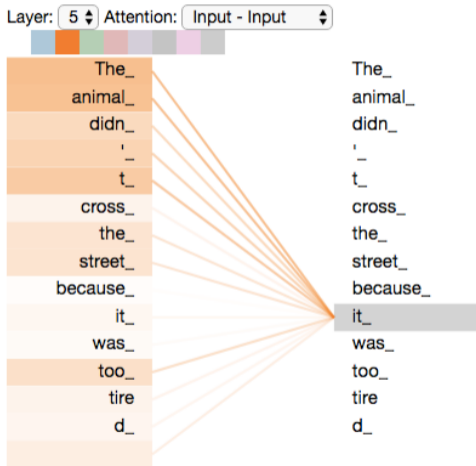


# Idea behind Self-Attention

Can you me help this sentence to translate  
↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑  
Kannst du mir helfen diesen Satz zu uebersetzen ?

Can you help me to translate this sentence  
↑ ↑ ↗ ↘ ↙ ↘ ↙ ↘  
Kannst du mir helfen diesen Satz zu uebersetzen ?

# Self-Attention as finding relevant words

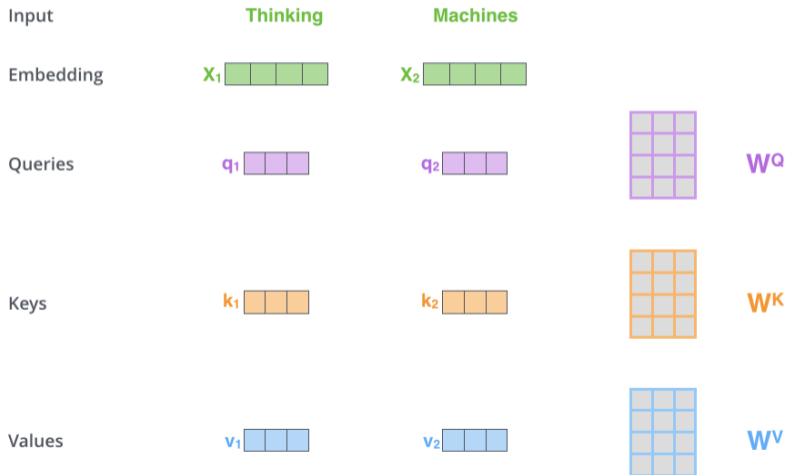


Try it yourself



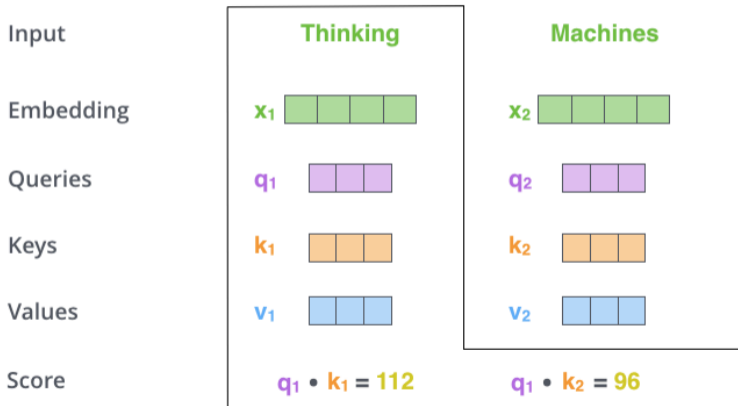
# Calculating Self-Attention

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



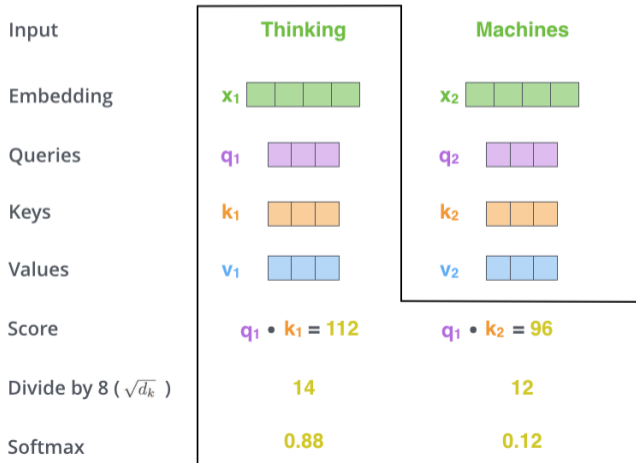
# Calculating Self-Attention

**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

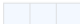
0.12

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

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


The diagram illustrates the calculation of the query matrix Q. It shows a green 2x4 input matrix X multiplied by a purple 4x3 weight matrix W<sup>Q</sup>, resulting in a purple 2x3 output matrix Q.

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


The diagram illustrates the calculation of the key matrix K. It shows a green 2x4 input matrix X multiplied by an orange 4x3 weight matrix W<sup>K</sup>, resulting in an orange 2x3 output matrix K.

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


The diagram illustrates the calculation of the value matrix V. It shows a green 2x4 input matrix X multiplied by a blue 4x3 weight matrix W<sup>V</sup>, resulting in a blue 2x3 output matrix V.

# Self-Attention in one equation

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

=

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

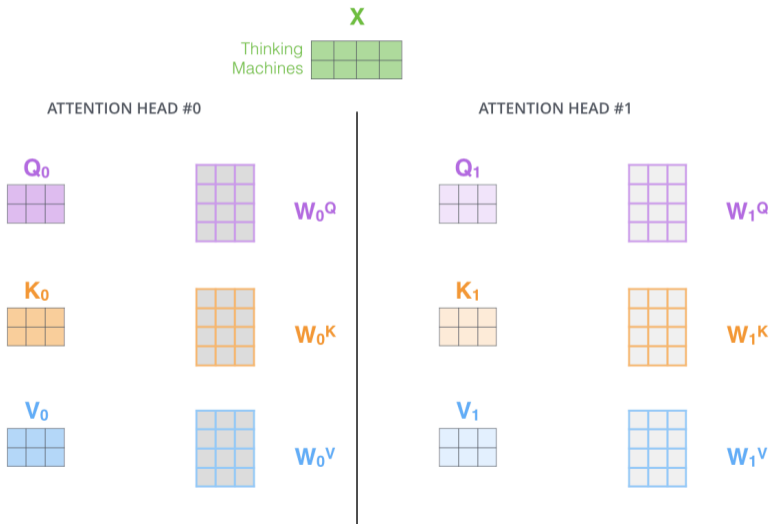
# Multi-head Attention

Motivation: Same words can have multiple meanings:

“She turned on the **light**” vs “The bag was very **light**”

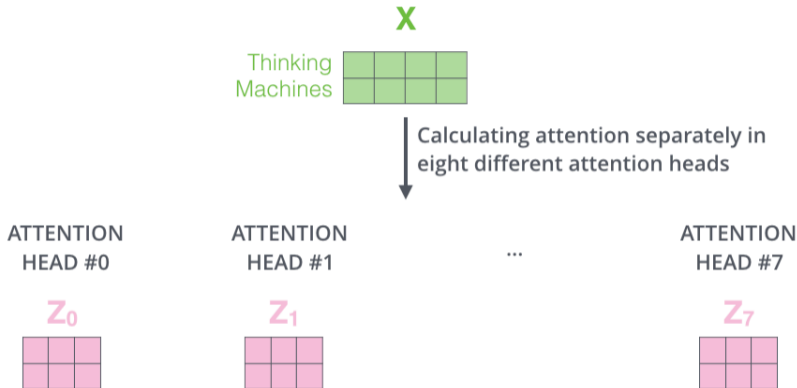
# Multi-head Attention

Motivation: Same words can have multiple meanings:





# 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

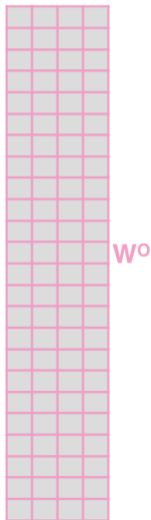


3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN



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

X



# Summary of Multi-head Self-Attention

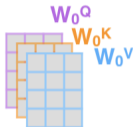
1) This is our input sentence\*

Thinking  
Machines

2) We embed each word\*



3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices



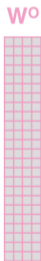
4) Calculate attention using the resulting  $Q/K/V$  matrices



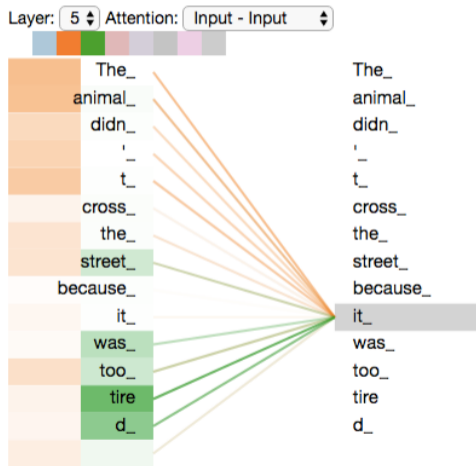
5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer



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

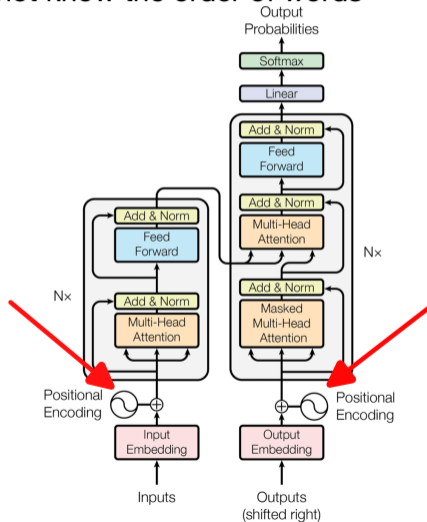


# Revisit the Visualization



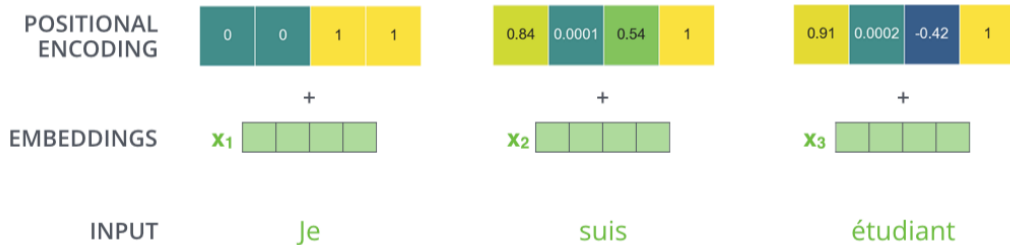
# Positional Encoding

Issue: The model does not know the order of words



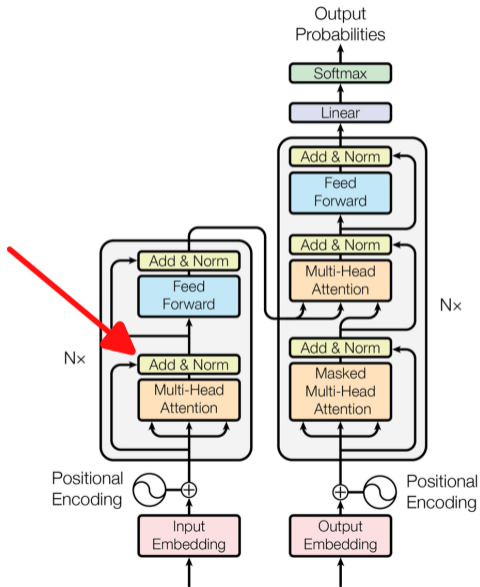
# Positional Encoding

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



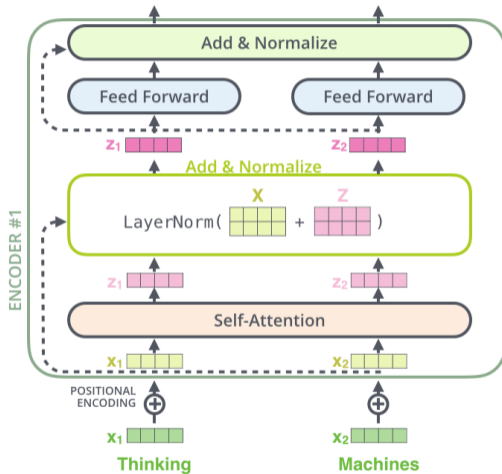


# Addition & Normalization

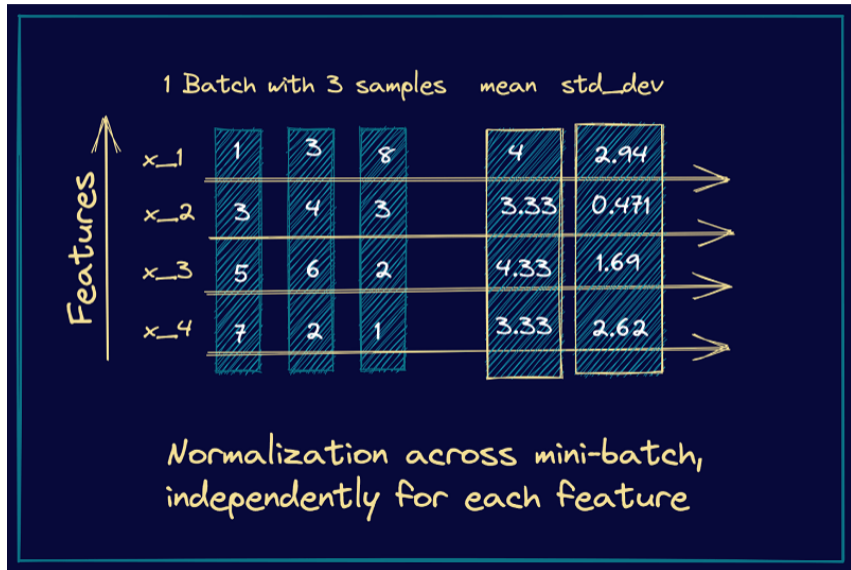




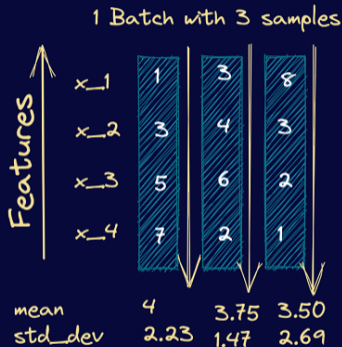
# Addition & Normalization



# Review: Batch normalization

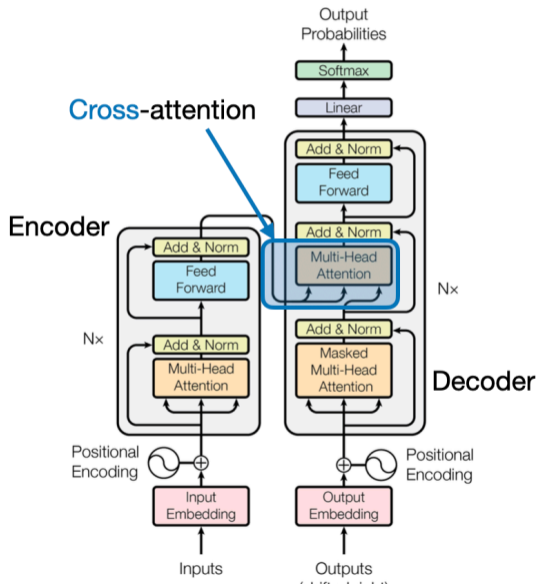


# Layer normalization

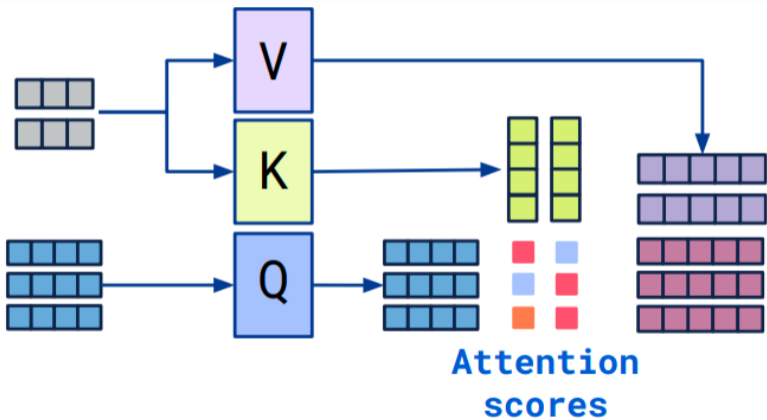


Normalization across features,  
independently for each sample

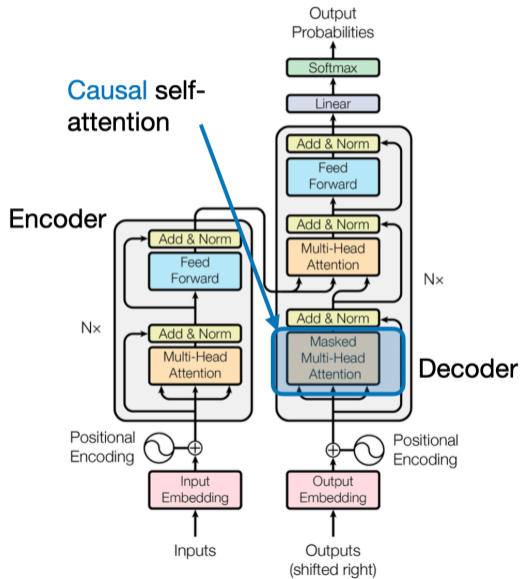
# The Decoder: Cross-Attention



# Cross-Attention



# Causal attention



# Causal attention

Attention weight between the tokens corresponding to “Life” and “short”

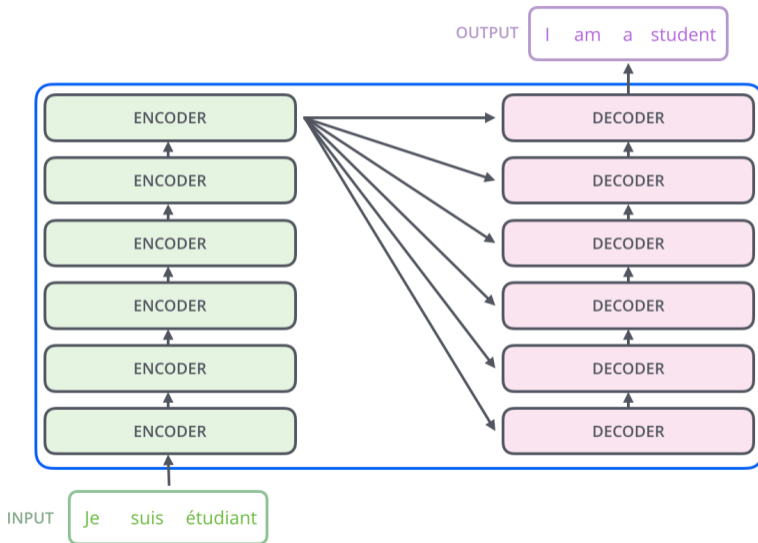
	Life	is	short	eat	desert	first
Life	0.17	0.13	0.18	0.16	0.15	0.18
is	0.03	0.68	0.02	0.08	0.14	0.02
short	0.19	0.06	0.25	0.14	0.11	0.23
eat	0.15	0.21	0.14	0.16	0.17	0.14
desert	0.13	0.27	0.11	0.16	0.18	0.12
first	0.19	0.02	0.31	0.11	0.07	0.27



Masked out future tokens are grayed out

	Life	is	short	eat	desert	first
Life	0.17					
is	0.03	0.68				
short	0.19	0.06	0.25			
eat	0.15	0.21	0.14	0.16		
desert	0.13	0.27	0.11	0.16	0.18	
first	0.19	0.02	0.31	0.11	0.07	0.27

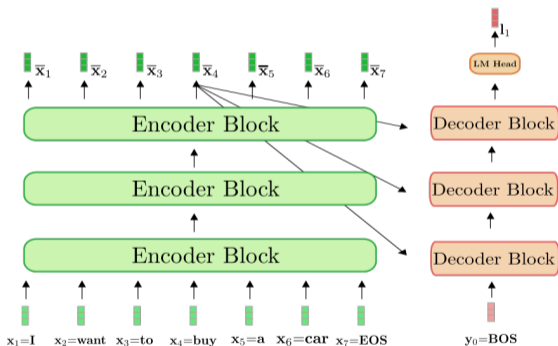
# Training: Teacher forcing



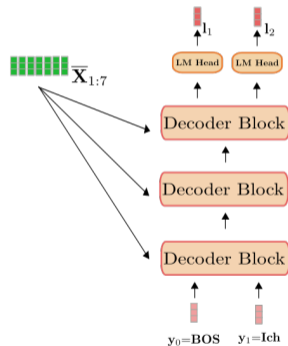


# Predictions

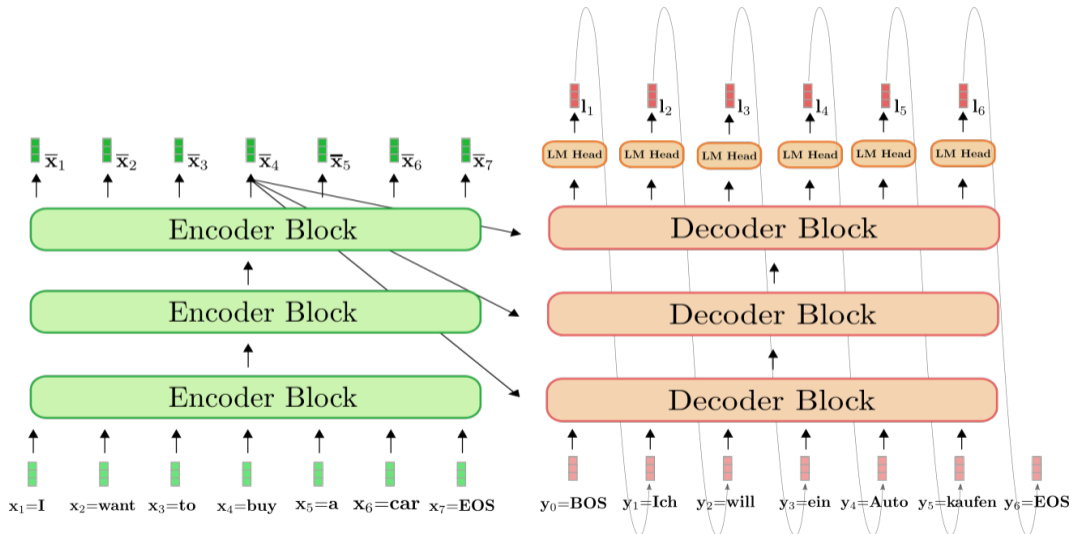
Step 1



Step 2



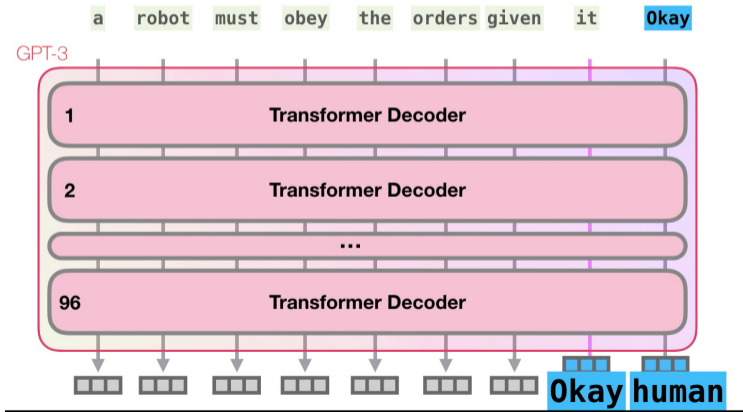
# Predictions



# **Transformer-based Language Models**

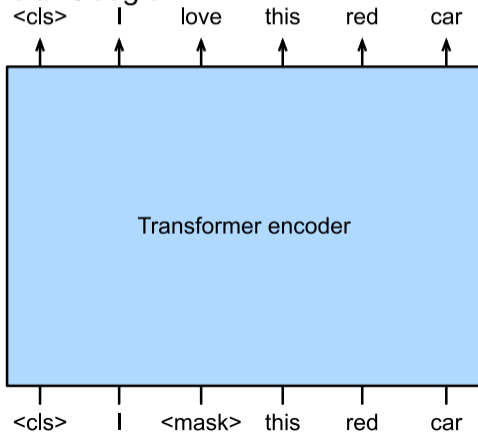
# GPT (Generative Pretrained Transformer)

Invented by researchers at OpenAI  
96 stacks of Transformer Decoder

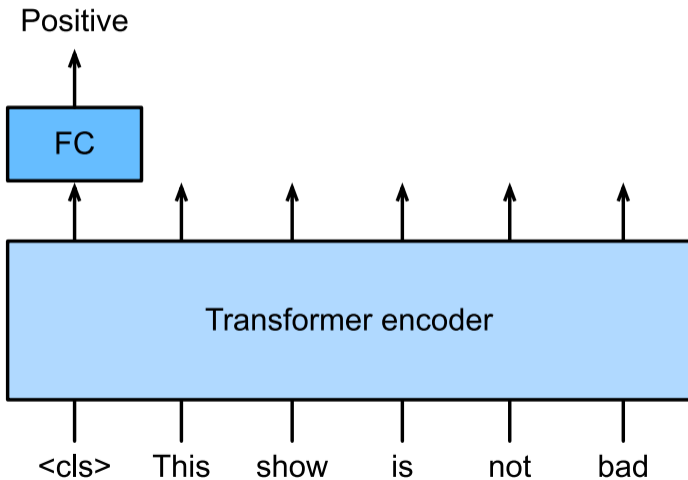


# BERT (Bidirectional Encoder Representations from Transformers)

Invented by researchers at Google AI



# Fine-tuning BERT



# ChatGPT

## Step 1

Collect demonstration data and train a supervised policy.

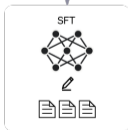
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3.5 with supervised learning.



## Step 2

Collect comparison data and train a reward model.

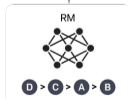
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



## Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.



The PPO model is initialized from the supervised policy.



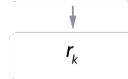
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



# 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>
- Sebastian Raschka. Understanding and Coding Self-Attention, Multi-Head Attention, Cross-Attention, and Causal-Attention in LLMs.  
<https://magazine.sebastianraschka.com/p/understanding-and-coding-self-attention>
- Patrick von Platen. Transformers-based Encoder-Decoder Models.  
<https://huggingface.co/blog/encoder-decoder>